

Distributed Proximal Policy Optimization for Contention-Based Spectrum Access

Akash Doshi and Jeffrey G. Andrews

Department of Electrical and Computer Engineering, University of Texas at Austin, TX 78712

Abstract—The increasing number of wireless devices operating in unlicensed spectrum motivates the development of intelligent adaptive approaches to spectrum access that go beyond traditional carrier sensing. We develop a novel distributed implementation of a policy gradient method known as Proximal Policy Optimization modelled on a two stage Markov decision process that enables such an intelligent approach, and still achieves decentralized contention-based medium access. In each time slot, a base station (BS) uses information from spectrum sensing and reception quality to autonomously decide whether or not to transmit on a given resource, with the goal of maximizing proportional fairness network-wide. Empirically, we find the proportional fairness reward accumulated by the policy gradient approach to be significantly higher than even a genie-aided adaptive energy detection threshold. This is further validated by the improved sum and maximum user throughputs achieved by our approach.

Index Terms—Medium access, proximal policy optimization, contention, deep reinforcement learning

I. INTRODUCTION

Spectrum sharing attempts to allow different transmitters to operate on the same allocated resource (spectrum/time) in a fair manner, while also providing high throughput. To alleviate spectrum constraints, 3GPP [1] standardized License Assisted Access (LAA) for LTE, and more recently released a study on 5G New Radio Unlicensed (NR-U) [2]. The approach adopted to access unlicensed spectrum in LAA/NR-U is known as Listen-Before-Talk (LBT) [1] [3] and requires each transmitter to perform a Clear Channel Assessment (CCA) before accessing spectrum i.e. a BS is allowed to transmit on a channel only if the energy level in the channel is less than the CCA threshold level for the duration of the CCA observation time [3]. However, a CCA threshold level based MAC decision – also referred to as an energy detect (ED) threshold in 5G NR – does not actually reflect the quality of reception (SINR) at the UE.

An optional collision reduction scheme known as Request-to-Send/Clear-to-Send (RTS/CTS) is supported by 802.11, but is known to inhibit potentially successful transmissions, and introduce significant additional overhead and latency [4]. Nearly all WiFi systems disable RTS/CTS. Recently, multi-agent reinforcement learning (RL) has been applied to design state-based policies that can improve the performance of unlicensed spectrum sharing [5]–[7]. Most recently, [8] presented a robust and scalable distributed RL design for radio resource management to mitigate interference. None of these papers thus far have attempted to model the asynchronous nature of the decisions made by the transmitters owing to contention. In [9], we developed a distributed deep RL spectrum sharing

algorithm incorporating contention-based medium access. It deployed Deep Q Networks (DQN) at each BS that sequentially decide whether or not to transmit, with the goal of maximizing proportional fairness (PF) network-wide. However, it suffered from slow training convergence and stability, and achieved a much smaller reward than a PF-based BS scheduler.

Policy gradient methods [10] are known to achieve significantly faster convergence than DQN algorithms, while also improving the reward earned by agents in a multi-agent environment. Consequently, in this paper, we design a novel distributed version of a recent policy gradient method known as Proximal Policy Optimization [11] to optimize medium access under the constraint of a contention-based access mechanism. We employ the paradigm of *centralized learning* with *decentralized execution*, such that each BS will decide whether and how to transmit based only on its own observations.

II. PROBLEM STATEMENT AND SYSTEM MODEL

We consider a downlink cellular deployment of N BSs, with a single UE scheduled per time slot per BS. The notation utilized henceforth is summarized in Table I. Assuming that the UE throughput $R_j[n]$ in each time slot n is approximated by the Shannon capacity $W \log_2(1 + \text{SINR}_j[n])$, the same UE is scheduled for reception for L consecutive time slots and each BS transmits at a constant power, the MAC algorithm at each BS has to decide whether or not to transmit to the UE in each time slot. We consider a simplified contention-based access mechanism in which each time slot is divided into a contention and data transmission period. At the start of the contention period consisting of N mini-slots, BS i draws a random counter $\theta_i \in \{0, \dots, N-1\}$, with the possibility that $\theta_i = \theta_j$ for $i \neq j$. The counter is decremented by 1 every mini-slot and when this counter expires, the BS ascertains if the channel is clear before transmitting a unique preamble for the remainder of the contention period, followed by data transmission, with the objective of each BS being to maximize the long-term data throughput seen by the UE. Mathematically, in [12], this is proved to be equivalent to

$$\max_{n \rightarrow \infty} \sum_{j=1}^N \log(\bar{X}_j[n]) \quad (1)$$

$$\text{where } \bar{X}_j[n] = (1 - 1/B)\bar{X}_j[n-1] + (1/B)R_j[n]. \quad (2)$$

