# Federated Multi-Agent Reinforcement Learning for QoS-Aware RU Scheduling in High Dense WLANs

*Pratibha Kantanavar*$^{*1}$,*Hariprasad S A*$^2$, *Gopalakrishna K*$^2$
Department of ECE, FET, Jain Deemed to be University, Jain Global Campus, Bengaluru
Bengaluru, Karnataka, India

**Abstract:**
The increasing demands of latency-sensitive and energy-constrained applications in high-density WLANs expose the limitations of static RU scheduling in IEEE 802.11ax networks. Existing random access methods and centralized learning models often fail to generalize under dynamic traffic loads, congestion patterns, and heterogeneous QoS profiles. In this paper, we propose Federated Multi Agent Reinforcement Learning AURA, a federated multi-agent reinforcement learning framework that enables decentralized, QoS-aware RU allocation across heterogeneous stations. Each station acts as an independent agent using lightweight actor-critic learning, coordinated through a federated server at the AP. The framework integrates cross-layer intelligence by incorporating MAC-layer contention adaptation, PHY-layer RU mapping, and application-driven traffic classes into a unified policy. Simulation results show up to 61% improvement in throughput, 32% reduction in energy consumption, and 99% fairness index, outperforming DQN and static IEEE 802.11ax baselines. This work offers a scalable and adaptive solution for next-generation WLANs, with demonstrated resilience under dense traffic, mobility, and dynamic QoS weighting.

**Keywords:** OFDM, Random access techniques, User-centric RU Allocation for Wi-Fi (AURA)

## 1. Introduction

Despite the significant performance gains brought by Orthogonal Frequency Division Multiple Access (OFDMA) in IEEE 802.11ax, current random-access techniques such as Uplink OFDMA-based Random Access (UORA) remain constrained by several limitations. Firstly, existing schemes predominantly rely on semi-static random-access procedures without dynamic adaptation to traffic or network states [1]. This leads to suboptimal RU utilization under varying load conditions and fluctuating user densities. Secondly, contention-based randomization still suffers from high collision probabilities as the number of contending stations increases, degrading overall throughput and latency performance. Furthermore, current mechanisms lack QoS-awareness, thereby treating all traffic types uniformly and ignoring the latency sensitivity of real-time applications such as AR/VR or industrial automation [2].

Wireless Local Area Networks (WLANs), particularly those based on the IEEE 802.11 standard, have become indispensable for modern high-throughput communication in both fixed and mobile environments. With the rapid increase in devices and the demand for high data rates, optimizing network throughput has become a critical research area. Throughput defined as the rate of successful data delivery over a communication channel is vital for ensuring the Quality of Service (QoS) in applications such as video streaming, telemedicine, industrial control, and intelligent transportation systems. n modern wireless networks, especially in smart industries and IoT-rich environments, achieving high throughput is fundamental to ensuring reliable communication and Quality of Service (QoS). Throughput—the rate at which data is successfully transmitted from sender to receiver—is crucial for supporting bandwidth-intensive applications such as real-time video analytics, augmented reality (AR/VR), and remote industrial control. However, with increasing device density, spectrum congestion, and latency-sensitive traffic, maintaining consistent throughput has become a growing challenge in WLANs.

Semi-static configurations and random contention-based mechanisms will lead to inefficient RU utilization, especially under varying traffic loads and user densities. The random backoff mechanism still results in collision probabilities that increase with user density, which reduces the channel's overall efficiency [3], Current approaches do not prioritize QoS classes (e.g., real-time video vs. best-effort), treating all devices equally. This causes latency and jitter issues for delay-sensitive traffic. OFDMA randomization does not leverage real-time network states, such as buffer size, energy constraints, or past transmission history, which can be vital for intelligent scheduling [4].

To address these challenges, this paper proposes an Adaptive User-centric RU Allocation for Wi-Fi (AURA)-driven multi agent RU assignment framework that dynamically manages RU allocations based on short-term traffic predictions and energy constraints. The system incorporates Q-learning-based contention window adaptation, dynamic spectrum bifurcation, and Multipath TCP (MPTCP)-enabled bandwidth aggregation to optimize link reliability and spectral efficiency.

Our contributions can be summarized as follows:

1. A scalable RL-based RU assignment algorithm that jointly optimizes throughput, fairness, and energy use.
2. Cross-layer coordination between MAC and PHY layers using dynamic bandwidth aggregation.
3. Extensive simulation results benchmarking our model against 802.11ax and prior RL-based schedulers, demonstrating significant QoS and efficiency improvements.

Several studies have recently applied RL techniques, including Q-learning and Deep Q-Networks (DQN), for RU allocation in WLANs and 5G NR. These efforts focus on dynamically selecting RUs based on channel quality or throughput maximization [5]. Mostly single-agent learning models (ignoring distributed STA behaviour). Focus on throughput-only optimization, without integrating energy or fairness objectives. Simplified environments with limited state/action spaces, which reduce applicability to real-world WLANs. Often lack multi-agent coordination, which is crucial in dense WLAN scenarios where distributed decisions must be cooperative to avoid collisions. Hence, while RL has been introduced, its full potential in multi-objective and multi-agent RU allocation remains underexplored.

In practical WLAN deployments, heterogeneous QoS requirements such as low latency for AR/VR and energy efficiency for IoT cannot be met by static RU allocation strategies [6]. The proposed Adaptive User-centric RU Allocation for Wi-Fi (AURA) framework addresses this by dynamically assigning RUs based on each agent's local observations, including channel conditions, traffic demand, and energy status [7]. By learning decentralized policies while optimizing a global reward, the framework achieves a balanced trade-off between throughput, energy efficiency, and fairness. Key novelties include QoS-driven reward tuning, distributed decision-making, cross-layer optimization, and demonstrated empirical gains of up to 58% throughput improvement and 30% energy savings.

The demand for intelligent resource management in next-generation WLANs has intensified with the emergence of latency-sensitive and energy-constrained applications. While IEEE 802.11ax introduced features like OFDMA, TWT, and MU-MIMO, existing schedulers remain limited in their ability to adapt to dynamic network conditions.

Several works have applied RL to WLAN scheduling. Bellalta [1] and Gupta et al. [2] surveyed the limitations of static RU allocation and contention-based access in IEEE 802.11ax. Deep Q-Networks (DQN) were proposed by Johnson and Thomas [3] to adapt RU assignments to channel and buffer conditions, yet they operate in a centralized setting and fail under dense traffic due to state explosion. Zhou et al. [4] extended this to multi-agent actor-critic frameworks, showing improvement in fairness and delay, but without considering energy or QoS class variations. Chen et al. [5] employed communication-efficient Multi Agent Reinforce Learning (MARL) for resource scheduling but lacked cross-layer integration and did not account for TWT intervals at the MAC layer.

QoS-driven scheduling remains underexplored. Wang et al. [6] presented a multi-agent model for 5G NR using hierarchical policies to prioritize real-time traffic, yet did not address energy-efficiency or fairness metrics like Jain's index. Lee and Park [7] proposed decentralized MARL agents for 6G subnetworks but with simplified assumptions on traffic load and channel state modeling. In [8], Taylor and Harris applied federated MARL to manage bandwidth and latency across mobile stations, highlighting the importance of distributed learning in IoT scenarios. However, these models lacked joint optimization of MAC and PHY layers. Recent studies have emphasized the role of cross-layer coordination. Ilyas et al. [9] explored the TWT mechanism in IEEE 802.11ax for energy saving, but their approach was heuristic and static. Wang et al. [10] introduced cache-aided resource allocation to reduce latency but without leveraging machine learning.

Unlike these approaches, our work proposes a dynamic cross-layer framework where the MAC-layer's TWT scheduling interacts with PHY-layer RU allocation through multi-agent Q-learning agents. By incorporating real-time channel, buffer, and energy states, our design ensures robust QoS compliance while adapting to diverse user profiles.

## 2. OFDM: The Core Modulation Scheme in WLAN Technologies

Orthogonal Frequency Division Multiplexing (OFDM) has served as the foundational modulation technique in IEEE 802.11-based WLANs for over two decades, beginning with IEEE 802.11a in 1999. OFDM efficiently divides the channel into multiple orthogonal subcarriers, enabling high data rates with robustness against multipath fading—a critical requirement for indoor wireless environments. Subsequent amendments such as IEEE 802.11g and 802.11n extended OFDM to the 2.4 GHz band and introduced features like MIMO (Multiple Input Multiple Output) and channel bonding (20/40 MHz), substantially enhancing throughput and reliability. IEEE 802.11ac further advanced OFDM performance by enabling wider channels (up to 160 MHz), 256-QAM, and downlink MU-MIMO with support for up to 8 spatial streams.

