

Approaches to Latency Management: Combining Hop-by-Hop and End-to-End Networking

A Technical Paper prepared for SCTE•ISBE by

Sebnem Ozer, Ph.D.

Senior Principal Architect
Comcast
1800 Arch St., Philadelphia, PA 19103
2152868890
Sebnem_Ozer@comcast.com

Carl Klatsky

Senior Principal Engineer, Product Development
Comcast
1800 Arch St., Philadelphia, PA 19103
215-286-8256
Carl_Klatsky@comcast.com

Dan Rice

VP
Comcast
1401 Wynkoop St Ste 300, Denver, CO 80202
720-512-3730
Daniel_Rice4@comcast.com

John Chrostowski

Executive Director
Comcast
1800 Arch St., Philadelphia, PA 19103
267-260-3695
John_Chrostowski@comcast.com

Table of Contents

Title	Page Number
1. Introduction.....	3
2. Latency Measurement.....	3
3. Latency Performance in Current and Emerging Network Architectures	10
4. End-to-end Support For Low Latency Services	17
5. Conclusion: Final Thoughts on Latency Management.....	20
Abbreviations	22
Bibliography & References.....	23
Acknowledgments	23

List of Figures

Title	Page Number
Figure 1 – Examples of Latency Measurement Methods	5
Figure 2 – TCP connection based latency measurement. Top: Internet RTT; Bottom: Client RTT	8
Figure 3 – Latency Under Load Measurement	9
Figure 4 – Top Left: Max DS and US LUL with suboptimal DS AQM settings; Bottom Left: Mean DS and US LUL with suboptimal DS AQM settings; Top Right: Max DS and US LUL with optimized DS AQM settings; Bottom Right: Mean DS and US LUL with optimized DS AQM settings	12
Figure 5 – Top: Max US LU; Bottom: Mean US LUL for different CM models and speed tiers.....	13
Figure 6 – Left: dslreports.com results with suboptimal DS AQM settings and D3.1 CM with BC with 250ms default target latency; Right: Optimized DS AQM settings and D3.1 CM with US AQM with 10ms default target latency	14
Figure 7 –Throughput (Mbps) values over time for each flow per CM HSD buffer size. Each flow has different e2e RTT. Green flow's server is the closest to the subscriber's home while purple flow's server is the farthest away. Top: 10ms target buffer size; Middle: 30ms target buffer size; Bottom: 50ms target buffer size.....	15
Figure 8 – End-to-end LL services support.....	17
Figure 9 – Low Latency Services Monitoring and Management for CCAP Systems.....	19
Figure 10 – Low Latency Services Monitoring and Management for Distributed Systems	19
Figure 11 – Low Latency Services Monitoring and Management Integrated within Data-driven and Knowledge-defined Architectures	20

List of Tables

Title	Page Number
Table 1 – Properties of Latency Measurement Methods	6
Table 2 – DOCSIS Networks Latency: RTT Between CM and CMTS	10
Table 3 – Architecture changes to support LL services	18

1. Introduction

MSO networks that deliver services between the source and destination ends consist of multiple hops, i.e. different network segments. Customer experience is shaped by many Quality of Service (QoS) features that may be defined per hop-by-hop and end-to-end views. Traditionally, speed performance has been addressed as the main contribution to Quality of Experience in an MSO network. However, latency and jitter have a significant impact on many current and emerging services. MSOs have started modeling, monitoring and managing a more complete performance concept, including speed, latency, jitter, packet loss, security and reliability. The importance of this approach has been reiterated in lockstep with a significant increase in traffic volumes, across a very different mix of residential services, and because of a sudden pandemic, as MSOs who adopted this approach were successful supporting additional traffic volumes. This paper will describe current work on latency measurement, optimization and management platforms with hop-by-hop and end-to-end features. We will present current achievements and results for lower latency systems in the cable industry, and ongoing work on new optimization techniques and big data analytics. Architectural examples will be provided for a data-driven and knowledge-based converged access network with low latency service assurance and agility. Finally, we will discuss roadmap items MSOs may adopt to provide low latency services within their 10G initiative.

2. Latency Measurement

Subscribers' Quality of Experience is a subjective concept affected by the Quality of Service along the path between the service endpoints. QoS metrics are defined from the system's perspective and can be measured and managed. QoE requires a multi-disciplinary approach to assess the user's perspective that may be affected by factors unique to the user and the user's interactions with the service. Mapping QoS to QoE is still evolving as new services and applications emerge. Performance metrics such as throughput (speed) have been regularly measured by MSOs, but speed is only one of the performance indicators. Different services and applications require different levels of speed, latency, jitter and packet loss. For example, depending on the online gaming type and platform, the impact of latency, jitter and packet loss on the gamer's experience and his/her lag perception may be different [1]. Customer experience with gaming, videoconferencing, VR/AR and many commercial services can only be assessed if all the corresponding QoS metrics -- and especially latency and jitter -- are also well modeled, monitored and managed. Fairness, availability, reliability and security are other factors affecting customer satisfaction.

MSOs need to assess and optimize the latency/jitter performance in their networks to improve their subscribers' experiences, and to ensure a strong competitive position, as intended by the 10G initiative. Upstream and downstream usage changed, due to the Covid-19 pandemic, and the steady increase of low latency services in residential networks is increasing the priority of latency management systems.

Figure 1 shows latency measurement (LM) methods widely used by operators. The LM-1 method is a passive measurement method described in more detail later in this section. It measures the round-trip time (RTT) of Transmission Control Protocol (TCP) handshakes to assess access and internet hops, as well as end-to-end latency values (including home, core networks and internet.) It has the advantage of using actual production traffic with no impact on the network. Due to TCP's ACK suppression and expediting implemented in CMs, it may not be directly mapped to UDP type latency. This latency is affected by network conditions and utilization levels (e.g. media access delay due to utilization in DOCSIS bonding groups, or Wi-Fi channel and queueing delay due to HSD service consumption.)

The LM-2 method is an active measurement for upstream (US) and downstream (DS) traffic latency under load (LUL). The test starts with TCP-based connection establishment only if the gateway utilization is low because of potential customer-facing performance impacts. This TCP-based connection latency within LM-2 (conn RTT in Figure 1 and Figure 3) reflects the performance in idle conditions. After the successful connection, LUL in both directions are measured sequentially. US LUL is round trip latency under US load while DS LUL is round trip latency under DS load and they reflect the queuing delay ranges in the CM and CMTS respectively. The LUL test may use iperf TCP as the load to utilize the bandwidth up to speed tier rates in DS and US directions. UDP pings (e.g. Netperf UDP_RR or iRTT) are transmitted concurrently. All the additional test traffic is excluded from any billing. UDP ping latency, jitter and packet loss metrics are collected along with TCP load throughput. Any bufferbloating issue can be detected and the performance of buffer control and active queue management techniques can be optimized with this test method. It can be measured at different endpoints to assess home Wi-Fi and access network performance.

The LM-3 method comprises well known ping measurements (both ICMP and DOCSIS pings). ICMP pings can be used when data and control planes overlap, but are not suitable for SDN networks where data and control planes are not the same. ICMP pings may be routed differently than TCP/UDP packets and may be processed with lesser priority. This latency is also affected by utilization levels similar to LM-1 method, however these are periodic synthetic ping data.

An overview of the properties of these measurement methods is provided in Table 1. LM-1 and LM-2 methods are newer techniques that MSOs started to integrate into their performance measurement platforms. More detailed information on LM-1 and LM-2 methods with test and measurement points that can be integrated are described below.

