

Scalable Joint Learning of Wireless Multiple-Access Policies and their Signaling

Mateus P. Mota^{*†}, Alvaro Valcarce^{*}, Jean-Marie Gorce[†],
^{*}Nokia Bell Labs, Nozay, France

Email: mateus.pontes_mota@nokia.com, alvaro.valcarce_rial@nokia-bell-labs.com

[†]National Institute of Applied Sciences, Lyon, France
Email: jean-marie.gorce@insa-lyon.fr

Abstract—In this paper, we apply an multi-agent reinforcement learning (MARL) framework allowing the base station (BS) and the user equipments (UEs) to jointly learn a channel access policy and its signaling in a wireless multiple access scenario. In this framework, the BS and UEs are reinforcement learning (RL) agents that need to cooperate in order to deliver data. The comparison with a contention-free and a contention-based baselines shows that our framework achieves a superior performance in terms of goodput even in high traffic situations while maintaining a low collision rate. The scalability of the proposed method is studied, since it is a major problem in MARL and this paper provides the first results in order to address it.

Index Terms—Multi-Agent Reinforcement Learning, Protocol Emergence, Wireless Communications.

I. INTRODUCTION

The goal of this paper is to explore a framework for jointly learning a channel access policy and its signaling policy for medium access control (MAC) in multiple-access scenarios. This study aims at proposing a general framework capable of producing application-tailored protocols which may lead to performance gains over more general purpose protocols.

It is expected that Artificial intelligence (AI) and machine learning (ML) will play a crucial role in 6G [1] in making the network more adaptable and self-upgradable, helping meeting the requirements while also making the network management and optimization simpler. One promising area in ML for achieving a more adaptable network system is reinforcement learning (RL). In particular, multi-agent reinforcement learning (MARL) has been used to emerge communication that allows a better cooperative behavior [2], [3]. The framework used in this paper leverages MARL to allow the network nodes to learn the channel-access policy and the communication needed to best collaborate with one another, thus also learning the signaling.

Related Work: RL has been used to develop channel access policies for the MAC in [4] and [5]. It has also been used to select which MAC protocol to use [6] or which blocks to use [7]. Differently from such works we propose to learn a channel access policy and its signaling. The idea of learning a given protocol and its signaling has already been addressed in a previous work [8], while in [9] we proposed the framework for emerging a MAC protocol in a multiple access scenario.

Contribution: This paper extends the previous one [9] in two ways:

- 1) Traffic model: By using a Poisson process, instead of limiting the total number of service data units (SDUs), making the new model more realistic. The Poisson process is used, for example, to model message arrivals in a packet data networks or the arrival of new telephone calls.
- 2) Scalability study: By evaluating the scalability both in terms of arrival rate as well as user equipments (UEs).

Since we propose to fully emerge a protocol for the base station (BS) and UEs, scalability may be an issue because the BS needs to communicate with all UEs.

This work is structured as follows. Section II describes the system model used and in Section III, we present a new framework allowing the emergence of MAC protocols with MARL. Finally, Section IV illustrates the performance of our algorithm with numerical results, where we compare the proposed solution with a baseline. The main conclusions are drawn in Section V.

II. SYSTEM MODEL

Consider a single cell with a BS serving L UEs operating according to a time division multiple access (TDMA) scheme, where each UE needs to deliver data to the BS. Each UE has a transmission buffer of capacity B MAC SDUs initially empty. The SDU arrival is modeled as a Poisson process with probability of arrival p_a . So, a new SDU is added to the buffer with probability p_a , until a maximum number T of steps is achieved. The average number of SDUs arriving at each UE's buffer in any given episode of duration T is then:

$$\lambda = p_a T \quad (1)$$

The network nodes can exchange information, using messages through the control channels. In the remainder of this paper, we refer to the UE MAC agent and the BS MAC agent as UE and BS, respectively.

The channel for the uplink data transmission is modeled as a packet erasure channel, where a transport block (TB) is incorrectly received with a probability referred to as transport block error rate (TBLER). The UEs use the same frequency resources on the uplink shared channel (UL-SCH), which leads to possible collisions. The downlink control messages (DCMs) and uplink control messages (UCMs) are transmitted over the downlink (DL) and uplink (UL) control channels, which

are assumed to be dedicated and error free, so without any contention or collision.

We assume that the sets of possible DL and UL control messages have cardinality D and U , respectively. For example, the DCMs in an DL control vocabulary of size $D = 4$ would have a bitlength Y_{DL} of $\log_2 D = 2$.