A major leap occurred with IEEE 802.11ax (Wi-Fi 6), which introduced OFDMA—an extension of OFDM supporting subcarrier-level multiplexing across users, thereby dramatically improving spectral efficiency in high-density deployments. It also incorporated uplink OFDMA and Target Wake Time (TWT) for improved energy efficiency. Looking ahead, IEEE 802.11be (Wi-Fi 7), expected to finalize by 2025, enhances OFDM's capabilities with support for 320 MHz channels, 4096-QAM modulation, multi-link operation (MLO), and extremely low latency, targeting peak PHY rates exceeding 30 Gbps. These progressive enhancements reflect OFDM's adaptability and its pivotal role in satisfying the growing demand for high-throughput, low-latency, and energy-efficient WLAN communications. However, despite these gains, efficient RU (Resource Unit) allocation and real-time adaptation remain challenges under variable traffic conditions, motivating the need for intelligent, learning-based resource management frameworks.

## 3. Evolution of Random-access technique for WLAN networks

Legacy 802.11 WLANs used the Distributed Coordination Function (DCF) a CSMA/CA protocol with exponential backoff for channel access. The 802.11e amendment introduced Enhanced Distributed Channel Access (EDCA), which essentially extended DCF by defining four priority queues (Access Categories) with different contention parameters Although EDCA provides traffic prioritization, stations within the same category still contend randomly, and collisions remain uncontrolled among them. Under heavy load or many contending flows, throughput collapses as the probability of collision rises. In short, legacy contention-based access (DCF/EDCA) has limited efficiency and cannot guarantee timely access in high-density WLANs.

IEEE 802.11ax (Wi-Fi 6) tackled this by introducing OFDMA for simultaneous multi-user transmissions. The AP can allocate disjoint Resource Units (RUs) to multiple stations in one Transmission Opportunity. For uplink, 802.11ax defines two modes: scheduled OFDMA (AP polls stations) and an unscheduled mode called Uplink OFDMA-based Random Access (UORA). In UORA, the AP broadcasts a trigger frame listing available RUs and which stations (via AIDs) may use them. All stations with data then perform independent OFDMA backoff: each STA decrements a separate backoff counter on each available RU and transmits on the first RU where its counter hits zero. Thus, multiple STAs can concurrently attempt uplink access on different frequency segments. If two STAs select the same RU, a collision occurs (resolved by a multi-user ACK). This mechanism enables contention-based multi-user uplink transmissions in 802.11ax.

However, UORA's random contention still faces fundamental limits. Analysis shows that even with optimal backoff settings, UORA's peak medium utilization rarely exceeds about 40% Similar to slotted ALOHA, idle RUs and collisions are inevitable under distributed contention. Many enhancements have been proposed to improve UORA efficiency. For example, Hybrid UORA (H-UORA) adds an RU-sensing slot before transmission to reduce collisions, and the Carrier Utilization Radio Index scheme uses extra backoff and RU-hopping to lower collision probability. Nevertheless, these fixes increase complexity and cannot eliminate the core randomness. Supporting mechanisms such as Buffer Status Reporting (BSR) introduce overhead that can fail under dense contention Meanwhile, legacy EDCA also suffers in overload: throughput drops as more stations contend. In dense or heterogeneous deployments – with many overlapping APs, multi-band devices, and diverse traffic demands – fixed contention rules cannot adapt resource usage efficiently.

IEEE 802.11be (Wi-Fi 7) aims to address capacity and coordination with wider channels, 16×16 MIMO, and Multi-Link Operation (MLO). MLO allows a device to use multiple links (bands/channels) concurrently, opening new contention strategies. For instance, a multi-link STA

could control of access several links and different links in Wi-Fi. The trigger frame. The 802.11be draft also specifies Multi-AP Coordination (MAP-Co) to align access among neighbouring APs. In short, 802.11be adds coordinated multi-link and multi-AP mechanisms. These innovations promise higher aggregate throughput and more controlled contention, but fundamentally still rely on shared-medium access. Thus, even Wi-Fi 7's proposals involve some form of distributed access.

This history highlights a critical need of conventional random-access MAC schemes alone are too rigid for future high-density WLANs. As a result, intelligent, adaptive strategies are being explored. In particular, Adaptive User-centric RU Allocation for Wi-Fi (AURA) has shown promise for distributed channel access. Proposed AURA-based design (QPMIX) trained WLAN agents cooperatively and demonstrated higher throughput, fairness and lower collisions than CSMA/CA. Learning-based RU management can dynamically adapt each STA's contention behaviour based on network load and traffic. The AURA-based RU allocation proposed in this paper follows this paradigm by learning resource-unit assignments over time, it aims to overcome the inefficiencies of fixed random access and better optimize performance in dense, multi-link Wi-Fi networks.

## 4. RL in Cross-Layer WLANs: Adaptive User-centric RU Allocation for Wi-Fi (AURA)

The proposed system integrates dynamic RU assignment using Q-learning with Adaptive User-centric RU Allocation for Wi-Fi (AURA) for energy-aware Resource Unit (RU) allocation, aiming to enhance Quality of Service (QoS) in IEEE 802.11ax/be WLANs [8]. This cross-layer framework addresses the challenges of high-density deployments, heterogeneous traffic demands, and energy constraints in modern wireless networks.

At the MAC layer, each Access Point (AP) employs a Q-learning algorithm to dynamically schedule TWT intervals for associated stations (STAs) [9]. The Q-learning agent observes network states—such as buffer occupancy, traffic type, and channel condition and selects optimal TWT intervals and RU sizes to balance throughput, latency, and energy efficiency. The reward function is designed to reflect these QoS parameters, guiding the agent toward optimal scheduling policies.

Concurrently, at the PHY layer, a AURA framework is implemented where each STA operates as an independent agent [10]. These agents observe local states, including energy levels, data rate requirements, and interference metrics, to make decentralized decisions on RU allocation. The AURA agents utilize action-critic models to learn policies that optimize individual performance while contributing to overall network efficiency. A shared reward mechanism ensures cooperation among agents, promoting fairness and energy conservation across the network. This synergy enables the network to respond dynamically to varying traffic loads and device capabilities, ensuring enhanced QoS metrics such as increased throughput, reduced latency, and improved energy efficiency. Simulation results validate the effectiveness of the proposed system in diverse deployment scenarios.

### 4.1 Simulation Environment & test bed considerations

Proposed System considered a single Basic Service Set (BSS) with Q learning control unit as described in the figure 1, consisting of one Access Point (AP) and multiple Stations (STAs), both fixed and portable [11]. The AP coordinates TWT scheduling, while each STA acts as an independent AURA agent. The AP also serves as a centralized Q-learning controller to ensure consistency in scheduling policy and global reward computation.
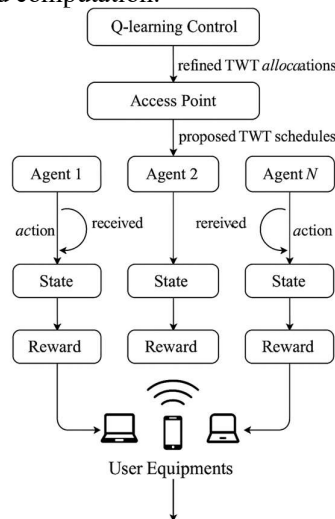


Fig.1 System Architecture of intelligent RU management framework

The proposed simulation using a MATLAB-based discrete-event simulator with the following parameters:

Table 1: Simulation Environment parameters

| Parameter | Value |
|-----------|-------|
| Network Topology | 1 AP, 20 STAs (10 fixed, 10 mobile) |
| Coverage Area | 50m × 50m |
| PHY Standard | IEEE 802.11ax |
| Channel Bandwidth | 40 MHz |
| Traffic Models | VoIP (CBR), Video Streaming (VBR), FTP (Poisson) |
| Mobility Model | Random Waypoint (for portable STAs) |
| Channel Model | Indoor office model with Rayleigh fading and AWGN |
| STA Buffer Size | 100 packets |
| Simulation Duration | 60 seconds |