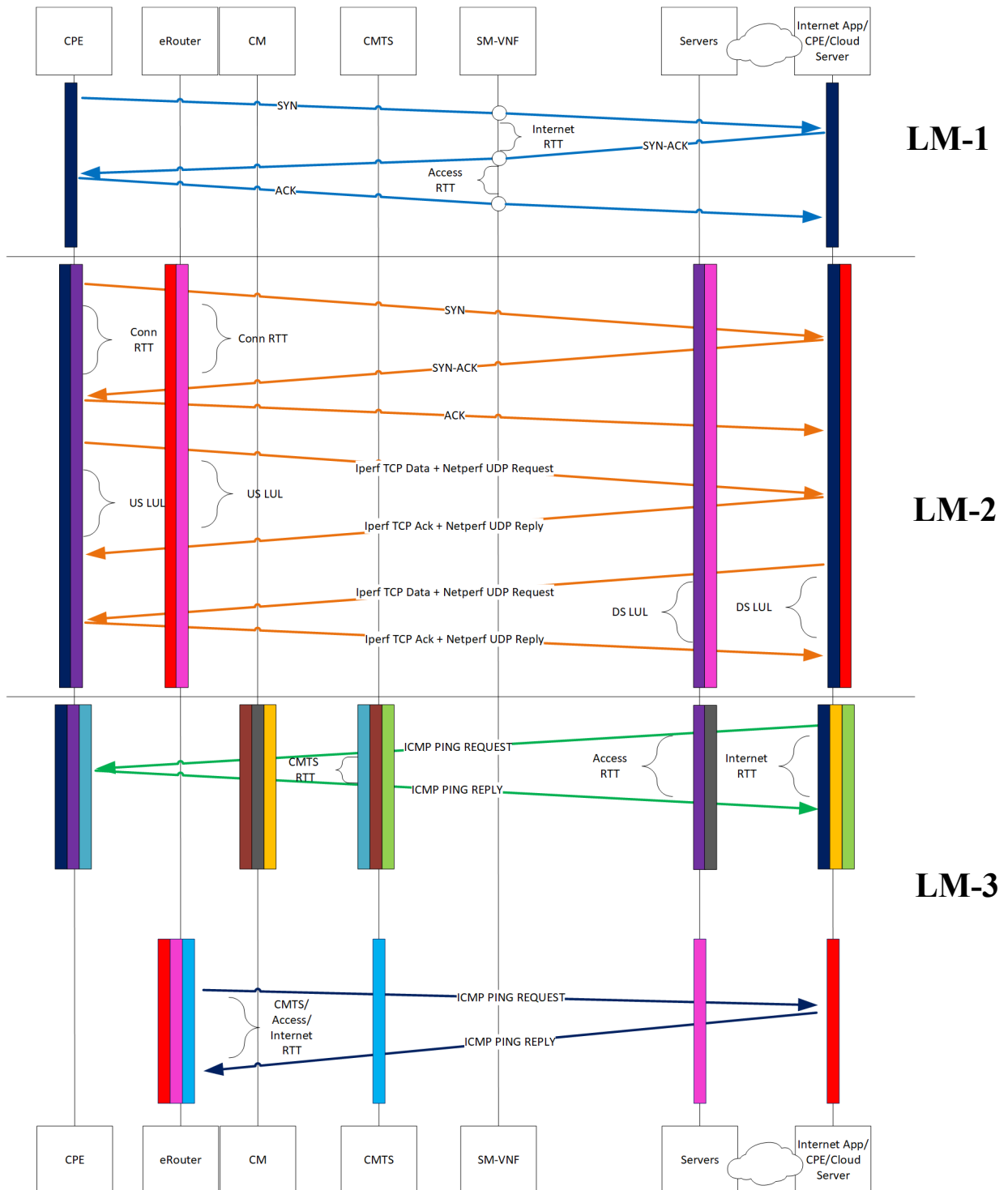


Figure 1 – Examples of Latency Measurement Methods

Table 1 – Properties of Latency Measurement Methods

Test/ Monitoring	Traffic End Point 1	Traffic End Point 2	Monitoring / Test Points	Traffic	Other metrics	Comments
LM-1 Passive Measurements	CPE	Internet	SM-VNF	TCP traffic	<ul style="list-style-type: none"> •TCP session info and retx count 	<ul style="list-style-type: none"> •No extra traffic & End-to-end approach •TCP traffic latency, no UDP latency •US latency is impacted by TCP suppression and expediting •Capability to measure per SF performance and home HSD utilization
LM-2 Active Measurements Under Load	CPE or gateway erouter	Netperf & Iperf Servers	Clients and server functionalities at endpoints	TCP iperf up to speed tier + Netperf UDP_RR / iRTT	<ul style="list-style-type: none"> •Speed •Jitter, packet loss 	<ul style="list-style-type: none"> •UDP ping traffic latency, TCP RTT •Should run only in idle times to not affect customer •Testing traffic must be excluded from billing
LM-3 Active Measurements Without load	CPE or gateway erouter or CM	CMTS or Servers	Ping Endpoints	ICMP Ping & DOCSIS MAC Ping	<ul style="list-style-type: none"> •Jitter, packet loss 	<ul style="list-style-type: none"> •ICMP/DOCSIS ping latency/jitter/packet loss •ICMP vs UDP/TCP diff (control plane) •Minimal extra traffic •ICMP is not directly applicable to SDN devices/vCM/vCMTS

The LM-1 method has been used in Comcast Networks with Software Defined Networking (SDN) and Virtual Network Functions (VNF) components. As SDN has evolved over the last decade, the telecommunications industry has seen the opportunity that SDN presents as a means to develop new product service offerings, and specifically the means to rapidly deploy VNF in support of new products. Comcast

has realized SDN as part of its Active Core and Access platforms. Comcast has also realized SDN as part of its Virtual Services Gateway (VSG) platform. The VSG platform uses commercial off-the-shelf (COTS) compute servers, supported with appropriate network interfaces, upon which various VNFs can be instantiated in support of new products and services.

The VSG platform is deployed in line with the production data traffic and thus a variety of VNFs are envisioned, covering usage metering, monitoring, telemetry, reporting, and traffic marking. Multiple VNFs can be instantiated as supported by the platform's compute and network capability. The initial hardware iteration of the VSG platform is based on a dual-socket x86 compatible motherboard, with twenty cores per CPU socket. The VSG platform hardware also includes dual-100G interfaces, for interconnection to peer network elements. This compute & networking configuration can support the initial set of VNFs envisioned and planned.

The first VNF that we deployed on the VSG platform was a usage metering function. By offloading the existing usage metering function from the CMTS, the function becomes more versatile. That's because changes to per-user profiles can be achieved via real-time configuration, as requested through the OSS / BSS.

The second VNF that we deployed was the telemetry function introduced above. This VNF monitors the anonymized TCP connections of connected subscribers and provides a collection of TCP statistics per subscriber. Within a configurable reporting interval, the VNF is able to track things like TCP session count, TCP connection duration, total session packet count, and TCP retransmission count, per direction. Two key statistics provided by this VNF, and shown in Figure 2, are TCP SYN-to-SYNACK RTT (Server RTT) and TCP SYNACK-to-ACK RTT (Client RTT).

