

# Multi-Access Point Coordination for Next-Gen Wi-Fi Networks Aided by Deep Reinforcement Learning

Lyutianyang Zhang, Hao Yin, Sumit Roy, *Fellow, IEEE*, Liu Cao

**Abstract**—Wi-Fi in the enterprise - characterized by overlapping Wi-Fi cells - constitutes the design challenge for next-generation networks. Standardization for recently started IEEE 802.11be (Wi-Fi 7) Working Groups has focused on significant medium access control layer changes that emphasize the role of the access point (AP) in radio resource management (RRM) for coordinating channel access due to the high collision probability with the distributed coordination function (DCF), especially in dense overlapping Wi-Fi networks. This paper proposes a novel multi-AP coordination system architecture aided by a centralized AP controller (APC). Meanwhile, a deep reinforcement learning channel access (DLCA) protocol is developed to replace the binary exponential backoff mechanism in DCF to enhance the network throughput by enabling the coordination of APs. First-Order Model-Agnostic Meta-Learning further enhances the network throughput. Subsequently, we also put forward a new greedy algorithm to maintain proportional fairness (PF) among multiple APs. Via the simulation, the performance of DLCA protocol in dense overlapping Wi-Fi networks is verified to have strong stability and outperform baselines such as Shared Transmission Opportunity (SH-TXOP) and Request-to-Send/Clear-to-Send (RTS/CTS) in terms of the network throughput by 10% and 3% as well as the network utility considering proportional fairness by 28.3% and 13.8%, respectively.

**Index Terms**—Wi-Fi 7, IEEE 802.11be, multi-AP coordination, channel access, proportional fairness, deep Q-learning

## I. INTRODUCTION

The rapid adoption of smartphones, tablets, and high-end mobile client devices has translated into the rapid growth of *network traffic flux* (measured in bits/s/Hz per unit area/volume). As tracked by Cisco Annual Internet Report [2], the number of Wi-Fi hotspots will grow four-fold, and the average mobile network connection speeds will triple from 2018 to 2023.

Multi-media streaming demands, e.g., 4K and 8K video [3], will stretch network capacity even beyond current generation (Wi-Fi 6) limits of 10 Gbps peak capacity. To address such technology bottlenecks, the IEEE 802.11 Working Group (WG) is standardizing the next generation of Wi-Fi—referred to as IEEE 802.11be (Wi-Fi 7) or Extremely High Throughput (EHT) networks.

Orthogonal frequency-division multiple access (OFDMA) adopted in 802.11ax has significantly enhanced the medium

This paper has been presented in part at the IEEE Vehicular Technology Conference (VTC) 2020-Fall [1].

Lyutianyang Zhang, Hao Yin, Sumit Roy, and Liu Cao are with Department of Electrical & Computer Engineering, University of Washington, Seattle, WA, USA (e-mail: {lyutiz, haoyin, sroy, liuca}@uw.edu). (*Corresponding author: Hao Yin*)

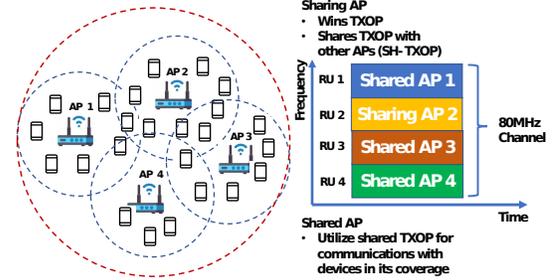


Fig. 1: Multi-AP Coordination with SH-TXOP: 4 APs share a 80 MHz universal channel. Each AP has its own coverage of stations (STAs) and operates on its own primary channel. Each resource unit (RU) represents a 20 MHz channel.

access control (MAC). OFDMA works on top of the legacy carrier-sense multiple access with collision avoidance (CSMA/CA) to provide extra features when access point (AP) contends for channel access, e.g., the use of trigger frame helps control the uplink transmission of stations [4] for greater efficiency in a single basic service set identifier (BSS-ID). The overlap of APs is defined as the intersection of their cells and operating frequency bands. In such a cluster of overlapping AP, channel access collision happens more frequently as the number of APs increases, especially when each AP has no prior knowledge of other APs' channel accessing policies. Thus, enabling some degree of collaboration among neighboring APs will permit more efficient utilization of the limited time and frequency resources, i.e., lower collision probability, higher network throughput. To this end, the next-generation standard 802.11be (Wi-Fi 7) introduces some additional features such as multi-AP coordination to further improve aggregate throughput in dense overlapping layout scenarios [5].

The emphasis of EHT WG is on *aggregate throughput in dense networking scenarios* and hence - building on the numerous physical layer (PHY) advances made in 802.11ac/ax - notably new control frames for coordination among APs that requires information exchange among the APs belonging to the coordinated AP set. Besides, the new 6 GHz bands opened up in the US and Europe are new green fields for the future Wi-Fi 7 standard, which gives more freedom in architecture and protocol design. In Fig.1, an example of the current 11be

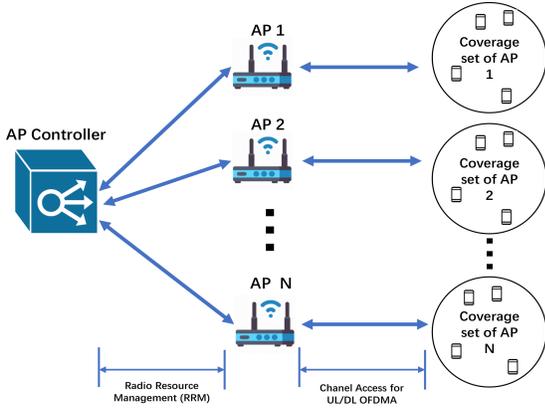


Fig. 2: Proposed Architecture for Multi-AP Coordination.

structure with Shared Transmission Opportunity (SH-TXOP) operation is introduced [6]. Four APs simultaneously operate on the shared 80-MHz bandwidth. Distributed coordination function (DCF) applied in SH-TXOP allows all APs to contend for channel access. In this example, AP 2 successfully gains the TXOP and becomes a sharing AP. AP 2 then collects information about the channel status and the traffic backlog from shared APs in the candidate shared AP set. Afterwards, the sharing AP will share the wide-band TXOP with the shared APs. Each AP must have a primary channel to operate in the dense overlapping network so that it can contend for TXOPs for the communication within its own coverage. Note that primary channels allocated to different APs are not necessarily the same. In such a scheme considering collaboration, there is no frequency overlap among four channels because each AP operates within its own allocated channel independently. The bandwidth of each channel occupied by a shared AP or sharing AP is 20 MHz.

#### A. Multi-AP Coordination Architecture and Related Work

DCF with CSMA/CA is a traditional MAC protocol for channel access in Wi-Fi networks. It has a long history of analysis using Markov models [7], [8] developed originally for a single (isolated) cell (e.g., single home networks) with saturated nodes.

In [9], a novel multi-AP coordination transmission scheme is proposed for 802.11be. TXOP in the proposed scheme is acquired after the AP completes its backoff procedure. Then, it performs a wide-band transmission if the secondary channel is idle during the point coordination function (PCF) inter-frame spacing (PIFS). This contention-based method for TXOP is similar to DCF Request-to-Send/Clear-to-Send (RTS/CTS) [8]. Then, the sharing AP sends an Announcement Trigger Frame (ATF) to the shared APs to allocate all 20 MHz channels, including the duration of TXOP and the scheduled channel information. This proposed scheme also aligns with Coordinated OFDMA (C-OFDMA). That is, each shared AP utilizes partial TXOP assigned by the sharing AP for its uplink/downlink (UL/DL) OFDMA transmission with its associated STAs.

As the backoff procedure for granting the wide-band TXOP is performed solely by the initiating AP, all responding APs

should terminate the entire transmission sequence before the TXOP duration, indicated by the received ATF. However, this leads to profound performance loss as the number of APs grows larger. In [9], only 2 APs operate on a 40 MHz channel. The backoff procedure not followed by ACK can lead to the following situation: more than one AP may think it has won the TXOP and start sending ATF, which will lead to a collision. Since there is no feedback such as ACK, in the end, the whole TXOP is wasted because more than one AP send different RU assignment in ATF and any AP will obtain confusing allocation scheme. As a result, the collision probability increases significantly with the increasing number of APs. Moreover, it is challenging to design a feedback mechanism such as ACK for the TXOP contention method. Responding to the sharing AP by all shared APs is a huge burden to the system performance because any failure reception of the ATF to any AP will lead to re-transmission.

