

Wi-Fi Access Latency Characterization

A technical paper prepared for presentation at SCTE TechExpo24

Lei Zhou
Principle Engineer II
Charter Communications
lei.zhou@charter.com

Table of Contents

Title	Page Number
1. Introduction.....	3
1.1. Latencies in Access Network.....	3
1.2. Active Queue Management.....	4
2. Wi-Fi Media Access Latencies.....	4
2.1. Overview of Wi-Fi Media Access Control.....	4
2.2. Enhanced Distributed Channel Access and Wi-Fi Multimedia.....	6
3. Characterization of Wi-Fi Latency.....	6
3.1. Test Methodology.....	6
3.2. Single-Station Latency.....	8
3.3. Wi-Fi Latency under Multiple Access Contention.....	10
3.4. Wi-Fi Latency under Multiple Access Contention.....	14
4. Conclusion.....	16
Abbreviations.....	17
Bibliography & References.....	17

List of Figures

Title	Page Number
Figure 1 - Wi-Fi Latency Test Setup.....	6
Figure 2 - Single-Station RTT-Load Plot for Station-1.....	9
Figure 3 - Single-Station RTT-Load Plot for Station-2.....	9
Figure 4 - Single-Station RTT-Load Plot for Station-3.....	10
Figure 5 - Multiple BE Station RTT-Load Plot for Station-1.....	11
Figure 6 - Multiple BE Station RTT-Load Plot for Station-2.....	11
Figure 7 - Multiple BE Station RTT-Load Plot for Station-3.....	12
Figure 8 - Multiple VI Station RTT-Load Plot for Station-1.....	12
Figure 9 - Multiple VI Station RTT-Load Plot for Station-2.....	13
Figure 10 - Multiple VI Station RTT-Load Plot for Station-3.....	13
Figure 11 - Three-Station WMM RTT-Load Plot for Station-1.....	14
Figure 12 - Three-Station WMM RTT-Load Plot for Station-2.....	15
Figure 13 - Three-Station WMM RTT-Load Plot for Station-1.....	15

List of Tables

Title	Page Number
Table 1 - System and Connection Information of Station Devices.....	7
Table 2 -Transport Parameters of Test Streams.....	8

1. Introduction

Low latency service is becoming a sought-after feature for access networks to improve user experiences of highly interactive applications such as gaming, video conferencing, virtual or augmented reality, and mission-critical computations. An overwhelming issue reported by users regarding experiences with those applications is the latency of the internet connection. An example of this is when a gamer playing a multi-player game is on a mission with co-players, and someone in their household starts a video streaming session. Another example is when a meeting participant in a real-time conversation starts a video sharing or file downloading session. Addressing the market sector of those latency-sensitive applications may open revenue opportunities for network operators.

From an end-to-end view, latency is the time that elapses between a user request and the completion of that request. When a user requests information from a remote host through an application, that request is processed locally into Internet Protocol (IP) packets. Then the packets are sent over the network to the remote host. There, the packets are processed, and a response is formed, starting the reply process for the return trip. Along the way, and in each direction, are network components known as switches, routers, protocol translators, transport and media changes. At each step, delays are introduced as the packets are buffered, processed and transmitted. These delays could add up to discernible waiting times for the user.

1.1. Latencies in Access Network

Focusing on access networks, latency could be attributed to three sources:

- Transmission latency is the time that it takes the transceivers of the communicating terminals to send/receive all the bits in an IP packet. It is determined by the packet size, the link speed (bandwidth), the modulation and coding scheme, and the physical distance between the two terminals over the communication media.
- Media access latency is the time that it takes the sending terminal to gain access to the communication channel. In most access networks, multiple terminals share a common channel through frequency or time divisions, and a centralized or distributed scheduler coordinates the channel access. A terminal must wait for its turn to start transmissions of its data. This waiting period may be of random length in contention-based media access schemes, such as in the Wi-Fi network.
- Queueing latency is the time that a packet must wait in a buffer before being taken by the transceiver. A Buffer is commonly implemented at the network interface of a terminal, which is crucial to smooth bursts and aggregate fragments of packet flows to achieve maximum bandwidth efficiency. When excessive numbers of packets that exceed the channel bandwidth for a terminal enter the buffer, a queue will build up, which will cause extra delays for any packets entering the buffer afterwards.

The three sources of latencies are not independent. The transmission latency and media access latency are part of the reasons that queues build up at the transmitter buffer. Queue build up also comes from congestion control protocols like TCP.

Reduction of latencies in access networks can target the three sources. More spectral resources can be allocated and advanced modulation technologies adopted, which increase the link speed in orders of magnitude. Examples of such improvements in Data Over Cable System Interface Specification (DOCSIS[®]) networks include mid/high-split, orthogonal frequency division multiplex (OFDM) and orthogonal frequency division multiple access (OFDMA) [1]. At the media access layer, an example in DOCSIS is proactive grant service (PGS) [2], which is a multiple access scheduling type offering to shorten media access delays by removing the request-grant cycle time. This paper will dive deeper into

the Wi-Fi media access latency in Section 2 and 3. First, we will provide a brief of an important technique that addresses the queuing latencies: Active Queue Management (AQM).

1.2. Active Queue Management

AQM is a solution to a salient phenomenon, buffer bloat [3], that is often a primary contributor of queuing latency. Buffer bloat results from the Transport Control Protocol (TCP). In seeking as much bandwidth as possible, TCP makes the transmitting host keep increasing its sending rate until it experiences packet loss at signals of missing acknowledgments. Then the transmitting host backs off by starting from a low rate or holding transmissions for some time until the buffer starts emptying. Lacking timely feedback of congestion situations, TCP tends to fill up the buffer at a bottleneck link and keep that buffer fully occupied for an extended time when a long session occurs. A full buffer is deprived of its capability to absorb traffic bursts and results in prolonged queuing delay.

AQM algorithms are designed to probabilistically mark the Explicit Congestion Notification (ECN) bits of ingress IP packets or drop the packet completely, based on estimates of the queue length and/or the average time that a packet spends waiting in the queue. In those algorithms, generally an increased probability of marking/dropping will be tuned to when the number of queued packets or the estimated queuing time is building up. The marking/dropping actions serve feedback to the end hosts that they should slow their data rates in response to the perceived congestions at the network link. With properly tuned parameters, AQM can reduce queuing latencies without sacrificing TCP throughputs.

Many AQM algorithms have been proposed, including Random Early Detection (RED) [4], Controlled Delay (CoDel) [5], Proportional Integral Controller Enhanced (PIE) [6], and various derivatives. DOCSIS 3.1 has adopted the PIE algorithm as the default AQM algorithm for cable modems [2] [7].

TCP itself is advancing in congestion control mechanisms making use of ECNs. The new TCP, TCP PRAGUE is dubbed Low Latency Low Loss Scalable Throughput (L4S) [8, 9]. It requires that network elements, including terminals and routers, be capable of marking ECNs to signal congestion, or in other words, AQM. The Internet Engineering Task Force (IETF) further recommends a dual-queue architecture that separates L4S and classic TCP packets into different transmission queues and applies different AQM rules to them. This difference of services allows L4S-supporting applications to enjoy low latencies without being interfered by traffic of classic TCP. DOCSIS 3.1 and DOCSIS 4.0 support the dual-queue architecture as Low Latency DOCSIS (LLD) [2].

