

Delay-aware Slicing and MAC Management using MCDA in IEEE 802.11 SD-RANs

Pedro H. Isolani^{*}, Daniel J. Kulenkamp[†], Johann M. Marquez-Barja[‡],
Lisandro Z. Granville[§], Steven Latré^{*}, and Violet R. Syrotiuk[†]

^{*}*Department of Computer Science, University of Antwerp - imec, Antwerp, Belgium*
{pedro.isolani, steven.latre}@uantwerpen.be

[†]*School of CIDSE, Arizona State University, Tempe, USA*
{dkulenk, syrotiuk}@asu.edu

[‡]*Department of Electronics - ICT, University of Antwerp - imec, Antwerp, Belgium*
johann.marquez-barja@uantwerpen.be

[§]*Institute of Informatics, Federal University of Rio Grande do Sul, Porto Alegre, Brazil*
granville@inf.ufrgs.br

Abstract—Advanced wireless services and applications are demanding lower latency and better reliability. In IEEE 802.11 networks, slicing abstractions at the Medium Access Control (MAC) layer can provide precise resource allocation and traffic isolation as a means to meet these performance requirements. In this paper, we propose a delay-aware approach for MAC management via airtime-based network slicing and user association using Multi-Criteria Decision Analysis (MCDA) in IEEE 802.11 Software-Defined Radio Access Networks (SD-RANs). To enable such an approach, we leverage Software-Defined Networking (SDN) to monitor queueing delay statistics at the Access Points (APs). We evaluate our approach in a testbed with two APs, controlled by SD-RAN and SDN controllers with four Stations (STAs) served by both Quality of Service (QoS) and Best-Effort (BE) flows dynamically. We compare our approach to a state-of-the-art user association approach. Our results show that in addition to load balancing flows across APs, our approach enhances the QoS delivery at runtime using slicing.

Index Terms—SDN, MAC management, network slicing, user association, IEEE 802.11 networks, SD-RAN, MCDA.

I. INTRODUCTION

Today’s advanced wireless services and applications are increasingly latency-sensitive [1]. Given the dynamic nature of the wireless environment, on-the-fly Medium Access Control (MAC) management is essential for efficient and reliable wireless communication [2]. To support traffic prioritization, IEEE 802.11e introduced the Enhanced Distributed Channel Access (EDCA) function [3]. EDCA defines access categories as priority levels, which determine how clients access the channel. Access categories with higher priority have more transmission opportunities than others and, once its clients reach the contention period, more time is reserved for transmission. However, EDCA does not ensure radio resource isolation among such access categories. Further, because EDCA is a distributed channel access method with no standard management interfaces, reliable connectivity cannot be guaranteed.

In contrast, network slices operate independently from one another and provide precise networking resources and traffic isolation among users and services [4]. The fact that slices can be dynamically instantiated, modified, and terminated implies a strong decoupling between software-based functions and the underlying network infrastructure. The 3GPP 5G architecture is embracing the Control/User Plane Split (CUPS) as one of the fundamental enablers for network programmability and End-to-End (E2E) network slicing [1]. In this context, Software-Defined Networking (SDN) has facilitated network innovation and simplified network management. Consequently, SDN-tailored systems are envisioned to ease the creation of logical and isolated networks via slice abstractions [5].

Current research efforts address network slicing and SDN to enhance resource utilization in IEEE 802.11 Radio Access Networks (RANs) [6]–[13]. There is extensive work on user association algorithms that focus on load balancing, mobility support, and fairness [14]–[16]. Such approaches often benefit from SDN to address resource management. Despite this work, deciding how to efficiently allocate, control, and manage users and slice resources remains challenging [17]. As well, latency-related constraints have not been considered.

In this paper, we propose a delay-aware approach for MAC management via airtime-based network slicing and user association using Multi-Criteria Decision Analysis (MCDA) in IEEE 802.11 Software-Defined Radio Access Networks (SD-RANs). By performing network slicing and network-triggered handovers, enhanced Quality of Service (QoS) delivery can be addressed. Besides, the ping-pong effect can be avoided. To enable such an approach, we monitor the queueing delay at the Access Points (APs). We evaluate our approach in a testbed with two APs, controlled by SD-RAN and SDN controllers, and four Stations (STAs) served by both QoS and Best-Effort (BE) flows dynamically. We compare our approach to a recent user association algorithm [16]. Our results show improved load balancing of flows across APs and QoS delivery.

II. RELATED WORK

Ensuring QoS in wireless networks is a longstanding research challenge that has only become more complex [18]. After the IEEE 802.11e amendment [3] established the foundations for traffic prioritization, many investigations started to focus on queue management [19]–[22]. Later, with the improvements provided by the IEEE 802.11n amendment [23], the focus shifted towards channel optimization and fairness [24]–[28]. Any solution requiring modifications to the driver (e.g., frame formats) becomes non-standard compliant.

Network slicing addresses precise infrastructure sharing, allowing reliable and improved QoS delivery [6]. Most consist of airtime-based Resource Allocation (RA) mechanisms for IEEE 802.11 network virtualization [7] [29] [30]. Airtime scheduling has been extensively studied as a means to overcome the well-known *IEEE 802.11 Performance Anomaly* [31]. Without such slicing capabilities, STAs equally share the available radio resources only if they experience the same or similar channel conditions. Otherwise, when an STA uses a lower bit rate, it results in performance degradation perceived by all.

Recent proposals [7]–[13] address network slicing in IEEE 802.11 networks. Richart et al. [7] propose a resource allocation mechanism to achieve infrastructure sharing and slicing on Wi-Fi APs. Later, Richart et al. [32] present an enhanced version of such scheduling with capacity limits, capable of achieving precise queueing delay for slices on an AP. In this manner, the aim is to achieve traffic isolation and precise network resources. However, the proposal was only assessed via simulation. Others [8]–[11] focus on practical implementations. However, runtime slice orchestration based on latency metrics is not addressed.

Coronado et al. [11] propose a framework that enables programmable and dynamic E2E network slicing over heterogeneous RANs. Deployed on a real-world testbed, slices and client traffic isolation have been evaluated through achieved throughput. However, given today’s stringent latency-related requirements, slice resource allocation requires further research. Vassilaras et al. [17] state that future wireless networks have to consider E2E latency requirements for a service chain where feasible slice embedding must also satisfy and guarantee the QoS requirements. Yet deciding how to efficiently control and manage such resources at runtime remains challenging.

On another front, user association algorithms focus on load balancing, mobility support, and fairness. They often take advantage of the SDN centralized view to control and manage the network [14]–[16]. By gathering the Received Signal Strength Indicator (RSSI) and throughput measurements, seamless handovers are triggered, improving the overall throughput of the network. However, latency-related metrics are not considered within their decision-making process.

In previous work [12] [13], we evaluate the impact of runtime slice reconfiguration on the E2E latency using ICMP. Subsequently, we integrate the queueing delay measurements into the formulation of a QoS optimization problem. In this work, we measure the queueing delay AP and use the results of

the SDN centralized monitoring to perform MAC management via network slicing and network-driven handovers and, hence, enhance the resource utilization and QoS delivered.