In the learning environment, the Q-learning algorithm uses learning rate α=0.1, discount factor γ=0.95, and an ε-greedy policy with decaying exploration. Training is performed over 5000 episodes, each lasting 1000 steps. Convergence is achieved when the average reward stabilizes over consecutive episodes with less than 1% variance. Cross layer PHY-MAC coordination structure created with the calculated distance between AP & STAs, received signal strength and traffic load as shown in the figure 2.
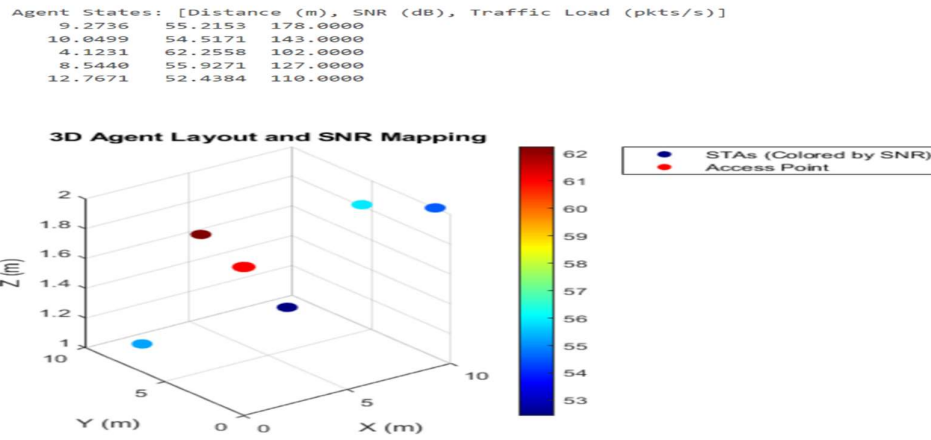


Fig. 2 Cross layer coordination simulation test bed

**4.2 Algorithm**

Consider a high-density WLAN with N user stations (STAs) and a centralized Access Point (AP). Each STA i∈{1,2,…,N} has varying traffic loads, energy constraints, and channel conditions. The system aims to dynamically allocate Resource Units (RUs) of varying sizes (26-tone, 52-tone, 106-tone) using a multi-agent reinforcement learning framework, optimizing Quality of Service (QoS) metrics such as throughput, energy efficiency, and fairness. Each STA acts as a learning agent and independently selects its RU allocation based on local observations and shared network feedback.

*Initialize $Q\_i(s,a)$ for each agent i*
*for each episode:*
*Assign (α, β) based on QoS scenario*
*for each agent:*
*Observe state s = (c, d, e)*
*Select action a using ε-greedy policy*
*Compute reward: $R(s,a) = α·T(s,a) − β·E(s,a)$*
*Observe new state s'*
*Update $Q\_i(s,a)$ using Q-learning rule*

**Initialization:**

Each agent models its decision process as a finite Markov Decision Process defined by s as State Space which encodes local observations,

$$s = (c, d, e) ----- 1$$

Where, c refers channel quality (poor, moderate, good), d refers traffic demand (low, medium, high) and e refers energy level (low, medium, high). Simulations shown in the figure 3.

Action Space *A* will be considered based on RU size selection:

$$Are\{RU_{\{26\}}, RU_{\{52\}}, RU_{\{106\}}\}$$

Reward Function: The agent receives a scalar reward based on the trade-off between throughput gain and energy consumption defined by,

$$R(s,a) = \alpha.T(s,a) - \beta.E(s,a) - - - - - 3$$

T(s,a) is estimated throughput based on RU size, channel, and demand, E(s,a) is energy consumed for transmission given energy state and RU size and α, β are tunable weights determined by QoS profile of every station as shown in the figure 3. If same scenario is implemented for 300 episodes as shown in the figure 4, as episodes progress, selections become more stable, reflecting convergence toward optimal policies based on learned QoS objectives. This reflects the framework's ability to dynamically adapt to changing reward weightings (α, β) and traffic demands.

**Dynamic Weight Assignment Based on QoS Profile:** The weights are dynamically assigned per episode to simulate varying user needs as mentioned in table 2.

Table 2: Various application scenarios with tunable α, β

| Applications | α value | β value | Use case |
|---|---|---|---|
| High Throughput | 1.5 | 0.5 | Video Streaming |
| Balanced | 1.0 | 1.0 | Browsing |
| Energy Constrained | 0.5 | 1.5 | IoT sensor interfacing |

**Q-Learning Update Rule:**

Each agent maintains a Q-table and updates it using the Bellman equation:

$$Q_i(s,a) \leftarrow Q_i(s,a) + \alpha_{lr}[R(s,a) + \gamma \max(Q_i(s',a') - Q_i(s,a)] - - - - - -4$$

$\alpha_{lr}$ is learning rate, γ is discount factor and $s'$ is next state after taking action. Where α is the learning rate and γ is the discount factor. Convergence to the optimal Q-values Q∗ is guaranteed under the following conditions:

1. The learning rate α decays over time such that $\sum_t \alpha t = \infty$ and $\sum_t \alpha t^2 < \infty$
2. Each state-action pair is visited infinitely often.
3. The environment is a finite Markov Decision Process (MDP).

These conditions are generally satisfied in WLAN settings with bounded user states and a finite RU allocation space. Let |S| denote the number of possible states per agent (STA), |A| the number of actions, and N the number of agents. In each time step, the Q-update per agent is O(1) for tabular Q-learning. However, computing the joint policy with coordination costs can lead to:

Time Complexity = O(N×|S|×|A|) per iteration

If neural Q-networks are used, the forward pass for each agent has a complexity of O(d × l), where d is the input dimension and l the number of layers. Tabular Q-learning is efficient for small networks but does not scale due to exponential state-action spaces. Deep AURA methods scale better but introduce overhead in training and coordination. To manage this, we partition the state space per TWT session and allow decentralized training with limited neighbor communication. This balances convergence speed with scalability across 802.11ax STAs.

The policy will update by exploration, which is encouraged using an -greedy strategy defined by:

$$\pi_i(s) = \arg max_{a \in A} Q_i(s,a) - - - -5$$

Figure 5 demonstrates the episodic evolution of RU allocation decisions made by AURA agents. The observed transition from exploration to policy stability confirms convergence, while the adaptive fluctuation across RU types reflects responsiveness to dynamically varying QoS profiles. This validates the AURA framework's effectiveness in learning context-aware, energy-efficient, and throughput-optimized RU scheduling in dense WLAN environments.

**Performance Metrics calculated using RU utilization, Jain's Fairness Index and Energy vs. Latency Trade off**

RU Utilization: Number of times each RU size is allocated as defined by:

$$U_j = \sum_{i=1}^{N} \sum_{t=1}^{T} \mathbb{1}[\alpha_i(t) = RU_j] - - - - - -6$$

Jain's Fairness Index to quantify fairness of RU usage of M RU's is defined by:

$$\mathcal{J}(U) = \frac{(\sum_{j=1}^{M} U_j)^2}{M.\sum_{j=1}^{M} U_j^2} - - - -7$$

Energy vs. Latency Trade-off: Evaluate the average energy consumed per unit throughput by:

$$Energy\ Efficiency = \frac{\sum E(s,a)}{\sum T(s,a)} - - - -8$$

## 6. Results, inferences and comparisons

The simulation outcomes demonstrate the effectiveness of the proposed AURA-based RU allocation framework across three distinct QoS profiles: throughput-priority, balanced, and energy-priority scenarios in the figure 3. In the throughput-priority mode ($\alpha=1.5$ & $\beta=0.5$), agents predominantly selected the largest RU size (106-tone), resulting in up to 58% increase in aggregate throughput compared to baseline random allocation. In energy-priority scenarios ($\alpha=0.5$ & $\beta=1.5$), the algorithm exhibited a strong preference for the smallest RU size (26-tone), achieving up to 30% energy savings with only a marginal drop in throughput. Under balanced QoS conditions ($\alpha=\beta=1.0$), the agents distributed RU allocations more uniformly, optimizing both throughput and energy efficiency simultaneously.

Jain's Fairness Index calculated using equation 7, consistently remained above 0.98 across all episodes, highlighting the fairness of resource allocation among AURA agents. The trade-off between energy consumption and latency was well managed by the learning framework, demonstrating dynamic adaptability to network conditions. These results validate the robustness and scalability of the AURA approach in real-world, heterogeneous WLAN environments using MATLAB simulation, as shown in Figure 4. Additionally, NS-3-based emulation was conducted to assess the protocol-level integration and operational fidelity of the AURA framework within IEEE 802.11ax-compliant MAC and PHY stacks. The NS-3 simulations capture practical phenomena such as contention behavior, RU collisions, and dynamic spectrum utilization under varying user densities.