While Proportional Fair (PF) scheduling would amount to an iterative BS scheduler computing the rate vector $\mathbf{R}^*[n]$ for

S_j	Signal Power	a_i	Action chosen = $\{0, 1\}$
I_j	Interference Power	o_i	Local observation
g_{ij}	Path gain with BS i	g'_{ij}	Path gain with BS j
R_j	Data Rate	π_i	Policy to choose a_i
\bar{X}_j	Average Rate	W	Channel Bandwidth

TABLE I: UE j and BS i Notation

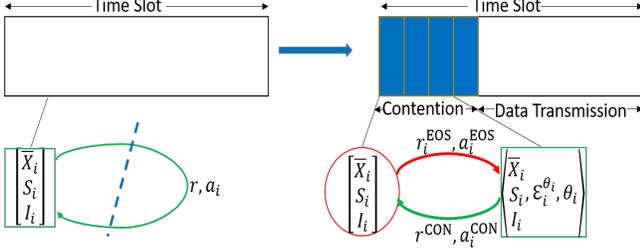


Fig. 1. The 2 state MDP at each agent capturing the actions taken and reward obtained on transitioning between the End-Of-Slot(EOS) and Contention(CON) states

every time slot n , such that

$$\mathbf{R}^*[n] = \arg \max_{\mathbf{R}[n]} \sum_{j=1}^N \frac{R_j[n]}{\bar{X}_j[n]}, \quad (3)$$

it requires a centralized controller and hence is not realizable in any practical decentralized deployment. In [9], we formulated (1) as a decentralized partially observable Markov decision process - *DEC-POMDP*- with the observation o_i received by BS i given by $o_i[n] = \langle \bar{X}_i[n-1], S_i[n-1], I_i[n-1] \rangle$ and a novel per-timestep reward structure given by $r[n] = \sum_{j=1}^N r_j[n] \forall n > 0$ and $r[0] = \sum_{j=1}^N \log(\bar{X}_j[0])$ where

$$r_j[n] = \log \left((1 - 1/B) \left(1 + \frac{R_j[n]}{(B-1)\bar{X}_j[n-1]} \right) \right). \quad (4)$$

We then incorporated contention by partitioning the 1-state MDP into 2 states, End-Of-Slot (EOS) and Contention (CON), as shown on the right in Fig. 1, with $\mathbf{o}_i^{\text{EOS}} = o_i$ and $\mathbf{o}_i^{\text{CON}} = \langle o_i, \mathcal{E}_i^{\theta_i}, \theta_i \rangle$, where $\mathcal{E}_i^{\theta_i} = \{\mathcal{E}_{ij}^{\theta_i}\}_{j \in [N]}$, such that $\mathcal{E}_{ij}^{\theta_i}$ is the energy measured at BS i due to an ongoing transmission between BS j and UE j . We have $r^{\text{CON}}[n] = r[n]$ and $a_i^{\text{CON}} = a_i$, while both r_i^{EOS} and a_i^{EOS} default to 0. The $\mathcal{E}_i^{\theta_i}$ vector served as a message exchanged between agents, and along with the addition of an LSTM layer to each neural network in the system, helped to overcome partial observability in a multi-agent environment [13] [14].