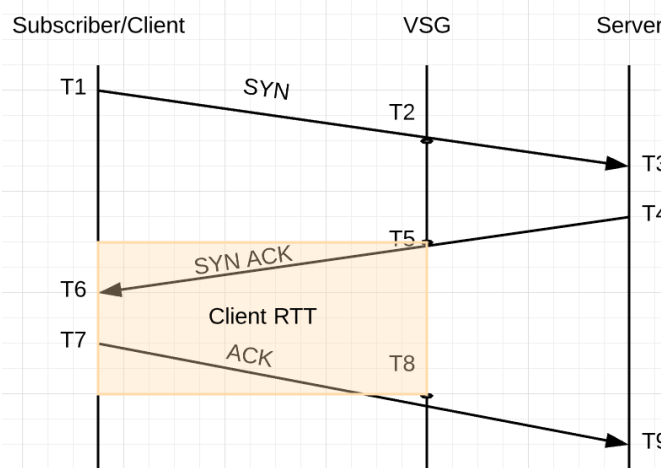
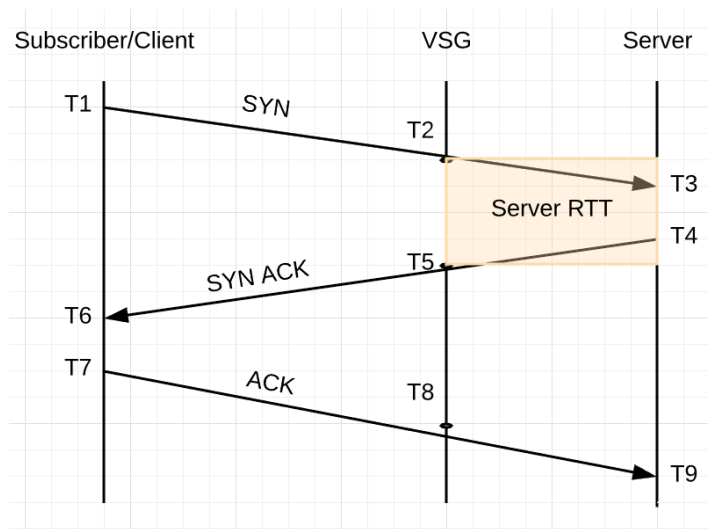


Figure 2 – TCP connection based latency measurement. Top: Internet RTT; Bottom: Client RTT

From the VSG's position in the network, looking at the TCP SYNACK-to-ACK timing (Client RTT) specifically gives a view into the performance of the access and home portion of the network, as the other network elements in the path to the customer device are the CMTS, cable modem and Wi-Fi Access Point. The per-subscriber metrics give an anonymized view into a specific customer's network experience. By looking at the aggregate subscriber metrics across a given service group or bonding group, we gain a view into the health of that service group or bonding group. Monitoring these metrics can serve as an early warning signal of the need for capacity expansion, even before traditional throughput and utilization metrics, thus providing a new way to plan for capacity expansion.

Providing insight into this timing serves also as a proxy for the customer experience. A part of how customers perceive the performance of the network is governed by the TCP 3-way handshake, to establish the TCP connection between the client and the server. The typical online application (Microsoft Office,

Google Docs, etc.) can open between 15-20 TCP connections. Delays in opening those TCP connections impact overall application performance, and thus impact the customer's overall experience.

While LM-1 helps the cable operators to assess the impact of typical home and serving group utilization levels on the network latency and jitter, LM-2 helps to assess the bufferbloating issues by measuring the latency and jitter at the highest home utilization levels.

The LM-2 method can be used within the MSO's network without requiring any HW test client in the subscriber's home as shown in Figure 3. In this case the test client is implemented within the eRouter firmware (e.g. RDK-B). The test can be initiated only when the subscriber's gateway (i.e. erouter+CM) has low utilization level not to degrade the subscriber's network performance. Hence, the connection RTT portion will be mostly ~10-20ms for current DOCSIS networks. This RTT includes media access delay in the US but queuing delay may not be high since the test is initiated when the home traffic volume is low.

Queuing delays in the CM and CMTS are the dominant factors in the US and DS LUL measurements respectively. UDP pings are impacted by queue building iperf TCP traffic since they are transmitted within the same service flow. LUL provides information on the latency variation subscribers may perceive when they use low latency applications. It reflects the impact of instantaneous bursts of queue building traffic (e.g. file download/upload, streaming) on the performance of the LL applications (e.g. gaming, videoconferencing).

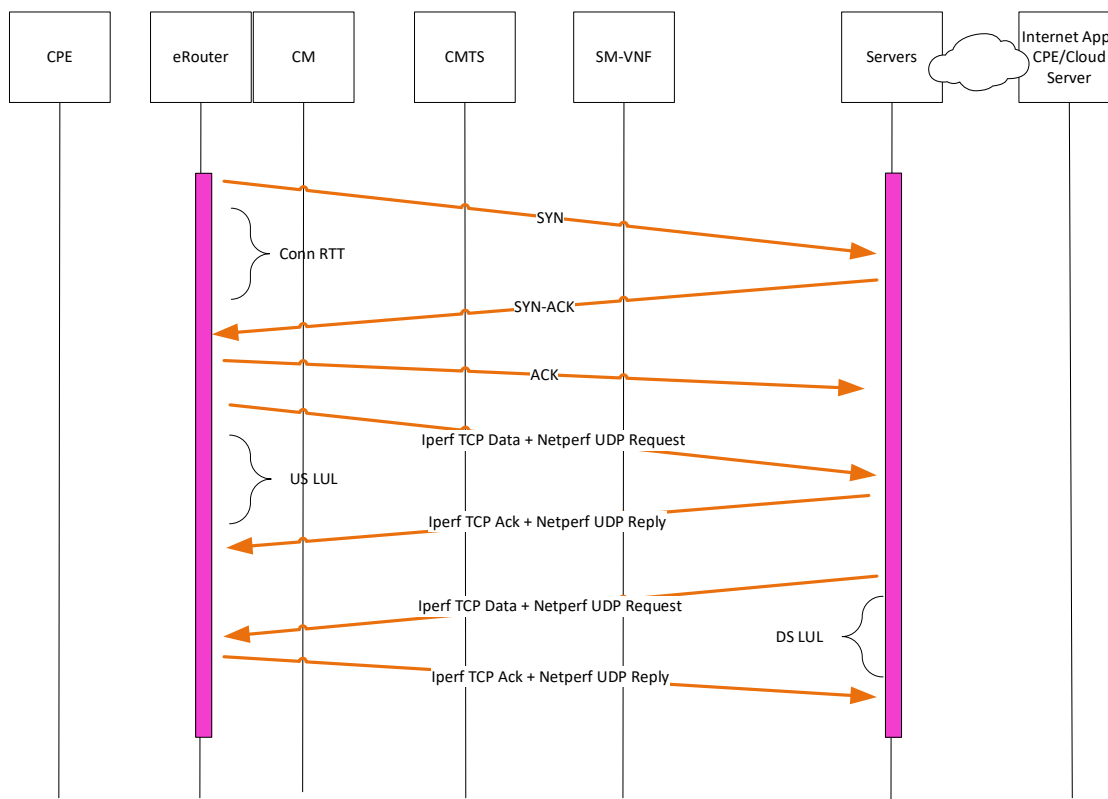


Figure 3 – Latency Under Load Measurement

3. Latency Performance in Current and Emerging Network Architectures

DOCSIS NETWORKS LATENCY

As summarized in Table 2 – DOCSIS Networks Latency: RTT Between CM and CMTS Table 2 [3], the latency performance of DOCSIS networks has been improving with each version. A major delay source is queuing [3], that has been managed by the buffer control (BC) and active queue management (AQM) introduced in the D3.0 and D3.1 specifications. However, these improvements are not adequate for low latency service requirements and for providing MSOs a competitive edge in 10G era. Note that although CMs have mandatory AQM specifications (i.e. the PIE algorithm), queue management at the CMTS varies based on the vendor and its firmware releases. Table 2 levels correspond to latency under US load. As a result, latency under DS load and bi-directional load may be higher than shown in Table 2 for existing implementations. For example, we can confirm US LUL results with ~10-20ms for ~90th percentile, but when we run a bi-directional load, and depending on the CMTS implementation and configuration, the US LUL results can be in the order of a few hundred of milliseconds. CMTSs may have large physical queues. If cable operators do not optimize the configurations for the BC and AQM features of older CMTS deployments, they may create high queueing delays and inefficient network utilization. In Figure 4, below, we show some examples with test cases and configurations to depict the performance changes.