2. Wi-Fi Media Access Latencies

Wi-Fi networks of the 802.11 standard [10] are a popular home networking technology that connect customer devices to the internet through the wireless medium. It is a critical factor impacting the customer experience of internet service latencies. The fluctuating radio carrier and interference levels in wireless medium and the contention-based MAC protocol of Wi-Fi networks can result in high and highly variable access latencies. As a Wi-Fi network becomes larger with more client devices connected to it, the access latency becomes dominant in the overall latency.

2.1. Overview of Wi-Fi Media Access Control

802.11 media access protocol is known as Distributed Coordinate Function (DCF) or its enhanced version, Enhanced Distributed Channel Access Function (EDCAF). DCF employs a carrier sensing and random backoff mechanism to coordinate the contending media access attempts from multiple Wi-Fi devices (stations or access points [AP]). Carrier sensing is the capability of a device to discern the idle or busy state of the channel to send data on. A busy channel means the channel is occupied by radio signal

transmissions, while an idle channel means that there are no transmissions on the channel. The random backoff randomizes the starting time of the transmissions from multiple devices to avoid collisions of concurrent attempts to access the shared channel. The DCF involves the following steps.

1. When a device has data to transmit, it senses the carrier continuously.
2. If the device detects the carrier being idle for a duration of Distributed Inter-Frame Space (DIFS), it immediately enters a random backoff. The backoff procedure is in the form of a countdown clock. At the beginning of the procedure, a clock is set with a value randomly chosen from a backoff window. The backoff window is denoted as CW_{min} and the backoff clock's initial value is a random number between 0 and $CW_{min} - 1$. The device keeps sensing the carrier during the backoff and the backoff clock counts down when the carrier is detected idle. The clock freezes any time the carrier is busy and reactivates when the carrier is idle for a duration of DIFS.
3. Once the backoff clock counts down to zero, the device transmits a packet rendered from the packet queue. After transmitting the data packet, the device waits for an acknowledgement (ACK) from the receiver. During the waiting period, the device continuously senses the carrier. If the ACK is received within a duration of Short Inter-Frame Space (SIFS), the data transmission is considered complete, and the device goes back to Step 1 if it has additional data in the transmission buffer.
4. If no ACK is received within SIFS or a transmission of another packet on the channel is detected by the carrier sensing, the data transmission in Step 3 is perceived as a failure and the device attempts to retransmit. The retransmission procedure starts by going back to the random backoff at Step 2 but with a backoff window of double the size of the initial backoff window, with an initial value randomly chosen from between 0 and $2 * CW_{min} - 1$. The retransmission repeats if Step 4 fails. At each repeat, the backoff window size doubles (i.e., the backoff clock's initial value for the k -th retransmission is randomly chosen from between 0 and $2^k * CW_{min} - 1$). The maximal number of retransmission attempts is seven by default, though it can be reconfigured.

From the DCF, the latency of a data packet transmitted through the Wi-Fi network includes the transmission of the modulated radio signals and the time spent on carrier sensing, random backoff, waiting for ACK, and retransmissions. We summarize these latency components into an expression below,

$$D = \sum_{k=0}^R \left(S_k + \text{DIFS} + B_k + F_k + \frac{\text{Packet_Size}}{\text{Phy_Rate}} + \text{SIFS} \right)$$

where R is the number of retransmissions and equal to 0 means no retransmissions, S_k is the time of carrier sensing before a DIFS is detected, B_k is the back off time, and F_k is the clock freezing time. S_k , B_k and F_k are all random variables. S_k and F_k are approximately of geometric distribution with a parameter equal to the probability of the channel being busy (aka, channel utilization). B_k is of uniform distribution with a range of 0 to $2^k * CW_{min} - 1$. The expectation of total latency in the above expression can be derived as

$$D = \frac{P}{1-P} \text{DIFS} + \left(2^R - \frac{1}{2} \right) CW_{min} + (R + 1) \frac{\text{Packet_Size}}{\text{Phy_Rate}} + \left(R + \frac{1}{2} \right) \text{SIFS}$$

where P is the channel utilization and R is the number of retransmissions.

2.2. Enhanced Distributed Channel Access and Wi-Fi Multimedia

EDCA is an enhancement to DCF by introducing Quality of Service (QoS) for different application data. Wi-Fi Multimedia (WMM) is the Wi-Fi Alliance specification that is based on 802.11 EDCA. EDCA defines four access categories (ACs): AC_BK (background), AC_BE (best effort), AC_VI (video), and AC_VO (voice). For each access category, there's an associated set of backoff window size CW_{min} values and Arbitration Inter Frame Spacing Numbers (AIFSN). AIFSN serves the same role as DIFS but is of different values for each access category. The net effect of using different AIFSN and CW_{min} values in carrier sensing and random backoff is a reduction in the average media access delay for high priority applications (mapped to AC_VI and AC_VO).

In EDCA framework, the devices of different access categories use carrier sensing and random backoff to compete for Transmission Opportunity (TXOP), which is of different value for each access category. AC_BK and AC_BE are assigned TXOP of lower values, dictating that they can send only one frame during their TXOP. AC_VI and AC_VO are assigned TXOP of larger values allowing them to send as many frames as possible within the TXOP duration. Larger TXOP gives high-priority applications more airtime which translates to higher link bandwidth.

3. Characterization of Wi-Fi Latency

In this section, the Wi-Fi media access latency in multi-station contention scenarios is investigated.

3.1. Test Methodology

The investigation is carried on a Wi-Fi network of one AP router and three station clients. The network configuration is illustrated in Figure 1. The AP router is a appropriate device supporting 802.11ax. The hardware and software information of the three stations is listed in Table 1. The network also includes a client computer that is connected to the AP router through ethernet. This client serves as a data endpoint for Wi-Fi tests and the test controller. The cable modem provides only control access to the Wi-Fi network elements via DOCSIS and is not part of the latency test. The AP router and the station devices are enclosed in an anechoic chamber (not shown in the figure) for radio isolation. Open-source tools are used for packet generations and latency measurements. Flent [11] is a network benchmarking tool, which wraps popular network performance test tools netperf [12] and iperf [13].