III. PROXIMAL POLICY OPTIMIZATION (PPO) FOR MEDIUM ACCESS DEC-POMDP

A. Proximal Policy Optimization

Given a state s and an action a , we have three terms associated with a typical single agent RL problem: $\pi(a|s; \Theta)$, $\mathcal{Q}^\pi(s, a)$ and $V^\pi(s)$. We denote by $\pi(a|s; \Theta)$ a policy parameterized by Θ that returns the probability of an agent selecting action a in state s . If an agent starts from state s , chooses action a and thereafter follows π , the expected reward

accumulated is represented by $\mathcal{Q}^\pi(s, a)$. Finally, $V^\pi(s)$ denotes the expected reward accumulated by an agent following π starting from state s . Policy-based model-free methods directly parameterize the policy $\pi(a|s; \Theta)$ and update Θ by performing gradient ascent on $J(\Theta) = \mathbb{E}[\gamma^n r[n]]$. The gradient is given by $\nabla_{\Theta} \log \pi(a|s; \Theta) \mathcal{Q}^\pi(a|s)$. To reduce the variance of this unbiased estimate of $\nabla_{\Theta} J(\Theta)$, a learnt baseline $b(s) \approx V^\pi(s)$ is subtracted [15] such that

$$\nabla_{\Theta} J(\Theta) = \nabla_{\Theta} \log \pi(a|s; \Theta) (\mathcal{Q}^\pi(a|s) - V^\pi(s)), \quad (5)$$

where $A(s, a) = \mathcal{Q}^\pi(a|s) - V^\pi(s)$ is known as the *advantage*. This approach can then be viewed as an actor-critic architecture, where actor refers to $\pi(a|s)$ while $V^\pi(s)$ is the critic. In actor-critic algorithms, both the actor and critic are represented by two separate neural networks, $\pi(a|s; \Theta)$ and $V(s; \vartheta)$, parameterized by Θ and ϑ respectively, and the following loss function is minimized

$$L(\Theta, \vartheta) = -J(\Theta) + (V(s, \vartheta) - V^{\text{target}}(s))^2. \quad (6)$$

A truncated version of generalized advantage estimation (GAE) [16] is traditionally utilized to compute $V^{\text{target}}(s_t)$ as

$$V(s_t; \vartheta) + \delta_t + (\gamma\lambda)\delta_{t+1} + \dots + (\gamma\lambda)^{T-t+1}\delta_{T-1}, \quad (7)$$

$$\text{where } \delta_t = r_t + \gamma V(s_{t+1}; \vartheta) - V(s_t; \vartheta). \quad (8)$$

Proximal Policy Optimization (PPO) [11] alters the loss function in (6) in two ways. Firstly, it replaces $J(\Theta) = \mathbb{E}[\pi(a|s; \Theta)A(s, a)]$ with

$$L^{\text{CLIP}}(\Theta) = \mathbb{E}[\min(r(\Theta)A, \text{clip}(r, 1 - \epsilon, 1 + \epsilon)A)], \quad (9)$$

$$\text{where } r(\Theta) = \pi(a|s; \Theta) / \pi(a|s; \Theta_{\text{old}}). \quad (10)$$

The motivation for this modified metric is to not reward excessively large policy updates, which is enforced via the clip function, while the ‘‘surrogate objective’’ $r(\Theta)A$ arises from a policy gradient approach known as trust region policy optimization (TRPO) [17], a precursor to PPO. The term $\pi(a|s; \Theta_{\text{old}})$ in (10) is a constant that corresponds to the evaluation of the policy π at the given (s, a) using the current weights Θ_{old} of the policy NN. Secondly, PPO adds an entropy bonus $S[\pi(\cdot|s; \Theta)]$ to ensure sufficient exploration. Consequently, the final PPO objective $L^{\text{PPO}}(\Theta, \vartheta)$ which is maximized each iteration is given by [11]

$$L^{\text{CLIP}}(\Theta) - c_1 (V(s, \vartheta) - V^{\text{target}}(s))^2 + c_2 S[\pi(\cdot|s; \Theta)]. \quad (11)$$