III. SYSTEM OVERVIEW

We consider a scenario where network slices are orchestrated by a logically centralized entity in an SDN-enabled network infrastructure. This entity has a global view of and control over all network resources. Multiple *tenants* (i.e., virtual operators or service providers) share the infrastructure and have their specific Service Level Agreements (SLAs). These SLAs are translated into QoS requirements for the network to support. To meet such requirements, we propose the use of network slicing. We focus on QoS within a slice as being a service, i.e., Quality of Service Slicing (QoSS), as defined by Richart et al. [7]. Figure 1 illustrates the system overview as a layered system model.

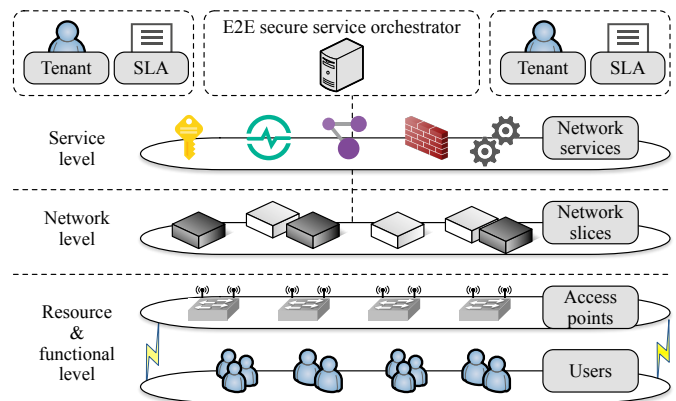


Fig. 1. System overview as a layered model.

The IEEE 802.11 RAN consists of a set APs responsible for delivering data from services down to users in the network, i.e., STAs. Each AP has resources to be shared and therefore to be managed properly. To represent the minimum chunk of wireless resources that can be assigned to a user, we use the *Resource Block abstraction* [33]. A resource block defines a Wi-Fi interface at a given AP, identified by the network interface identifier (i.e., MAC address), operating channel (e.g., 1, 6, 11), and the type of channel (e.g., High Throughput (HT) 20MHz, Very High Throughput (VHT) 40MHz). In our system, each AP has only one interface and therefore only one resource block; henceforth, resource blocks are synonymous with APs. Figure 2 shows a simplified queue structure along with the data traffic flow within an AP.

At the APs, frames from slices are classified into queues based on the definition of the traffic rules (e.g., Open-Flow rules) and are dequeued following the Airtime Deficit Weighted Round Robin (ADWRR) scheduling algorithm [11]. In ADWRR, a portion of airtime, or *quantum* Q^s , is allocated to each slice s in each transmission round. When a larger value for Q^s is assigned to a slice s more radio resources are allocated to s ; this provides a mechanism to support services with stricter performance requirements. Nevertheless,

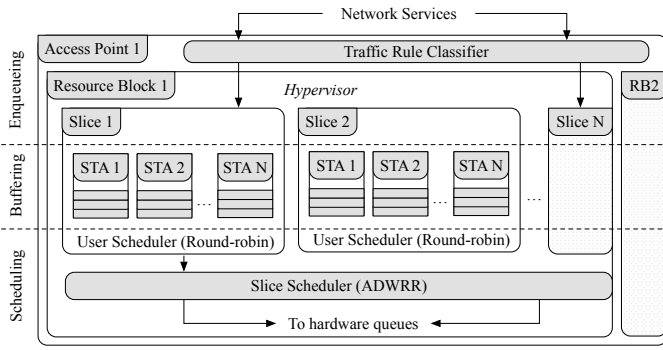


Fig. 2. Simplified slice queue structure and data traffic flow in an AP.

it is important to notice that inactive traffic rules do not cause any performance degradation to the system.

To perform network-triggered handovers and airtime-based network slicing, our approach relies on the 5G-EmPOWER platform¹, which includes an SD-RAN controller called 5G-EmPOWER controller, a backhaul controller implementation of the SDN controller Ryu², and an agent that runs at each AP. The IEEE 802.11 interface of the AP is set monitor mode for radio measurements collection. We extended the SD-RAN controller to allow flow demands and QoS requirements to be informed. Thus, the SD-RAN controller can calculate the expected throughput and verify the QoS. The network intelligence is implemented at the SD-RAN controller, which communicates with the APs at the data plane through its southbound interface using a persistent Transmission Control Protocol (TCP) connection.

IV. DELAY-AWARE SDN-BASED APPROACH

In this section, we present our delay-aware approach for MAC management via airtime-based network slicing and user association using MCDA in IEEE 802.11 SD-RANs. We apply MCDA whenever we want to decide which APs STAs should be assigned according to the high-level objectives of load balancing while considering the delay constraints of QoS flows. Besides user association, we complement our approach with network slicing performed at the IEEE 802.11 SD-RAN. We describe the details of our approach next.

A. Formulating the Load Balancing Problem using MCDA

The IEEE 802.11 RAN consists of a set B of APs, responsible for delivering n services to a set T of STAs. Each service is instantiated in a slice, with S^b denoting the slices of AP b . Each STA t is therefore served by a subset $\alpha \subseteq S^b$ of slices.

We select six criteria for MCDA to evaluate for an AP b : (i) the overall channel load θ^b in B/s; (ii) the total measured dequeuing rate $\mu^b = \sum_{s \in S^b} \mu^s$, where μ^s is the dequeuing rate of slice s ; (iii) the total expected throughput $\mu_{\text{EXP}}^b = \sum_{s \in S^b} \sum_{t \in T} \mu_{\text{EXP}}^{s,t} \cdot t^b$, where $\mu_{\text{EXP}}^{s,t}$ is the expected throughput for STA t in slice s given by its dequeuing rate

while t^b evaluates to true if STA t is associated with b ; (iv) the total measured queueing delay $D^b = \sum_{s \in S^b} D^s$, where D^s is the average queueing delay in slice s ; (v) t_{RSSI}^b , the RSSI perceived at b from STAs within range; and (vi) the indicator variable t^b . The first four criteria are minimized to avoid resource overuse. We use (iv) to avoid APs with a high number of active or overflowing queues. This reduces the chance of a Network Interface Card (NIC) overload and channel saturation. The last two criteria are maximized to improve the chances of using higher data rates, and of fewer connection disruptions.

The weight of each criterion depends on the flow type, either QoS or BE. We use the Analytic Hierarchy Process (AHP) [34] to inform our selection of weights for each flow type, and then tune the resulting weights to avoid the ping-pong effect. Another consideration is that we want BE flows to be more likely to undergo handovers than QoS ones because handovers are detrimental to delay. The MCDA criteria and the resulting weights by flow type (\mathcal{W}_{BE} and \mathcal{W}_{QoS}) are listed in Table I.

TABLE I
MCDA CRITERIA, OBJECTIVES, AND WEIGHTS FOR AP b

Criterion	Objective	\mathcal{W}_{BE}	\mathcal{W}_{QoS}	Description
θ^b	MIN	0.10	0.10	Overall channel load of b .
μ^b	MIN	0.15	0.10	Measured dequeuing rate of b .
μ_{EXP}^b	MIN	0.40	0.20	Overall expected throughput of b .
D^b	MIN	0.15	0.20	Measured avg queueing delay of b .
t_{RSSI}^b	MAX	0.10	0.20	Measured RSSI from STA t of b .
t^b	MAX	0.10	0.20	True if STA t is associated with b .