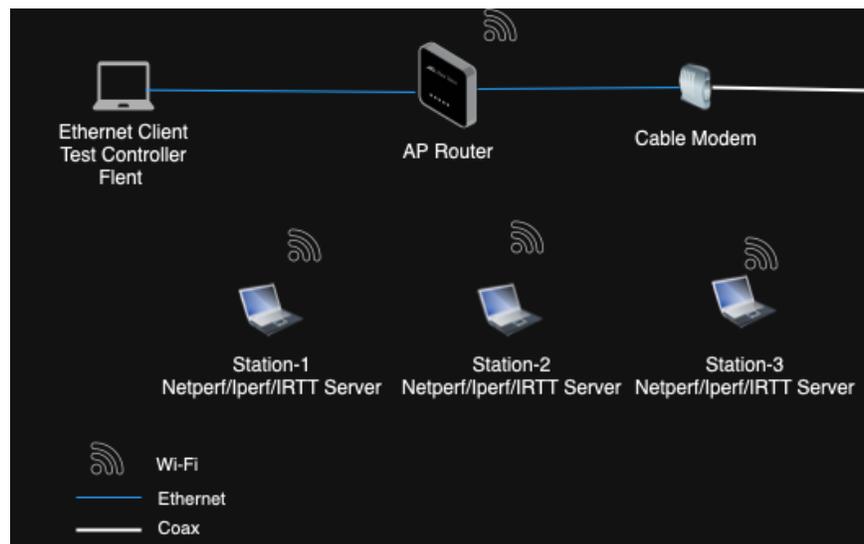


Figure 1 - Wi-Fi Latency Test Setup

Table 1 - System and Connection Information of Station Devices

	Station-1	Station-2	Station-3
System Info			
Platform	Macbook Pro 18 (M1)	Macbook Pro 17 (M1)	Dell Precision 5550 (Core i7-10850H)
OS	MacOS 12.5	MacOS 12.5	Ubuntu 20.04.3 LTS
Wi-Fi	BCM 4387	BCM 4387	AX201
Wi-Fi Connection Info			
Phy Mode	802.11ax	802.11ax	802.11ac
Channel	44	44	44
RSSI/Noise	-17 dBm/-83 dBm	-20 dBm / -90 dBm	-23 dBm/
MCS	11	11	11
NSS	2	2	2
Phy-Rate	1200 Mbps	1200 Mbps	1200 Mbps

Round trip time (RTT) of User Datagram Protocol (UDP) packets between a station and the ethernet client is used as the metric of the Wi-Fi latency. The RTT is measured under specified UDP data rates to characterize the Wi-Fi latency under multiple access contentions. The UDP RTT test method requires that UDP packets be sent from each station to the ethernet client and bounced back. This method makes the multiple access traffic load deviate from (higher than) the specified UDP rates because Wi-Fi uplink and downlink transmissions are sharing the same radio frequency channel and the bounced UDP packets worsen the multiple access contention. To avoid the complication of decoupling downlink and uplink multiple access, two UDP test streams are generated from each station: one with the specified data rates provides the traffic load on the uplink, and the other of negligible rate (50 Kbps) will be bounced for RTT measurement. Since the UDP stream for RTT measurements has minimum impact to the channel access, the interference of the downlink to the uplink is minimized; and the sampled RTT closely approximates two times the uplink multiple access latency under the specified traffic load. The parameters of the two streams are listed in Table 2. Other configurations such as Differentiated Services Code Point (DSCP) marking (for WMM AC mapping purpose) are the same for the two streams.

Table 2 -Transport Parameters of Test Streams

	Protocol	Packet Size	Intended Load
Stream 1 (Load)	UDP	1024 Bytes	100 – 1000 Mbps
Stream 2 (RTT Sample)	UDP	1024 Bytes	50 Kbps

The Wi-Fi multiple access latency characterization test includes generating the two UDP streams from one, two or three stations simultaneously. We run the tests for duration of 1 minute and 5 unique runs are performed. The following subsections will present the test results.

3.2. Single-Station Latency

The first set of results, as shown in Figure 2, Figure 3 and Figure 4, are plots of UDP RTT versus traffic load for each of the three stations when they monopolize the Wi-Fi network individually. In each figure, eight plots are presented, depicting the 50 percentile and 99 percentile value of the measured RTT of the latency-sampling stream using the four WMM AC. The maximal throughput values annotated on the figures are the goodput of the loading UDP stream.

This data primarily lays the baseline of the Wi-Fi latency performance for the three stations. The different implementations of the 802.11 protocol stack affect the Wi-Fi latency characteristics can be inferred. While Station-1 and Station-2 show consistencies in latency and latency variation (measured roughly by the difference between 50 percentile and 99 percentile latency values) before they reach maximum throughput, Station-3 shows significantly higher 99 percentile latency and latency variations. Station-3 also shows high latency in low traffic load regions, especially when no load is present. This behavior is determined to be a result of the packet aggregation feature of the Wi-Fi chip – that is – the Wi-Fi transmitter buffer holds multiple data packets and transmits them on one PDU. Packet aggregation algorithms usually aggregate data bursts arrived within a time window and transmit in a TXOP period. The data rate may affect how much bursts will be aggregated within the time window. Therefore, the latency for a higher data load may be slightly lower than a that for a lower one at certain load range.