Table 2 – DOCSIS Networks Latency: RTT Between CM and CMTS

		When Idle	Under Load	95-99 th Percentile
DOCSIS 3.0 Early Equipment		~10ms	~1000ms	~1000ms
DOCSIS 3.0 w/ Buffer Control		~10ms	~100ms	~100ms
DOCSIS 3.1 w/ Active Queue Management		~10ms	~10ms*	~100ms
Low Latency DOCSIS 3.1	Dual Queue	<10ms	<10ms	<10ms
	PGS	~1ms	~1ms	~1ms

- Latency under DS load and bi-directional load may be ~50-100ms with current CMTS AQM implementations.

In the examples shown in Figure 4, the pre- and post-optimization settings for the DS AQM are shown. As seen in Figure 4, ~1-2 seconds of DS LUL is observed for pre-optimization, while post-optimization reduces this latency to the order of ~100-200ms. The suboptimal settings cause higher queueing delay caused by bursty queue-building traffic, as CMTSs may have large physical queues. Latency is reduced with the optimized settings and speed tests results show that there is no degradation in throughput. Note that outliers, e.g. some unexpected high latency results, due to test failures are not removed in these examples. When a D3.1 CM's US AQM is disabled, the default BC is 250ms as defined in the earlier D3.1 specifications, with US LUL around 250-300ms, which can be seen in the graphs (red CMs). This is also displayed in Figure 5 that shows the test US LUL results per speed tier rates and CM models, which are configured for different BC and AQM settings. The results show that when appropriate features are not enabled or configured with optimized settings, high US LUL may be observed. For example, a D3.1 CM with US AQM disabled and default BC settings (i.e. 250ms) has higher US LUL compared to D3.1CMs with US AQM enabled and D3.0 CMs with lower buffer sizes.

Figure 6 displays dslreports.com test results for the same CMTS and CM with and without optimized AQM settings. In this case, the test is done from a laptop connected via wired Ethernet to the CM. Optimized AQM settings improve LUL results (grade for bufferbloat) while still meeting throughput requirements.

These examples show that cable operators should measure LUL and audit CM AQM and BC configurations to make sure optimized settings are deployed.



Figure 4 – Top Left: Max DS and US LUL with suboptimal DS AQM settings; Bottom Left: Mean DS and US LUL with suboptimal DS AQM settings; Top Right: Max DS and US LUL with optimized DS AQM settings; Bottom Right: Mean DS and US LUL with optimized DS AQM settings

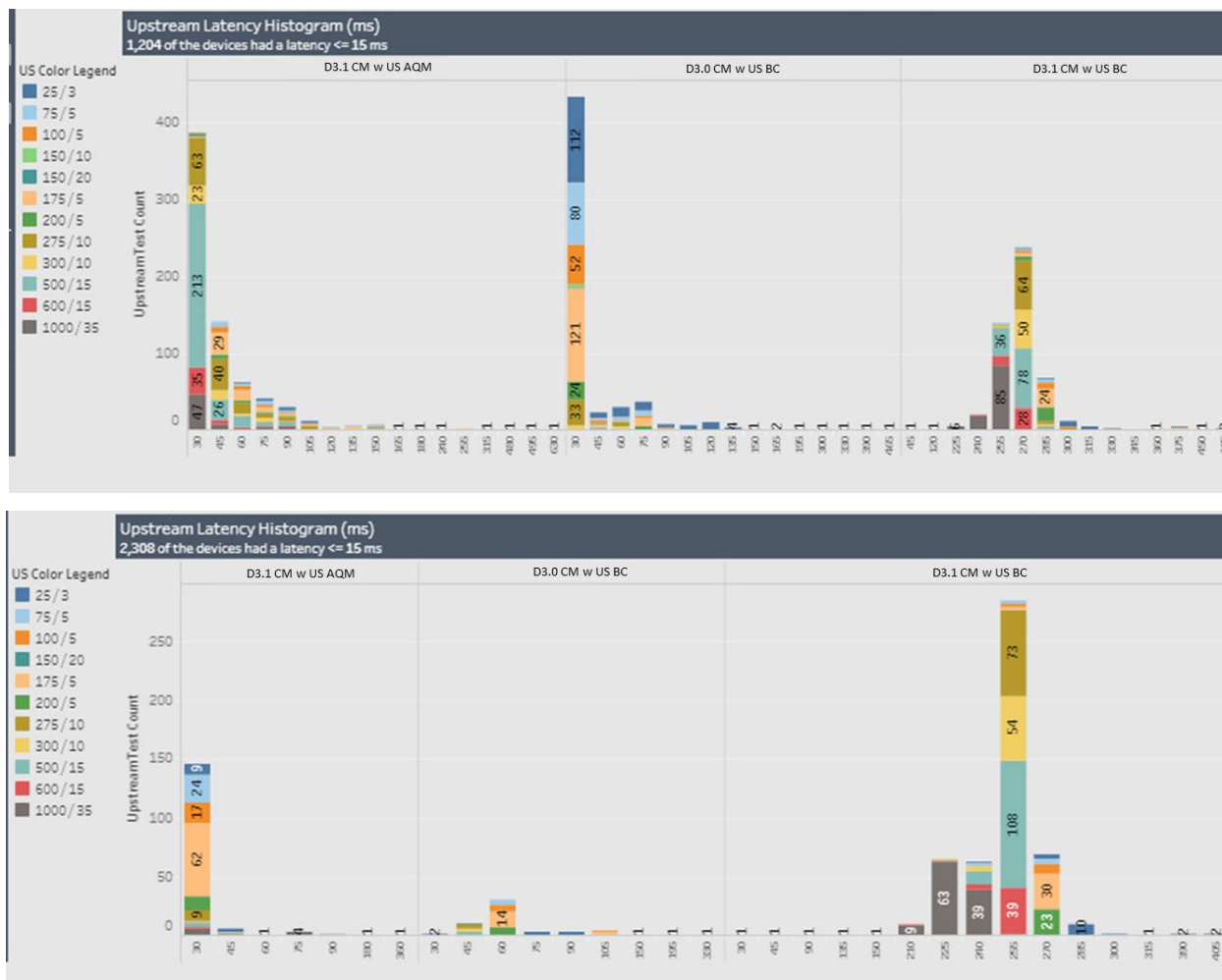


Figure 5 – Top: Max US LU; Bottom: Mean US LUL for different CM models and speed tiers