Several guidelines for choosing the appropriate method to solve an MCDA problem are given in [35]. Given that our problem has quantitative weights, a quantitative scale of comparisons, no uncertainty in the decision problem, and the decision problem is characterized by a complete ranking, we select the Technique for Order of Preference by Similarity to Ideal Solution (TOPSIS) [36]. TOPSIS ranks the alternative solutions by minimizing the distance to the positive ideal solution and the farthest geometric distance from the negative ideal solution.

B. Monitoring Queueing Delay

To implement our MCDA approach, several extensions were required on the APs and also at the controllers. To obtain the queueing delay statistics, we extended the agent implementation (running on APs) with a frame-tracking functionality. This enabled calculation of the average time spent by frames on each slice/queue of an AP. To make use of the new statistics, we extended both the *OpenEmpower* protocol (used for the controller and AP communication) and the SD-RAN controller. Several handler apps at the controller periodically request and calculate the needed metrics. The metrics are maintained at the controller, and the Simple Moving Average (SMA) and Simple Moving Median (SMM) of the last ten measurements are calculated. Finally, the two management apps we implement for our user association and network slicing algorithms can make use of the queueing delay statistics.

¹<https://github.com/5g-empower/5g-empower.github.io>

²<https://osrg.github.io/ryu/>

Algorithm 1 User Association Algorithm

Input:

```
1: every                                ▷ configuration loop interval (20 sec used)
2:  $\mathcal{C}, \mathcal{W}_{\text{QoS}}, \mathcal{W}_{\text{BE}}$         ▷ set of MCDA criteria and weights
3:  $\forall s \in S^b : D_{\text{QoS}}^s$           ▷ max queueing delay of each slice  $s$ 
4:  $\forall s \in S^b : \mu_{\text{EXP}}^{s,t}$       ▷ expected dequeuing rate of each STA  $t$ 

5: loop every
6:   for each  $b \in B$  do                ▷ iterate over all APs
7:      $b_{\text{STATS}} \leftarrow \text{GETRBSTATS}(b)$ 
8:   for each  $t \in T$  do                ▷ iterate over all STAs
9:      $\mu_{\text{EXP}}^t \leftarrow \text{GETSTAEEXPECTEDLOAD}(t)$ 
10:     $\mathcal{W}^t \leftarrow \text{GETSTAWEIGHTS}(t)$ 
11:   for each  $b \in B$  do                ▷ iterate over all APs
12:      $\mu_{\text{EXP}}^b = \text{GETRBEXPECTEDLOAD}(b, \mu_{\text{EXP}}^t)$ 
13:     if  $t^b = \text{true}$  then  $\mu_{\text{EXP}}^b = \mu_{\text{EXP}}^b - \mu_{\text{EXP}}^{s,t}$ 
14:      $\text{TOPSIS.ALTERNATIVE}(\mathcal{C}, \mathcal{W}^t, \mu_{\text{EXP}}^b, b_{\text{STATS}})$ 
15:      $b_{\text{BEST}} \leftarrow \text{TOPSIS.BESTALTERNATIVE}()$ 
16:     if  $t^{b_{\text{BEST}}} \neq \text{true}$  then
17:        $\text{DOHANDOVER}(t, b_{\text{BEST}})$ 
18:
19: function  $\text{GETSTAWEIGHTS}(t)$ 
20:   for each  $s \in S^b$  do                ▷ iterate over all slices of an AP
21:     if  $D_{\text{QoS}}^s$  then return  $\mathcal{W}_{\text{QoS}}$ 
22:   return  $\mathcal{W}_{\text{BE}}$ 
23:
24: function  $\text{GETSTAEEXPECTEDLOAD}(t)$ 
25:    $\mu_{\text{EXP}}^t \leftarrow 0$ 
26:   for each  $s \in S^b$  do                ▷ iterate over all slices of an AP
27:      $\mu_{\text{EXP}}^t = \mu_{\text{EXP}}^t + \mu_{\text{EXP}}^{s,t}$ 
28:   return  $\mu_{\text{EXP}}^t$ 
29:
30: function  $\text{GETRBEXPECTEDLOAD}(b, \mu_{\text{EXP}}^t)$ 
31:    $\mu_{\text{EXP}}^b \leftarrow 0$ 
32:   for each  $t \in T$  do                ▷ iterate over all STAs
33:     if  $t^b = \text{true}$  then
34:        $\mu_{\text{EXP}}^b = \mu_{\text{EXP}}^b + \mu_{\text{EXP}}^{s,t}$ 
35:   return  $\mu_{\text{EXP}}^b$ 
```

C. Using MCDA in the User Association Algorithm

Our user association algorithm is given in Algorithm 1. At a high level, this algorithm solves the load balancing problem formulated in §IV-A. To do such a thing, this algorithm periodically decides to which APs STAs should be assigned using an MCDA method and then perform the handovers accordingly. To make the handover decisions, the algorithm first gathers the monitored statistics from all APs (line 7). The expected throughput of each AP depends on how STAs are distributed and their demands (line 12). The expected throughput of t is calculated based on its active flows. Once obtained, the handover decision for STA t can be made.

To avoid the ping-pong effect, we subtract the expected throughput of an STA t from the overall expected throughput of the AP with which it is connected (line 13). This prevents the expected throughput of t from affecting its own handover decisions. Next, TOPSIS solves the MCDA problem returning b_{BEST} , the highest-ranked AP according to the criteria and weights (line 15). An STA undergoes a handover only if is not associated with its top-ranked AP b_{BEST} (line 17). At each

Algorithm 2 Network Slicing Algorithm

Input:

```
1: every                                ▷ configuration loop interval (5 sec used)
2:  $\forall s \in S^b : D_{\text{QoS}}^s$           ▷ max queueing delay of each slice  $s$ 
3:  $Q_{\text{MIN}}, Q_{\text{MAX}}$                 ▷ min, max quantum (10 us, 12 000 us used)
4:  $Q_{\text{INC}}, Q_{\text{DEC}}$                 ▷ increase, decrease factors (10%, 30% used)

5: loop every
6:   for each  $b \in B$  do                ▷ iterate over all APs
7:      $\text{RECONFIGURE}(b, \text{REQUIREMENTSMET}(b))$ 
8:
9:   function  $\text{REQUIREMENTSMET}(b)$ 
10:    for each  $s \in S^b$  do                ▷ iterate over all slices of an AP
11:      if  $D_{\text{QoS}}^s$  then
12:        if  $D^s > D_{\text{QoS}}^s$  then return  $Q_{\text{DEC}}$ 
13:      return  $Q_{\text{INC}}$ 
14:
15:   function  $\text{RECONFIGURE}(b, Q_{\text{FACTOR}})$ 
16:    for each  $s \in S^b$  do                ▷ iterate over all slices of an AP
17:      if  $D_{\text{QoS}}^s == \emptyset$  then
18:         $Q^s \leftarrow \text{GETCURRENTQUANTUM}(s)$ 
19:         $Q_{\text{NEW}}^s \leftarrow Q^s \cdot Q_{\text{FACTOR}}$ 
20:        if  $Q_{\text{NEW}}^s > Q_{\text{MAX}}$  then  $Q_{\text{NEW}}^s \leftarrow Q_{\text{MAX}}$ 
21:        if  $Q_{\text{NEW}}^s < Q_{\text{MIN}}$  then  $Q_{\text{NEW}}^s \leftarrow Q_{\text{MIN}}$ 
22:        if  $Q_{\text{NEW}}^s \neq Q^s$  then  $b.\text{SETSLICE}(Q_{\text{NEW}}^s)$ 
```

loop, the worst-case execution of this algorithm must iterate over the set of STA T , the set of APs B , and the set of slices of each AP S^b , where $b \in B$. In the worst case, the execution time is $\mathcal{O}(|T| \cdot |B| \cdot |S^b|)$.