A novel system architecture is thereby proposed for multi-AP coordination in 802.11be to decrease the collision probability of channel access, as is shown in Fig.2. The centralized AP controller (APC) implements channel configuration, i.e., assigning primary channels to all APs with consideration of proportional fairness (PF). Under this system architecture, APs do not need to contend for the wide-band TXOP, and ACK is much easier to design. The core function provided by APC is called radio resource management (RRM) [10], which automatically monitors traffic, capacity, and reliability of operating APs. RRM can periodically reconfigure the 802.11 networks for best efficiency by performing functions such as radio resource monitoring and dynamic channel assignment. Each AP contends for TXOP on the assigned primary channel which is further utilized for the communications with its associated STAs by UL/DL OFDMA.

The suffering from the traditional DCF characterized by collision probability in dense overlapping networks also prompts the applications of state-of-the-art machine learning techniques. Deep reinforcement learning (DRL) is a machine learning technique that enables an agent to take actions in an environment aiming to maximize the cumulative rewards. Reinforcement learning (RL) principles have shown potential on optimizing resource allocation in various aspects in wireless communication, see [11]–[14]. Nonetheless, the application of DRL in wireless networks must be made wisely, as the network utility unavoidably oscillates due to the unstable nature of RL [15]. Moreover, in the overlapping multi-AP network, the network throughput is affected by the different channel qualities, the number of users, etc. Therefore, the designed AP coordination algorithm should also account for those factors. Adversarial RL-based method [16] is proposed as the solution to single-band multi-AP coordination in 11be. Deep Q-network (DQN) is investigated as an enhancement (higher utility) for CSMA in heterogeneous networks in which more than one multiple access protocol coexist [17], [18]. However, this work assumes that all stations run on the same frequency band. As a more ambitious study, multiple-agent deep learning multiple access with imperfect transmission feedback [19] is studied. The main objective in [19] is to recover the lost transmission feedback between AP and STA due to imperfect

channel condition. In [20], multi-channel access for multiple STAs and one AP based on deep reinforcement learning is investigated. This paper concentrates on that all mobile users utilize improved channel access approaches to communicate with one AP. Our work differs from [19], [20] by considering the throughput maximization problem for multiple overlapping and collaborating APs following 11be standards in which AP is in charge of channel access for TXOP. In [1], we consider a single channel access problem for the communication between only one AP and multiple STAs. As an extended version to [1], we propose in this paper a novel system architecture that consists a centralized APC with a PF solution to multi-channel allocation problem.

By contrast, our work also differs from the above research work by investigating a new multi-band multi-AP coordination network with APC based on IEEE 802.11be, and our proposed protocol simultaneously enables TXOP contentions at different frequency bands.

### B. Contribution

This paper proposes a novel coordinated multi-AP architecture and a corresponding channel access mechanism aligning with IEEE 802.11be to maximize the aggregate network throughput while preserving fairness among APs. The major contributions of this paper are listed as below:

- We propose a multi-AP system with APC as well as formulate a dynamic resource allocation and channel access optimization problem. The resource allocation process is considered as a Markov decision process (MDP). We choose the previous observation of channels and actions as the state, transmission at a channel as the action, and successful/unsuccessful transmissions at multiple channels as positive/negative rewards.
- Deep reinforcement learning channel access (DLCA) protocol is proposed. For each AP in the coordinated multi-AP set, DLCA is deployed to contend for channel access. The first AP winning the contention gains the TXOP on its primary channel. The First-Order Model-Agnostic Meta-Learning (FOMAML) is then applied to DLCA to enhance the overall performance. We also develop a greedy algorithm to maintain PF among APs.
- Simulation results show that the performance of DLCA protocol is verified to have strong stability and outperform baselines such as SH-TXOP and RTS/CTS in terms of the network throughput as well as the network utility in dense overlapping Wi-Fi networks.

The paper is organized as follows: Sec II introduces the proposed system model aligning with IEEE 802.11be protocol. Next, in Sec III, DLCA plus greedy algorithm with FOMAML are combined and developed as the DLCA protocol for this system. Sec. IV contains a suite of performance evaluations for the proposed DLCA protocol. The comparison between DLCA protocol and baselines is implemented to demonstrate the robustness and efficiency of our proposed DLCA protocol.

TABLE I: Main Acronyms

TXOP	Transmission Opportunity
SH-TXOP	Shared Transmission Opportunity
STA	Station
RU	Resource Unit
DIFS	Distributed Inter-Frame Space
SIFS	Short Inter-Frame Space
RTS/CTS	Request-to-Send/Clear-to-Send

## II. SYSTEM MODEL AND PROBLEM FORMULATION OF DLCA

In this section, we firstly introduce the multi-AP network with APC. Then the DCF RTS/CTS is introduced and the novel DLCA is proposed. We formulate each AP's dynamic resource allocation and channel access optimization problem as MDP.

### A. Proposed Multi-AP network

In 802.11be system, the aggregation of 5 and 6 GHz spectrum allows simultaneous operation on different bands or channels (orthogonal frequency resource allocation). Our proposed network model aligning with 11be protocol is shown in Fig.2. The coordinated multi-AP set is defined as  $\mathcal{N} = \{1, \dots, n, \dots, N\}$ , and  $\mathcal{F} = \{1, \dots, f, \dots, F\}$  denotes the available orthogonal channel set. Each AP can be allocated with a different channel from the available channel sets. Suppose, at the  $t^{th}$  contention for TXOP, AP  $n$  observes the channel states of its allocated primary channel (the  $f^{th}$  channel of 20 MHz), which yields the observation vector  $o_t^n(f) \in \{0, 1\}$  where 0 and 1 denote the IDLE and BUSY channel state, respectively. Action vector is denoted as  $a_t^n(f) \in \mathcal{A} \triangleq \{0, 1\}$  where  $a_t^n(f) = 1$  represents AP  $n$  contends for the  $f^{th}$  channel at  $t^{th}$  contention for TXOP and  $a_t^n(f) = 0$  represents AP  $n$  does not contend for the  $f^{th}$  channel at  $t^{th}$  contention for TXOP. A successful transmission occurs if a sole AP occupies a TXOP. In the following section, the action and observation vector are simplified as  $a_t^n$  and  $o_t^n$ , respectively, because the  $f^{th}$  channel is implicitly linked to AP  $n$  after APC has made the channel allocation decision.

### B. Channel Access Mechanism

This section briefly introduces three packet mode channel access protocols that can be potentially utilized in our proposed multi-AP network with APC, including our proposed DLCA protocol. They are described as follows:

- *Distributed coordination function (DCF) basic* [8]: Suppose an AP wants to occupy the TXOP on its primary channel. It waits until the channel is sensed idle for a distributed inter-frame space (DIFS). Then, a backoff process is initiated. Backoff intervals are slotted, and the discrete backoff time is uniformly distributed in the range  $[0, W - 1]$ , where  $W$  is defined as the contention window size, and  $CW_{min}$  represents the minimum contention window. The backoff counter is utilized for AP to decide whether to access channel at the current time slot. The backoff counter value is initialized by uniformly choosing an integer from the range  $[0, W - 1]$ . Then,

it is decremented by one at the end of each idle slot. Note that the backoff counter will be frozen when a packet transmission is detected on the channel and will be reactivated until the channel is sensed idle again for a DIFS period. The AP contends for TXOP when its backoff counter reaches zero. The ACK follows after the completion of the TXOP unless collision happens. If unsuccessful TXOP happens, contention window size  $W$  is doubled after each unsuccessful transmission, up to a maximum value  $CW_{max} = 2^m CW_{min}$ , where  $m$  represents the largest times the contention window size can be doubled.

- *DCF Request-to-Send/Clear-to-Send (RTS/CTS)*: AP transmits a short frame of RTS to APC after the backoff counter is decremented to zero. When APC detects an RTS frame, it responds, after a short inter-frame space (SIFS), with a CTS frame. The AP can only occupy the TXOP of its primary channel if the CTS frame is correctly received. The RTS and CTS frames also carry the information of the TXOP duration to be transmitted. This information can be heard by any listening AP, which can then update a network allocation vector (NAV) containing the information that the duration of the channel being busy. Therefore, an AP can suitably delay further transmission by detecting just one frame among the RTS and CTS frames and thus avoid collision. The major difference between basic and RTS/CTS is that RTS/CTS will send a RTS and decide to contend for the TXOP only after a CTS is received from the APC. DCF basic, on the other hand, contends for TXOP without sending a RTS. Hence, when collision happens, DCF basic wastes a whole TXOP duration while DCF RTS/CTS only wastes a RTS/CTS duration. It is noteworthy that TXOP can take up to 8.16 ms and RTS/CTS only take up to 0.4 ms. Hence, DCF RTS/CTS mechanism is very effective in terms of network throughput, especially for large data load in TXOP, as it reduces the average number of wasted time slots involved in the contention process [8].
- *Deep reinforcement learning channel access (DLCA)*: In the DLCA protocol, APC periodically assigns the primary channel to each AP considering PF (as formulated in section III-C). Then, each AP senses the channel and obtains an observation from its primary channel environment, indicating the channel is BUSY or IDLE. Based on the observed results, each AP implements inference regarding the next action utilizing its trained deep Q learning model to maximize the network throughput in its coverage (as formulated in section III-A). It is noteworthy that once the AP employing DLCA decides to contend for TXOP, the similar protocol to DCF RTS/CTS in Fig.3 is followed. The only difference is that DLCA has no backoff time, and it transmits RTS as long as it determines to contend for TXOP, which makes the Coordinated-OFDMA [3] possible because each contention process has a constant duration and can be visualized as a time slot (summation of RTS/CTS, TXOP length, and ACK). For each time slot, either one of APs wins the TXOP, or all APs that send the RTS do not receive CTS from APC, leading to

Time\_Out which indicates that RTS was not approved by APC.