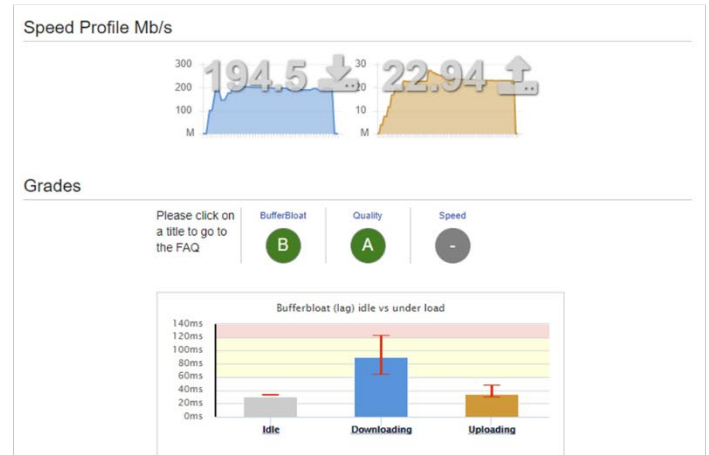
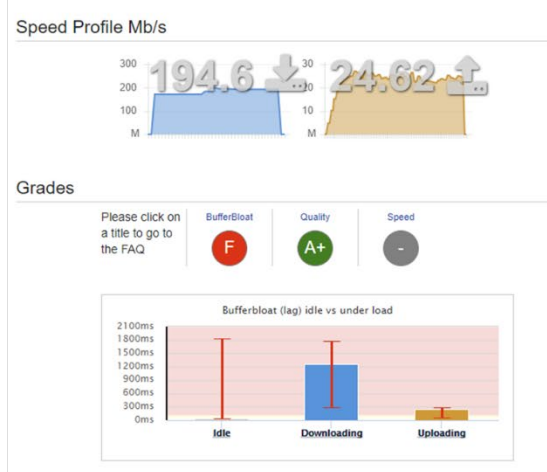


Figure 6 – Left: dsreports.com results with suboptimal DS AQM settings and D3.1 CM with BC with 250ms default target latency; Right: Optimized DS AQM settings and D3.1 CM with US AQM with 10ms default target latency

Depending on the buffer control and AQM, other factors, such as DOCSIS configurations, speed tier rates, additional RTT over the data path (e.g. additional delay at the home or northbound of the MSO's network), transport and congestion control protocols, duration of flows, rate adaptation schemes and the mix of traffic and utilization levels, all can affect the system performance. For example, Figure 7 shows that buffer sizes in a D3.0 CM affect how much throughput each concurrent flow within the subscriber's HSD SF gets, with different end-to-end RTTs. The graphs show US throughput for each flow within an HSD service of 10Mbps US speed tier rate. Although it may be desirable for a flow with a smaller RTT (e.g. an edge computing service) to have better performance, starving flows with long end-to-end RTTs should be avoided. In this example, 10, 30 and 50 ms target buffer sizes are set for the BC. 4 flows are destined to different servers, green flow's server is the closest to the subscriber's home while the purple flow's server is the farthest away. MSOs may improve their network with new features while optimizing configurations for the earlier versions of DOCSIS components.

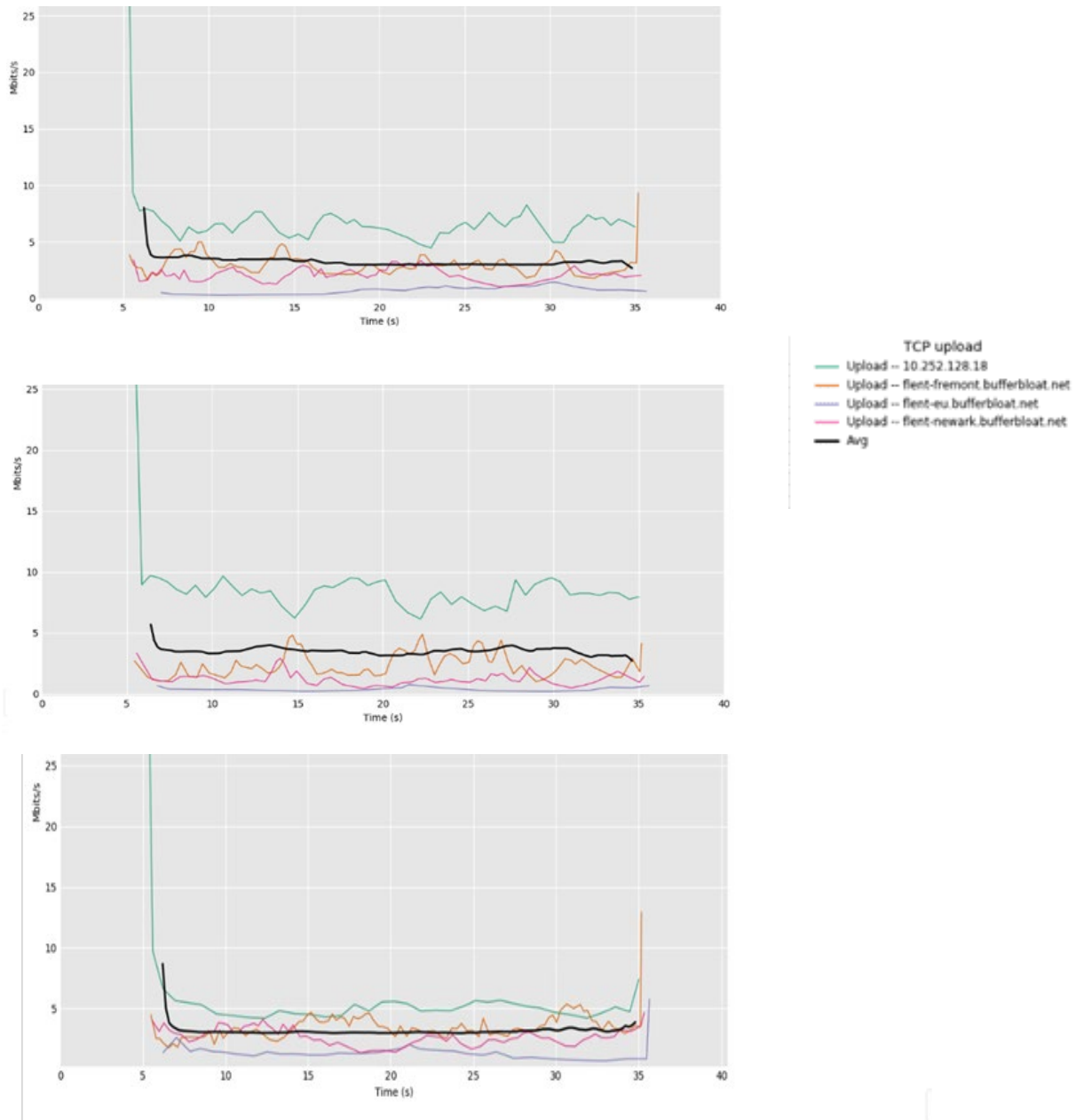


Figure 7 –Throughput (Mbps) values over time for each flow per CM HSD buffer size. Each flow has different e2e RTT. Green flow’s server is the closest to the subscriber’s home while purple flow’s server is the farthest away. Top: 10ms target buffer size; Middle: 30ms target buffer size; Bottom: 50ms target buffer size

As described in Section 3, new D3.1 Low Latency DOCSIS (LLD) features with dual queue AQM and reduced bandwidth allocation map (MAP) interval time will enable non-queue-building low latency traffic to have an RTT of <10ms for the 95-99th percentile. This will improve gaming experiences, because gaming control traffic can be classified as NQB LL traffic and can be transmitted with much better jitter characteristics. Work is ongoing as it relates to scalable congestion control algorithms, so that streaming services can be also transmitted as NQB LL traffic. This area requires more research and development to analyze traffic characteristics and their impact on other service flows.