D. Network Slicing Algorithm

Algorithm 2 is responsible for adapting the network slice configurations at runtime. Based on the maximum queueing delay threshold and the quantum adjustments, the network slicing algorithm aims to satisfy the QoS requirements of the QoS flows by reallocating resources from the BE slices to the QoS slices. Periodically, this algorithm checks, for each AP, whether the requirements of all QoS slices are met. When all requirements of an AP are met, the quantum is increased by a factor of Q_{INC} (line 13), releasing resources until all slices share the AP equally. Otherwise, the quantum is decremented by a factor of Q_{DEC} (line 12), leaving more resources for the QoS-constrained slices. Q_{MIN} and Q_{MAX} are thresholds that prevent traffic in BE slices from being blocked and from exceeding a maximum quantum configuration, respectively. A new quantum Q_{NEW}^s is set for a slice on an AP only when it differs from its current one.

It is important to emphasize that inactive traffic rules, i.e., slices, do not cause any performance degradation to the system. The limit of resources in which a slice might utilize only proceeds when the AP or its channel is saturated and the remaining resources must be shared with other active slices. At each loop, the worst-case execution of this algorithm must iterate over the set of the set of APs B and the set of slices of each AP S^b , where $b \in B$. In the worst case, the execution time is $\mathcal{O}(|B| \cdot |S^b|)$.

V. PERFORMANCE EVALUATION

A. Methodology and Workload

Figure 3 shows the real testbed used to evaluate our approach. It is made up of a single computer hosting the both controllers, two APs, and four STAs. The APs are based on the PC Engines APU2D4 (x64) processing board, equipped with one Qualcomm Atheros AR958x 802.11 a/b/g/n each. The STAs are Raspberry Pis 4 Model B+ with 802.11b/g/n/ac.

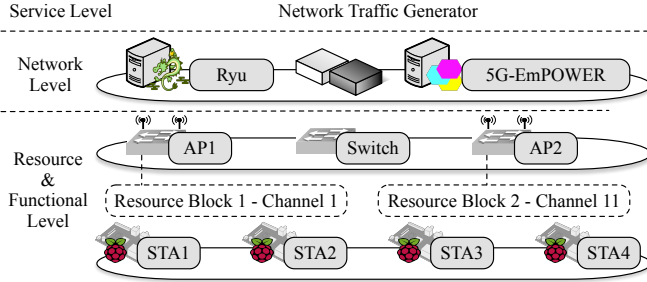


Fig. 3. Testbed deployment scenario.

To replicate the scenario where both APs perceive similar RSSIs from STAs and therefore ping-pong effect is more susceptible to occur, APs and STAs are positioned about 2 meters apart from one another. APs are set to operate on non-overlapping channels (1 and 11). The supported Modulation and Coding Scheme (MCS) rate indices are from 0 to 7 because the STAs operate in the 2.4 GHz band. Our experiments were conducted in a closed office environment with little to no external interference. We generate five User Datagram Protocol (UDP) flows with each flow representing a different service in the network. For each flow, a dedicated slice with the default quantum Q^s of 12 000 us is instantiated. Downstream traffic is generated from the controllers' host towards the four STAs, with the workload parameters of Table II. To avoid static flow rates and arrival times, we generated the flows following the Poisson distribution with MGEN³ having a fixed frame size of 1024 bytes. The experiments were run for five minutes.

TABLE II
WORKLOAD PARAMETERS USED DURING EXPERIMENTATION

Event	Time (sec)	Flow	STA	$\mu_{EXP}^{s,t}$	D_{QoS}^s
1	10	QoS 1	1	2 Mbps	30 ms
2	40	BE 1	2	4 Mbps	N/A
3	70	BE 2	2	4 Mbps	N/A
4	100	BE 3	3	4 Mbps	N/A
5	130	BE 1	2	0 Mbps	N/A
		BE 2	2	0 Mbps	N/A
6	160	BE 1	2	15 Mbps	N/A
		BE 2	2	15 Mbps	N/A
		BE 3	3	15 Mbps	N/A
7	190	QoS 2	4	2 Mbps	50 ms

The SD-RAN controller periodically polls for monitoring statistics on both APs and calculates the SMA and SMM of the last measurement window. The SMM is used to avoid the

masking effect in the presence of outliers. The polling frequency and window size are set to 1 sec and 10 measurements and re-configuration loops are given in Algorithms 1 and 2.

B. Results and Analysis

Figure 4 illustrates the STAs association status as well as the throughput and queueing delay per slice for a representative experiment. Vertical dotted lines mark the events (see Table II). Horizontal dotted lines mark the QoS requirements. Figures 4a and 4b show to which AP each STA is assigned during the experiment and, therefore, when handovers occurred. We begin the experiment with the STAs associated with AP 1 and because there are no active flows, only the RSSI and channel load measurements distinguish the ideal AP, they remain assigned to AP 1. (Figure 4a). After the first flow starts (at second 10), the overall expected throughput and measured throughput and queueing delay of the AP increase. This causes other STAs (which are not receiving data) to be handed over to AP 2 (Figure 4b). At event 3, AP 2 has more resources being used (two 4 Mbps flows versus one 2 Mbps flow on AP 1). This causes STAs 3 and 4 to be reassigned to AP 1. Figures 4c and 4d present the dequeuing rate per slice.

At event 5, when BE flows 1 and 2 are stopped, STA 3 using slice BE 3 is reassigned while it had an active flow. In this case, because the APs are operating on different channels, the Channel Switch Announcement (CSA) mechanism is triggered. CSA is defined by the IEEE 802.11h amendment to enable APs to announce switching to a new channel before their transmission begins on that channel. Beacon messages containing the CSA information are sent to the STA before it switches to the new channel. This allows STAs, which support CSA, to transit to the new channel with minimal downtime. In our experimentation, we assume STAs do not support CSA, causing even more performance degradation when handovers are performed. Indeed, the flow that was reassigned was running on a BE slice while the first flow, running on a QoS slice (QoS 1), had no connection disruptions. Figures 4e and 4f present the queueing delay per slice.

At event 6, the three BE flows increase their throughput to 15 Mbps. This event was introduced to determine if QoS can be ensured with higher demands in the network. When this event occurs, the queueing delay for slice QoS 1 was not satisfied for a short period (i.e., above the 30 ms requirement when Address Resolution Protocol (ARP) messages are exchanged) because the BE slices are using as much as possible from their airtime available. This increased the delay of other slices on the same AP, therefore, slice adaptations were required.