At each time step t , the BS can send one control message to each UE and each UE can send one control message to the BS while being able to send data protocol data units (PDUs) through the UL-SCH. Furthermore, the UEs can also delete a SDU from the buffer at each time step.

We define the cellwide goodput G (in SDUs/TTIs) as the number of MAC SDUs received by the BS per unit of time. SDUs received by the BS several times are only counted once:

$$G = \frac{N_{RX}}{T} \quad (2)$$

where N_{RX} represents the number of unique SDUs received. Since the BS can only receive at most one SDU per time step $N_{RX} \leq 1$, the maximum cellwide goodput on average can be calculated as:

$$G_{\max} = \min(p_a L, 1) \quad (3)$$

The collision rate Γ is the number of steps in which a collision happened divided by the total number of time steps:

$$\Gamma = \frac{N_c}{T}. \quad (4)$$

where N_c represents the total number of time steps in which at least two SDUs collided.

III. EMERGING A MAC PROTOCOL WITH MARL

A. MARL Formulation

We formulate the problem defined above as a MARL cooperative task, where the MAC layers of the network nodes (UEs and BS) are RL agents that need to learn how to communicate with each other to solve an uplink transmission task. In addition, the UE agents need to learn when to send data through the UL-SCH and when to delete an SDU, in other words, to learn how to correctly manage the buffer. In order to decide how to act, an agent needs to consider the messages received from the other agents. In addition, the UEs also take into account their buffer status when taking actions, while the BS takes into account the state of the UL-SCH, i.e idle, busy or collision-free reception.

We model this problem as a decentralized partially observable Markov decision process (Dec-POMDP) [10], augmented with communication. A Dec-POMDP for n agents is defined by the global state space \mathcal{S} , an action space $\mathcal{A}_1, \dots, \mathcal{A}_n$, and an observation space $\mathcal{O}_1, \dots, \mathcal{O}_n$ for each agent. In Dec-POMDP, an agent observation does not fully describe the environment state. All agents share the same reward and the action space of each agent is subdivided into one environment action space and a communication action space. The communication action represents the message sent by an agent and it does not affect the environment directly, but it may be passed to other agents. In this work, the agent state x_i may comprise not only

the agent's current observation, but also previous observations, actions and received messages.

We use the following notations:

- o_t^u : Observation received by the u^{th} UE at time step t .
- o_t^b : Observation received by the BS at time step t .
- n_t^u : The UCM sent from the u^{th} UE at time step t .
- m_t^u : The DCM sent to the u^{th} UE at time step t .
- a_t^u : Environment action of the u^{th} UE at time step t .
- x_t^u : Agent state of the u^{th} UE at time step t .
- x_t^b : Agent state of the BS at time step t .

Observations: The observation $o_t^u \in \{0, \dots, B\}$ is a integer representing the number of SDUs in the buffer of the UE u at that time t . Similarly, the observation o_t^b received by the BS is a discrete variable with $L + 2$ possible states:

$$o_t^b = \begin{cases} 0, & \text{if the UL-SCH is idle} \\ u, & \text{if the UL-SCH is detected busy with a} \\ & \text{single PDU from UE } u, \text{ correctly decoded} \\ L + 1, & \text{non-decodable energy in the UL-SCH} \end{cases} \quad (5)$$

where $u \in \{0, \dots, L\}$.

Actions: The environment action $a_t^u \in \{0, 1, 2\}$ is interpreted as follows:

$$a_t^u = \begin{cases} 0: & \text{do nothing} \\ 1: & \text{transmit the oldest SDU in the buffer} \\ 2: & \text{delete the oldest SDU in the buffer} \end{cases} \quad (6)$$

We highlight that the DCM and UCM messages, m and n , are communication actions that the agents select while also being information available to the other agent's state as received message.

The agent state at time step t is a tuple comprising the most recent k observations, actions and received messages:

- UE u : $x_t^u = (o_t^u, \dots, o_{t-k}^u, a_t^u, \dots, a_{t-k}^u, n_t^u, \dots, n_{t-k}^u, m_t^u, \dots, m_{t-k}^u)$
- BS: $x_t^b = (o_t^b, \dots, o_{t-k}^b, \mathbf{n}_t, \dots, \mathbf{n}_{t-k}, \mathbf{m}_t, \dots, \mathbf{m}_{t-k})$, with \mathbf{n} and \mathbf{m} containing the messages from all the UEs.