In Fig.3, we can have either DCF RTS/CTS or DLCA as the packet mode channel access method. If the DCF basic method is utilized, the transmission result will only be known to the AP until the TXOP duration is finished. The data load within the TXOP is much larger than the traditional scenario in which the DCF basic method is applied, leading to intolerable performance loss [8].

### C. Maximum Achievable Data Rate

Each AP's action  $a_t^n$  is directly related to the throughput of its coverage. The more TXOPs each AP gains, the higher throughput is reached. However, the number of successful TXOP contention is not the only factor to the throughput. The channel conditions between each AP and its associated STAs on different frequency bands are different, and they also vary over time. In this paper, we consider that the overall spectral efficiency of AP on its primary channel can be obtained by taking the average of the individual spectral efficiencies between AP  $n$  and all the associated STAs. The channel spectral efficiency is assumed to be known by APC and denoted as  $C_t^n \in \mathcal{R}^F$  (bit/s/Hz). The spectral efficiency  $C_t^n(f)$  represents the maximum data that AP  $n$  can achieve at time slot  $t$  on channel  $f$ .

### D. Access Point Model

Define a Markov decision process (MDP) for an AP over a finite state space  $\mathcal{S} \in \{0, 1\}^{2L} \times \mathcal{Z}$ , where  $L$  denotes the state size. The finite state space  $\mathcal{S}$  is a set that contains concatenations of observation vectors, action vectors, and  $c_t^n \in \mathcal{Z}$  that represents the total number of APs contending for TXOPs in the AP  $n$ 's operating frequency channel (including AP  $n$  itself), i.e.,  $\mathbf{s}_{t+1}^n = [a_{t-L+1}^n, o_{t-L+2}^n, \dots, a_t^n, o_{t+1}^n, c_t^n]$ . The transition function  $\delta(\mathbf{s}_t^n, \mathbf{s}_{t+1}^n; a_t^n)$  denotes the probability that the state  $\mathbf{s}_t^n$  transfers to the state  $\mathbf{s}_{t+1}^n$  after taking action  $a_t^n$ .  $r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) \in \mathcal{R}$  denotes the reward of AP  $n$  at  $t^{\text{th}}$  TXOP contention results from its state-action-state pair  $(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n)$ . The accumulated discounted reward  $R_t^n \in \mathcal{R}$  for AP  $n$  can be expressed as

$$R_t^n = \sum_{k=0}^{\infty} \gamma^k r(\mathbf{s}_{t+k}^n, a_{t+k}^n, \mathbf{s}_{t+k+1}^n), \quad (1)$$

where  $\gamma \in (0, 1]$  is a discounting factor. The policy of AP  $n$   $\pi(n) : \mathcal{S} \rightarrow \mathcal{A}$  is assumed to be stationary, and the decision of the policy only depends on the current state. Hence, each AP aims to solve the following problem:

$$\operatorname{argmax}_{\pi(n)} \mathbf{E}_{\delta} [R_t | \mathbf{s}_t^n = \mathbf{s}, a_t^n = a, \pi(n)], \quad (2)$$

which is the objective of each AP and the expectation with respect to the transition probability function  $\delta$  is denoted as  $E_{\delta}[\bullet]$ . Since each AP simultaneously takes actions on its primary channel where each action is associated with an objective (maximization of the accumulated reward), this is overall a multi-agent problem. The calculation of the reward corresponding to various state-action pairs is detailed in the next section.

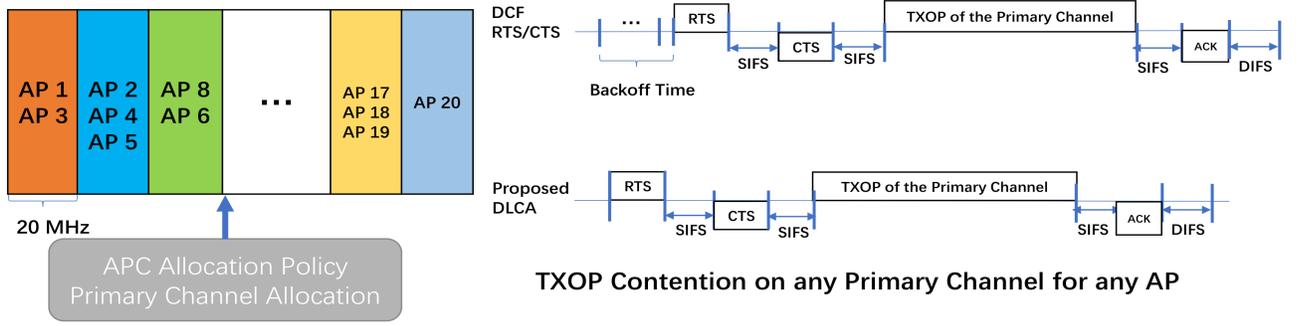


Fig. 3: The primary channels of APs are assigned by the APC.  $F$  channels from 5 and 6 GHz bands are available for IEEE 802.11be. Then APs in each 20 MHz band contends for TXOP.

### E. System Reward

Every action made by an AP has corresponding feedback (CTS/Time\_Out). The system reward is calculated as follows:

$$r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) = \sum_{l=0}^{L-1} \eta^l y_{t-l}^n, \quad (3)$$

where  $y_{t-l}^n = 1$  denotes the successful feedback for the action  $a_{t-l}^n$  and  $y_{t-l}^n = -1$  for an unsuccessful contention, and  $\eta \in [0, 1]$  is a factor such that the more recent action is, the more weight it will have in the system reward. The overall reward estimation algorithm is shown in Appendix A.

## III. DLCA PROTOCOL

In this section, we develop the following steps: a) Q-learning satisfying the Bellman optimality condition is introduced; b) we introduce deep reinforcement learning in which a deep Q network is utilized as a model for the action-value function; c) First-Order Model-Agnostic Meta-Learning (FOMAML) is applied to enhance the convergence rate and the stability of the deep Q network; d) The greedy algorithm considering PF is proposed for multi-AP coordination.

### A. Q-learning

In this section, we introduce standard Q-learning and  $\epsilon$ -greedy policy as the foundation for the following deep Q-learning. As defined in the above section, each action  $a_t^n$  transfers the current state  $\mathbf{s}_t^n$  of AP  $n$  to  $\mathbf{s}_{t+1}^n$  with reward  $r_{t+1}^n$ . The action-value function of AP  $n$  is denoted as follows:

$$J(\mathbf{s}_t^n, \{A_t^n\}) \triangleq \mathbf{E}_\delta \left[ \sum_{k=0}^{\infty} \gamma^k r(\mathbf{s}_{t+k}^n, A_{t+k}^n, \mathbf{s}_{t+k+1}^n) \mid \mathbf{s}_t^n \right], \quad (4)$$

where the action-value function  $J(\mathbf{s}_t^n, a_t^n) : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  outputs the accumulated reward with respect to the state  $\mathbf{s}_t^n$  and the corresponding action  $a_t^n$ . The optimal value function is defined as:

$$V^*(\mathbf{s}_t^n) = \max_{A_t^n} J(\mathbf{s}_t^n, A_t^n), \quad (5)$$

where  $V^*(\mathbf{s}_t^n)$  can be further written as the Bellman optimality equation:

$$V^*(\mathbf{s}_t^n) = \max_{a_t^n \in \mathcal{A}} \sum_{\mathbf{s}_{t+1}^n \in \mathcal{S}} \delta(\mathbf{s}_t^n, \mathbf{s}_{t+1}^n; a_t^n) [r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) + \gamma V^*(\mathbf{s}_{t+1}^n)], \quad (6)$$

where  $\delta(\mathbf{s}_t^n, \mathbf{s}_{t+1}^n; a_t^n)$  represents the transition probability from state  $\mathbf{s}_t^n$  to  $\mathbf{s}_{t+1}^n$  after taking action  $a_t^n$ . The optimal Q-function can then be expressed as the follows:

$$Q^*(\mathbf{s}_t^n, a_t^n) = \sum_{\mathbf{s}_{t+1}^n \in \mathcal{S}} \delta(\mathbf{s}_t^n, \mathbf{s}_{t+1}^n; a_t^n) \{r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) + \gamma V^*(\mathbf{s}_{t+1}^n)\}, \quad (7)$$

in which the Q-function is a fixed point of a contraction operator  $\mathcal{H}$  [21], defined for a generic function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  as the follows:

$$(\mathcal{H}Q)(\mathbf{s}_t^n, a_t^n) = \sum_{\mathbf{s}_{t+1}^n \in \mathcal{S}} \delta(\mathbf{s}_t^n, \mathbf{s}_{t+1}^n; a_t^n) \{r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) + \gamma \max_{a_{t+1}^n \in \mathcal{A}} Q(\mathbf{s}_{t+1}^n, a_{t+1}^n)\}. \quad (8)$$