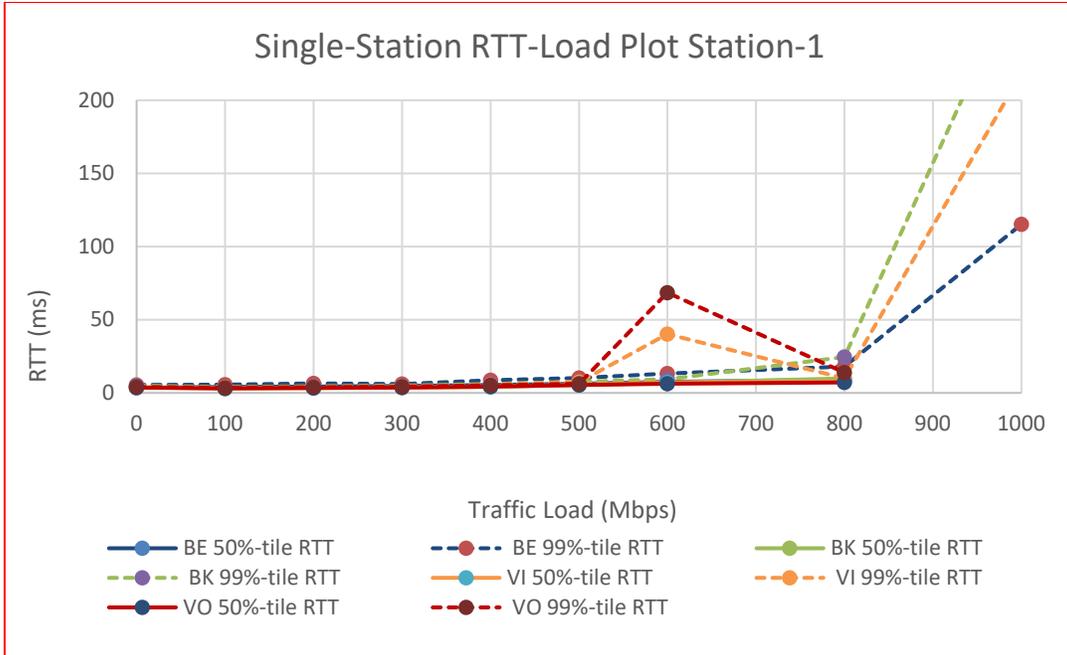


Figure 2 - Single-Station RTT-Load Plot for Station-1

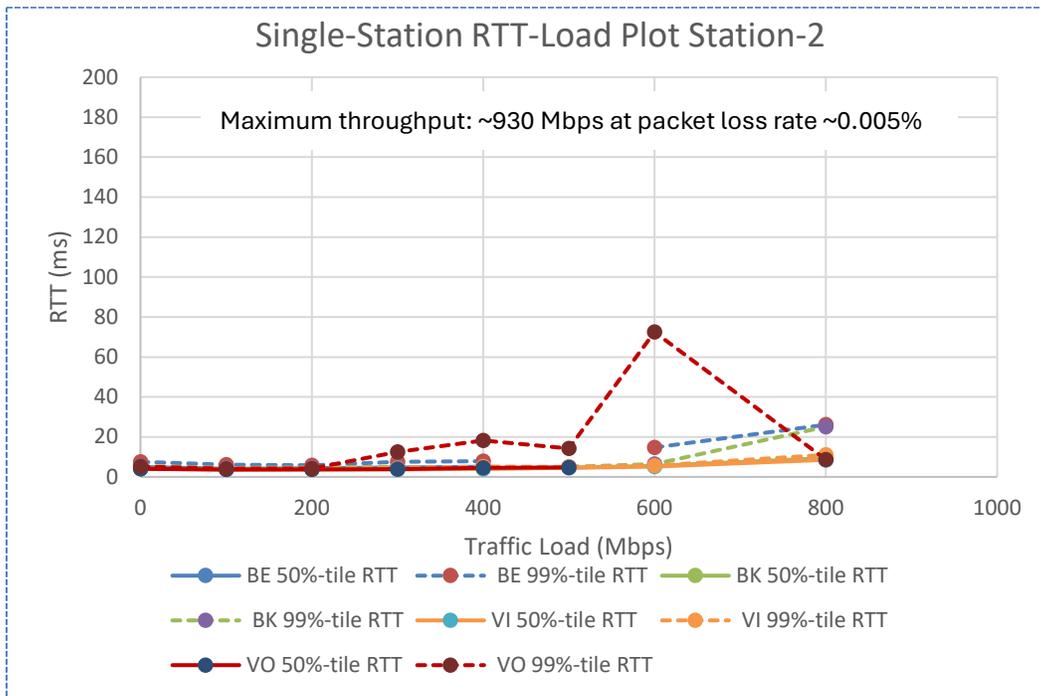


Figure 3 - Single-Station RTT-Load Plot for Station-2

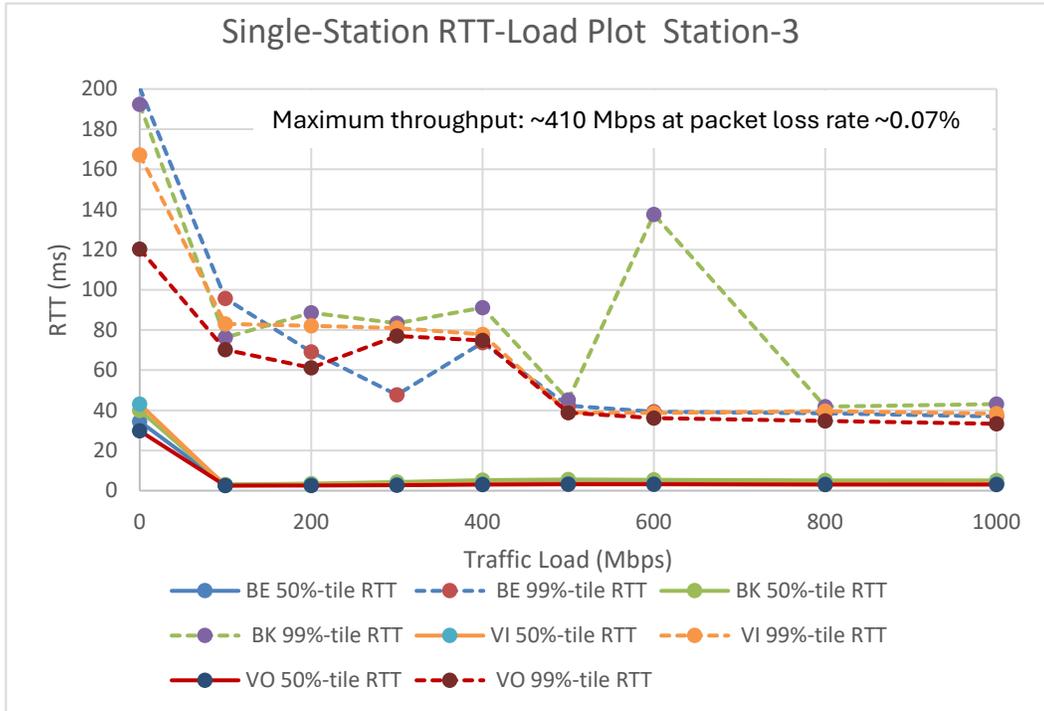


Figure 4 - Single-Station RTT-Load Plot for Station-3

3.3. Wi-Fi Latency under Multiple Access Contention

In this set of tests, UDP streams are generated from two or all three stations simultaneously, and all streams are of the same AC and bandwidth. Figure 5 through Figure 10 are plots of UDP RTT versus traffic load for each station in such multiple access contention scenarios. Figure 5, Figure 6 and Figure 7 are of the case when both or all stations are of BE, and Figure 8, Figure 9 and Figure 10 are of the case when both or all stations are of VI. In the two-station scenarios, Station-1 and Station-2 are transmitting. In the three-station scenario, all three stations are transmitting. Proper single-station plots are reproduced in each figure for comparisons. Note that traffic load axes in the figures are per station. In other words, the aggregated network loads are two or three times values represented on these axes.

Station-1 and Station-2 show robustness against multiple access contentions in the sense that the RTT does not increase if the aggregated network load does not exceed the capacity. Station-3 is more vulnerable to multiple access contentions. Its RTT increases linearly with the number of contending stations.