B. Adapting PPO to a Medium Access DEC-POMDP

We now define N loss functions $L(\Theta_i^{\pi \text{CON}}, \vartheta_i^{V \text{CON}}, \vartheta_i^{V \text{EOS}})$ corresponding to each BS i , by adapting (11) to the 2-state MDP presented in Fig. 1, as follows

$$L^{\text{PPO}}(\Theta_i^{\pi \text{CON}}, \vartheta_i^{V \text{CON}}) - c_3 L^{\text{VF}}(\vartheta_i^{V \text{EOS}}), \quad (12)$$

at each BS i , where

$$L^{\text{VF}}(\vartheta_i^{V \text{EOS}}) = (V_i^{\text{EOS}}(\mathbf{o}_i^{\text{EOS}}) - V_i^{\text{target, EOS}}(\mathbf{o}_i^{\text{EOS}}))^2, \quad (13)$$

while a similar expression for $L^{\text{VF}}(\vartheta_i^{V \text{CON}})$ is already part of $L^{\text{PPO}}(\Theta_i^{\pi \text{CON}}, \vartheta_i^{V \text{CON}})$. Note that the only key changes from

(11) are that we have added a term for training V_i^{EOS} and the input to the EOS and CON NNs will be $\mathbf{o}_i^{\text{EOS}}$ and $\mathbf{o}_i^{\text{CON}}$ respectively, instead of the full system state s . To overcome this partial observability, an LSTM layer is introduced in V_i^{CON} , V_i^{EOS} and π_i^{CON} at every BS i . Now, to compute $V_i^{\text{target,CON}}$ and $V_i^{\text{target,EOS}}$, we first observe that (8), when applied to the EOS-CON transition yields

$$\delta_{i,n}^{\text{CON}} = r[n] + \gamma^{\frac{1}{2}} V_i^{\text{EOS}}(\mathbf{o}_i^{\text{EOS}}[n+1]) - V_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}}[n]) \quad (14)$$

$$\delta_{i,n}^{\text{EOS}} = \gamma^{\frac{1}{2}} V_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}}[n]) - V_i^{\text{EOS}}(\mathbf{o}_i^{\text{EOS}}[n]). \quad (15)$$

Note that the factor of $\gamma^{\frac{1}{2}}$ in (14) and (15) is simply meant to keep the overall discount factor to γ in one time step. Substituting (14) and (15) into (7) and replacing T by the episode length L , we obtain

$$V_i^{\text{target,EOS}} = V_i^{\text{EOS}}(\mathbf{o}_i^{\text{EOS}}) + \delta_{i,n}^{\text{EOS}} + (\gamma^{\frac{1}{2}} \lambda) \delta_{i,n}^{\text{CON}} + \dots + (\gamma^{\frac{1}{2}} \lambda)^{L-n+1} \delta_{i,L-1}^{\text{EOS}} \quad (16)$$

$$V_i^{\text{target,CON}} = V_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}}) + \delta_{i,n}^{\text{CON}} + (\gamma^{\frac{1}{2}} \lambda) \delta_{i,n+1}^{\text{EOS}} + \dots + (\gamma^{\frac{1}{2}} \lambda)^{L-n+1} \delta_{i,L-1}^{\text{CON}}. \quad (17)$$

Note that we will denote $V_i^{\text{target,CON}}(\mathbf{o}_i^{\text{CON}}) - V_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}})$ in (17) as \hat{A}_i^{CON} , an estimate of A_i^{CON} . This will be utilized for computing $L^{\text{CLIP}}(\Theta_i^{\pi^{\text{CON}}})$ via (9).

C. Generating an episode

An episode refers to a collection of L consecutive time slots. At the beginning of time slot n , a random counter θ_i is drawn for each BS i . Each π_i^{CON} outputs two probabilities corresponding to the actions 0 and 1, with

$$a_i = \arg \max_{a \in A_i} \pi_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}})[a]. \quad (18)$$

While generating an episode during the training of the algorithm, we simply sample the action randomly from the probability distribution outputted by $\pi_i^{\text{CON}}(\mathbf{o}_i^{\text{CON}})$ [18].