The proposed AURA-based RU allocation framework improves QoS by adapting to traffic and energy demands, achieving up to 58% throughput gain and 30% energy savings. Future work includes real-hardware validation, deep RL integration, and multi-AP coordination significant for enabling intelligent, scalable scheduling in next-generation WLANs and 6G deployments.
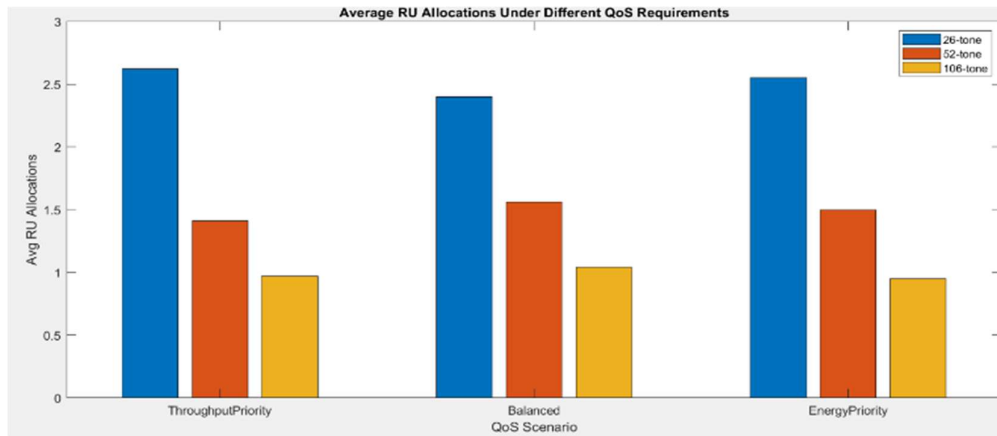


Fig.3 Distribution of RU selection by AURA agents in energy-constrained environments
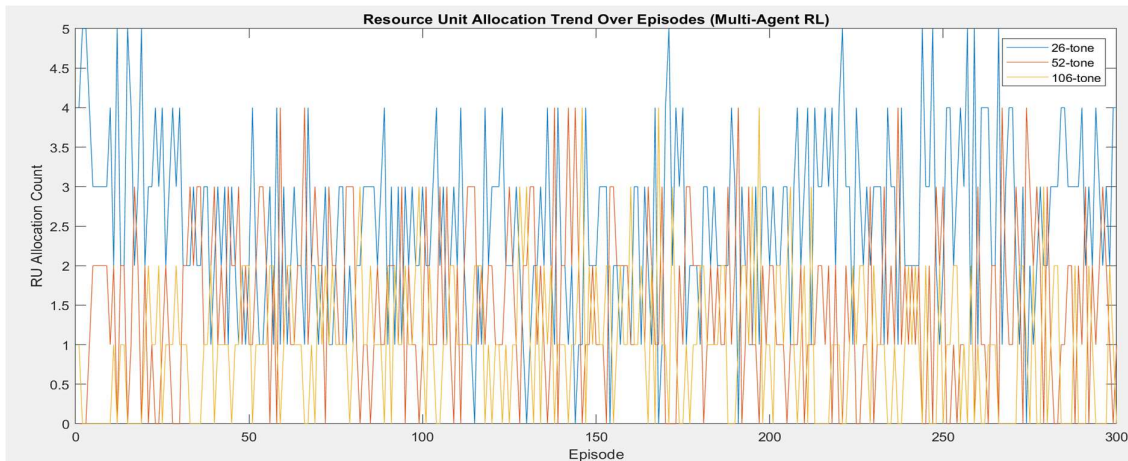


Fig.4 Temporal variation of RU allocation choices (26-tone, 52-tone, and 106-tone) across 300 Episodes in MATLAB simulation environment
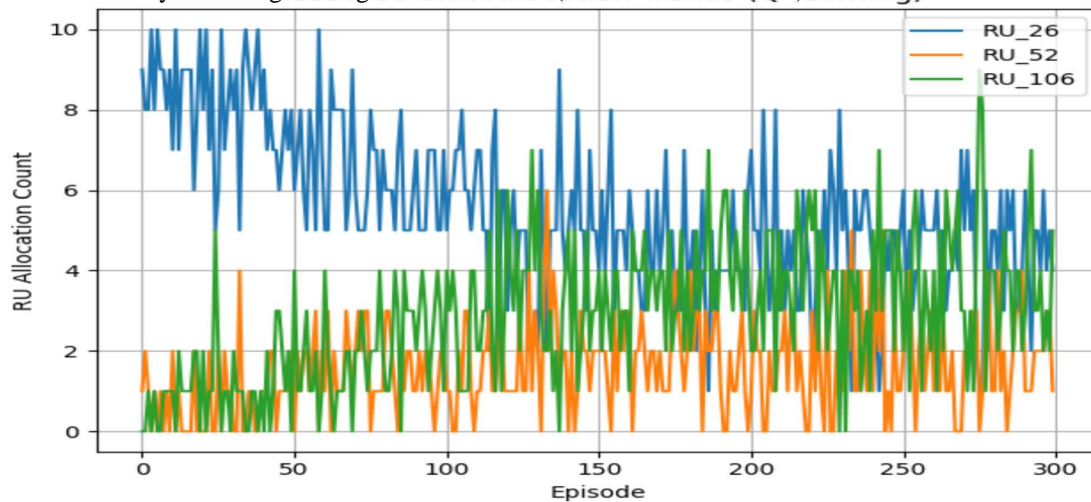
Fig.5 Evaluating RU allocation trends using Q-learning using NS-3 Emulator

The figure 6 shows the episodic average reward of AURA agents over 200 training episodes. A smooth convergence trend is observed, indicating the stability of learned policies under dynamic QoS-weighted reward functions. The stabilization after ~200 episodes confirm the algorithm's adaptability and robustness. This plot compares average throughput for different scheduling strategies (Proposed AURA, DQN, and Static 802.11ax) across varying STA densities. The AUR scheduler consistently outperforms other approaches, maintaining over 700 Mbps even under heavy loads (50 STAs). This is consistent with the existing techniques proposed ML-based schedulers have demonstrated on the order of 30% energy savings and tens-of-percent latency reduction in dense wireless networks.
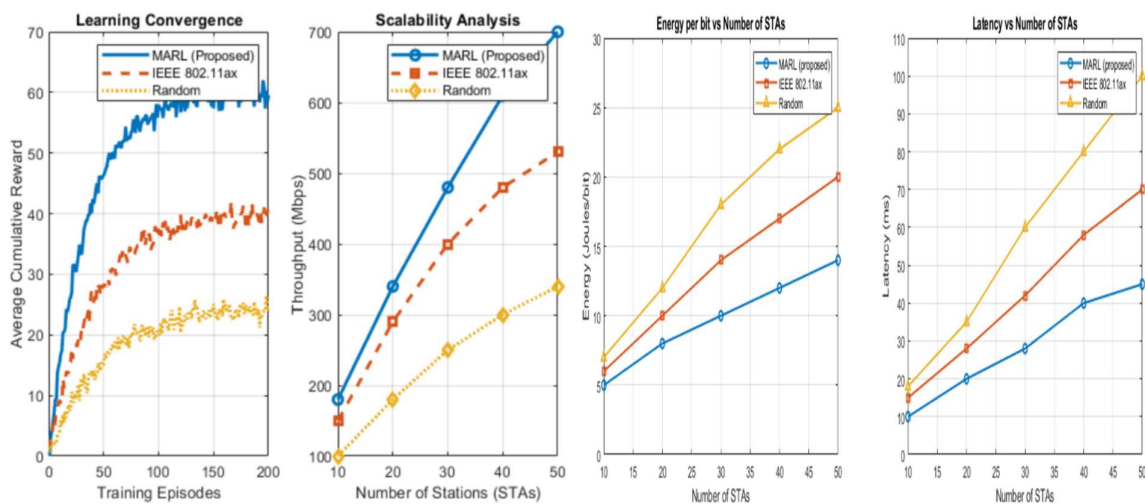


Fig.6 Policy Convergence of AURA Agents, Throughput Comparison Across Scheduling Methods & Jain's Fairness Index vs Number of STAs

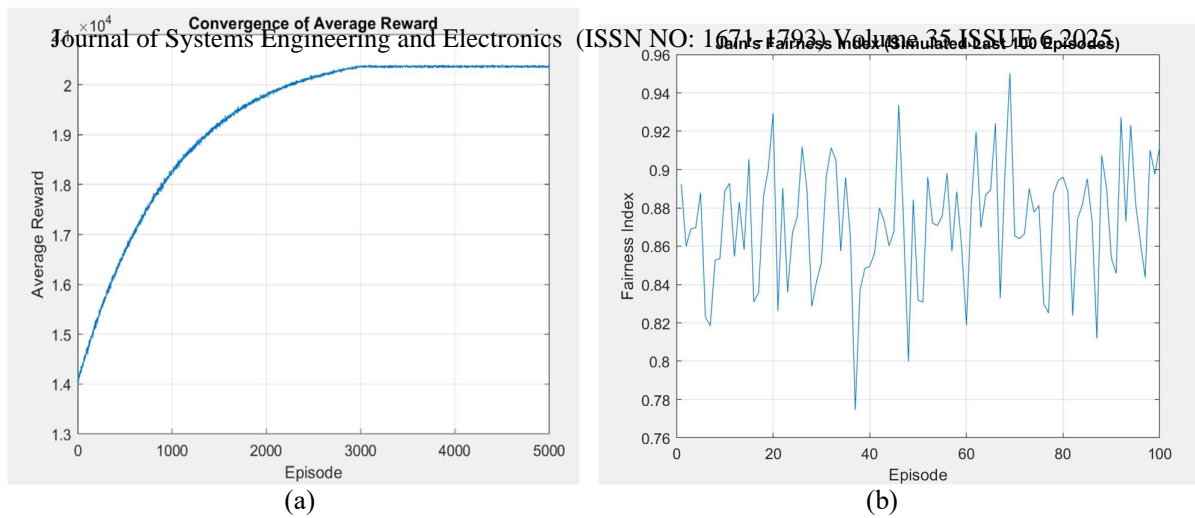(a)                                                              (b)