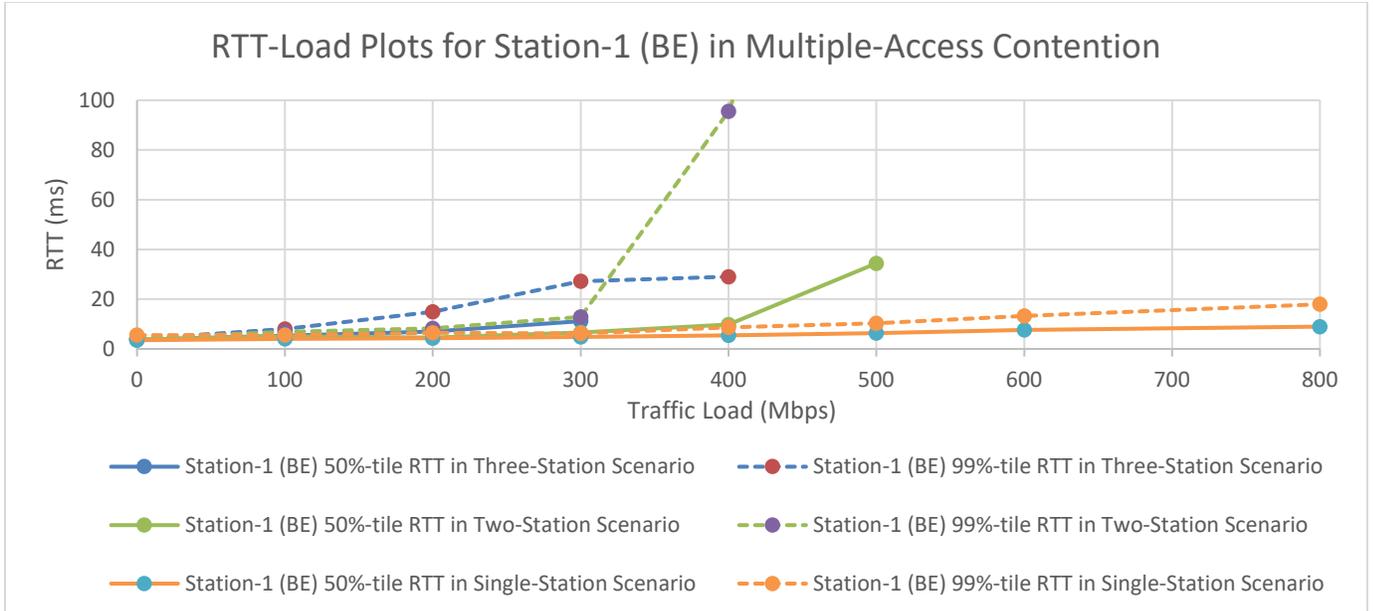


Figure 5 - Multiple BE Station RTT-Load Plot for Station-1

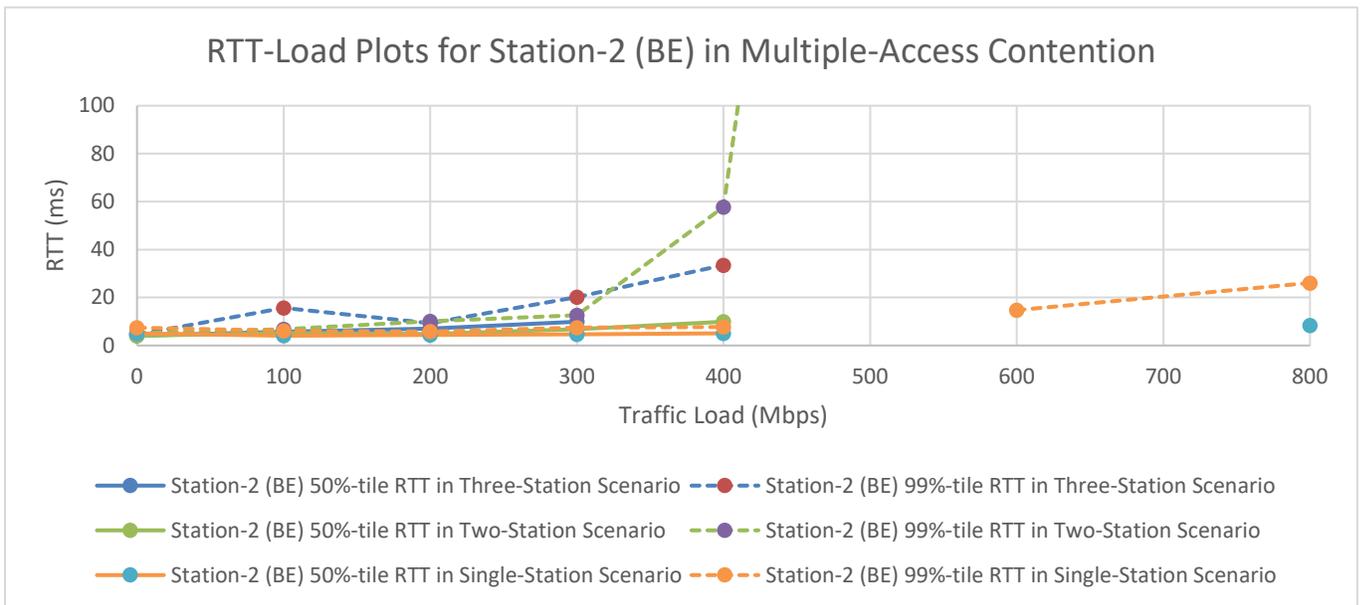


Figure 6 - Multiple BE Station RTT-Load Plot for Station-2

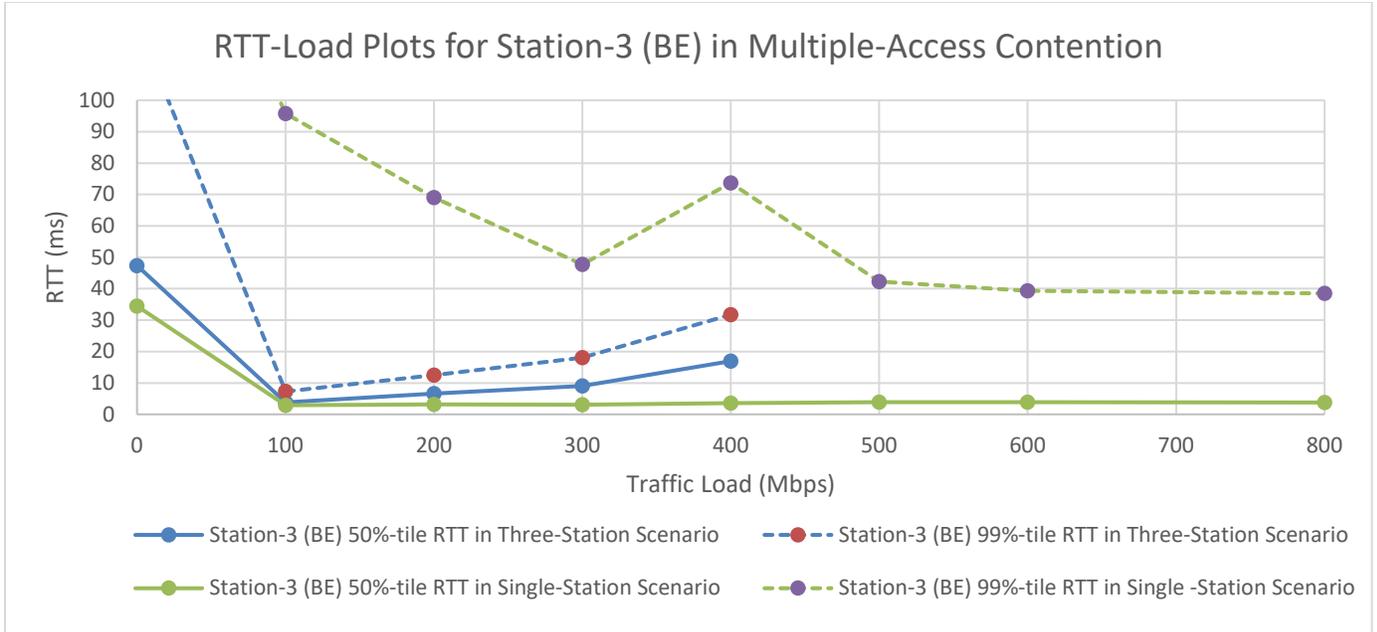


Figure 7 - Multiple BE Station RTT-Load Plot for Station-3

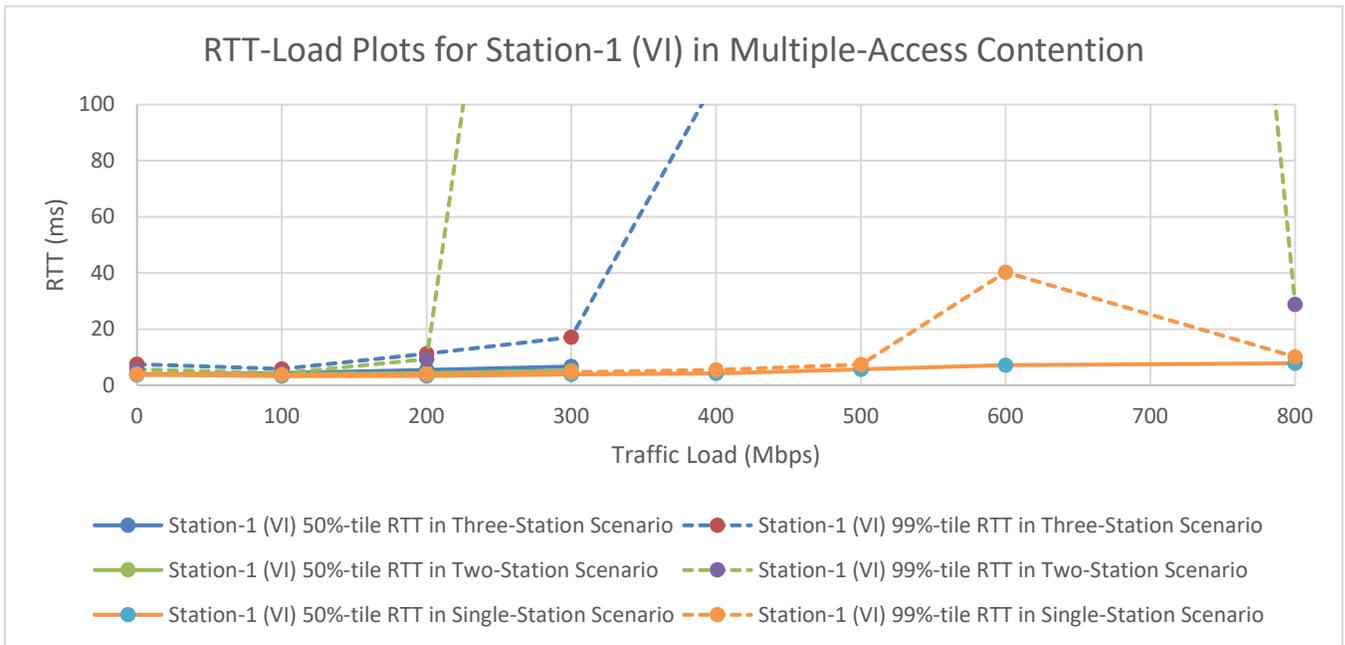


Figure 8 - Multiple VI Station RTT-Load Plot for Station-1

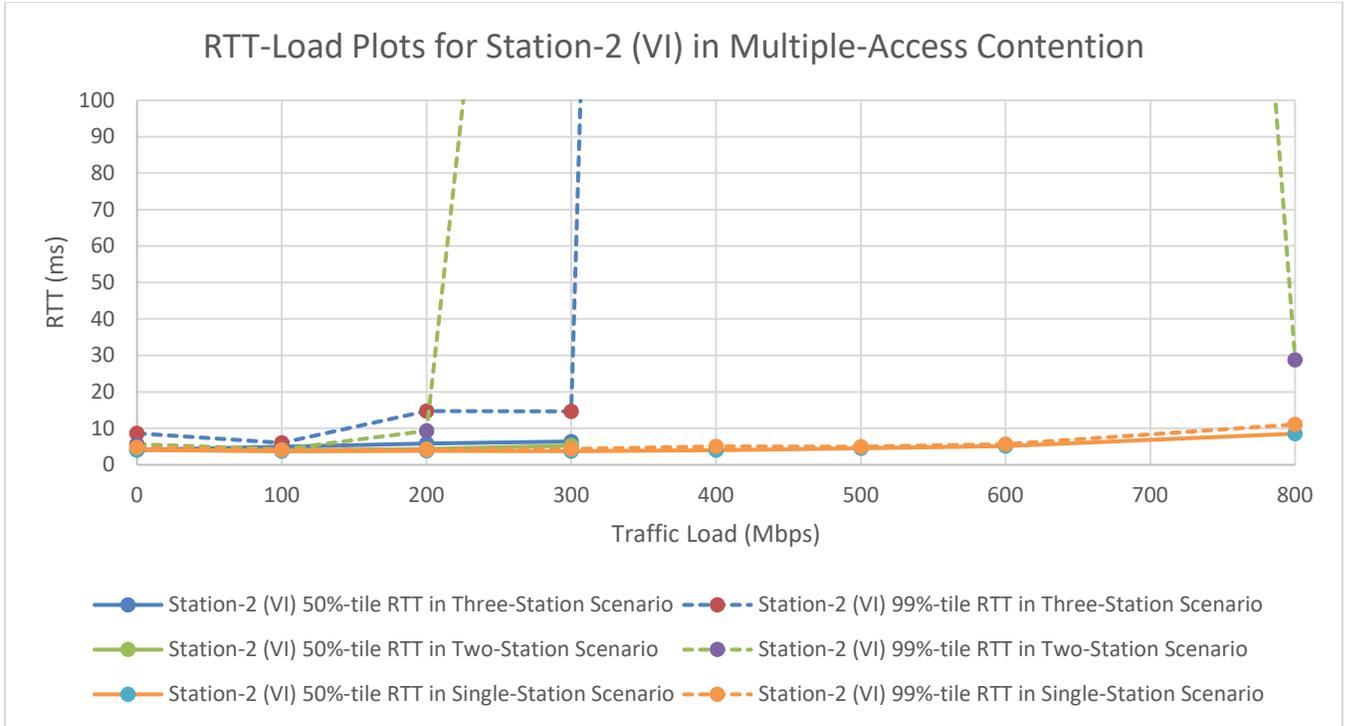


Figure 9 - Multiple VI Station RTT-Load Plot for Station-2

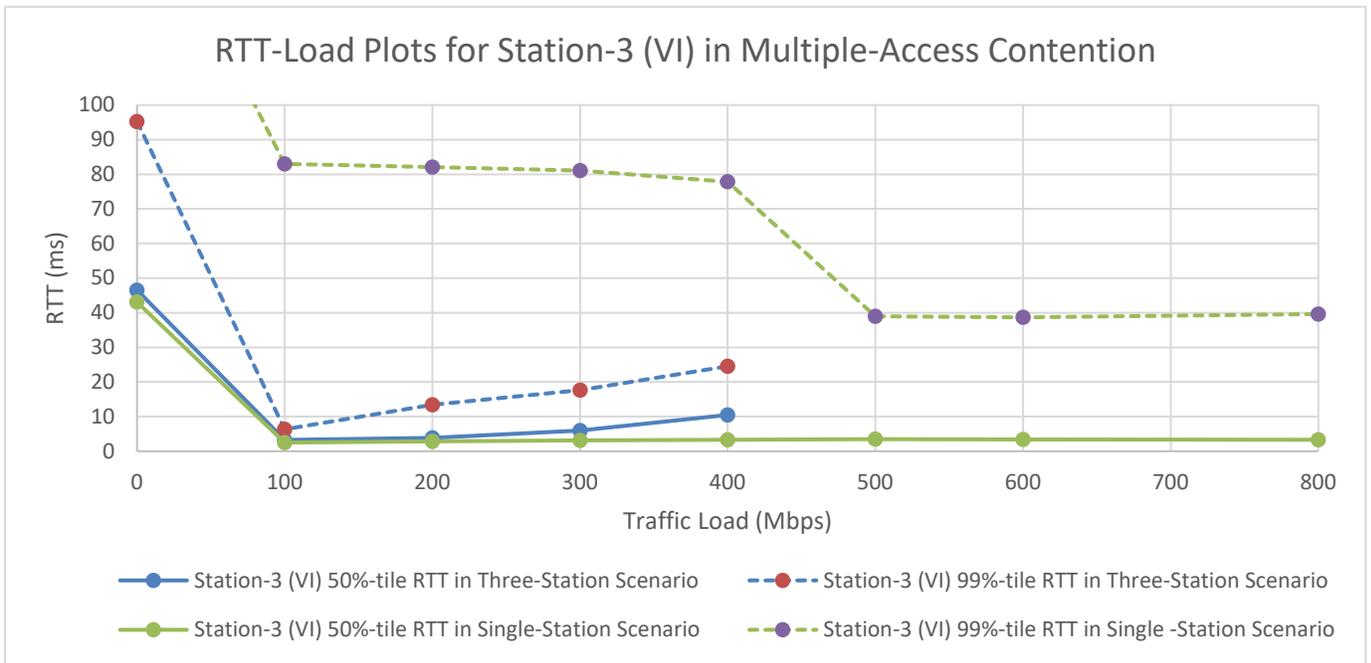


Figure 10 - Multiple VI Station RTT-Load Plot for Station-3

3.4. Wi-Fi Latency under Multiple Access Contention

This set of tests characterizes the effect of WMM QoS to the Wi-Fi latency. In the tests, UPD streams are generated from two or all three stations simultaneously but of different WMM AC. The streams from Station-1 are marked BE, Station-2 is marked VI, and Station-3 is marked VO. The plot of the RTT versus the per-station traffic load for each station is presented in Figure 11, Figure 12 and Figure 13.

The test results do not support WMM QoS as a means of prioritizing traffic, as higher priority AC (Station-3) does not necessary offer lower RTT compared to other ACs contending for the channel, or even in the case that all stations are of the same AC.