Figure 5 shows the quantum values for all slices during the experiment. The Q^s of all BE slices on the AP were reduced according to Algorithm 2 which released resources to the QoS slice. Specifically, in the second 211, the Q^s values for the two BE slices dropped by 30%, allowing the delay requirements of the QoS slice to be satisfied. Later, resources were released to the BE gradually, with Q^s values increased by a 10% factor until it reaches the maximum threshold. At event 7, a new QoS flow is introduced. Since there are 3 active flows assigned to

³<https://www.nrl.navy.mil/itd/ncs/products/mgen>

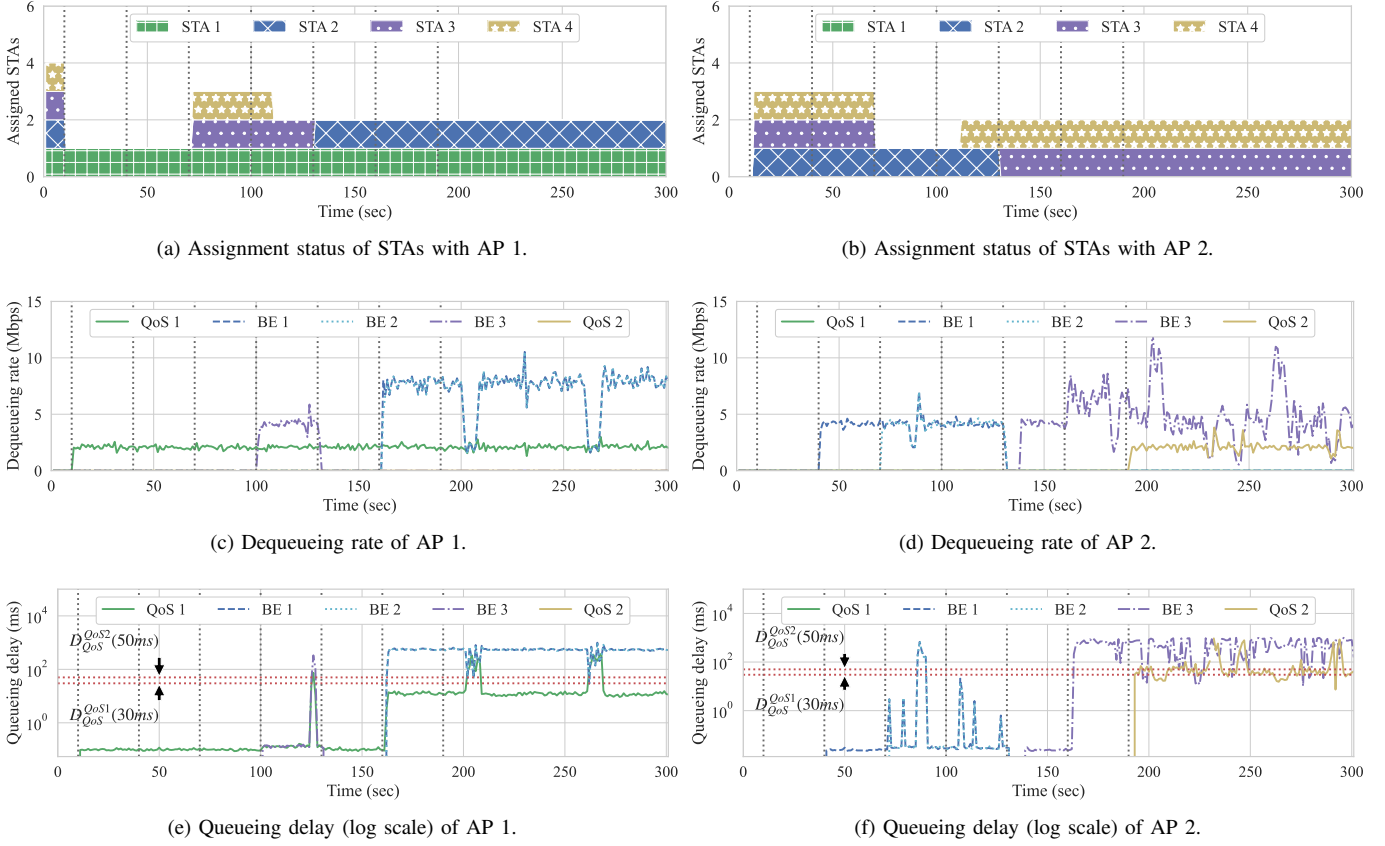


Fig. 4. STA assignment status, dequeuing rate, and queuing delay per slice over the experiment time span running our proposed approach.

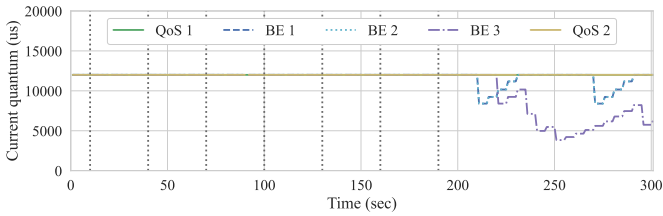


Fig. 5. Current quantum of the slices over the experiment time span.

AP 1 and its overall queuing delay is currently the highest, the AP selected for this new flow is AP 2. As Figure 4f shows, the new QoS flow is introduced at second 190 and its queuing delay of 50 ms is guaranteed for most of the time with analogous slice re-configurations on AP 2.

Overall, during this experiment, the channel load was insignificant, being on average 0.5 Kbps with a standard deviation of 0.2 Kbps on each of channels 1 and 11. During their active periods, QoS flow 1 and 2 had averages very close to the expected dequeuing rate (2 Mbps) while BE 1, 2, and 3, from event 6, had around 7 Mbps each. In this case, more than half of the frames are waiting within the queues until the NIC is ready and it is their turn to transmit. Therefore, such BE flows experience queuing delays of almost 1 sec as presented in Figures 4e and 4f.

C. Comparative Study

We now compare our approach to a state-of-the-art approach for user association. Gómez et al. [16] proposed an algorithm that improves the aggregated goodput compared to traditional approaches based on signal strength. Their proposal uses three indicators: (i) average RSSI of an AP, (ii) AP load, and (iii) channel occupancy. Indicator (i) refers to the average of the uplink RSSIs for all STAs connected to the APs. The second (ii) represents the load of the APs, while the third (iii) represents the channel occupancy in which APs are operating.

Gómez et al. claim that the ping-pong effect can be reduced by using the maximum \max_i , minimum \min_i , and median η_i values for an indicator i . When $\max_i - \min_i > \eta_i$, the best AP is determined and handovers might be performed. To determine the best AP, for each STA, the algorithm selects the AP which has the maximum product of the average RSSIs times the load of the AP plus channel occupancy. However, the ping-pong effect might occur when APs perceive similar RSSIs. For example, assuming two APs perceiving similar RSSIs from STAs and an uneven load on APs or channel occupancy. In this case, at each reconfiguration, STAs will be reassigned to the AP with a lower load or the one operating in the less busy channel. Consequently, AP loads and the channel occupancy will be swapped and the same behavior will be observed, causing STAs to ping-pong.