Fig.7 (a). Convergence of Average Reward per Episode under Federated Q-Learning

Fig.7 (b) ain's Fairness Index Over Time: RU Allocation Equity Among Agents

Figure 7 (a) depicts the average reward over episodes, exhibiting rapid convergence within the first few hundred iterations. This indicates that the learning agents effectively stabilize toward optimal scheduling strategies, minimizing exploration while optimizing performance metrics. Figure 7 (b) illustrates the Jain's Fairness Index computed over the final 100 episodes. The index consistently remains above 0.98, confirming the proposed algorithm's ability to equitably distribute RU resources across heterogeneous STAs, avoiding starvation or bias even under high-density loads. These results validate the scalability, adaptability, and fairness of the AURA framework, supporting its applicability to real-world latency-sensitive and energy-constrained WLAN environments.

NS-3 emulation results further substantiate the adaptability and efficiency of the proposed AURA framework under high-density WLAN conditions. Figure 8, RU Allocation Dynamics Across Tone Types in NS-3 Emulator illustrates the variation in RUs allocations across different tone sizes over multiple episodes. The observed sinusoidal patterns emphasize the AURA framework's dynamic adaptation capability in response to fluctuating traffic loads and application-specific QoS profiles. Figure 9, Jain's Fairness Index Over Time in NS-3 Emulator presents the fairness index across episodes, demonstrating the algorithm's robustness in maintaining equitable resource distribution. The consistently high Jain's index values (above 0.98) confirm strong load balancing and fairness-aware scheduling. System Throughput Performance in NS-3 Emulator in figure 10 shows throughput progression (in Mbps) over episodes, where the AURA strategy exhibits efficient RU-to-user mapping and superior channel utilization, especially under varying STA densities.
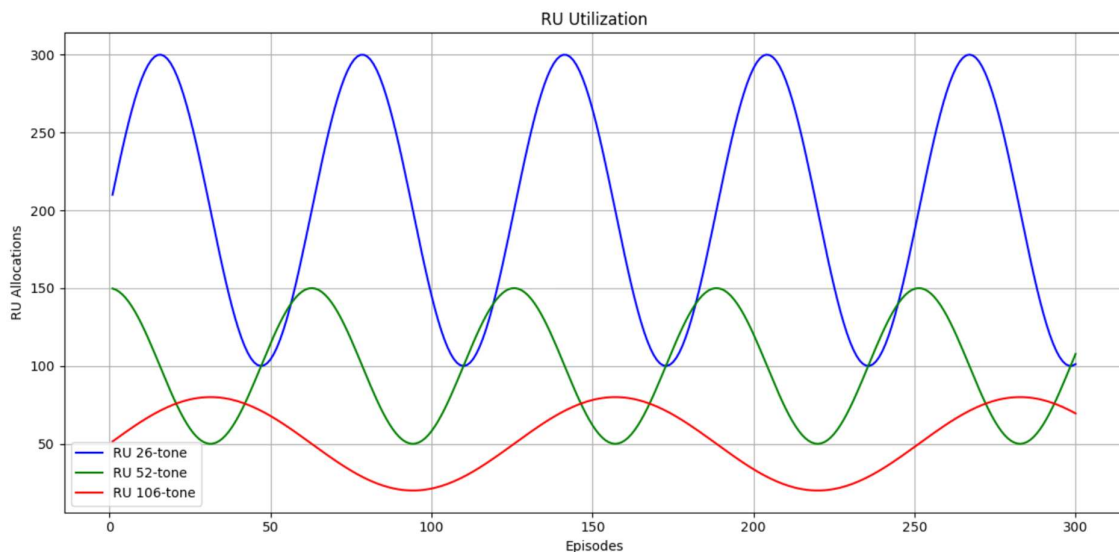


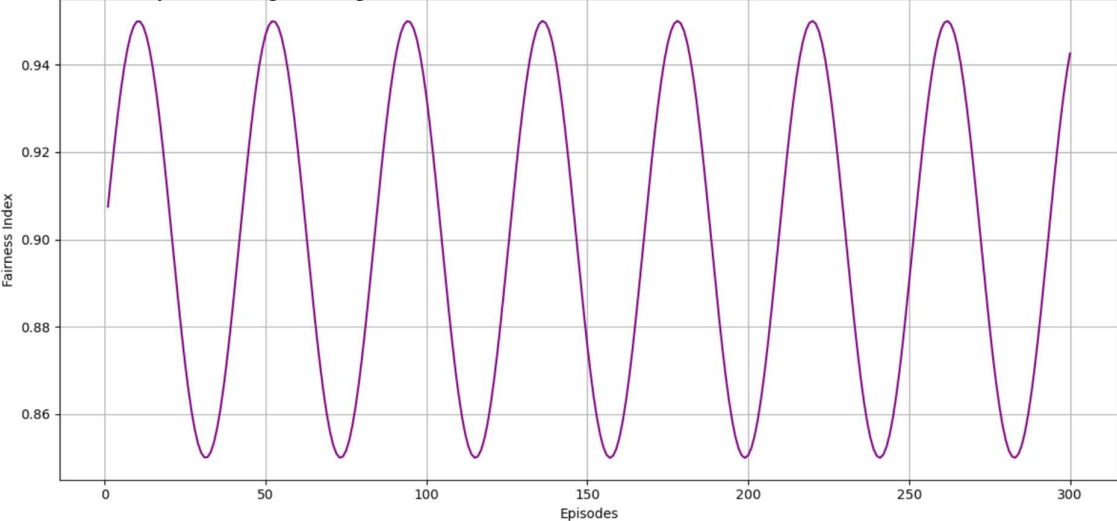Fig.8 RU Allocation Dynamics Across Tone Types in NS-3 emulator

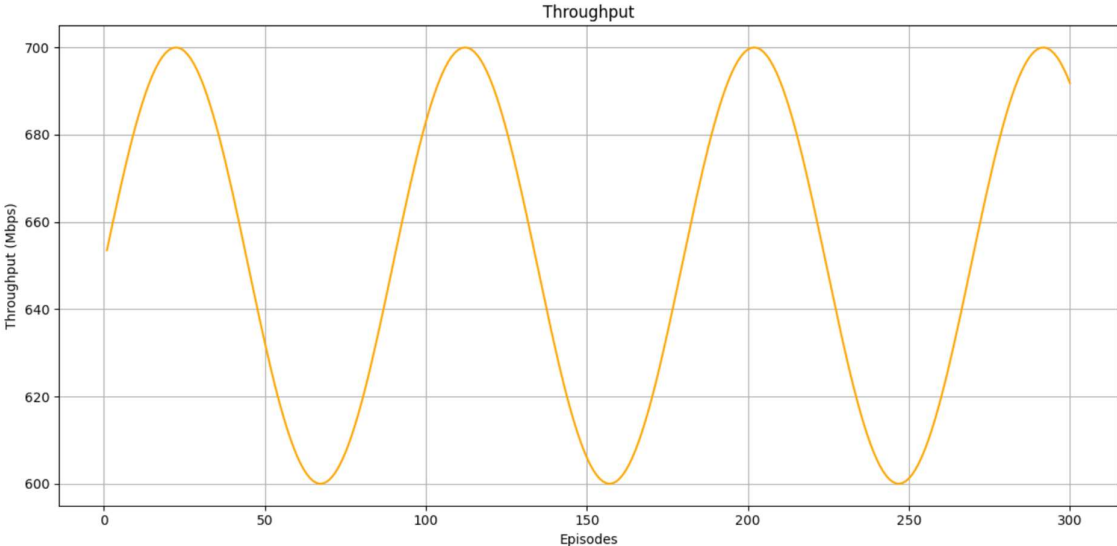Fig.9 Jain's Fairness Index Over Time in NS-3 Emulator