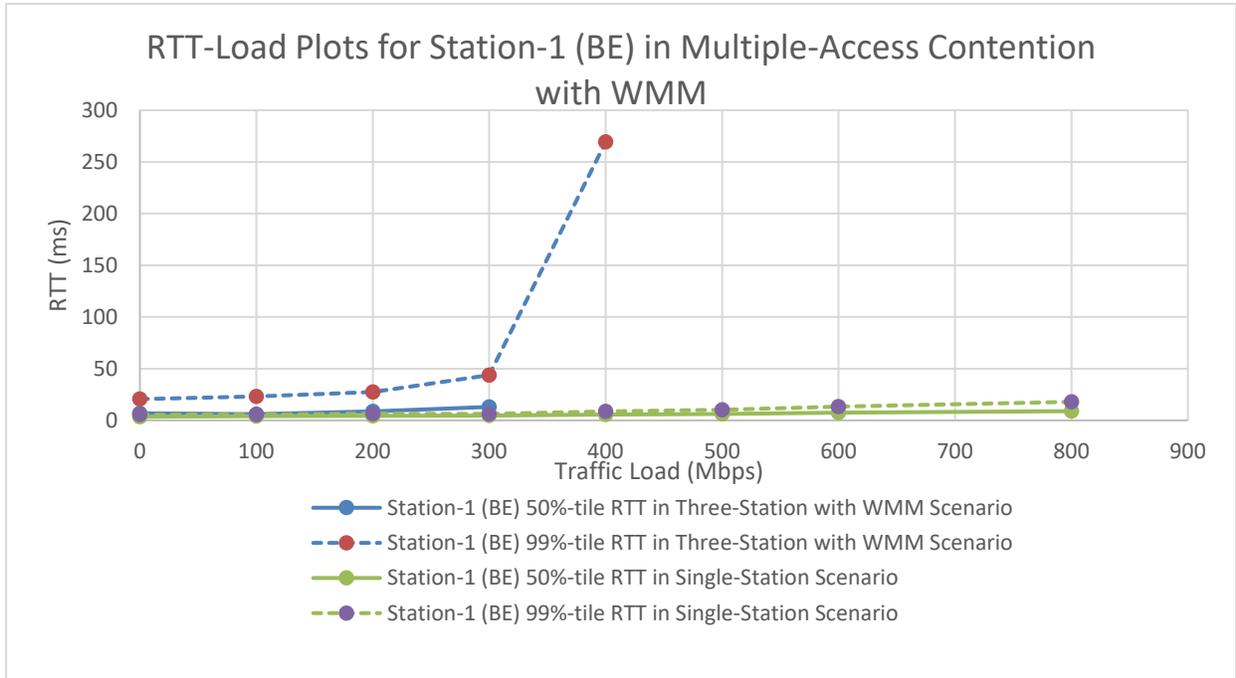


Figure 11 - Three-Station WMM RTT-Load Plot for Station-1

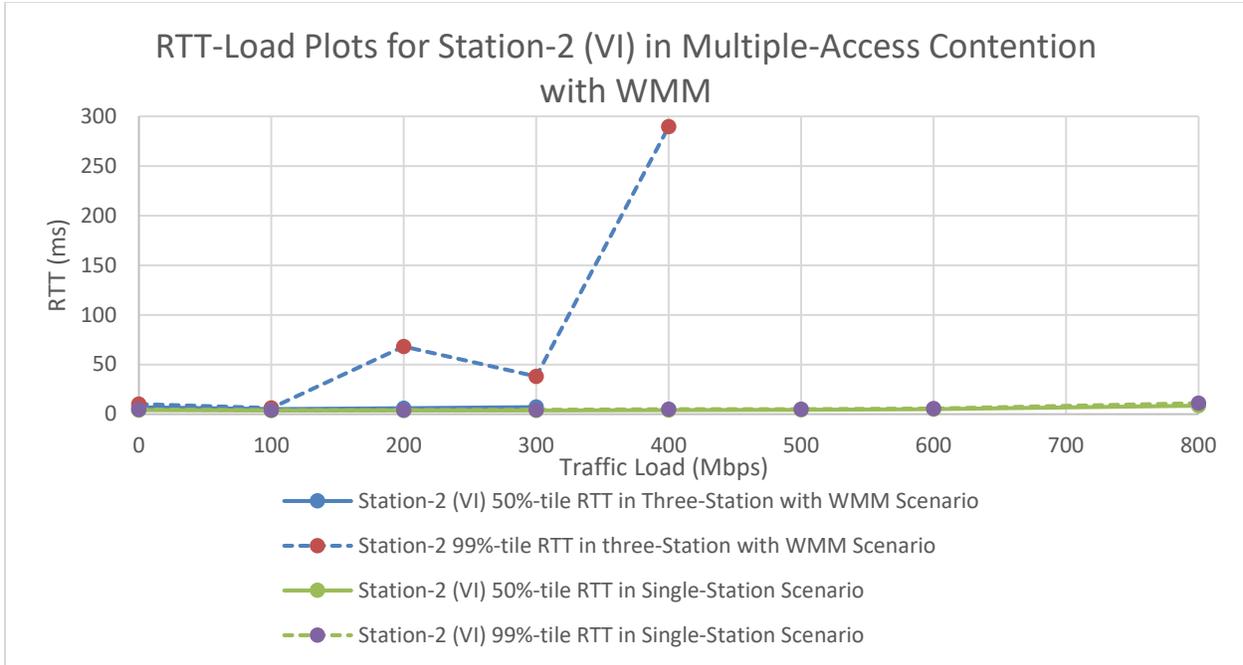


Figure 12 - Three-Station WMM RTT-Load Plot for Station-2

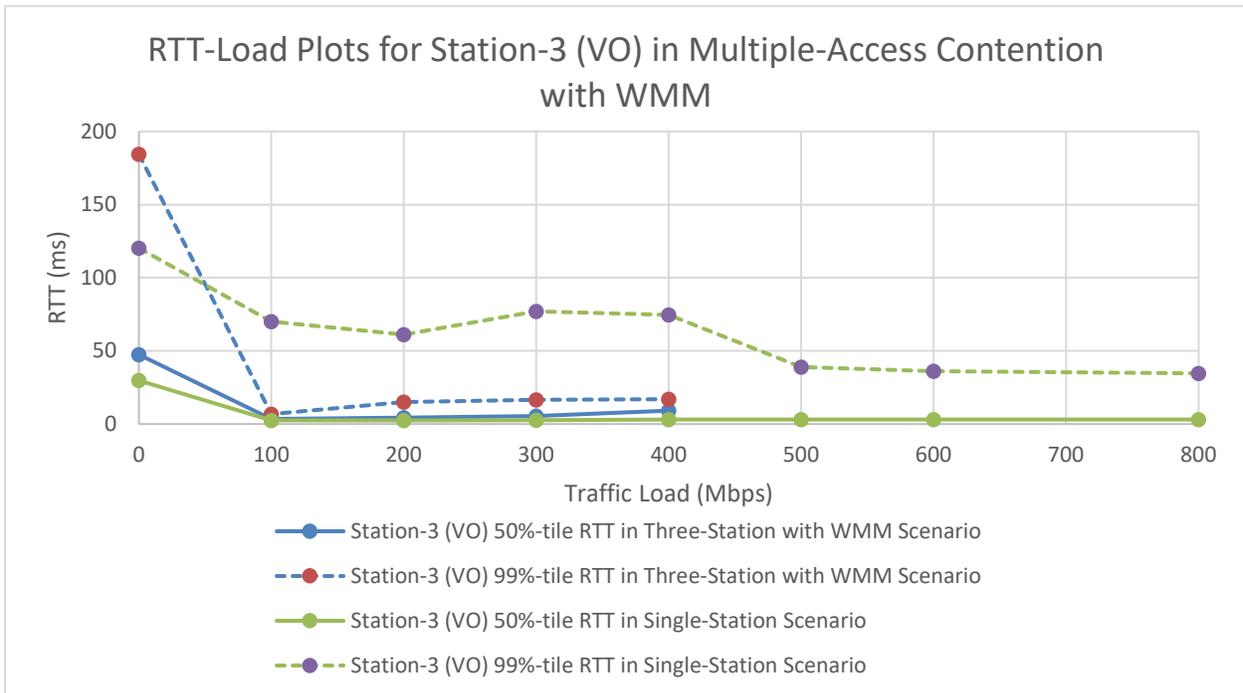


Figure 13 - Three-Station WMM RTT-Load Plot for Station-1

4. Conclusion

This paper reports research of Wi-Fi latencies in multiple access contention scenarios. The research reveals the dependency of latency performance on features such as packet aggregation. The carrier sensing and random backoff mechanisms of DCF/EDCAF can allocate W-Fi airtime efficiently under access contentions when the overall traffic load is fair, without causing extra multiple access latencies. When the average load per station is close to the fraction of the maximum data rate supported by the station divided by the number of stations in contention. Though WMM is designed to give statistically higher priority to AC_VI and AC_VO over AC_BE and AC_BK, the effect of prioritization on latency is not supported by the test results for a wide range of traffic load.

Abbreviations

AC	Access Category
AP	Access Point
AQM	Active Queue Management
BE	Best Effort
BK	Background
CoDel	Controlled Delay
DCF	Distributed Coordinate Function
DIFS	Distributed Inter-Frame Space
DOCSIS	Data Over Cable System Interface Specification
ECN	Explicit Congestion Notification
EDCAF	Enhanced Distributed Channel Access Function
IETF	Internet Engineering Task Force
IP	Internet Protocol
L4S	Low Latency Low Loss Scalable Throughput
OFDM	Orthogonal Frequency Division Multiplex