We assume the episode ends when a maximum number of steps T is reached. The reward given at each time step is:

$$r_t = \begin{cases} +\rho, & \text{if a new SDU was received by the BS} \\ -\rho, & \text{if an UE deleted a SDU that has} \\ & \text{not been received by the BS} \\ 0, & \text{else,} \end{cases} \quad (7)$$

where ρ is a positive integer. This choice of reward is possible by leveraging the centralized training and decentralized execution (CTDE). During the centralized training, a centralized reward system can be used to observe the buffers of the BS and UEs in order to assign the reward. For wireless systems, centralized training can be achieved in a simulation environment as well as a testbed, i.e. a server farm.

B. Training Algorithm

The proposed RL solution is based on the multi-agent deep deterministic policy gradient (MADDPG) algorithm [11]. This

algorithm is well suited to partially observable environments when strong coordination is needed, due to its centralized critic architecture.

In MADDPG, each agent has an actor network that depends only on its own agent's state in order to learn a decentralized policy μ_i with parameters θ_i . Each agent also has a centralized critic network that receives the agent states and actions of all agents in order to learn a joint action value function $Q_i(x, a)$ with parameters φ_i , where $x = (x_1, x_2, \dots, x_n)$ is a vector containing all the agents' states and $a = (a_1, a_2, \dots, a_n)$ contains the actions taken by all of the agents. The critic networks are only used during the centralized training.

The critic network parameters φ are updated by minimizing the loss given by the temporal-difference error

$$L^i := \mathbb{E}_{x, a, r, x' \sim \mathcal{D}} \left[(y^i - Q_i(x, a_1, \dots, a_n; \varphi_i))^2 \right] \quad (8)$$

where \mathcal{D} denotes the experience replay buffer in which the transition tuples (x, a, r, x') are stored, Q' and μ' represent the target critic network and the value of the target actor network, with parameters θ' and φ' , respectively, and y^i is the temporal-difference target, given by

$$y^i := r + \gamma Q'_i(x', a'_1, \dots, a'_n; \varphi'_i) \Big|_{a'_k = \mu'_k(x_k)} \quad (9)$$

where γ is the discount factor. The actor network parameters θ are updated using the sampled policy gradient

$$\nabla_{\theta_i} J = \mathbb{E}_{x, a \sim \mathcal{D}} \left[\nabla_{a_i} Q_i(x, a) \nabla_{\theta_i} \mu_i(x_i) \mid a_i = \mu_i(x_i) \right]. \quad (10)$$

The target networks parameters are updated as

$$\varphi'_i \leftarrow \tau \varphi_i + (1 - \tau) \varphi'_i \quad (11)$$

$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i \quad (12)$$

where $\tau \in [0, 1]$ is the soft-update parameter.

Architecture: The actor and critic networks have the same architecture; a fully connected multilayer perceptron (MLP) with two hidden layers, of 64 neurons each. The activation function of all hidden layers is the rectified linear unit (ReLU). In order to improve training of our MADDPG solution, we make use of parameter sharing [2] for similar network nodes, in this case the UEs. Since UE index is not included in the agent's state, any policy that uses the agent's identity is not capable of effectively solving the task due to the parameter sharing, because it would lead to collisions.

Similarly to the original work [11], we use the Gumbel-softmax [12] trick to soft-approximate the discrete actions to continuous ones. The Gumbel-softmax reparameterization also works to balance exploration and exploitation. The exploration-exploitation trade-off is controlled by the temperature factor ζ .

After training finishes, we have successfully trained a population of N_{rep} protocols. We then select the protocol that performed the best at any point during training across all different protocols, i.e the historically best protocol. This selection step can be seen as a "survival of the fittest" approach because only one protocol of the population of N_{rep} is chosen going forward.

TABLE I
SIMULATION PARAMETERS

Parameter	Symbol	Value
Number of UEs	L	[2, 3, 4, 5]
Size of transmission buffer	B	20
Avg. number of SDUs per UE	λ	[2, 4, 6, 8, 10, 12]
SDU arrival probability	p_a	[0.083, 0.16, 0.25, 0.33, 0.41, 0.5]
Transport block error rate	TBLER	10^{-1}
DCM vocabulary size	D	3
UCM vocabulary size	U	2
Duration of episode (TTIs)	T	24
Reward function parameter	ρ	3
Number of training episodes	N_{train}	100k
Number of evaluation episodes	N_{eval}	500
Number of test episodes	N_{test}	5000
Number of randomized repetitions	N_{rep}	8