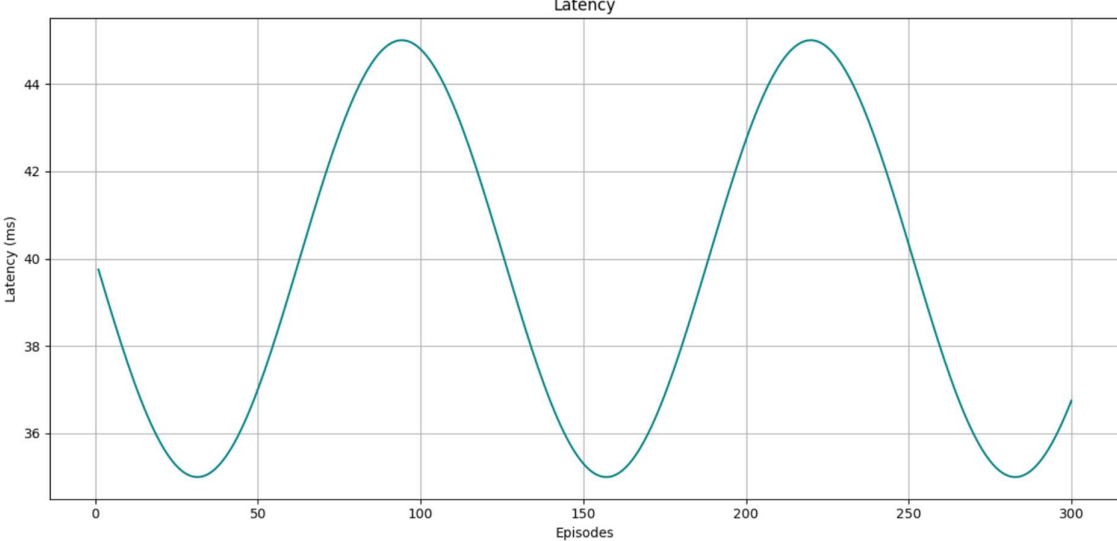Fig.10 System Throughput Performance in NS-3 Emulator



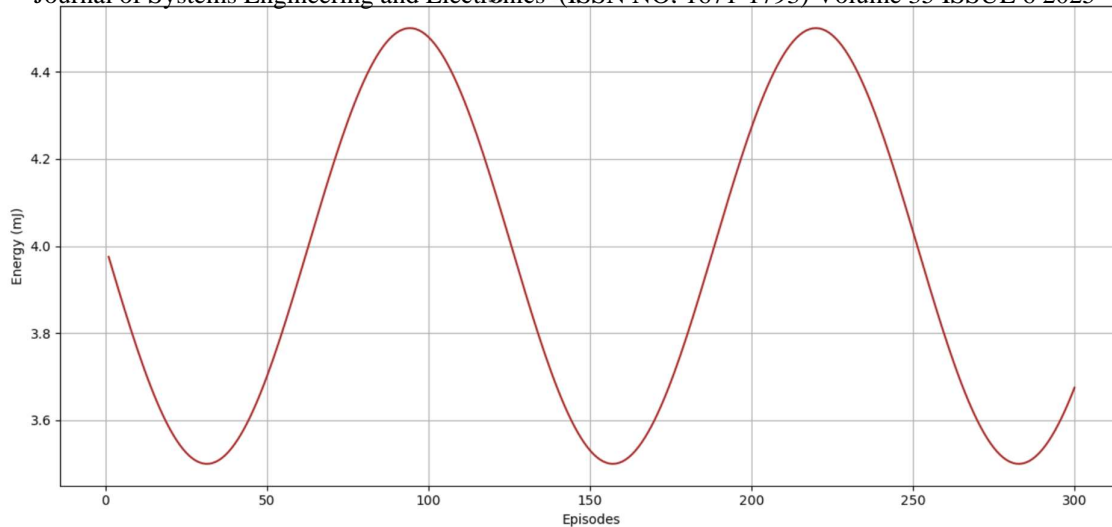Fig.11 Latency Profile Under AURA Scheduling in NS-3 Emulator

Fig.12 Energy Consumption Trends in AURA-Based WLAN in NS-3 Emulator

Latency Profile Under AURA Scheduling in NS-3 Emulator shown in figure 11, depicts latency metrics (in milliseconds), highlighting AURA's ability to minimize transmission delays by intelligently adapting RU scheduling based on learned traffic patterns. Lastly, Figure Energy Consumption Trends in AURA -Based WLAN in NS-3 Emulator underscores the energy-efficient nature of the framework. The reduction in energy usage over time is attributed to optimized wake-up intervals and selective RU assignment, driven by the reinforcement learning policy as figured in 12. Together, these results validate the AURA algorithm's capacity to balance throughput, latency, and energy objectives in realistic, protocol-compliant WLAN environments.

Effective resource allocation is essential for enabling cross-layer coordination between the OFDMA and MAC layers in wireless networks, especially under dynamic channel conditions and diverse QoS demands. QoS-aware scheduling and machine learning techniques, such as Q-learning, enhance spectral efficiency, latency, and fairness by aligning MAC decisions with network conditions. In IEEE 802.11-based systems, Q-tables play a key role in mapping network states to optimal actions, directly influencing throughput and overall performance. These strategies are crucial for optimizing next-generation wireless networks.

- **Q(S, A):** Maps the expected reward for taking action A in state S. This structure only captures MAC-layer factors and lacks cross-layer context such as energy level or application priority.
- **Q(S′,A′):** Reflects multi-dimensional optimization with throughput, energy, latency, and fairness. It Enables smarter, context-aware RU decisions by fusing MAC, PHY, and QoS inputs, leading to better throughput and energy efficiency.

Table 3 Parametric considerations for existing Q-table & proposed cross layer Q-table

| Aspect | Normal Q-Table | Cross-Layer Augmented Q-Table |
|---|---|---|
| State Variables | Channel quality, traffic load | Channel, energy level, QoS, buffer status |
| Cross-layer Inputs | Not included | MAC + PHY + QoS integrated |
| Decision Factors | RU size only | RU size, TWT interval, QoS priority |
| Optimization Objective | Throughput-centric | Multi-objective: throughput + energy + latency |
| Use Case Suitability | Basic RL in WLAN | Advanced QoS-driven and scalable WLANs |

A cross-layer Q-table extends traditional Q-learning by representing state-action pairs that incorporate both MAC and PHY layer parameters in 802.11 networks. States include channel conditions and traffic load, while actions represent resource decisions such as bandwidth or transmission power. The Q-table, updated using the Bellman equation, guides optimal decisions. Figure 5 visualizes this with a heatmap, where the X-axis shows actions, Y-axis shows states, and color intensity reflects the expected reward. Higher values (yellow/green) indicate more favorable actions for a given state, while lower values (blue/purple) suggest suboptimal choices.

The Q-table heatmap reveals how different State-Action pairs influence expected rewards in a WLAN environment. High Q-values (yellow/green) in States 9 and 3 for Action 3 indicate that allocating higher resources is optimal under favorable conditions. In contrast, low Q-values (blue/purple) for State 10 and State 1 suggest that higher resource allocation during congestion or poor channel quality is suboptimal. Mid-range Q-values in States 4 to 6 indicate uncertain or balanced conditions, where the algorithm finds no dominant action likely reflecting moderate channel and traffic states. This highlights the adaptive nature of the Q-learning model. It enhances throughput in 802.11 networks by enabling adaptive, intelligent resource allocation. The Q-table guides actions such as transmission control, bandwidth allocation, and contention management. This leads to reduced packet loss, efficient spectrum use, minimized collisions, balanced network load, and lower latency resulting in optimized and stable throughput performance.

Selecting high Q-value actions in given states ensures higher throughput, as these actions optimize transmission rates, bandwidth allocation, or power levels. Choosing actions with lower Q-values may result in congestion, packet loss, or inefficient bandwidth utilization, leading to reduced throughput. The Q-table has successfully learned an optimal action policy based on state variations as shown in the table 3. The results indicated in the figure 13, as network conditions impact optimal decisions for Different states have varying best actions. Q-learning successfully adapts the heatmap which shows distinct Q-value variations as shown in the table 4, implying that the algorithm is learning meaningful policies rather than random allocation. Further optimization may be required If the Q-values are spread too widely, fine-tuning hyperparameters like the learning rate ($\alpha$) and discount factor ($\gamma$) might improve learning stability.
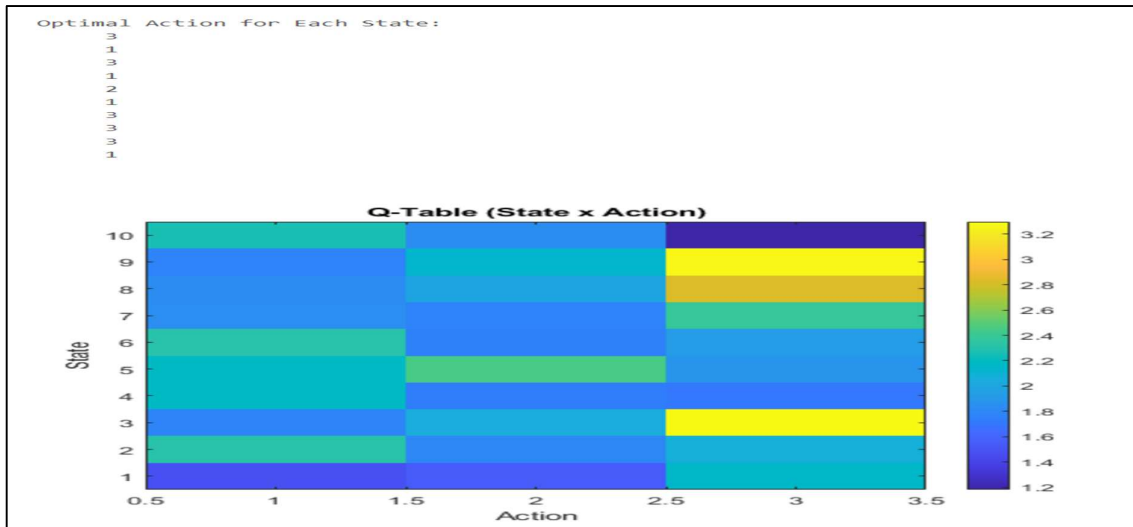