OFDMA	Orthogonal Frequency Division Multiple Access
QoS	Quality of Service
PIE	Proportional Integral Controller enhanced
PGS	Proactive Grant Service
RED	Random Early Detection
RTT	Round Trip Time
SIFS	Short Inter-Frame Space
TCP	Transport Control Protocol
TXOP	Transmission Opportunity
UDP	User Datagram Protocol
VI	Video
VO	Voice
Wi-Fi	
WMM	Wi-Fi Multi-Media

Bibliography & References

- [1] CableLabs, Data-Over-Cable Specifications DOCSIS 3.1 MAC and Upper Layer Protocols Interface Specifications, 2013.
- [2] CableLabs, Data-Over-Cable Specifications DOCSIS 3.1 Physical Layer Specifications, 2013.
- [3] [Online]. Available: <https://www.bufferbloat.net/projects/bloat/wiki/Introduction/>.

- [4] S. a. J. V. Floyd, "Random Early Detection (RED) Gateways for Congestion Avoidance," *IEEE/ACM Transactions on Networking*, 1993.
- [5] K. a. J. V. Nichols, "Controlling Queue Delay," *Communications of the ACM*, vol. 55, no. 7, pp. 42-50, 2012.
- [6] IETF RFC8033, *Proportional Integral Controller Enhanced (PIE): A Lightweight Control Scheme to Address the Bufferbloat Problem*, 2017.
- [7] IETF RFC8034, *Active Queue Management (AQM) Based on Proportional Integral Controller Enhanced (PIE) for Data-Over-Cable Service Interface Specifications (DOCSIS) Cable Modems*, 2017.
- [8] IEEE Standards, *IEEE Standard for Information Technology--Telecommunications and Information Exchange between Systems - Local and Metropolitan Area Networks--Specific Requirements - Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications*, 2024, 2021, 2013, 2009, 2003, 1999 .