Consider as an example $N = 3$ with BS 0, 1 and 2 being allocated counter values $\langle \theta_0, \theta_1, \theta_2 \rangle = \langle 2, 0, 1 \rangle$ in time slot n . Since $\theta_1 = 0$, BS 1 goes first and measures the energy from ongoing transmissions to compute $\mathcal{E}_1^{\theta_1}$. It senses no other BS's transmitting ($\mathcal{E}_1^{\theta_1} = [0, 0, 0]$), and in combination with $\bar{X}_1[n-1]$, $S_1[n-1]$ and $I_1[n-1]$ of the UE it serves, it determines $a_1[n]$ using the policy given in (18). Let us assume it chose to transmit (transmission is not a given simply because $\mathcal{E}_1^{\theta_1} = [0, 0, 0]$). BS 2 is scheduled next, detects BS 1 is transmitting such that $\mathcal{E}_2^{\theta_2}$ is non-zero and π_2^{CON} instructs it not to transmit. Finally BS 0 also detects a non-zero $\mathcal{E}_0^{\theta_0}$, but chooses to transmit. Note that while the training procedure, elaborated in Section IV, will require training V_i^{CON} , V_i^{EOS} and π_i^{CON} , testing the learnt policy using (18) only requires π_i^{CON} .

Once all the BS's have taken an action a_i , the action vector \mathbf{a} ($\langle 1, 1, 0 \rangle$ in this example) and $\{g_{ij}\}$ are used to calculate the reward $r[n]$ and the updated average rates $\bar{\mathbf{X}}[n]$. These determine the observations $\mathbf{o}_i^{\text{EOS}}$ for the next time slot. In the

next time slot $n+1$, a new counter θ'_i is drawn at each BS i and the process repeated.

IV. SIMULATION DETAILS

The performance metric is the expected cumulative reward $\sum_{n=0}^L \gamma^n r[n]$, with $r[n]$ given by (4). Note that for $\gamma \rightarrow 1$, $\sum_{n=0}^L \gamma^n r[n] \rightarrow \sum_{j=1}^N \log(\bar{X}_j[L])$. In each iteration, N_{batch} episodes are generated by N_{batch} π_i^{CON} (actors) at each BS i acting in parallel. An overview of the training procedure is presented in Algorithm 1, while the simulation parameters are summarized in Table II.

Policy gradient methods in multi-agent environments typically exhibit very high variance and perform poorly in absence of both stationarity and the Markov property. Consequently, to stabilize the training and improve the learnt policy, we incorporate a *decentralized actor centralized critic* approach, first proposed in [20]. The motivation behind this approach is to use extra information to ease training, so long as this information is not used at test time i.e. centralized training with decentralized execution. In Algorithm 1, we observe that only π_i^{CON} is required for generating an episode i.e. at test time. Hence, we change the input to both V_i^{EOS} and V_i^{CON} by replacing $\mathbf{o}_i^{\text{EOS}}$ with \mathbf{s}^{EOS} at each BS i . While we defined $\mathbf{o}_i^{\text{EOS}}[n+1] = \langle \bar{X}_i[n], S_i[n], I_i[n] \rangle$, we have $\mathbf{s}^{\text{EOS}}[n+1] = \langle \bar{\mathbf{X}}[n], \mathbf{S}[n], \mathbf{I}[n] \rangle$ i.e. it will contain the average rate, signal and interference power of all UEs in the previous time slot. Hence the input to V_i^{CON} will be $\langle \mathbf{s}^{\text{EOS}}, \mathcal{E}_i^{\theta_i}, \theta_i \rangle$, while the input to π_i^{CON} remains $\mathbf{o}_i^{\text{CON}}$.

In order to have a fair comparison with the DQN algorithm from [9], we make two changes in the implementation of distributed DQN. Firstly, we alter the input to Q_i^{EOS} to \mathbf{s}^{EOS} in place of $\mathbf{o}_i^{\text{EOS}}$. Secondly, in each iteration, N_{batch} episodes are generated using the current Q_i^{EOS} and Q_i^{CON} , in place of the replay memories D^{EOS} and D^{CON} utilized in [9] that added one episode generated using the current NNs and removed the oldest episode every iteration. Finally, we will also compare with the PF, ED and Adaptive ED baselines. ED allows a BS to transmit only if $\sum_{j=1}^N \mathcal{E}_{ij}^{\theta_i} < E_0$. We employ $E_0 = -72$ dBm [1]. Adaptive ED finds the ED threshold that maximizes $\sum_{n=0}^L \gamma^n r[n]$ for the given configuration of UEs from a set of ED thresholds ranging from -22 to -92 dBm.