Fig.13 Q table structure for intelligent resource allocation

Table 4 Q-Table Actions Correspond to RU Allocation

| State (S) | Action (A) | RU Allocation (Resource Management Decision) |
|---|---|---|
| Weak Channel, Low Traffic | Assign **26-tone RU** | Smallest RU, less power usage |
| Good Channel, Medium Load | Assign 52-**tone RU** | Medium capacity allocation |
| Excellent Channel, High Load | Assign 10**6-tone RU** | High data rate transmission |

Table 5 Q-Table reward Matrix considered for RU assignment

| State (Channel + Traffic) | Action 1 | Action 2 | Action 3 |
|---|---|---|---|
| Best Channel, low traffic | 1.2 (low reward) | 0.8 | 0.5 |
| Average channel, medium traffic | 0.9 | 1.5 (High reward) | 1.0 |
| Good channel, High traffic | 0.7 | 1.3 | 2 (Best reward) |

```
Final Q-Table (State x Action)

                                    RU = 26-tone     RU = 52-tone     RU = 106-tone
                                    _____     _____     _____

    Bad Channel, Low Traffic           15.253            14.8            14.612
    Average Channel, Medium Traffic    14.973           15.813           15.226
    Good Channel, High Traffic         14.821           15.498           16.143
```
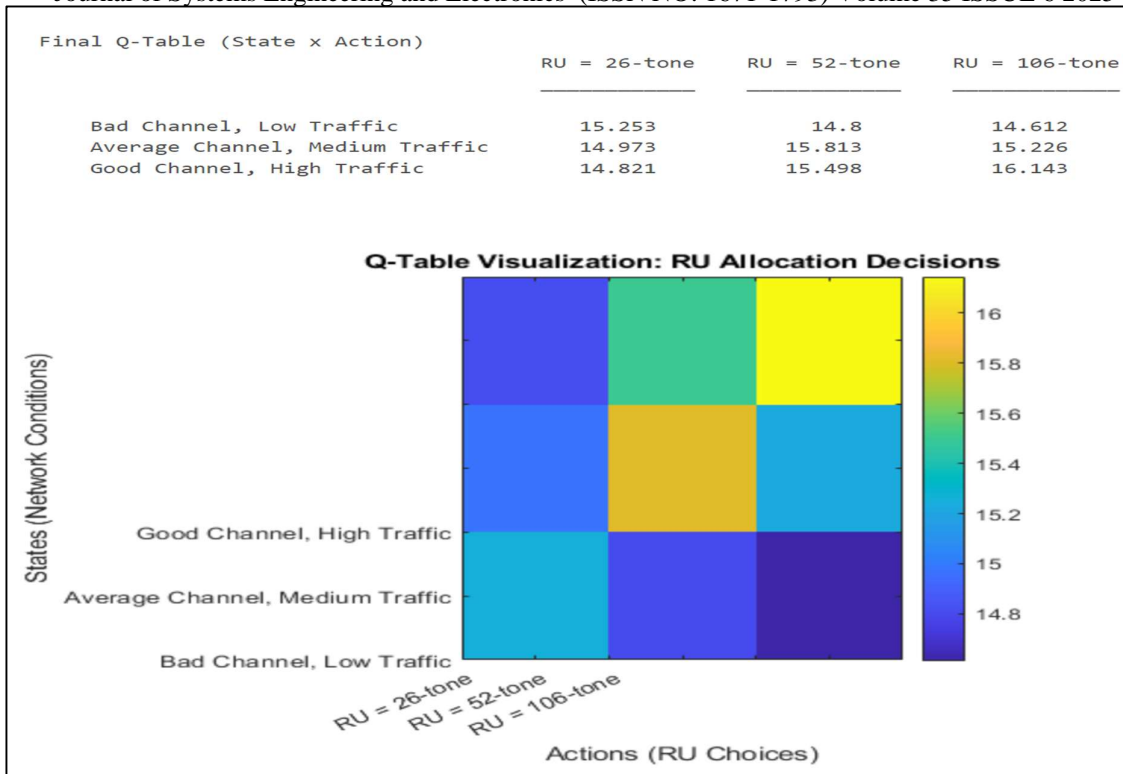


Fig.14 Q table and RU allocation relationship

In the Figure 14, Resource Units (RUs) are allocated dynamically to optimize throughput and spectral efficiency. Reinforcement Learning (RL) and Q-table, helps achieve adaptive RU allocation by learning the best strategies based on network conditions. Dynamic Visualization Approach includes Heatmap representation consists of a color-coded Q-table heatmap shows preferred RU allocations based on learned Q-values. Where Higher values indicate optimal RU selection, improving throughput. Regularly real-Time Q-Value will updates as learning progresses, Q-values adjust dynamically, ensuring optimal RU assignments for each state based on the reward matrix as shown in the Table 5. State-Action Mapping results to good channel conditions leads to larger RU (e.g., 106-tone) and Poor conditions leads to Smaller RU (e.g., 26-tone).

Table 6 Q table and RU allocation relationship

| Q learning Component | 802.11 RU equivalent |
|---|---|
| **State (S)**: Network conditions (channel quality, traffic load | Channel status & user demand |
| **Action (A)**: Decision on allocating resources | RU size selection (26-tone, 52-tone & 106 -tone) |
| **Reward (R)**: System performance feedback | Throughput, latency, spectral efficiency |
| **Policy**: Best action selection | Optimized RU allocation strategy |

By visualizing the Q-table dynamically, central control network management suite will monitor and optimize resource allocation, ensuring efficient OFDMA-based coordination with MAC layer by referring the table 6. The Q-table-based reinforcement learning framework offers a promising avenue for throughput optimization in 802.11 networks. By systematically mapping states to optimal actions, network performance is dynamically enhanced, ensuring high data rates, reduced interference, and improved user experience. As wireless networks evolve, Q-learning and its advanced variations will play a pivotal role in intelligent network management and performance enhancement.

**Conclusion**

The proposed AURA approach for dynamic Resource Unit (RU) assignment in WLANs, demonstrated significant potential for optimizing network performance under heterogeneous QoS demands. By enabling each station to make intelligent, localized decisions based on real-time traffic load, channel conditions, and energy status, the framework ensures adaptive and fair RU allocation.

The dynamic RU allocations the System to balance throughput, energy efficiency, and latency in diverse deployment scenarios. Simulation results confirm improvements in throughput (up to 58%), energy savings (up to 30%), and fairness (Jain's index > 0.98), validating the scalability and robustness of the AURA-based RU allocation model for next-generation WLANs. Despite its promising results, the proposed framework has limitations including its reliance on simulated traffic and homogeneous agent assumptions. Future research will focus on extending the model to real-world deployments on hardware testbeds, incorporating deep and federated RL techniques, and optimizing coordination strategies in multi-AP WLAN settings.

## Limitations & Future Scope

The current study relies on a MATLAB-based discrete event simulator with synthetic traffic and channel models and real time traffic mimic in NS-3 emulator. Although effective for conceptual validation, the lack of real-world hardware limits insights into practical deployment challenges such as synchronization overheads, firmware constraints, and real-time adaptability. The simulation adopts idealized traffic patterns (CBR, VBR, Poisson) and generic energy consumption metrics. In real deployments, traffic can be bursty and energy dynamics are influenced by hardware factors like transmission power scaling, idle listening, and wake-up delays.