The LLD architecture is premised on the concept of separating Queue Building (QB) traffic from Non-Queue Building (NQB) traffic. An example of QB traffic is the typical file transfer. Online gaming control traffic is an example of NQB traffic. Current network deployments carry both QB & NQB traffic, transiting in the same service flow. When the service flow is unsaturated, this is not a problem. But when the service flow is at capacity, it is quite possible that NQB traffic gets stuck behind QB traffic. When this occurs, the NQB traffic incurs latency delays while the queue drains out the QB traffic and the NQB traffic awaits its transmit opportunity.

The LLD architecture attempts to alleviate this issue by separating QB & NQB traffic into separate sub-service flows, each with their own queue, as part of an Aggregate Service Flow (ASF). The available bandwidth of the overall ASF is still the same as what the customer has purchased, but the NQB & QB queues are drained in a manner that allows NQB traffic to be rapidly dispatch from the queue, while also ensuring that the QB queue receives adequate transmit opportunities.

Proactive grant scheduling is another D3.1 LLD feature, targeting a ~1ms RTT for the 99th percentile. Network efficiency for PGS must be analyzed for possible use cases.

Wi-Fi NETWORKS LATENCY

Although avoiding Wi-Fi can result in more deterministic latency, jitter and packet loss, subscribers run low latency services such as gaming and videoconferencing over Wi-Fi all the time, because it is simply more convenient than wired Ethernet cables. Features such as WMM and AQM [4],[5] aim to improve latency and jitter for low latency services in Wi-Fi networks. Similar to DOCSIS networks, the tradeoffs between latency/jitter and throughput and fairness must be balanced (e.g. the tradeoff between efficiency with frame aggregation, vs faster channel access).

Although significant improvements may be observed with correct queue management and channel access control, factors such as outside interference and high concurrent utilization have been limiting the Wi-Fi performance for low latency services. The need for a more deterministic quality of service has been supported by Wi-Fi 6 (802.11ax) with MU-MIMO and OFDMA, enabling both DS and US increased simultaneous communications [7]. Removing contention along with deterministic QoS features (e.g. multi-user EDCA) that can provide low latency, jitter and overhead will be key to supporting low latency services.

In addition, Wi-Fi 6E, with its additional 1200 MHz in the 6 GHz spectral range, would be suitable especially for MDU type environments, but can be also used with multiple routers in large houses in the

future. Although 6 GHz is used by other networks, and more testing is required for range-speed-latency characterization, all of these improvements point to an optimistic future for low latency services. Most residential applications will not need Gbps symmetrical throughput anytime soon, but having this kind of speed both on the Wi-Fi and DOCSIS networks enables new services, including commercial and industrial services. Enterprise trials [6] have proven 2 Gbps speeds with consistent connections and latency around 2ms.

CORE NETWORKS LATENCY

Game developers and providers, gaming router developers and network optimization companies have been developing products to optimized routing/tunneling to best gaming servers by measuring RTTs and/or establishing private networks. SDN-WAN has been applied for low latency services to intelligently shift traffic and dynamically adjust to network and traffic changes. Delivery at the edge by linking up data centers with ISPs' last mile networks has been proven to reduce latency and jitter significantly and improve fairness among online gamers.

MSOs have a unique position to provide the best end-to-end performance for gamers by applying low latency features at the access and home networks as well as at their edge core and peer routing platforms.

4. End-to-end Support For Low Latency Services

In addition to support LL service performance requirements, MSOs need to provide architecture changes for LL service classification, marking, service integration and assurance. Table 3 summarizes the main features for those architecture changes, and Figure 8 illustrates them for the corresponding components.

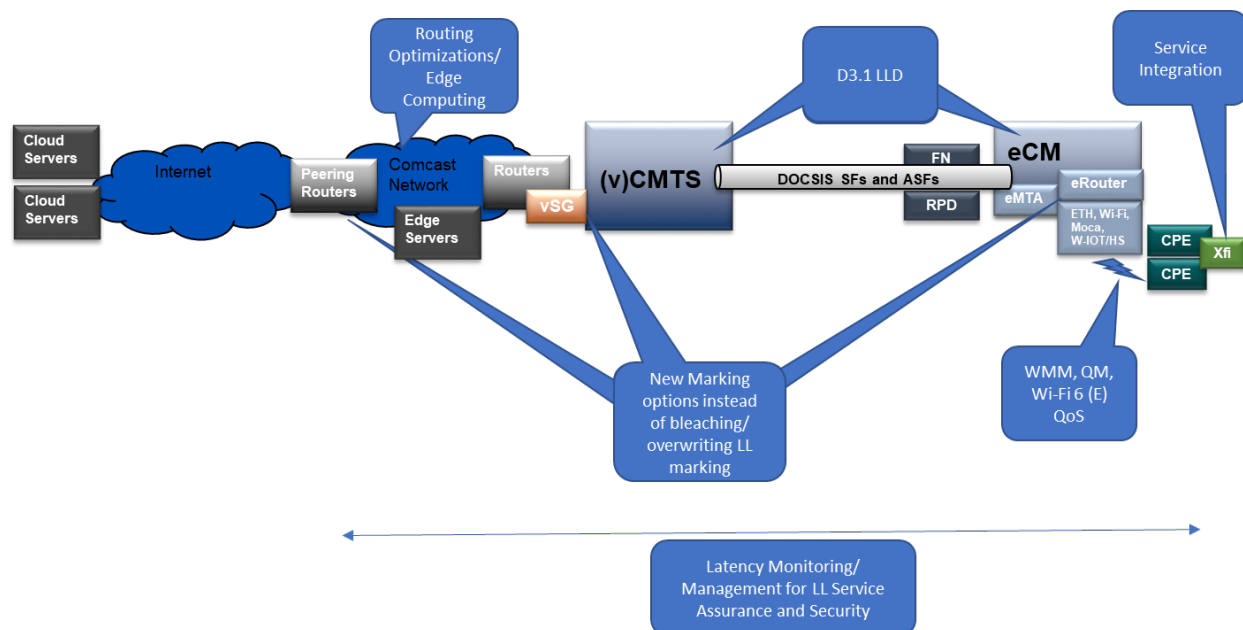


Figure 8 – End-to-end LL services support

Table 3 – Architecture changes to support LL services

Current Architecture

- No classification and marking in untrusted network segments
- Single queue for all HSD applications and scheduler (in DOCSIS and most Wi-Fi networks)
- Lower US speeds; Wi-Fi contention and unreliable QoS
- Single performance measurement for all services within HSD SF
- Single network optimization technique for all services within HSD SF in DOCSIS and Wi-Fi networks
- Loose coupling in design, development, testing, and operations among network segments
- Limited business models for HSD apps

Target Architecture

- End-to-end classification with new security measures
- Dual/multiple queue and weighted scheduler for NQB-LL and other applications
- > Gbps DOCSIS and Wi-Fi with deterministic latency/jitter for LL services
- Differentiated performance measurement for NQB-LL and other applications
- Network optimization techniques per traffic requirements within HSD service flows
- Orchestrated end-to-end design, development, testing and operations for LL services
- Flexible business models for LL services

Marking options are being discussed in IETF WGs such as Low Latency, Low Loss, the Scalable (L4S) proposal and the Differentiated Services Code Point (DSCP) marking proposal[8][9]. Until these approaches are adopted widely, MSOs may support LL NQB traffic classification by changing their current marking options and taking new security measures. In this sense, another VNF envisioned for the VSG platform is in support of the Low Latency DOCSIS architecture.