TABLE II
TRAINING ALGORITHM PARAMETERS

Parameter	Symbol	Value
Memory length	k	3
Replay buffer size		10^5
Batch size		1024
Number of neurons per hidden layer		{64, 64}
Interval between updating policies		96
Optimizer algorithm		Adam
Learning rate	α	10^{-3}
Discount factor	γ	0.9
Policy regularizing factor		10^{-3}
Gumbel-softmax temperature factor	ζ	1
Target networks soft-update factor	τ	10^{-3}

IV. RESULTS

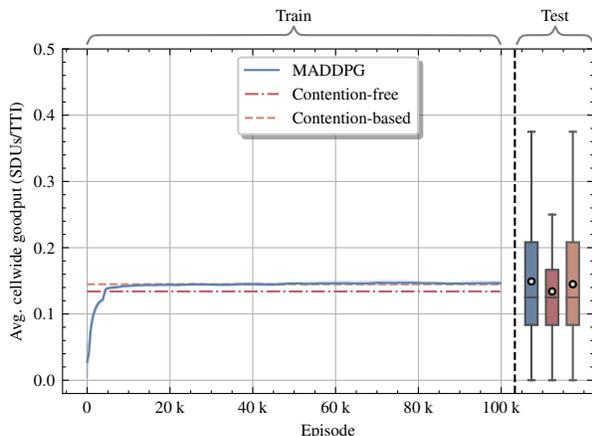
A. Simulation Procedure and Parameters

The transmission buffer of each user starts empty and the SDU arrival probability is p_a for each UE. The system is trained for a fixed number of episodes N_{train} . During training, we evaluate the policy on a fixed set of N_{eval} evaluation episodes with disabled exploration and disabled learning in order to assess the current performance of the communication protocol. The protocol that performed the best on the evaluation episodes during the whole training procedure is selected and its performance is assessed in N_{test} episodes with exploration and learning disabled. This whole procedure represents a single training repetition. We evaluate a total of N_{rep} repetitions, each with a different random seed. A summary of the main simulation parameters is provided in Table I, while the parameters of the MADDPG and deep deterministic policy gradient (DDPG) algorithms are listed in Table II.

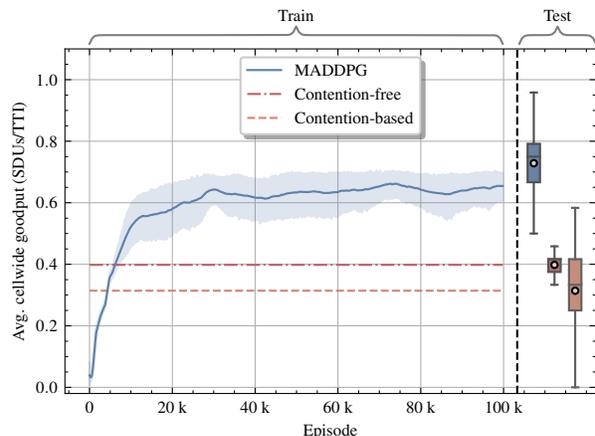
B. Baseline Solutions

We compare the proposed solution with a contention-free (i.e. BS-controlled, scheduled) and a contention-based (i.e. grant-free) baseline.

In the contention-free protocol, the UE sends a scheduling request (SR) if its transmission buffer is not empty and it only transmits if it has received a scheduling grant (SG). Similarly,



(a) Avg. Number of SDUs per UE: $\lambda = 2$. Arrival rate: $p_a = 0.83$.



(b) Avg. Number of SDUs per UE: $\lambda = 12$; Arrival rate: $p_a = 0.5$.

Fig. 1. Goodput comparison during the training procedure. Number of UEs: $L = 2$; TBLER = 10^{-1} .

it only deletes a TB from the transmission buffer after the reception of an acknowledgement (ACK). At each time step, the BS receives zero or more SRs. It then chooses one of the requesters at random to transmit in the next time-step, sending a SG to the selected UE. However, if the UE had made a successful data transmission simultaneously with an SR, the BS will send an ACK to this UE and its SR is ignored.

In the contention-based protocol, each UE transmits with probability p_t if its transmission buffer is not empty. Similarly to the contention-free baseline, the UE only deletes a TB after the reception of an ACK. At each time step, the BS sends an ACK to a UE if it received a TB from the UE. For each experiment, the transmission probability chosen is the one that performs better in terms of goodput.