In the case of model-free reinforcement learning, the above Q-function is impossible to obtain since the transition probability is unknown. Hence, the Q-learning algorithm searches the optimal Q-function with samplings from the episodes of the MDP. Then, the Q-learning algorithm utilizes the following updating rule:

$$Q(\mathbf{s}_t^n, a_t^n) \leftarrow Q(\mathbf{s}_t^n, a_t^n) + \beta \{r(\mathbf{s}_t^n, a_t^n, \mathbf{s}_{t+1}^n) + \gamma \max_{a_{t+1}^n \in \mathcal{A}} Q(\mathbf{s}_{t+1}^n, a_{t+1}^n) - Q(\mathbf{s}_t^n, a_t^n)\}, \quad (9)$$

where the learning rate is denoted by  $\beta$ . While each AP updates  $Q(\mathbf{s}_t^n, a_t^n)$ , it also makes decisions based on  $Q(\mathbf{s}_t^n, a_t^n)$ , i.e., choosing the action corresponding to the largest Q-value. For the  $\epsilon$ -greedy policy, the optimal action is given by

$$a_t^n = \begin{cases} \operatorname{argmax}_{a_t^n} Q(\mathbf{s}_t^n, a_t^n), & P = 1 - \epsilon. \\ \text{random action}, & P = \epsilon, \end{cases} \quad (10)$$

where  $\epsilon$  denotes the probability of choosing random action. The greedy policy helps the Q-learning policy to search for more possibilities of actions randomly. It can help the policy converge faster and prevent the policy from being stuck at a

sub-optimum. Q-learning is proven to converge to the optimum action-values with probability 1 so long as all actions are repeatedly sampled in all states, and the action-values are represented discretely [22].

### B. Gradient descent in Deep Q-Learning

The traditional Q-learning algorithm can be applied to solve for the optimal policy. However, traditional Q-learning is impractical if the dimension of action-state space is large, i.e., the curse of dimensionality [23]; thus, the well-known DQN is proposed in [15] to approximate the action-state Q-value function and the neural network used to achieve the approximation is called Q neural network (QNN).

Each AP is equipped with a QNN which outputs the approximated Q-value  $\{Q(s_t^n, a_t^n; \theta^n) | a_t^n \in \mathcal{A}\}$  given the input state  $s_t^n$  and action  $a_t^n$ . The optimal policy is to choose the action with the largest Q-value. Unlike the tabular update for Q-learning in Eq (9), the QNN in deep Q-learning can be trained by minimizing prediction errors of  $Q(s_t^n, a_t^n; \theta^n)$  at each AP and time slot, where  $\theta^n$  denotes the trainable QNN weights on AP  $n$ . After the reward is obtained, the state transfers to  $s_{t+1}^n$ . The pair  $(s_t^n, a_t^n, r(s_t^n, a_t^n, s_{t+1}^n), s_{t+1}^n)$  then forms a single training sample for QNN and is stored in training set  $\mathcal{D}^n$ . Please note that we will sample training data  $d_s^n$  from the training set  $\mathcal{D}^n$  for each update in the training process. Next, we define the prediction loss function of QNN as

$$L(\theta^n) = (v - Q(s_t^n, a_t^n; \theta^n))^2, \quad (11)$$

where  $Q(s_t^n, a_t^n; \theta^n)$  is the output of QNN at time slot  $t$  and the approximate value function is defined as

$$v = r(s_t^n, a_t^n, s_{t+1}^n) + \gamma \max_{a_{t+1}^n} Q(s_{t+1}^n, a_{t+1}^n; \theta^n), \quad (12)$$

in which the second term  $\gamma \max_{a_{t+1}^n \in \mathcal{A}} Q(s_{t+1}^n, a_{t+1}^n; \theta)$  is obtained by searching the maximum output of QNN with respect to the selection of action  $a_{t+1}^n$  given  $s_{t+1}^n$ . Then, we can update the trainable QNN weights using the semi-gradient algorithm [24] as below:

$$\theta^n \leftarrow \theta^n + \rho [v - Q(s_t^n, a_t^n; \theta^n)] \nabla Q(s_t^n, a_t^n; \theta^n), \quad (13)$$

where  $\nabla$  is the gradient with respect to  $\theta^n$ . Moreover, each AP is employed with deep Q-learning to search for the optimal QNN. However, this inevitably results in a performance loss, for some APs cannot avoid learning aggressive policies to maximize the contention benefits for themselves, and some APs learn conservative policies to avoid collision, especially when AP networks are densely overlapping. Hence, in our proposed protocol, each APs sends its QNN weights to APC that takes the average of the weights of all QNNs as follows:

$$\theta_g = \frac{1}{N} \sum_{n=1}^N \theta^n, \quad (14)$$

where  $\theta_g$  is denoted as global weight. Minimizing the above equation is equivalent to minimize the summation of all loss functions of QNNs from all APs, i.e.,

$$\sum_{n=1}^N L(\theta^n) = \sum_{n=1}^N (v - Q(s_t^n, a_t^n; \theta^n))^2. \quad (15)$$

---

### Algorithm 1: DLCA Algorithm.

---

**Data:**  $\theta^n$ ,  $s_0^n$ , and  $t = 0$ .

**while**  $t \geq 0$  **do**

**if**  $\text{mod}(t, T)$  is not  $T - 1$  **then**

$$a_t^n = \begin{cases} \text{argmax}_{a_t^n} Q(s_t^n, a_t^n), & P = 1 - \epsilon. \\ \text{random action}, & P = \epsilon \end{cases}$$

    Obtain  $r(s_t^n, a_t^n, s_{t+1}^n)$  according to Algorithm 3 and store the tuple  $(s_t^n, a_t^n, r_t^n, s_{t+1}^n)$  to the training batch  $\mathcal{D}$ ;

**if** *Update* **then**

        1. Sample  $d_s^n$  transitions from  $\mathcal{D}^n$ .

        2. Calculate the target value as follows:

$$v = r_t^n + \gamma \left( \max_{a_{t+1}^n} Q(s_{t+1}^n, a_{t+1}^n; \theta^n) \right)_Q.$$

        3. For each randomly sample tuple in the training batch with  $d_s$  samples, update  $\theta^n$  with the following gradient descent method:

$$\theta^n \leftarrow \theta^n - \rho \nabla_{\theta^n} L(\theta^n).$$

**else**

        4. Each AP send its QNN weights to APC that obtains the global QNN  $\theta_g = \frac{1}{N} \sum_{n=1}^N \theta^n$ .

        5. Sample  $d_s^g$  from  $\{\mathcal{D}^1, \dots, \mathcal{D}^N\}$  and update  $\theta_g$  with the following gradient descent method:

$$\theta_g \leftarrow \theta_g +$$

$$\frac{\rho}{N} \sum_{n=1}^N [v - Q(s_t^n, a_t^n; \theta^n)] \nabla Q(s_t^n, a_t^n; \theta^n).$$

        6. Send  $\theta_g$  back to each AP:  $\theta^n \leftarrow \theta_g$ .

$t = t + 1$ ;

---

In the meantime, minimizing the above equation implicates the following gradient descent for the summed loss function:

$$\theta_g \leftarrow \theta_g + \frac{\rho}{N} \sum_{n=1}^N [v - Q(s_t^n, a_t^n; \theta^n)] \nabla Q(s_t^n, a_t^n; \theta^n). \quad (16)$$

Hence, the global QNN weight  $\theta_g$  is equivalently obtained by iterations over the sampled data batch  $d_s^g$  collected from all APs' local data set  $\{\mathcal{D}^1, \dots, \mathcal{D}^N\}$ , which enables a faster convergence rate and lower loss function value. It is noteworthy that  $\theta_g$  is sent back to all APs from APC after global gradient descent completes according to Eq (16). Then,  $\theta_g$  replaces the previous QNN weights for future training and inference. This method is called First-Order Model-Agnostic Meta-Learning (FOMAML), which can further enhance all APs' models with non-IID local data [25].