This architecture can improve the customer experience by having QB & NQB traffic serviced, separately as addressed in Section 3. The challenge for the network operator then becomes how to separately identify NQB from QB traffic as it traverses the network. Assuming that business rules provide guidance on what qualifies as QB & NQB traffic, the data packets from each type of traffic must have some unique identifier, so that the CMTS can direct each type of traffic into the correct sub-service flow within the ASF. Here, a VNF can be deployed onto the VSG which includes rule mapping, to re-mark the data packets of the two types of traffic. This could be through changing the packet's DSCP value, based on other packet values, such as source or destination IP address, port number, protocol number, etc. A VNF supporting this re-marking can be dynamically updated to support changing business rules on the classification of QB vs NQB traffic, and re-mark (or not) accordingly. Thus, by providing a platform to instantiate this VNF, the VSG platform can play a part in implementing the overall LLD architecture live on the network.

VNFs such as a VSG platform are also crucial for end-to-end latency management systems. Figure 9 illustrates an example monitoring tool that collects hop-by-hop information such as network (including AQM, BC, and QoS MIBs), utilization and latency test results. A VSG or similar VNF can detect anomalies or changes that require correction actions. This information can also be used for performance

prediction and business intelligence purposes. Such a system may have limitations due to SNMP based polling scales, different monitoring systems with varying collection times and lack of efficient VNF to network segment mapping. These limitations are especially crucial for latency management since instantaneous changes in the network and traffic conditions may create lag spikes that are hard to analyze. New, push-based telemetry systems, that MSOs have been deploying for SDN/NFV enabled distributed architectures, can overcome these limitations, and are shown in Figure 10. Monitoring agents may be distributed and linked to local VNFs to have aggregated data processed in a cloud-based end system.

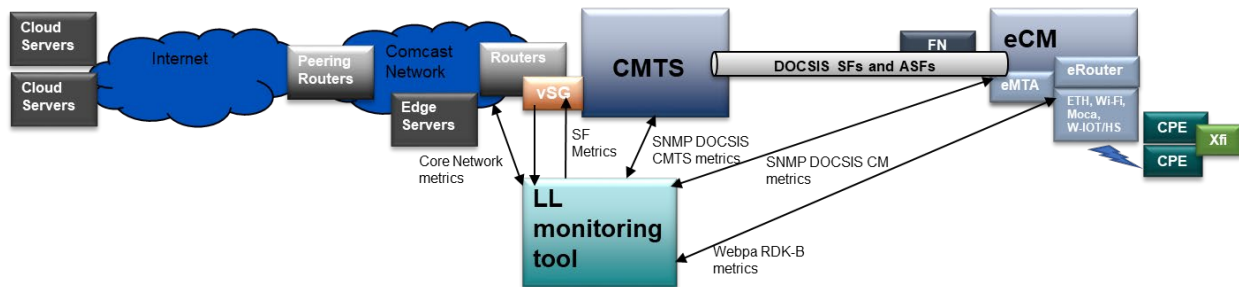


Figure 9 – Low Latency Services Monitoring and Management for CCAP Systems

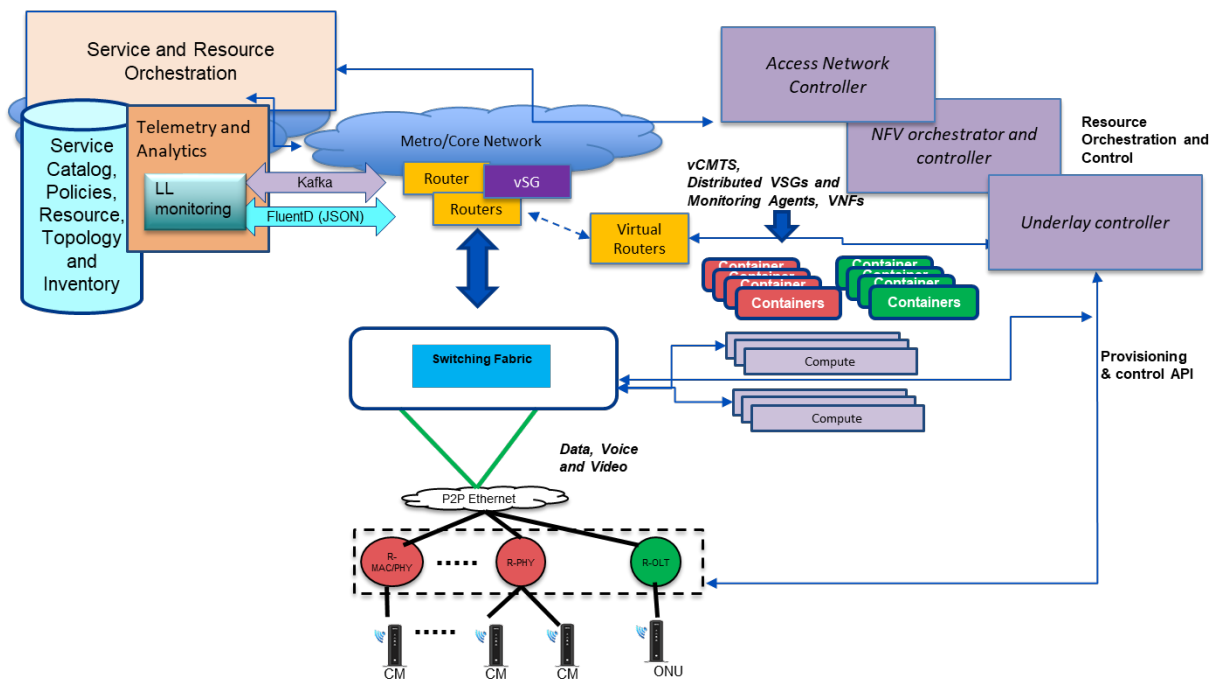


Figure 10 – Low Latency Services Monitoring and Management for Distributed Systems

Figure 11 displays an architecture where resources may be managed and configured based on monitoring and prediction systems control [2]. LL monitoring and management can then be part of the data center as an orchestration function. Today, even some simple configurations may not be available, because of a lack of provisioning flexibility. For example, buffer sizes in older modems may be expressed in bytes, and setting them based on speed tiers may require the coordination of service class names and boot file settings. There may be other configuration parameters that need to be coordinated as well. Flexible configurations and operations may seem like features for the farther future, but today MSOs already deploy SDN/NFV enabled distributed systems and new telemetry platforms. Configuration flexibility is already an integral part of such systems. The next section discusses new architectures where low latency services may be integrated.