C. Results

1) *Learning Performance*: We first analyze the performance over the training procedure, comparing the proposed solution with the baselines in Fig. 1. The solid lines in Figs. 1a and 1b show the average performance in the evaluation episodes during the training and the shaded areas represent the 95% confidence interval (CI). After assessing the performance on the last N_{eval} evaluation episodes, we select the best performing repetitions for each solution in terms of average goodput to compare using boxplots of the test episodes.

The main conclusions we can draw from Fig. 1 are:

- In the lower arrival rate showed in Fig. 1a, the proposed solution seems to learn a protocol that performs like a contention-based one. This conclusion is supported by the similar box plots on the test episodes.
- In higher arrival rates showed in Fig. 1b, the proposed solution drastically outperforms both baselines, which indicates it learns a completely different protocol.
- The contention-free baseline shows a better performance on low arrival rates, but when the arrival probability increases the contention-free baseline outperforms it.

2) *Scalability*: In this set of results, we analyze the scaling capabilities of the proposed solution across two dimensions,

the number of UEs and the SDU arrival rate. The performance is evaluated on N_{test} test episodes by comparing the average goodput and collision rate achieved when changing across one dimension while the other is fixed. For the MADDPG solution, we also show the 95% CI across randomized repetitions. The upper bound shows the maximum average goodput when all the SDUs are received.

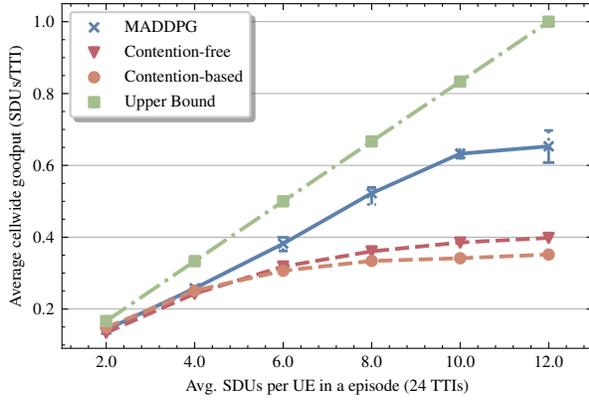
The proposed framework is capable of producing protocols that outperform both baselines in terms of goodput when the arrival rate increases while maintaining a low collision rate, as shown in Fig. 2. Also, the CI increases when the arrival rate increasing, indicating that in more difficult conditions there's a bigger variability in the emerged protocol.

Scalability to growing numbers of UEs is proving challenging, as shown in Fig. 3. The proposed framework consistently outperforms both solutions for up to four UEs and has similar performance to the contention-free solution on average for five UEs, but it is unable to scale as well as it does when scaling with traffic. Although the proposed solution achieves a lower collision rate than the contention-based solution, it seems unable to effectively deal with more UEs while avoiding collisions, which can explain why the cellwide goodput drops when increasing the number of UEs.

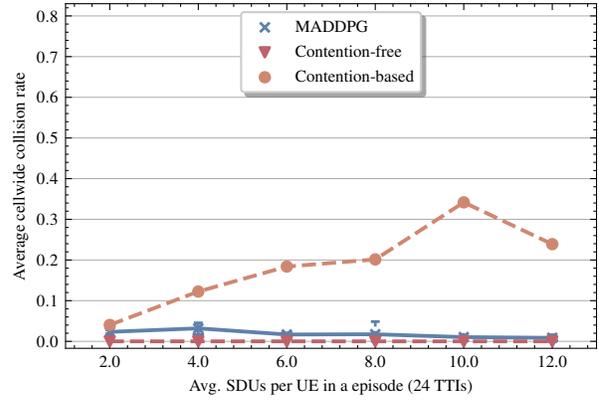
V. CONCLUSIONS AND PERSPECTIVES

We have applied a framework to emerge a MAC protocol and have demonstrated through simulations that cooperative MARL augmented with communication provides an original approach to emerge a protocol by jointly learning the channel access policy and its signaling. The results indicate the capabilities of the MARL framework to produce protocols that outperform the baselines. In addition, the results illustrate the capabilities of the framework to adapt to different arrival rates and to different number of UEs.

In our future works we will propose extensions to deal with even more UEs. We will also propose comparisons of different MARL algorithms, and we will evaluate accurately the impact of the vocabulary sizes used in the control channels.