**Remark 1.** In Algorithm 1, samples are collected from the local data set in step 1. Then, in step 2, target value is calculated. Gradient descent method is implemented in step 3 based on the collected samples and the calculated target value. Step 4, 5, and 6 represent the FOMAML method and is triggered once every  $T$  local training loops. It is noteworthy that the  $Q$ -value in step 2 can be obtained by looping the action corresponding to the largest  $Q$ -value with time complexity of  $\mathcal{O}(F)$ . The gradient in machine learning is normally computed using the back-propagation method [26] as a numerical solution with time complexity of  $\mathcal{O}(FM)$ . Algorithm 1 is typically executed in batch mode - such that QNN update occurs once per batch to reduce computation load. FOMAML is only triggered every period of  $T$  to reduce the communication overhead between the APC and AP. The global QNN  $\theta_g$  is trained on APC using global information gathered by APC, i.e.,  $\{d_s^1, \dots, d_s^N\}$ .

### C. AP Coordination: Greedy Algorithm

In the above section, each AP runs with a deep Q-learning algorithm independently. However, the channel on which each AP should run is not described. In this section, the APC policy that allocates channels to all APs considering PF is proposed.

Denote  $\phi_t^n(f)$  as the instantaneous proportional achievable data rate for AP  $n$  at time slot  $t$  at channel  $f$ . We assume the block fading channel condition to explore the convergence property of our proposed algorithm, i.e.,  $C_t^n(f) = C^n(f)$  is assumed to be constant over multiple TXOP slots. Then we have  $\phi^n(f) = \frac{C^n(f)}{\mathbf{n}(f)}$ . Denote  $\mathbf{x}_t^n(f)$  as the actual data rate of AP  $n$  at time slot  $t$  at channel  $f$ , the allocation scheme for AP  $n$  can then be expressed as follows:

$$f^* = \operatorname{argmax}_f P_t^n(f), \quad (17)$$

where

$$P_t^n(f) = \frac{\phi^n(f)}{\tilde{D}_t^n} \quad (18)$$

in which

$$\tilde{D}_t^n = \left(1 - \frac{1}{t}\right) \tilde{D}_{t-1}^n + \frac{\mathbf{1}^T \mathbf{x}_t^n}{t}, \quad (19)$$

and  $\tilde{D}_t^n$  represents the average throughput of AP  $n$  up to time slot  $t$ . Note that only one element in  $\mathbf{x}_t^n \in \mathcal{R}^F$  is non-zero. After allocating one channel to an AP, the scheduler updates the ratio  $P_t^n(f)$  for the next AP's allocation. The proposed greedy algorithm considering PF [27] on APC is specified in Algorithm 2. In Algorithm 2, each AP is allocated with a channel in the while-loop. The average throughput of AP  $n$  up to time slot  $t$  is calculated. Then, the channel is chosen to maintain the current PF. The greedy algorithm can guarantee the asymptotic PF, and the corresponding proof is shown in Appendix B.

Round Robin (RR) and PF have been developed as two common scheduling strategies. Among those, PF is widely considered in wireless networks. Different AP has different average throughput due to the number of previously gained TXOPs and various channel spectral efficiency. The greedy

---

### Algorithm 2: Greedy Algorithm Considering PF on APC.

---

**Data:**  $x_0^n, \tilde{D}_0^n, C^n(f)$ , and  $n = 0$ .

$\mathbf{n}(f) = 0$  for all  $f$ ;

**while**  $n \leq N$  **do**

$$\tilde{D}_t^n \leftarrow \left(1 - \frac{1}{t}\right) \tilde{D}_{t-1}^n + \frac{\mathbf{1}^T \mathbf{x}_t^n}{t}.$$

$$f^* \leftarrow \operatorname{argmax} \left\{ \frac{\phi^n(f)}{\tilde{D}_t^n} \right\}$$

where

$$\phi^n(f) = \frac{C^n(f)}{\mathbf{n}(f)}.$$

$n \leftarrow n + 1$ ;

$\mathbf{n}(f) \leftarrow \mathbf{n}(f) + 1$ ;

Allocate AP  $n$  to  $f^{\text{th}}$  channel;

---

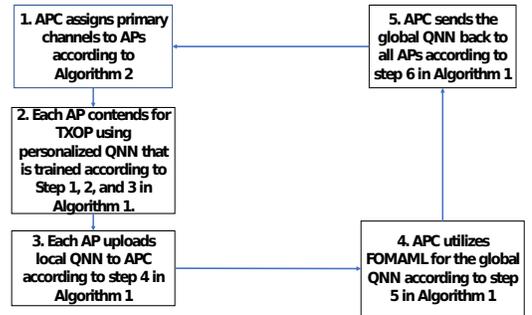


Fig. 4: Flow Chart of DLCA Protocol: DLCA + Greedy Algorithm + FOMAML.

algorithm considering PF exploits these variations by allocating the primary channel to the AP with the best conditions for the upcoming TXOP slot. As a design approach, this approach is superior to RR. In the end, the proposed DLCA protocol is shown in Fig.4.

**Remark 2.** In Algorithm 2, the computation complexity of the PF scheduling method is  $\mathcal{O}(NF)$ . Similar to FOMAML, The greedy algorithm is only triggered every period of  $T$  on APC to reduce the communication burden between the APC and AP. The information of instantaneous data rate  $\mathbf{x}_t^n$ , average data rate  $\tilde{D}_t^n$ , and spectral efficiency  $C^n(f)$  can be exchanged between APs and APC through the wired connection such as Light Weight Access Point Protocol (LWAPP), wireless connection such as DCF. FOMAML and the greedy algorithm triggered with a suitable period has negligible impact on the network throughput as long as the QNN is of lightweight.

## IV. PERFORMANCE EVALUATION

In this section, we present simulation results for a dense overlapping network that implements our DLCA protocol. The performance comparison between the DLCA protocol, SH-TXOP, and DCF RTS/CTS for the overlapping network characterizes the superiority of our DLCA protocol. In the

end, we evaluate proportional fairness and stability achieved by the greedy algorithm in the DLCA protocol.

#### A. Throughput in Multi-AP and Multi-band Networks

With the FCC opening up the 6 GHz [28] band for unlicensed use for 5G wireless networks, joint operation in 5 and 6 GHz is feasible with orthogonal sub-channels with a bandwidth of 20 MHz. Our simulations consist of a fully overlapping multi-AP network using 5 GHz and 6 GHz bands for a total of  $F$  sub-channels of 20 MHz. The TXOP slot is granted to the shared APs as a multiple of  $32 \mu s$ , and the maximum amount of time granted is 8.16 ms. The TXOP is thus set as 8.16 ms. In the simulation, AP is assumed to be operating in the saturation mode, i.e., it is always necessary for AP to gain TXOPs because AP needs to communicate with its associated STAs in common Wi-Fi networks continuously. The average value of the spectral efficiency on each channel is 40 Mbps [3], and the spectral efficiency is assumed to be an uniform distribution, i.e.,  $C_t^n(f) \sim U[1, 3]$  bit/s/Hz. FOMAML mechanism and the greedy algorithm are triggered with the period of  $T = 100$  ms. We conduct Monte-Carlo simulations with 100 independent trials and then take the average of the results.

We firstly consider SH-TXOP without APC [9] as the baseline. In the SH-TXOP protocol, APs share the universal frequency band rather than operating on different sub-channels. After a certain AP wins the wide-band TXOP, it starts to share the TXOP on each 20 MHz sub-channel to other APs based on round robin (RR) method, indicated by Announcement Trigger frame (ATF) that contains the information such as the channel allocation scheme for shared APs. For example, if there are 16 APs and 4 channels and AP 1 wins the TXOP, then AP 1 allocates the first sub-channel as the primary channel for itself. Next, it allocates the second sub-channel to AP 2, and etc. Each AP contends for the wide-band TXOP using IEEE 802.11 DCF basic [8]. The Backoff Window size of DCF basic is  $CW_{min} = 32$  and  $m = 6$ . Note that DCF basic method can be utilized in this scenario. Since the packet size in the process of gaining the sharing opportunity of TXOP is negligible, DCF basic is applicable in such a scenario. However, DCF basic is well known to perform worse than DCF RTS/CTS [8]. To further enhance the DCF RTS/CTS, a model is proposed in [29] that increases the network throughput by optimizing the initial backoff window size for DCF RTS/CTS. For this simulation, all APs are equally allocated to a fixed primary channel at the beginning. For example, if there are 16 APs and 4 channels, then AP 1–4 are allocated to the first channel as their primary channel, and AP 5–8 are allocated to the second channel as their primary channel and so on. Next, each AP contends for TXOP in its primary channel using DCF RTS/CTS with the optimized initial backoff window size [29]. This method in the simulation is called RTS/CTS. The aggregate network throughput  $x$  of the entire network with  $N$  APs for this method is expressed as follows:

$$x = \sum_{n=1}^N x^n = \sum_{n=1}^N \frac{z^n LU}{L + \delta}, \quad (20)$$

TABLE II: System Parameters for Multi-AP Networks

Parameters	Value
slot time ( $\mu s$ )	50
SIFS ( $\mu s$ )	28
DIFS ( $\mu s$ )	128
PHY Header ( $\mu s$ )	20
TXOP ( $\mu s$ )	640
CTS_Timeout ( $\mu s$ )	300
ACK_Timeout ( $\mu s$ )	300
Headers (Bytes)	36
ACK (Bytes)	14 + PHY Header
RTS (Bytes)	20 + PHY Header
CTS (Bytes)	14 + PHY Header
ATF (Bytes)	16 + PHY Header

where  $L$  is the information bits in one packet,  $\delta = \delta_0 U$  stands for the protocol overhead in the unit of bits. The channel bit rate  $U$  can be further written as the product of the sub-channel bandwidth and the link spectral efficiency. The above equation is further explained in Appendix C. The system parameters are summarized in Table II.