V. SUMMARY OF RESULTS

We consider 4 BSs lying at corners of a rectangle of breadth 20 m and length 20 m in *Layout 1* (L1) and 60 m in *Layout 2* (L2). As the rectangle length is increased, for most choice of 4 UE's, the inter-BS energies \mathcal{E}_i will more accurately reflect the quality of reception. This is because the separation between UE's from different BS's reflects the inter-BS separation more accurately as we move from L1 in Fig. 2a to L2 in Fig. 2b.

The validation curve is shown for Layout 1 and 2 in Fig. 3a and 3b respectively for both the DQN and PPO methods, along with the constant benchmarks provided by the PF and ED baselines. It is obtained by evaluating the trained models obtained after every 50 iterations on 15 randomly sampled configurations and averaged over 20 realizations of

N	4
Layout	InH-Office [2]
Noise PSD	-174 dBm/Hz
Bandwidth W	20 MHz
(UE, BS) Noise Figure	(9, 5) dB
Fading Coefficient α	0.1
Smoothing Window B	10
Center frequency f_c	6 GHz

(a) Data Generation Parameters

Initial Learning Rate η	L1, L2, L3: $(4, 4, 2) \times 10^{-4}$
Learning Rate Decay	L1 & L2: 0.85 / 500 updates L3: 0.5 / 250 updates
Optimizer	Adam [19]
N_{batch}	8
Training Iterations	800
ϵ, γ, L	0.2, 1 - 1e-6, 2000
$ state_h_{i,n} $ for (PPO, DQN)	(128, 256)

(b) RL Training Parameters

TABLE II: Simulation Parameters

Algorithm 1: Spectrum Sharing Proximal Policy Optimization

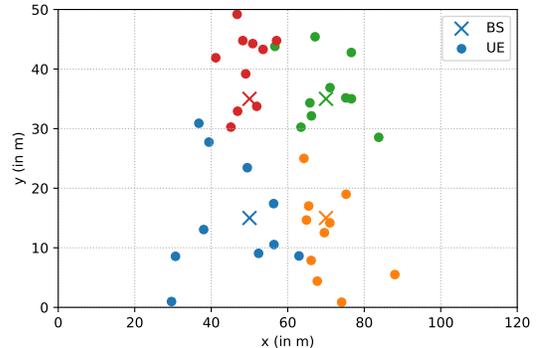
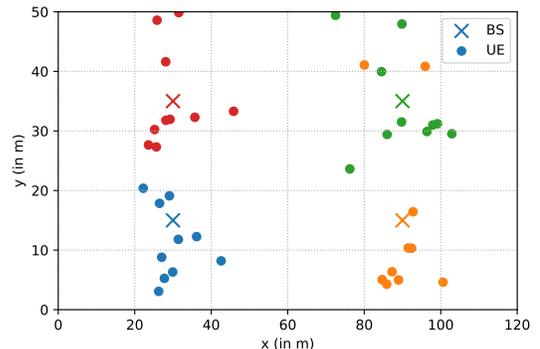
```

for iteration = 1, 2, ... do
  for actor = 1, 2, ...,  $N_{\text{batch}}$  do
    Generate an episode of  $L$  time slots as detailed
    in Section III-C.
    In each time slot, each BS  $i$  chooses to transmit
    with probability  $\pi_i^{\text{CON}}[1]$ 
  end
  for  $i = 1, 2, \dots, N$  do
    Compute  $V_i^{\text{target, EOS}}, V_i^{\text{target, CON}}$  and  $\hat{A}_i^{\text{CON}}$ 
    using (16) and (17) at each time for all actors.
    Perform 1 epoch of gradient ascent with batch
    size  $N_{\text{batch}} \times L$  on (12) to update weights of
     $\pi_i^{\text{CON}}, V_i^{\text{CON}}$  and  $V_i^{\text{EOS}}$ .
  end
end

```

each configuration. Two key observations can be made from the PF and ED baselines: firstly, as we move from L1 to L2, both baselines accumulate a larger cumulative reward. This is because the increasing separation between UEs from different BSs allows more BSs to transmit simultaneously. Secondly, a single standardized threshold of -72 dBm cannot provide the same degree of fairness in different scenarios. In fact, for L1, -72 dBm is a very pessimistic threshold that ends up primarily switching off all the BSs, hence it has not even been plotted. On the other hand, the RL PPO algorithm consistently outperforms even the adaptive ED threshold for all three layouts. More importantly, the RL PPO algorithm tends to always converge faster to the optimal solution than DQN, has a more stable training curve and even significantly outperforms DQN in some instances e.g. Layout 2.