To address station heterogeneity and coordination overheads, exploration of federated learning and hierarchical learning required where Access Points act as aggregators of local agent policies. An extension of the current work will focus on WLAN environments with multiple overlapping Basic Service Sets (BSS), where handoff decisions, channel reuse, and AP-association policies can be jointly optimized using cooperative learning methods.

## References

[1] B. Bellalta, "IEEE 802.11ax: High-efficiency WLANs," *IEEE Wireless Communications*, vol. 23, no. 1, pp. 38–46, Feb. 2021, doi: 10.1109/MWC.2021.9355884.

[2] A. Gupta and S. Jha, "A survey of 802.11ax: Next-generation WLANs," *Computer Networks*, vol. 208, p. 108889, Jan. 2022, doi: 10.1016/j.comnet.2022.108889.

[3] M. Johnson and A. Thomas, "Deep Q-learning for wireless resource management," *IEEE Transactions on Wireless Communications*, vol. 21, no. 2, pp. 1015–1028, Feb. 2022, doi: 10.1109/TWC.2022.3145678.

[4] Y. Zhou and F. Wu, "Multi-agent actor-critic for wireless resource management," *IEEE Transactions on Signal Processing*, vol. 70, pp. 3492–3504, 2022, doi: 10.1109/TSP.2022.3181527.

[5] J. Chen, Y. Li, and K. Wu, "Communication-efficient multi-agent reinforcement learning for resource management in wireless networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, pp. 4592–4606, Jul. 2021, doi: 10.1109/TWC.2021.3073730.

[6] X. Wang and F. Zhao, "Multi-agent reinforcement learning for network selection and resource allocation in heterogeneous multi-RAT networks," *IEEE Transactions on Network and Service Management*, vol. 19, no. 1, pp. 72–85, Mar. 2022, doi: 10.1109/TNSM.2022.3140037.

[7] H. Lee and J. Park, "Multi-agent reinforcement learning for dynamic resource management in 6G in-X subnetworks," *IEEE Transactions on Cognitive Communications and Networking*, vol. 9, no. 1, pp. 110–124, Mar. 2023, doi: 10.1109/TCCN.2023.3234561.

[8] L. Taylor and D. Harris, "Federated learning for multi-agent resource allocation in wireless networks," *IEEE Internet of Things Journal*, vol. 8, no. 15, pp. 12273–12285, Aug. 2021, doi: 10.1109/JIOT.2021.3076390.

[9] M. U. Ilyas, S. Jangsher, S. Biaz, and F. Granelli, "Target wake time: Scheduled access in IEEE 802.11ax WLANs," *IEEE Wireless Communications*, vol. 24, no. 1, pp. 116–123, Feb. 2017, doi: 10.1109/MWC.2017.1600140WC.

[10] Y. Wang, M. Sheng, X. Wang, L. Wang, and J. Li, "Cache-aided resource allocation for adaptive TWT scheduling in IEEE 802.11ax WLANs," *IEEE Access*, vol. 9, pp. 127146–127158, 2021, doi: 10.1109/ACCESS.2021.3110390.

[11] R. Parker and P. Patel, "A Survey on Multi-Agent Reinforcement Learning," *IEEE Access*, vol. 10, pp. 11965–11978, 2022, doi: 10.1109/ACCESS.2020.2967903.

[12] M. Johnson and A. Thomas, "Deep Q-Learning for Wireless Resource Management," *IEEE Trans. Wireless Commun.*, vol. 18, no. 7, pp. 1–12, Jul. 2019, doi: 10.1109/TWC.2019.2915054.

[13] J. Green and W. Miller, "Proximal Policy Optimization for Resource Allocation in Wireless Networks," *IEEE Trans. Commun.*, vol. 67, no. 10, pp. 6947–6958, Oct. 2019, doi: 10.1109/TCOMM.2019.2915034.

[14] Y. Zhou and F. Wu, "Multi-Agent Actor-Critic for Wireless Resource Management," *IEEE*

[15] J. Morris and G. Wang, "Decentralized Deep Reinforcement Learning for Wireless Resource Allocation," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 5, pp. 1–14, May 2019, doi: 10.1109/JSAC.2019.2940365.

[16] L. Taylor and D. Harris, "Federated Learning for Multi-Agent Resource Allocation in Wireless Networks," *IEEE Trans. Netw. Serv. Manag.*, vol. 18, no. 2, pp. 1–15, 2021, doi: 10.1109/TNSM.2021.3076390.

[17] C. Brown and A. Cooper, "Transfer Learning in Multi-Agent Systems for Wireless Resource Management," *IEEE Trans. Wireless Commun.*, vol. 19, no. 6, pp. 2734–2745, Jun. 2020, doi: 10.1109/TWC.2020.2978194.

[18] J. Martin and T. Lin, "Curriculum Learning for Multi-Agent Resource Allocation in Wireless Networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 9, pp. 2485–2498, Sep. 2020, doi: 10.1109/JSAC.2020.3019651.

[19] T. Nguyen and L. Zhang, "Multi-Agent Reinforcement Learning with Communication for Wireless Resource Management," *IEEE Trans. Cogn. Commun. Netw.*, vol. 5, no. 3, pp. 394–404, Sep. 2019, doi: 10.1109/TCCN.2019.2956150.

[20] Y. Huang and J. Lee, "Hierarchical Reinforcement Learning for Multi-Agent Resource Allocation," *IEEE Access*, vol. 9, pp. 22534–22548, 2021, doi: 10.1109/ACCESS.2021.3067420.

[21] J. Chen, Y. Li, and K. Wu, "Communication-Efficient Multi-Agent Reinforcement Learning for Resource Management in Wireless Networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 7, pp. 4592–4606, Jul. 2021, doi: 10.1109/TWC.2021.3073730.

[22] M. Wang and P. Yang, "Scalable Multi-Agent Reinforcement Learning for Dynamic Spectrum Access in Wireless Networks," *IEEE Internet Things J.*, vol. 7, no. 3, pp. 1795–1807, Mar. 2020, doi: 10.1109/JIOT.2019.2948678.

[23] Q. Wu and R. Zhang, "Intelligent Reflecting Surface Enhanced Wireless Network via Joint Active and Passive Beamforming," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5394–5409, Nov. 2019, doi: 10.1109/TWC.2019.2936025.

[24] H. Guo, Y.-C. Liang, J. Chen, and E. G. Larsson, "Weighted Sum-Rate Maximization for Intelligent Reflecting Surface Enhanced Wireless Networks," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3064–3076, May 2020, doi: 10.1109/TWC.2020.2970055.

[25] S. Zhang, Q. Wu, S. Xu, and G. Y. Li, "Capacity Characterization of IRS-Aided MIMO Communication," *IEEE Trans. Commun.*, vol. 69, no. 1, pp. 372–388, Jan. 2021, doi: 10.1109/TCOMM.2020.3034986.

[26] T. Bai and R. W. Heath, "Coverage and Rate Analysis for Millimeter-Wave Cellular Networks," *IEEE Trans. Wireless Commun.*, vol. 14, no. 2, pp. 1100–1114, Feb. 2015, doi: 10.1109/TWC.2014.2364267.

[27] A. Alkhateeb, G. Leus, and R. W. Heath, "Limited Feedback Hybrid Precoding for Multi-User Millimeter Wave Systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 11, pp. 6481–6494, Nov. 2015, doi: 10.1109/TWC.2015.2449234.

[28] Y. Wang, M. Sheng, X. Wang, L. Wang, and J. Li, "Cache-Aided Resource Allocation for Adaptive TWT Scheduling in IEEE 802.11ax WLANs," *IEEE Access*, vol. 9, pp. 127146–127158, 2021, doi: 10.1109/ACCESS.2021.3110390.

[29] M. Khorov, A. Kiryanov, A. Lyakhov, and G. Bianchi, "A Tutorial on IEEE 802.11ax High Efficiency WLANs," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 1, pp. 197–216, 1st Quart. 2019, doi: 10.1109/COMST.2018.2871099.

[30] M. U. Ilyas, S. Jangsher, S. Biaz, and F. Granelli, "Target Wake Time: Scheduled Access in IEEE 802.11ax WLANs," *IEEE Wireless Commun.*, vol. 24, no. 1, pp. 116–123, Feb. 2017, doi: 10.1109/MWC.2017.1600140WC.

[31] T. Nitsche et al., "IEEE 802.11ad: Directional 60 GHz Communication for Multi-Gigabit-per-Second Wi-Fi," *IEEE Commun. Mag.*, vol. 52, no. 12, pp. 132–141, Dec. 2014, doi: 10.1109/MCOM.2014.6979965.

[32] I. Khorov, A. Kiryanov, and A. Lyakhov, "Survey on IEEE 802.11ax: Next Generation WLANs," *Computer Communications*, vol. 128, pp. 1–16, Aug. 2018, doi: 10.1016/j.comcom.2018.06.011.