We utilize PyTorch [30] to train QNN for DLCA. The simulations are conducted on a server with a CPU (Intel Core i7-9700k) and a GPU (NVIDIA GeForce GTX 2080Ti) in Python language. QNN is constructed by  $h = 5$  fully connected layers with 64 neurons in each layer, which is illustrated in Fig.5. The operation of the QNN starts by taking the state vector as the input. Then, it outputs two Q values corresponding to action - transmission and action - wait respectively. AP decides to transmit if the corresponding Q value is larger and wait otherwise. Table III lists the hyper-parameter of the deep Q-learning. ReLU (Rectified Linear Unit) defined as

$$f(x) = x^+ = \max(0, x) \quad (21)$$

is utilized as the activation function to the input of each neuron in the QNN. Unlike other activation functions such as the sigmoid function, ReLU can help the QNN avoid the vanishing gradient issue because the gradient of  $f(x)$  when  $x > 0$  is always a constant. Therefore, choosing ReLU can prompt faster learning process and better performance.

TABLE III: Hyper-Parameters of QNN

Parameters	Value
State size	40
Batch size	32
Learning rate $\rho$	0.001
$\gamma$ in Eq (12)	0.9
$\epsilon$ in Eq (10)	0.05
$\eta$ in Eq (3)	0.5
Step size of $\lambda$	0.1
Gradient descent step size $u$	0.25

The simulation result of the network throughput is shown in Fig.6. We validate the precision of the RTS/CTS model by showing that the simulation results are close to their corresponding theoretical results. Three curves related to the DLCA protocol are plotted respectively. The curve labeled with DLCA reflects the simulation of running distributed deep reinforcement learning on each AP without primary channel allocation from the APC. In this case, each AP only has a fixed primary channel allocation similar to the RTS/CTS simulation.

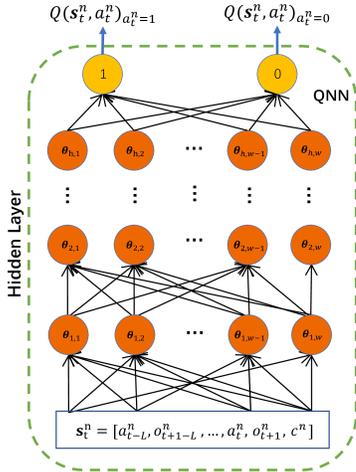


Fig. 5: QNN: Fully-connected Neural Network.

On the other hand, DLCA + greedy method represents the simulation of running distributed deep reinforcement learning on each AP with primary channel allocation decision from APC. Note that APC does not update QNN globally in DLCA + greedy method. In Fig.6(b), when the number of APs is small, the performance of DLCA and DLCA + greedy method is better than SH-TXOP. The reason is that when the number of APs is  $N = 8$  for both DLCA methods, each AP has its own exclusive primary channel and no collision happens at all, which can be observed from Fig.7. However, 8 APs have to contend for the wide-band TXOP for SH-TXOP. The average slots wasted in the collision for DLCA and DLCA + greedy method are zero, not to mention that SH-TXOP has more average IDLE slots per TXOP than DLCA + greedy. As the number of APs increases, SH-TXOP outperforms DLCA and DLCA + greedy. This is because some APs can develop very aggressive TXOP contention policy without the supervision from FOMAML. Note that each AP only wants to maximize its own total reward. Hence, it is possible that APs assigned to one channel are all aggressive and we can view this training process as Prisoner's Dilemma; that is, if one AP does not develop aggressive TXOP contention policy, then it has no chance to get any TXOP forever. Therefore, it is likely that DLCA will lead to high collision probability. However, SH-TXOP and DCF RTS/CTS have backoff counter to avoid collision probability if collision happens. Hence, as the number of APs increases, average slots wasted in collisions and average IDLE slots per TXOP all increases for DLCA and DLCA + greedy in Fig.7 and 8, which leads to severe performance loss. Although both RTS/CTS method and DLCA + greedy + FOMAML perform well, DLCA + greedy + FOMAML provides higher network throughput than RTS/CTS method by 3% for the total number of AP ranging from 8 to 56. This is because the overhead remains in RTS/CTS and the backoff window takes up time slots without sending any data packet. One can observe the RTS/CTS method has more average collision slots and IDLE slots from Fig.7 and 8. In Fig.6(c), for the case of 16 channels with 8 APs, only half of the channel resources are utilized. Hence, the network throughput grows

linear with increasing number of APs at the beginning. In Fig.6, DLCA + Greedy + FOMAML outperforms SH-TXOP by 10% when the number of APs is  $N = 56$  in average of three cases. The advantage of FOMAML can also be demonstrated in Fig.9, during the training process, the network throughput in both methods have large variance in the initial learning phase. However, DLCA + greedy + FOMAML has faster convergence rate and smaller variance. This can be attributed to the fact that each AP's self-training only reaches local optimality, emphasizing the important role of FOMAML as the global optimizer that achieves the necessary AP coordination for throughput maximization.

### B. PF in Multi-AP and Multi-band Networks

In the above section, aggregate network throughput is simulated. However, aggregate network throughput does not reflect the throughput of each AP, leading to a potential issue that the maximum throughput can always be achieved by having the same AP holding the channel, and no fairness exists at all. Hence, fairness must be guaranteed so that each AP in the network can utilize the TXOP for UL/DL communication with associated STAs. The greedy algorithm has been proven to enable PF among APs asymptotically previously. This section implements simulations to study a network utility metric that describes PF and network throughput. Meanwhile, the stability of our proposed algorithm is also investigated.

To be consistent with the notations in Appendix.B, the network utility is defined as  $\sum_{n=1}^N \log(\bar{D}^n)$ , where  $\bar{D}^n$  is the average data rate of AP  $n$ . We simulate the network utility of DLCA + greedy + FOMAML, SH-TXOP, and RTS/CTS for comparison. In Fig.10, DLCA protocol performs better than RTS/CTS by 13.8% and SH-TXOP by 28.3% when the number of APs equals  $N = 56$ . In  $\sum_{n=1}^N \log(\bar{D}^n)$ , the  $\log$  term punishes the AP that has low throughput. Therefore, the network utility demonstrates joint network throughput and PF. We next show the stability of DLCA + greedy + FOMAML in terms of AP's PF ratio, which is expressed as

$$b^n = \frac{\bar{D}^n}{\phi^n}. \quad (22)$$

According to Eq (28), the closer  $b^n$  of each AP is to each other, the better PF is achieved. In Fig.11, the PF ratios of 3 APs converge to 0.87, 0.84, and 0.81 respectively after 30000 steps. Then, we exchange the value of spectral efficiency  $C_t^n(f)$  between AP 1 and AP 3 at step 30000 to demonstrate the stability. After a sudden change at step 30000, three curves experience drastic oscillation. Then, we can observe that the curves of AP 1 and AP 3 converge again eventually. This result indicates that the fairness broken by a sudden network change can be restored very quickly. Hence, our greedy algorithm is shown to be robust and efficient.