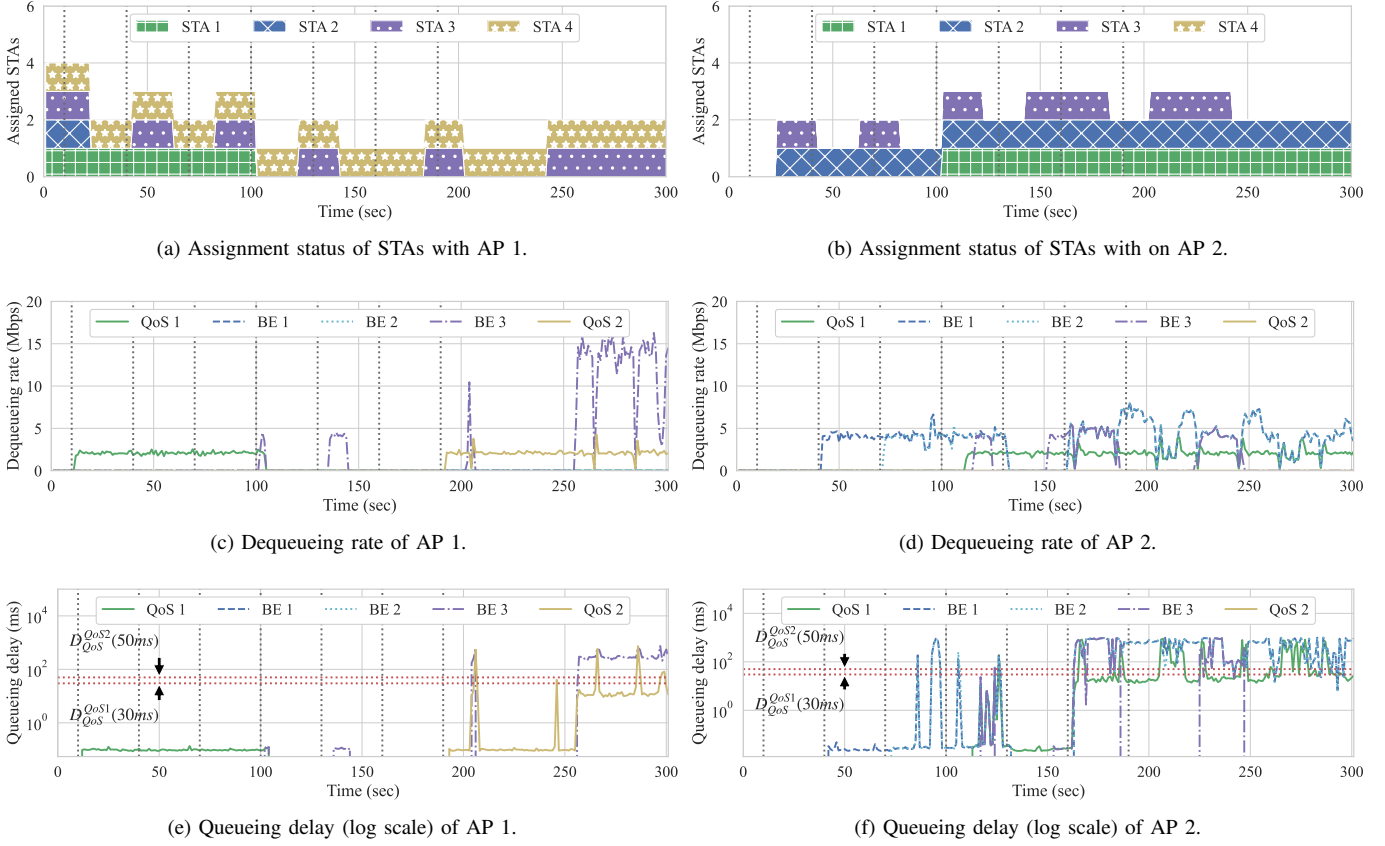


Fig. 6. STA assignment status, dequeuing rate, and queuing delay per slice over the experiment time span running the approach from Gómez et al. [16].

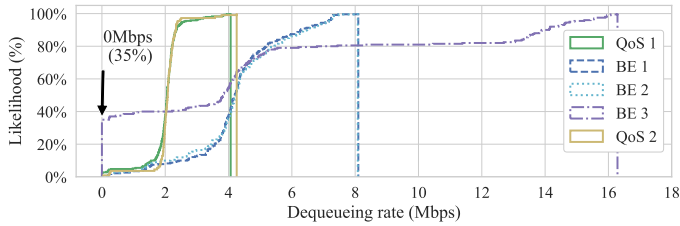
We evaluate Gómez et al. using the same configuration loop interval and workload parameters presented in Table II. Figure 6 illustrates the STA association status, dequeuing rate, and queuing delay per slice during a 5 min experiment. As before, the vertical/horizontal dotted lines mark the events and QoS requirements, respectively (see Table II). Figures 6a and 6b show to which AP each STA is assigned during the experiment. As can be seen, the number of handovers triggered during this experiment is higher, especially for STA 3. The total number was 12 while only 8 where performed using our approach. Besides, with our approach, only one handover occurred with an STA that had an active flow and the flow was a BE flow. On the other hand, using the algorithm presented in Gómez et al., a connection disruption happened for a QoS flow. At the second 100, STA 1 was moved from AP 1 to AP 2 while the QoS 1 flow was active. Because the APs are operating on different channels and hence the STA has to switch its channel, no data was received for about 8 seconds.

As expected, the ping-ping effect might occur for STAs perceived by multiple APs with similar RSSI values. In this experiment, STA 3 has suffered 9 handovers where 6 happened when BE 3 flow was active. The impact of the handovers as well as the expected dequeuing rate and delay requirements for QoS flows were not considered. Our proposal, on the other hand, considers weighted criteria so that STAs running BE

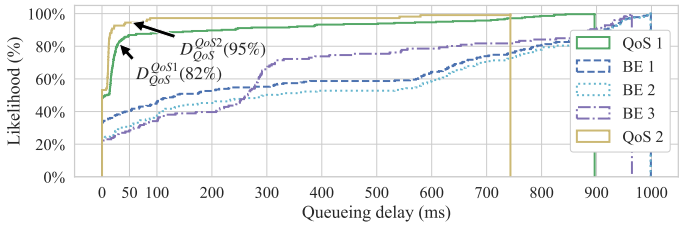
flows are more likely to suffer handovers while STA running QoS flows tend to remain connected to the same AP. Figure 7 shows the Cumulative Distribution Function (CDF) of both dequeuing rate and queuing delay per slice running the flows of both approaches. As we can see, due to the excessive handovers performed by STA 3 running Gómez et al. BE 3 flow has not received data for 35% of the experiment duration while it is no more than 5% of the time with our approach.

Figures 7b and 7d present the likelihood of the queuing delay of each slice using the CDF. The plot shows that QoS 1 flow, which was active for most of the experiment time, had its delay below its requirement D_{QoS}^{QoS1} with the probability of 95% while the probability was 82% with the other approach. This is due to the network slicing algorithm that interactively adapts a slice's airtime according to the delay requirements. For the QoS 2 flow, our approach could only maintain its delay lower than D_{QoS}^{QoS2} with a 68% probability compared to Gómez et al. that had 95% probability. However, with closer examination, we can see that at around the second 100 mark, the BE 3 flow has stopped on AP 1 and this has favored QoS 2 flow that was running along with it.