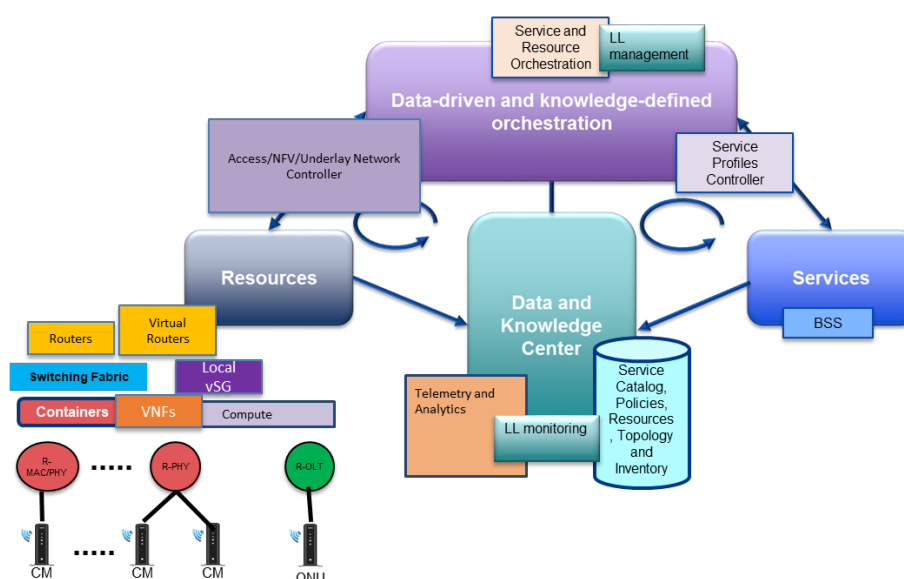


Figure 11 – Low Latency Services Monitoring and Management Integrated within Data-driven and Knowledge-defined Architectures

5. Conclusion: Final Thoughts on Latency Management

Although service assurance with performance management has been always the main driving force in designing access network architectures, a unified platform with an end-to-end orchestration approach has not been fully adopted, largely because of design limitations in many operators' networks. Recent changes in MSO network and service architectures provide the building blocks for such a unified platform, including:

- *Networking improvements:* New features in both wireline (e.g. DOCSIS/PON) and wireless (e.g. Wi-Fi) networking improve both the customer experience and network efficiency significantly. These features include frequency split and extended spectrum in DOCSIS, FDX, higher rate PON technologies and coherent optics, 802.11ax (Wi-Fi 6), low latency and distributed architectures for DOCSIS, Wi-Fi and mobile networks. Although some of these technologies aim at a specific hop or segment of the MSO network, initiatives like low latency networks target end-to-end improvements.
- *Data-driven networks and Monitoring:* Changes in telemetry, e.g. adapting push-based and cloud-hosted telemetry, enable data-driven network and service architectures. Both hop-by-hop, end-to-end latency and other performance metrics can be collected for overall latency management, troubleshooting, operations and planning purposes. Low latency services such as gaming will benefit from different latency measurement approaches, including concurrent and multi-hop measurements.
- *Software Defined Networking:* SDN enables centralized orchestration and coordination of the distributed controllers in different network and service segments. Dynamic and flexible configurations will help low latency services such as gaming, for example, by avoiding separate configurations of home and access network components. Some of the traditional latency measurement techniques assume that the control and data planes overlap, which wouldn't be the case for SDN networks. On the other hand, SDN enables an end-to-end data path view, with associated capabilities and monitoring that may be easily controlled for hop-by-hop and end-to-end measurements.
- *Network Function Virtualization:* MSOs have been introducing new VNFs over the control and data paths, with innovative functionalities in the areas of subscriber and service flow management that can help the differentiation of services per their traffic requirements. Virtualization in access networks may help to integrate new queueing and scheduling functionalities -- while special design requirements need to be considered for low latency services, as these designs may introduce additional latency not found in purpose-built, hardware-based architectures.
- *Knowledge-defined Networking:* Advances in the application of machine learning (ML) techniques to MSO networks and services open new doors for better performance prediction and management with self-optimizing capabilities. Recent advances in proactive network management (PNM) in DOCSIS and Wi-Fi networks can be extended for low latency services. The advantage of a knowledge-defined network is the ability to apply multi-hop PNM for end-to-end service assurance. In addition to smart networks and operations, MSOs have been using knowledge-defined systems for smart homes and customer interfacing platforms, which facilitate new low latency service offerings while assessing the customer experience for these services.
- *Cloud-based applications vs edge computing:* MSO integration of both cloud and edge computing based applications, depending on the service requirements, will enable low latency service providers to select the best architecture for the optimized customer experience. For example, cloud-based game providers can optimize network peering or consider edge computing based on performance, hardware and cost requirements.
- *Open source products and standards:* Many standardization efforts like the low-latency work within CableLabs and IETF enable support in a larger ecosystem. MSO use of open source products, and flexible integration of third party services in their platform, ensure the compliancy with regulations and policies as well.
- *Security:* Security, in terms of network and customer privacy protection, has become a vital item to be integrated into the design, instead of a later add-on feature. Low latency services require new classification and marking design, which requires new security elements. A proactive design approach is pivotal to a secure end-to-end solution, even while each network segment can have its own security feature (e.g. queue protection in D3.1 LLD specs).

- *Accessibility:* Similar to security in proactive inclusions, accessibility was a core part of the initial architecture design process for MSOs, instead of a later add-on feature. Low latency services may include strategies that offer control options to subscribers. MSOs that have an established framework to incorporate accessibility requirements early in the design can easily integrate low latency services by meeting every subscriber's needs.

Abbreviations

AR	Augmented reality
ASF	Aggregate service flow
AQM	Active Queue Management
BC	Buffer Control
CM	Cable modem
CMTS	Cable modem termination system
COTS	Commercial off-the-shelf
DOCSIS	Data over cable service interface specification
DS	Downstream
DSCP	Differentiated Services Code Point
IETF	Internet Engineering Task Force
ISBE	International Society of Broadband Experts
LL	Low Latency
LLD	Low Latency DOCSIS
LM	Latency Measurement
LUL	Latency Under Load
MAP	Bandwidth Allocation MAP
ML	Machine learning
MIB	Management Information Base
MU-MIMO	Multi-User Multi-Input Multi-Output
NFV	Network Functions Virtualization
NQB	Non-queue-building
PNM	Proactive Network Management
QoE	Quality of experience
QoS	Quality of service
QB	Queue-building
PIE	Proportional Integral Enhanced
RTT	Round-trip time
SCTE	Society of Cable Telecommunications Engineers
TCP	Transmission Control Protocol
SDN	Software Defined Networking
US	Upstream
VNF	Virtual Network Function
VR	Virtual Reality
VSG	Virtual Subscriber Gateway
WG	Working group
WMM	Wi-Fi MultiMedia

Bibliography & References

- [1] *Supporting The Changing Requirements For Online Gaming*, K. Scott Helms, SCTE-ISBE Workshop 2018
- [2] *The Future of Operations: Building a Data-Driven Strategy*, Sebnem Ozer, Sinan Onder, and Nagesh Nandiraju, “SCTE Journal on Network Operations, 2018
- [3] *Low Latency DOCSIS: Overview And Performance Characteristics*, Greg White, Karthik Sundaresan and Bob Briscoe, SCTE-ISBE Workshop 2019.
- [4] <https://www.cablelabs.com/10g/latency#:~:text=Our%20Low%20Latency%20DOCSIS%20technology,slowing%20down%20all%20other%20data.>
- [5] https://www.bufferbloat.net/projects/bloat/wiki/What_can_I_do_about_Bufferbloat/
- [6] <https://wballiance.com/wbas-first-phase-of-wi-fi-6e-trials-shows-the-massive-potential-of-wi-fi-in-the-6ghz-band/>
- [7] *The Importance of Wi-Fi 6 Technology For Delivery of Gbps Internet Service*, David John Urban, SCTE-ISBE Workshop 2019.
- [8] <https://datatracker.ietf.org/doc/draft-ietf-tsvwg-l4s-arch/>
- [9] <https://datatracker.ietf.org/doc/draft-ietf-tsvwg-nqb/>

Acknowledgments

The authors would like to thank and acknowledge all those who helped to make this paper possible. This paper includes several tests done and dashboards created by Aaron Tunstall, Sarulatha Subbaraj, Soomin Cho, Peifong Ren, Ray Hammer, Lei Zhou and Joe McHale.