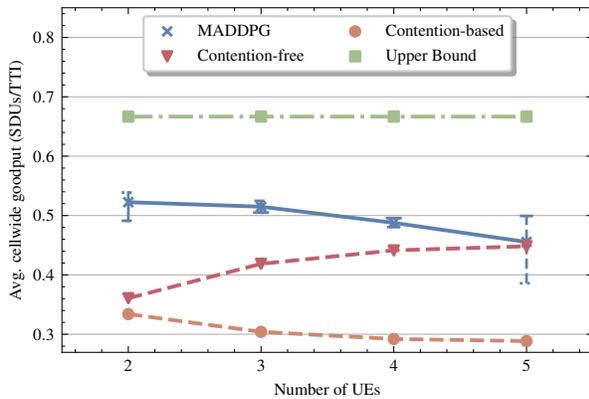


(a) Goodput when increasing p_a .

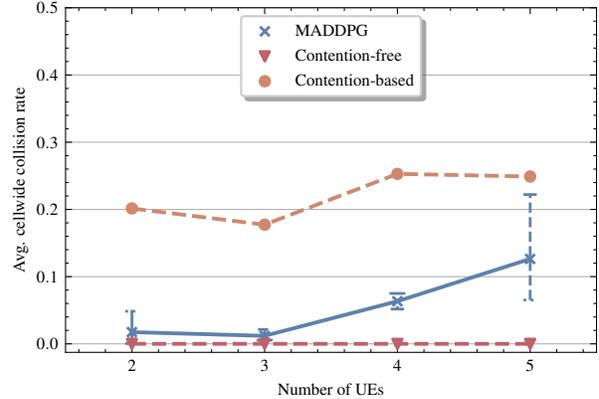


(b) Collision rate when increasing p_a .

Fig. 2. Scaling the arrival rate while maintaining the number of UEs fixed: $L = 2$ UEs.



(a) Goodput when increasing the number of UEs.



(b) Collision rate per number of UEs.

Fig. 3. Scaling the number of UEs while maintaining the average number of SDUs in the cell per episode fixed: $L\lambda = 16$.

Additionally, interpretability will be investigated to better understand the key for improvements used by our RL based algorithms.

ACKNOWLEDGMENT

The work of Mateus P. Mota is funded by Marie Skłodowska-Curie actions (MSCA-ITN-ETN 813999 WIND-MILL).

REFERENCES

- [1] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, "The roadmap to 6G: AI empowered wireless networks," *IEEE Communications Magazine*, vol. 57, no. 8, pp. 84–90, 2019.
- [2] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, "Learning to communicate with Deep multi-agent reinforcement learning," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 2145–2153.
- [3] S. Sukhbaatar, A. Szlam, and R. Fergus, "Learning multiagent communication with backpropagation," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 2252–2260.
- [4] H. Dutta and S. Biswas, "Towards multi-agent reinforcement learning for wireless network protocol synthesis," in *2021 International Conference on Communication Systems NETWORKS (COMSNETS)*, 2021, pp. 614–622.

- [5] Z. Guo, Z. Chen, P. Liu, J. Luo, X. Yang, and X. Sun, "Multi-agent reinforcement learning based distributed channel access for next generation wireless networks," *IEEE Journal on Selected Areas in Communications*, pp. 1–1, 2022.
- [6] A. Gomes, D. F. Macedo, and L. F. Vieira, "Automatic mac protocol selection in wireless networks based on reinforcement learning," *Computer Communications*, vol. 149, pp. 312–323, 2020.
- [7] H. B. Pasandi and T. Nadeem, "Towards a learning-based framework for self-driving design of networking protocols," *IEEE Access*, vol. 9, pp. 34 829–34 844, 2021.
- [8] A. Valcarce and J. Hoydis, "Towards joint learning of optimal MAC signaling and wireless channel access," *IEEE Transactions on Cognitive Communications and Networking*, pp. 1–1, 2021.
- [9] M. P. Mota, A. Valcarce, J.-M. Gorce, and J. Hoydis, "The emergence of wireless mac protocols with multi-agent reinforcement learning," in *2021 IEEE Globecom Workshops (GC Wkshps)*. IEEE, 2021, pp. 1–6.
- [10] F. A. Oliehoek, M. T. Spaan, and N. Vlassis, "Optimal and approximate q-value functions for decentralized pomdps," *Journal of Artificial Intelligence Research*, vol. 32, pp. 289–353, 2008.
- [11] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Advances in neural information processing systems*, 2017, pp. 6379–6390.
- [12] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with Gumbel-Softmax," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=rkE3y85ee>