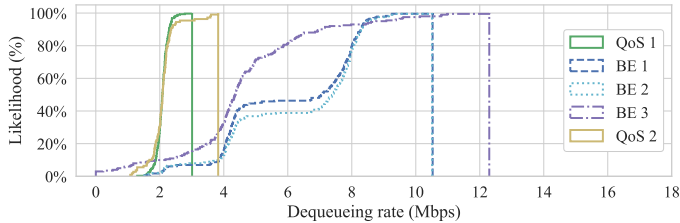
To get an overall picture, we ran 30 experiments using both approaches. Then, we computed the average queuing delay and the average overall throughput achieved per slice at the end of each experiment. In this manner, we can evaluate, besides



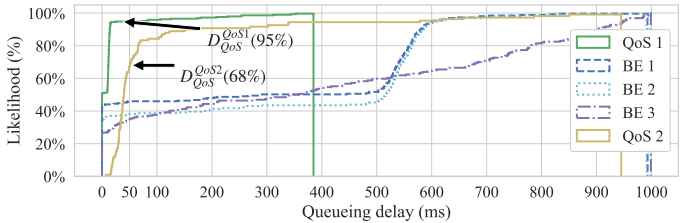
(a) Dequeuing rate per slice with the approach from Gómez et al. [16].



(b) Queuing delay per slice with the approach from Gómez et al. [16]

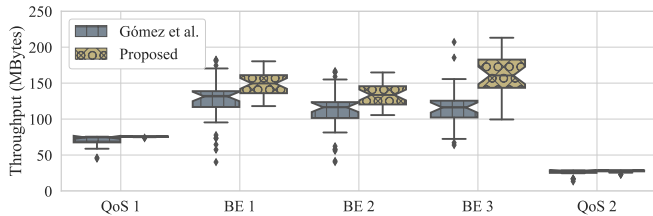


(c) Dequeuing rate per slice with our proposed approach.

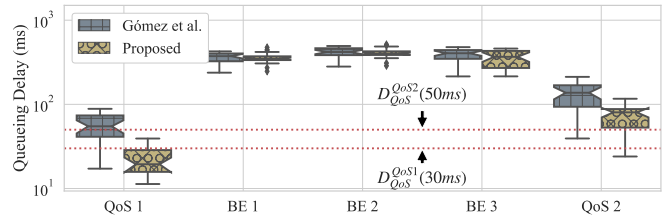


(d) Queuing delay per slice with our proposed approach.

Fig. 7. CDF for the dequeuing rate and queuing delay for the approach from Gómez et al. [16] and for our approach.



(a) Overall throughput of slices.



(b) Queuing delay (log scale) of slices.

Fig. 8. Box-and-whiskers plots showing the overall throughput and average queuing delay of slices for 30 experiments.

the overall throughput achieved per slice, whether the queuing delay requirements were maintained (on average) during the experiment time span. Figure 8 presents the overall throughput and queuing delay per slice using both approaches. Figure 8a shows that our approach has less variability and fewer outliers, not to mention higher average overall throughput, especially for the BE flows. Most importantly, the average queuing delay of QoS 1 flow—which remained active for most of the experimentation—was not guaranteed using the approach from Gómez et al. (see Figure 8b).

Differently, our approach maintained the QoS 1 flow below its requirements. Both of the 1st and 3rd quartiles (interquartile range) were below 30 ms, meaning that about 75% of the averages were below its queuing delay requirements (D_{QoS}^{QoS1}). However, for the QoS 2 flow, in both approaches, fewer than 25% percent of averages were below its requirements (D_{QoS}^{QoS2}). Despite the averages hiding disparities, our slicing algorithm needs some time to react and enhance the QoS.

VI. CONCLUSION

In this paper, we proposed a delay-aware approach for MAC management via airtime-based network slicing and user association using MCDA in IEEE 802.11 SD-RANs. To perform

load balancing in our network, we use the TOPSIS method to decide which APs STAs are assigned and therefore slices are allocated. Six criteria were used for the decision-making with different weights defined for the QoS-constrained and BE services. Unlike existing work in load balancing, our approach periodically analyzes the queuing delay of slices and, by performing airtime-based slice re-configurations at runtime, enhances the QoS when needed. Through experimentation in real hardware, the results show that our approach is capable of performing runtime load balancing and ensures QoS even with higher demands. As future work, we plan to deploy our solution on a large scale testbed.

ACKNOWLEDGMENT

This research received partial funding from The European Union's Horizon 2020 Research and innovation program, under grant agreement No. 826284 (ProTego), from The U.S. National Science Foundation NeTS Award 1813451, and was also supported by the FLEXNET project: "Flexible IoT Networks for Value Creators" (Celtic 2016/3), in the Eureka Celtic-Next Cluster.