In the end, we conclude that the design of multi-AP network with APC has a higher upper limit than the multi-AP network without APC in terms of the network throughput. The choice of RTS/CTS or DLCA + greedy + FOMAML is constrained by the hardware and energy cost. For the AP with limited power constraint, one can choose RTS/CTS with lower power

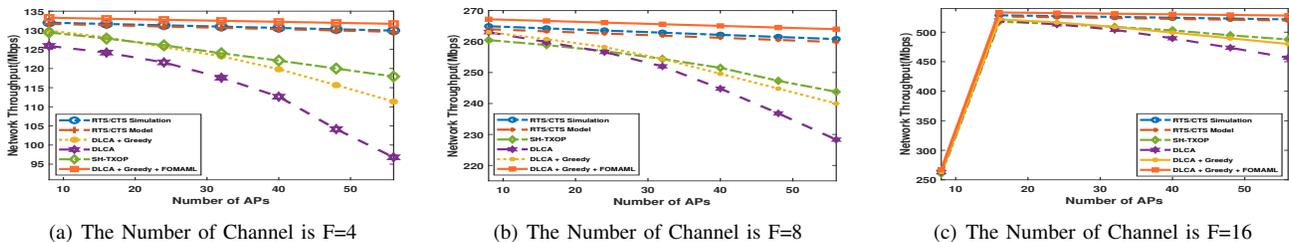
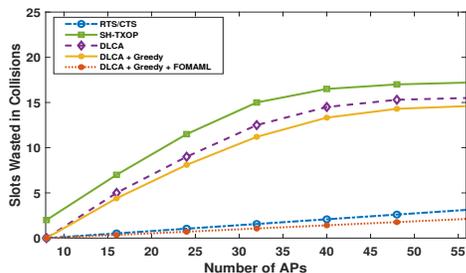
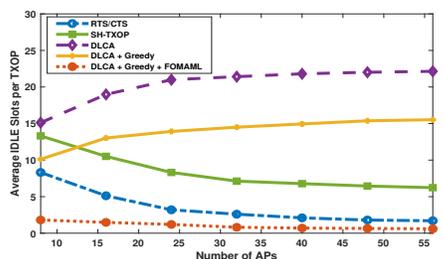


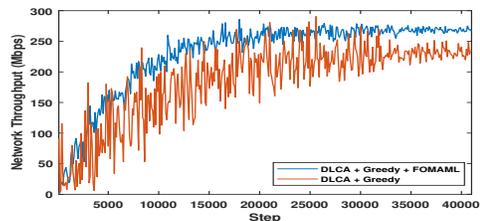
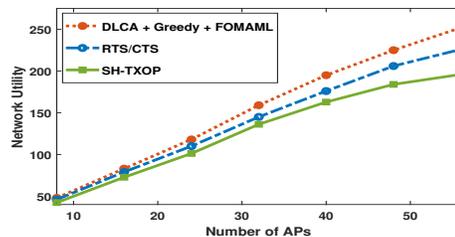
Fig. 6: Network Throughput vs Number of APs.

Fig. 7: Average Collision vs Number of APs. The Number of Channel is  $F = 8$ .Fig. 8: Average IDLE vs Number of APs. The Number of Channel is  $F = 8$ .

consumption but the performance loss, especially the lack of PF, might lead to an unsatisfied user experience. On the other hand, if the power budget is high enough and each AP is able to run light-weight QNN with suitable CPU or GPU, AP can reach higher network throughput, and proportional fairness among APs can be also guaranteed.

## V. CONCLUSION

In this paper, we propose enhancements to the RRM architecture for dense overlapping Wi-Fi networks that align with the proposed coordinated AP operation in Wi-Fi 7 (802.11be). Specifically, we develop a novel multi-AP coordination system architecture with DLCA protocol. The proposed protocol considers not only the network throughput maximization but also the proportional fairness among APs. The performance of DLCA related algorithms is then evaluated via simulations and compared with SH-TXOP protocol and RTS/CTS as benchmarks. The numerical results show that DLCA outperforms SH-TXOP and state-of-the-art RTS/CTS with an optimized initial back-off window in terms of network throughput and

Fig. 9: Convergence: Throughput vs Training Steps. The Number of Channel  $F = 8$ , The Number of APs  $N = 16$ . Each time step is one round of gradient descent in Eq. (13).Fig. 10: Network Utility ( $\sum_{n=1}^N \log(\bar{D}^n)$ ) vs Number of APs. The number of channel is  $F = 8$ .

network utility. Moreover, convergence rate and stability are also demonstrated in the simulation.

This paper studies the fully overlapping dense Wi-Fi networks. In our future work, we will investigate partially overlapping dense Wi-Fi networks. In such a case, optimization for dynamic AP coordination set, frequency reuse in different coordination set for 802.11be, and coexistence with other protocols will be considered.

## APPENDIX

### A. Reward Estimation Method

Algorithm 3 is a modified version of Monte Carlo method in [24] that aims to estimate reward and help Q-learning converges faster. For AP  $n$  operating in  $f^{th}$  frequency channel, we initialize the total reward to be zero in step 1. Then, if the current action  $a_t^n(f) = 1$ , we use while loop to find all feedback of the transmission action (action value is equal to 1) in the state vector. Suppose one of the feedback is ACK, we add one to the weighted total reward value in step 2 since it is a successful transmission. Otherwise, we minus one to punish the weighted total reward value in step 3. Suppose the current action  $a_t^n(f) = 0$ , then the total reward value is 1 if

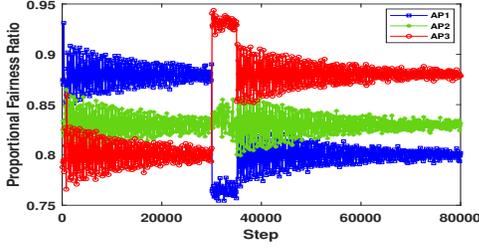


Fig. 11: PF ratio  $b^n$  vs Time Step with Greedy Algorithm. PF ratios of 3 APs out of  $N = 18$  APs on  $F = 8$  channels are depicted. Each time step is one slot time in Tab. II.

busy channel is sensed in step 4 since the AP successfully avoids a potential collision. The total reward value is set to  $-1$  as a punishment if idle channel is sensed in step 5.

---

**Algorithm 3:** Reward estimation method.

---

**Data:**  $f$ ,  $s_t^n$ , feedback for all actions in  $s_t^n$ .

**Result:**  $r(s_t^n, a_t^n, s_{t+1}^n)$

1.  $r(s_t^n, a_t^n, s_{t+1}^n) \leftarrow 0$ ;

**if**  $a_t^n(f) == \{1\}$  **then**

$l \leftarrow L$ ;

**while**  $l \geq 0$  **do**

**if** the feedback of  $a_{t-l}^n(f) = 1$  is ACK **then**

            2.  $r(s_t^n, a_t^n, s_{t+1}^n) \leftarrow \eta \times r(s_t^n, a_t^n, s_{t+1}^n) + 1$

**else**

            3.  $r(s_t^n, a_t^n, s_{t+1}^n) \leftarrow \eta \times r(s_t^n, a_t^n, s_{t+1}^n) - 1$

$l \leftarrow l - 1$ ;

**else**

**if**  $o_{t+1}^n(f) == \{1\}$  **then**

        4.  $r(s_t^n, a_t^n, s_{t+1}^n) \leftarrow 1$

**else**

        5.  $r(s_t^n, a_t^n, s_{t+1}^n) \leftarrow -1$

---

### B. Proof of Asymptotic Proportional Fairness

We show that the greedy algorithm considering PF maximizes the aggregate throughput while guaranteeing the asymptotic PF on one channel. Denote  $p_t^n(f)$  as the probability of the AP  $n$  at time slot  $t$  being assigned with channel  $f$ . Assume the spectral efficiency does not change within  $t$  slots, then we have

$$\bar{D}_t^n = \sum_f \sum_t p_t^n(f) \phi^n(f). \quad (23)$$

Therefore, the network utility maximization problem [27] is

$$\max \sum_n \log \left( \sum_f \sum_t p_t^n(f) \phi^n(f) \right) \quad (24)$$

$$\text{s.t.} \sum_n \sum_f p_t^n(f) \leq 1, p_t^n(f) \geq 0, \forall n, f, t. \quad (25)$$

Note that the above problem is a convex problem since log function with composition of an affine function still preserves concavity. Applying Lagrange multipliers, we obtain the following:

$$\sum_n \log \left( \sum_t \sum_f p_t^n(f) \phi^n(f) \right) - \sum_t \lambda_t \left( \sum_n \sum_f p_t^n(f) - 1 \right) \quad (26)$$

Taking the derivative w.r.t.  $p_t^n(f)$ , we obtain the optimal solution as

$$\frac{\phi^n}{\bar{D}_t^n} - \lambda^* = 0 \quad \text{if } p_t^n(f) > 0, \quad (27)$$

Asymptotically, the PF algorithm helps the AP network to reach the PF, i.e.,

$$\lim_{t \rightarrow \infty} \frac{\bar{D}_t^1}{\phi^1} = \dots = \lim_{t \rightarrow \infty} \frac{\bar{D}_t^n}{\phi^n}. \quad (28)$$

Hence, allocation method shown in Eq. (17) follows the optimal condition of the IEEE 802.11 network utility maximization problem with PF.

### C. Performance of Multi-AP IEEE 802.11 RTS/CTS Networks

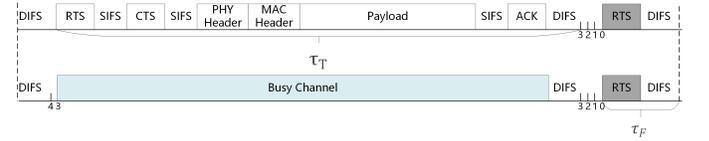


Fig. 12: Graphic illustration of successful transmission and collision in DCF networks with the RTS/CTS access mechanism.