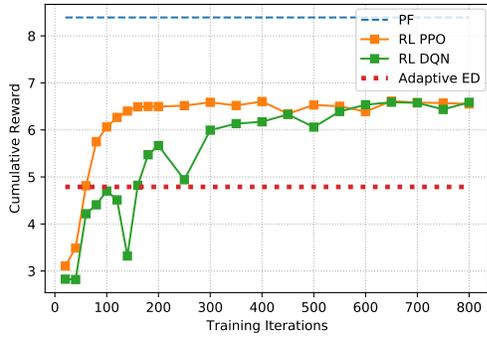
For every realization of each UE configuration, we compute the sum rate $W \sum_{j=1}^N \bar{X}_j[L]$ and max rate $W \max_j \bar{X}_j[L]$ obtained using the RL PPO algorithm at the end of L time-steps. The sum and max rate, averaged over all realizations and configurations, and evaluated using the trained model obtained after every 600 iterations, are plotted for L1 and L2 in Fig. 4a and 4b respectively, along with the corresponding PF and Adaptive ED baselines. Consistent with the higher cumulative rewards earned by the RL algorithm, RL PPO achieves a sum rate at least equal to the adaptive ED algorithm, but always

(a) Layout 1: $l = 20$ (b) Layout 2: $l = 60$ Fig. 2. Two layouts of 4 BS's at the corners of a $l \times 20$ m rectangle.

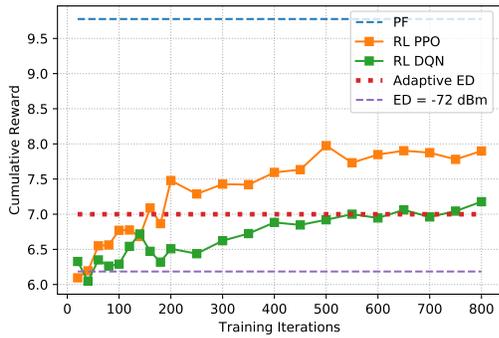
manages to increase the gap between the maximum and sum rate, hence providing a greater degree of fairness than the adaptive ED algorithm.

VI. CONCLUSIONS & FUTURE DIRECTIONS

The distributed PPO algorithm designed in this paper jointly utilized the information from LBT-based spectrum sensing at the BS along with the average rate, signal and interference power seen by the UE it serves to determine whether a BS will transmit in the designated time slot. Consequently, it was found to significantly outperform a configuration adaptive ED threshold, and also achieve improved UE throughputs. With a view to the design of a learning based BS, the framework developed in this paper has the potential to be

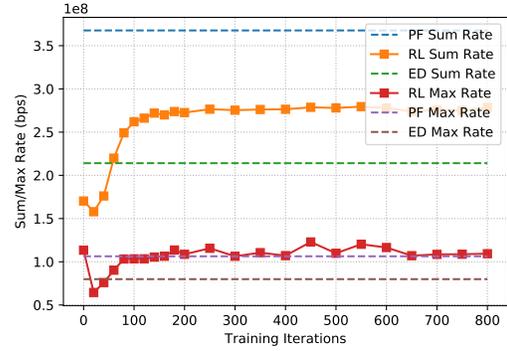


(a) Layout 1

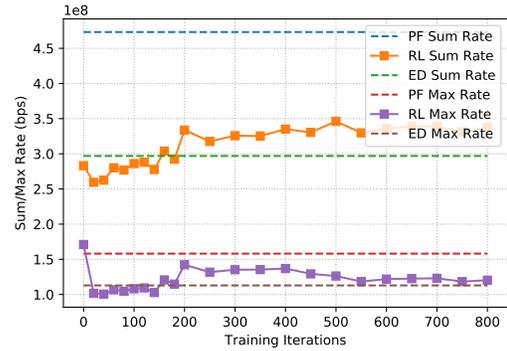


(b) Layout 2

Fig. 3. Cumulative Reward evaluated on Validation Set v/s Training Iterations for *Layout 1* and *Layout 2*