REFERENCES

- [1] 5GPPP Architecture Working Group. (2017) View on 5G Architecture. [Online]. Available: <https://5g-ppp.eu/wp-content/uploads/2018/01/5G-PPP-5G-Architecture-White-Paper-Jan-2018-v2.0.pdf>
- [2] P. H. Isolani, M. Claeys, C. Donato, L. Z. Granville, and S. Latré, "A Survey on the Programmability of Wireless MAC Protocols," *IEEE Communications Surveys Tutorials*, pp. 1–1, 2018.
- [3] Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specification. Amendment 8: Medium Access Control (MAC) Quality of Service (QoS), ANSI/IEEE Std 802.11e, LAN/MAN Standards Committee of the IEEE Computer Society Std., 2005.
- [4] 3GPP, "Study on management and orchestration of network slicing for next generation network," 3GPP, Tech. Rep. TR 28.801 V15.0.0T, 2017.
- [5] E. Coronado, S. N. Khan, and R. Riggio, "5G-EmPOWER: A Software-Defined Networking Platform for 5G Radio Access Networks," *IEEE Transactions on Network and Service Management*, vol. 16, no. 2, pp. 715–728, June 2019.
- [6] M. Richart, J. Baliosian, J. Serrat, and J. Gorricho, "Resource Slicing in Virtual Wireless Networks: A Survey," *IEEE Transactions on Network and Service Management*, vol. 13, no. 3, pp. 462–476, Sep. 2016.
- [7] M. Richart, J. Baliosian, J. Serrati, J. Gorricho, R. Agüero, and N. Agouline, "Resource allocation for network slicing in WiFi access points," in *2017 13th International Conference on Network and Service Management (CNSM)*, Nov 2017, pp. 1–4.
- [8] T. Høiland-Jørgensen, M. Kazior, D. Täht, P. Hurtig, and A. Brunstrom, "Ending the Anomaly: Achieving Low Latency and Airtime Fairness in WiFi," in *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. Santa Clara, CA: USENIX Association, 2017, pp. 139–151. [Online]. Available: <https://www.usenix.org/conference/atc17/technical-sessions/presentation/hoilan-jorgesen>
- [9] J. J. Aleixendri, A. Betzler, and D. Camps-Mur, "A practical approach to slicing Wi-Fi RANs in future 5G networks," in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*, 2019, pp. 1–6.
- [10] T. Høiland-Jørgensen, P. Hurtig, and A. Brunstrom, "PoliFi: Airtime Policy Enforcement for WiFi," in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE Press, 2019, p. 1–6. [Online]. Available: <https://doi.org/10.1109/WCNC.2019.8885440>
- [11] E. Coronado, R. Riggio, J. Villalón, and A. Garrido, "Lasagna: Programming Abstractions for End-to-End Slicing in Software-Defined WLANs," in *2018 IEEE 19th International Symposium on "A World of Wireless, Mobile and Multimedia Networks" (WoWMoM)*. IEEE, 2018, pp. 14–15.
- [12] P. H. Isolani, N. Cardona, C. Donato, J. Marquez-Barja, L. Z. Granville, and S. Latré, "SDN-based Slice Orchestration and MAC Management for QoS delivery in IEEE 802.11 Networks," in *2019 Sixth International Conference on Software Defined Systems (SDS)*, June 2019, pp. 260–265.
- [13] P. H. Isolani, N. Cardona, C. Donato, G. A. Pérez, J. M. Marquez-Barja, L. Z. Granville, and S. Latré, "Airtime-Based Resource Allocation Modeling for Network Slicing in IEEE 802.11 RANs," *IEEE Communications Letters*, vol. 24, no. 5, pp. 1077–1080, 2020.
- [14] E. Coronado, R. Riggio, J. Villalón, and A. Garrido, "Wi-balance: Channel-aware user association in software-defined Wi-Fi networks," in *NOMS 2018 - 2018 IEEE/IFIP Network Operations and Management Symposium*, 2018, pp. 1–9.
- [15] E. Zeljković, N. Slamnik-Kriještorac, S. Latré, and J. M. Marquez-Barja, "ABRAHAM: Machine Learning Backed Proactive Handover Algorithm Using SDN," *IEEE Transactions on Network and Service Management*, vol. 16, no. 4, pp. 1522–1536, 2019.
- [16] B. Gómez, E. Coronado, J. Villalón, R. Riggio, and A. Garrido, "User Association in Software-Defined Wi-Fi Networks for Enhanced Resource Allocation," in *Proc. of IEEE WCNC*, Seoul, South Korea, 2020.
- [17] S. Vassilaras, L. Gkatzikis, N. Liakopoulos, I. N. Stiakogiannakis, M. Qi, L. Shi, L. Liu, M. Debbah, and G. S. Paschos, "The Algorithmic Aspects of Network Slicing," *IEEE Communications Magazine*, vol. 55, no. 8, pp. 112–119, Aug 2017.
- [18] M. Agiwal, A. Roy, and N. Saxena, "Next Generation 5G Wireless Networks: A Comprehensive Survey," *IEEE Communications Surveys Tutorials*, vol. 18, no. 3, pp. 1617–1655, thirdquarter 2016.
- [19] H. Luo and M. Shyu, "An Optimized Scheduling Scheme to Provide Quality of Service in 802.11e Wireless LAN," in *2009 11th IEEE International Symposium on Multimedia*, Dec 2009, pp. 651–656.
- [20] P. Serrano, A. Banchs, P. Patras, and A. Azcorra, "Optimal Configuration of 802.11e EDCA for Real-Time and Data Traffic," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 5, pp. 2511–2528, Jun 2010.
- [21] W. L. Pang, D. Chieng, and N. N. Ahmad, "Adaptive Priority Sliding Admission Control and Scheduling Scheme for DCF and EDCA WLANs," *Wireless Personal Communications*, vol. 70, no. 1, pp. 295–321, May 2013. [Online]. Available: <https://doi.org/10.1007/s11277-012-0695-2>
- [22] E. Charfi, C. Gueguen, L. Chaari, B. Cousin, and L. Kamoun, "Dynamic frame aggregation scheduler for multimedia applications in IEEE 802.11n networks," *Transactions on Emerging Telecommunications Technologies*, vol. 28, no. 2, p. e2942, 2017.
- [23] Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications. Amendment 5: Enhancements for Higher Throughput, ANSI/IEEE Std 802.11n, LAN/MAN Standards Committee of the IEEE Computer Society Std., 2009.
- [24] M. G. Sarret, J. S. Ashta, P. Mogensen, D. Catania, and A. F. Cattoni, "A Multi-QoS Aggregation Mechanism for Improved Fairness in WLAN," in *2013 IEEE 78th Vehicular Technology Conference (VTC Fall)*, Sep. 2013, pp. 1–5.
- [25] D. Kim and S. An, "Throughput enhancement by Dynamic Frame Aggregation in multi-rate WLANs," in *2012 19th IEEE Symposium on Communications and Vehicular Technology in the Benelux (SCVT)*, Nov 2012, pp. 1–5.
- [26] S. V. Azhari, O. Gurbuz, and O. Ercetin, "QoS based aggregation in high speed IEEE802.11 wireless networks," in *2016 Mediterranean Ad Hoc Networking Workshop (Med-Hoc-Net)*, June 2016, pp. 1–7.
- [27] B. Maqhat, M. Dani Baba, R. A. Rahman, and A. Saif, "Performance analysis of fair scheduler for A-MSDU aggregation in IEEE802.11n wireless networks," in *2014 2nd International Conference on Electrical, Electronics and System Engineering (ICEESE)*, Dec 2014, pp. 60–65.
- [28] S. Seytnazarov and Y. Kim, "QoS-Aware Adaptive A-MPDU Aggregation Scheduler for Voice Traffic in Aggregation-Enabled High Throughput WLANs," *IEEE Transactions on Mobile Computing*, vol. 16, no. 10, pp. 2862–2875, Oct 2017.
- [29] K. Nakauchi, Y. Shoji, and N. Nishinaga, "Airtime-based resource control in wireless LANs for wireless network virtualization," in *2012 Fourth International Conference on Ubiquitous and Future Networks (ICUFN)*, July 2012, pp. 166–169.
- [30] K. Guo, S. Sanadhya, and T. Woo, "WiFi: Virtualizing WLAN Using Commodity Hardware," in *Proceedings of the 9th ACM Workshop on Mobility in the Evolving Internet Architecture*, ser. MobiArch '14. New York, NY, USA: ACM, 2014, pp. 25–30. [Online]. Available: <http://doi.acm.org/10.1145/2645892.2645893>
- [31] M. Heusse, F. Rousseau, G. Berger-Sabbatel, and A. Duda, "Performance anomaly of 802.11b," in *IEEE INFOCOM 2003. Twenty-second Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE Cat. No.03CH37428)*, vol. 2, March 2003, pp. 836–843 vol.2.
- [32] M. Richart, J. Baliosian, J. Serrat, J. L. Gorricho, and R. Agüero, "Slicing With Guaranteed Quality of Service in WiFi Networks," *IEEE Transactions on Network and Service Management*, vol. 17, no. 3, pp. 1822–1837, 2020.
- [33] R. Riggio, M. K. Marina, J. Schulz-Zander, S. Kuklinski, and T. Rasheed, "Programming Abstractions for Software-Defined Wireless Networks," *IEEE Transactions on Network and Service Management*, vol. 12, no. 2, pp. 146–162, June 2015.
- [34] K. D. Goepel, "Implementation of an Online Software Tool for the Analytic Hierarchy Process (AHP-OS)," *International Journal of the Analytic Hierarchy Process*, vol. 10, no. 3, Dec. 2018. [Online]. Available: <https://www.ijahp.org/index.php/IJAHp/article/view/590>
- [35] J. Wątróbski, J. Jankowski, P. Ziemia, A. Karczmarczyk, and M. Ziolo, "Generalised framework for multi-criteria method selection," *Omega*, vol. 86, pp. 107–124, 2019.
- [36] C.-L. Hwang and K. Yoon, "Methods for multiple attribute decision making," in *Multiple attribute decision making*. Springer, 1981, pp. 58–191.