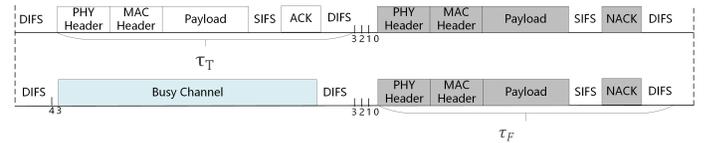


Fig. 13: Graphic illustration of successful transmission and collision with DCF basic mechanism.

The aggregate network data rate  $x$  is the average number of information bits successfully transmitted per second, which is the sum of the average number of information bits that AP  $n$  successfully transmits per second  $x^n$

$$x = \sum_{n=1}^N x^n, \quad (29)$$

The duration to transmit a packet consisting of  $L$  information bits is given by  $\frac{L}{U} + \delta_0$ , where  $U$  is the channel bit rate, and  $\delta_0$  is the protocol overhead in seconds. With the RTS/CTS mechanism and DCF basic illustrated in Fig. 12 and 13 respectively,  $\delta_0$  is given by

$$\begin{aligned} \delta_0^{RTS} &= \frac{\text{RTS} + \text{CTS} + \text{ACK}}{U_b} + \text{Header} + \text{DIFS} + 3 \times \text{SIFS} \\ \delta_0^{\text{basic}} &= \frac{\text{ACK}}{U_b} + \text{Header} + \text{DIFS} + \text{SIFS}, \end{aligned} \quad (30)$$

where  $U_b$  denotes the basic rate. In the end, the data rate  $x^n$  of each BSS  $i$  can be expressed as follows:

$$x^n = \frac{Lz^n}{\frac{L}{U} + \delta_0} = \frac{z^n LU}{L + \delta},$$

where  $\delta = \delta_0 U$  stands for the protocol overhead in the unit of bits. The channel bit rate  $U$  can be further written as the product of the channel bandwidth of each AP and the link spectral efficiency. The details of derivation of  $z^n$  is related to parameters of  $\tau_T$  and  $\tau_F$  [29]. As for the throughput of Multi-AP IEEE 802.11 DCF basic network without APC,  $\tau_T^{DCF} = \tau_T - (RTS + CTS + 2 \times SIFS)$ .  $\tau_F^{DCF}$  is equal to  $\tau_T^{DCF}$  in basic DCF since there is no RTS/CTS and the collision leads to a waste of whole packet time instead of RTS/CTS with short duration, which is the reason why RTS/CTS can enhance the system throughput.

## REFERENCES

- [1] L. Zhang, H. Yin, Z. Zhou, S. Roy, and Y. Sun, "Enhancing WiFi multiple access performance with federated deep reinforcement learning," in *IEEE 92nd Vehicular Technology Conference (VTC2020-Fall)*.
- [2] Cisco, "Cisco Annual Internet Report (2018–2023)," *White Paper*, Cisco System Inc., Mar. 2020.
- [3] D. Lopez-Perez, A. Garcia-Rodriguez, L. Galati-Giordano, M. Kasslin, and K. Doppler, "IEEE 802.11be Extremely High Throughput: The Next Generation of Wi-Fi Technology Beyond 802.11ax," *IEEE Communications Magazine*, vol. 57, no. 9, pp. 113–119, Sept. 2019.
- [4] E. Khorov, A. Kiryanov, A. Lyakhov, and G. Bianchi, "A Tutorial on IEEE 802.11ax High Efficiency WLANs," *IEEE Communications Surveys Tutorials*, vol. 21, no. 1, pp. 197–216, Sep. 2019.
- [5] J. Liu, T. Pare, Y. Seok, J. Wang, F. Hsu, and J. Yee, "Multi-AP Enhancement and Multi-Band Operations," Mediatek Inc., Tech. Rep. IEEE 802.11-18/1155r1, Jun. 2018.
- [6] S. Naribole, W. B. Lee, K. Srinivas, R. Duan, and A. Ranganath, "Shared TXOP Protocol," Samsung Inc., Tech. Rep. IEEE 802.11-20/0277r1, Mar. 2020.
- [7] B. Kwak, N. Song, and L. Miller, "Performance Analysis of Exponential Backoff," *IEEE/ACM Transactions on Networking*, vol. 13, no. 2, pp. 343–355, Apr. 2005.
- [8] G. Bianchi, "Performance Analysis of the IEEE 802.11 Distributed Coordination Function," *IEEE Journal on Selected Areas in Communications*, vol. 18, no. 3, pp. 535–547, Mar. 2000.
- [9] W. Ahn, "Novel Multi-AP Coordinated Transmission Scheme for 7th Generation WLAN 802.11 be," *Entropy*, vol. 22, no. 12, p. 1426, 2020.
- [10] L. Sequeira, J. L. de la Cruz, J. Ruiz-Mas, J. Saldana, J. Fernandez-Navajas, and J. Almodovar, "Building an SDN Enterprise WLAN Based on Virtual APs," *IEEE Communications Letters*, vol. 21, no. 2, pp. 374–377, 2017.
- [11] Y. Wang, X. Li, P. Wan, and R. Shao, "Intelligent Dynamic Spectrum Access Using Deep Reinforcement Learning for VANETs," *IEEE Sensors Journal*, vol. 21, no. 14, pp. 15 554–15 563, Feb. 2021.
- [12] H. Song, L. Liu, J. Ashdown, and Y. Yi, "A Deep Reinforcement Learning Framework for Spectrum Management in Dynamic Spectrum Access," *IEEE Internet of Things Journal*, vol. 8, no. 14, pp. 11 208–11 218, Jan. 2021.
- [13] W. Ahsan, W. Yi, Z. Qin, Y. Liu, and A. Nallanathan, "Resource Allocation in Uplink NOMA-IoT Networks: A Reinforcement-Learning Approach," *IEEE Transactions on Wireless Communications*, vol. 20, no. 8, pp. 5083–5098, Mar. 2021.
- [14] H. Ding, F. Zhao, J. Tian, D. Li, and H. Zhang, "A Deep Reinforcement Learning for User Association and Power Control in Heterogeneous Networks," *Ad Hoc Networks*, vol. 102, p. 102069, May 2020.
- [15] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski et al., "Human-level Control through Deep Reinforcement Learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [16] Y. Kihira, Y. Koda, K. Yamamoto, T. Nishio, and M. Morikura, "Adversarial Reinforcement Learning-based Robust Access Point Coordination against Uncoordinated Interference," *arXiv preprint arXiv:2004.00835*, 2020.
- [17] Y. Yu, T. Wang, and S. Liew, "Deep-reinforcement Learning Multiple Access for Heterogeneous Wireless Networks," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1277–1290, Mar. 2019.
- [18] Y. Yu, S. C. Liew, and T. Wang, "Non-uniform Time-step Deep Q-network for Carrier-sense Multiple Access in Heterogeneous Wireless Networks," *IEEE Transactions on Mobile Computing*, Apr. 2020.
- [19] —, "Multi-Agent Deep Reinforcement Learning Multiple Access for Heterogeneous Wireless Networks with Imperfect Channels," *IEEE Transactions on Mobile Computing*, pp. 1–1, Feb. 2021.
- [20] X. Ye, Y. Yu, and L. Fu, "Multi-channel Opportunistic Access for Heterogeneous Networks Based on Deep Reinforcement Learning," *IEEE Transactions on Wireless Communications*, pp. 1–1, July 2021.
- [21] M. G. Bellemare, G. Ostrovski, A. Guez, P. Thomas, and R. Munos, "Increasing the action gap: New operators for reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [22] C. Watkins and P. Dayan, "Technical Note: Q-Learning," *Machine learning*, vol. 8, no. 3-4, pp. 279–292, May 1992.
- [23] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep Reinforcement Learning: A Brief Survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, Nov. 2017.
- [24] R. Sutton and A. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [25] A. Fallah, A. Mokhtari, and A. Ozdaglar, "On the convergence theory of gradient-based model-agnostic meta-learning algorithms," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 1082–1092.
- [26] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, MIT press, 2016.
- [27] S. Rayadurgam and Y. Lei, *Communication Networks: An Optimization, Control, and Stochastic Networks Perspective*, Cambridge University Press, 2013.
- [28] Federal Communications Commission (FCC), "Unlicensed Use of the 6 GHz Band," *Docket No. 17-183*, Apr. 2020.
- [29] Y. Gao, L. Dai, and X. Hei, "Throughput Optimization of Multi-BSS IEEE 802.11 Networks With Universal Frequency Reuse," *IEEE Transactions on Communications*, vol. 65, no. 8, pp. 3399–3414, May 2017.
- [30] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., "Pytorch: An Imperative Style, High-performance Deep Learning Library," in *Advances in Neural Information Processing Systems*, 2019, pp. 8026–8037.