(a) Layout 1



(b) Layout 2

Fig. 4. Max and Sum UE Rate evaluated on Validation Set v/s Training Iterations for *Layout 1* and *Layout 2*

applied in a variety of decision-making problems, including adaptive modulation and coding, beam selection for scheduling, coordinated scheduling and channel selection in the frequency domain. In essence, these extensions tremendously increase the dimensionality of the output action space, from a simple transmit Yes/No decision to a choice of MCS, beamformer, user and subcarrier.

REFERENCES

- [1] 3GPP, "Feasibility Study on Licensed-Assisted Access to Unlicensed Spectrum," 3rd Generation Partnership Project (3GPP), Technical Report (TR) 36.889, Jul. 2015, version 13.0.0.
- [2] —, "Study on NR-based access to unlicensed spectrum," 3rd Generation Partnership Project (3GPP), Technical Report (TR) 38.889, Dec. 2018, version 16.0.0.
- [3] ETSI, "Broadband Radio Access Networks (BRAN); 5GHz high performance RLAN," European Telecommunications Standards Institute (ETSI), European Standards (EN) 301 893, Jul. 2014, version 1.7.2.
- [4] J. L. Sobrinho, R. De Haan, and J. M. Brazio, "Why RTS-CTS is not your ideal wireless LAN multiple access protocol," in *Proc., IEEE Wireless Networking and Comm. Conf.*, vol. 1, Mar. 2005, pp. 81–87.
- [5] J. Lundén, V. Koivunen, S. R. Kulkarni, and H. V. Poor, "Reinforcement learning based distributed multiagent sensing policy for cognitive radio networks," in *IEEE Intl. Symposium on Dynamic Spectrum Access Networks (DySPAN)*, May 2011, pp. 642–646.
- [6] O. Naparstek and K. Cohen, "Deep multi-user reinforcement learning for distributed dynamic spectrum access," *IEEE Trans. on Wireless Communications*, vol. 18, no. 1, pp. 310–323, Nov. 2018.
- [7] M. Tonnemacher, C. Tarver, J. Cavallar, and J. Camp, "Machine Learning Enhanced Channel Selection for Unlicensed LTE," in *IEEE Intl. Symposium on Dynamic Spectrum Access Networks (DySPAN)*, Nov. 2019, pp. 1–10.
- [8] N. Naderializadeh, J. Sydir, M. Simsek, and H. Nikopour, "Resource management in wireless networks via multi-agent deep reinforcement learning," *IEEE Trans. on Wireless Communications*, Jan. 2021.
- [9] A. Doshi, S. Yerramalli, L. Ferrari, T. Yoo, and J. G. Andrews, "A Deep Reinforcement Learning Framework for Contention-Based Spectrum Sharing," *IEEE Journal on Sel. Areas in Communications*, vol. 39, no. 8, pp. 2526–2540, Jun. 2021.
- [10] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, Oct. 2018.
- [11] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv:1707.06347*, Jul. 2017.
- [12] B. Wang and D. Zhao, "Scheduling for long term proportional fairness in a cognitive wireless network with spectrum underlay," *IEEE Trans. on Wireless Communications*, vol. 9, no. 3, pp. 1150–1158, Mar. 2010.
- [13] M. Hausknecht and P. Stone, "Deep Recurrent Q-learning for Partially Observable MDPs," in *AAAI Fall Symposium Series*, Sep. 2015.
- [14] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *Adv. in Neural Info. Process. Systems*, May 2016, pp. 2137–2145.
- [15] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3–4, pp. 229–256, May 1992.
- [16] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *Proc. ICLR*, Jun. 2015.
- [17] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, Jun. 2015, pp. 1889–1897.

- [18] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [19] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, Dec. 2014.
- [20] R. Lowe, Y. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Adv. in Neural Info. Process. Systems*, vol. 30, Jun. 2017.