



VIRTUAL EXPERIENCE
OCTOBER 11-14



How Network Topology Impacts Rf Performance: A Study Powered By Graph Representation Of The Access Network

A Technical Paper prepared for SCTE by

Maher Harb

Director of Data Science, Next Generation Access Network
Comcast
maher_harb@comcast.com

Karthik Subramanya

Senior Engineer, Next Generation Access Network
Comcast
karthik_subramanya@comcast.com

Ramya Narayanaswamy

Senior Manager, Next Generation Access Network
Comcast
ramya_narayanaswamy@cable.comcast.com

Sanket Walavalkar

Executive Director, Next Generation Access Network
Comcast
sanket_walavalkar@comcast.com

Dan Rice

VP, Next Generation Access Network
Comcast
daniel_rice4@comcast.com

Table of Contents

Title	Page Number
1. Introduction.....	4
2. Constructing the Graph.....	4
3. Visualizing the Graph.....	6
4. Amplifier Cascade Length Analysis.....	9
5. Conclusion.....	15
Abbreviations	15
Bibliography & References.....	16

List of Figures

Title	Page Number
Figure 1 - Schematic representation of some of the main entities captured in the graph database.	5
Figure 2 - Screen capture from Graphistry—an interactive graph map that allows browsing attributes of the different vertices and edges by hovering over the element with the mouse tip.	6
Figure 3 - Example of a basic graph visual for a relatively small size node. The tree trunk (brown circle) corresponds to the CMTS, and the leaves (gray circles) correspond to the cable modems/IP devices. The pathway connecting the two traverses multiple physical and logical elements as shown here.	7
Figure 4 - Example of graph visual in which the thickness of the link (edge) is weighted by number of cable modems connected through the link. The label on each vertex represents the total number of cable modems that are hierarchically located underneath the vertex.	8
Figure 5 - Example of a graph visual in which the color of the cable modem symbol indicates the upstream SNR on the 25-to-40 dB scale shown in the adjacent color bar. The one cable modem colored gray had missing telemetry data for that polling time sample.	9
Figure 6 - Visual example of a large node in which devices, billing addresses, and drops were removed for clarity. Two paths from CMTS to tap are highlighted in purple. One path traverses a single amplifier and the other traverses 5 amplifiers. Devices attached to the latter are expected to experience larger amplifier distortions in the downstream path.	10
Figure 7 - Distribution of amplifier cascade length across devices (top panel) and distribution of total number of amplifiers across nodes (bottom panel). This data was generated by querying the graph database.	11
Figure 8 - The top panel shows the relationship between RxMER and amplifier cascade length for different aggregating percentiles (25th, 50th, and 75th). There's a clear downward trend, as highlighted by the linear best-of-fit lines. The bottom panel shows the corresponding modulation distribution. Once again, the trend is visible by the decreasing ratio of 4096-QAM as the amplifier cascade length increases.	13
Figure 9 - Linear Fit to the RxMER vs. amplifier cascade length data (blue line with gray confidence band). The horizontal black line, at 38 dB, is the current threshold in PMA for assigning a 4096-QAM modulation.	14
Figure 10 - Linear Fit to the US SNR vs. total number of amplifiers (blue line with gray confidence band). The horizontal black line at ~20 dB is the current threshold in PMA for assigning a 64-QAM modulation.	15



UNLEASH THE
POWER OF LIMITLESS
CONNECTIVITY
VIRTUAL EXPERIENCE
OCTOBER 11-14



List of Tables

Title	Page Number
Table 1 - Correlations between the node features and device telemetry.	12

1. Introduction

We have recently embarked on a project with the aim of capturing all the building blocks of the access network, their relationships, and their properties in a graph representation encompassing vertices and edges. This representation is to be available in a high performance and scalable graph database that allows access to the data through application programming interface (API) endpoints and in batch. The graph database mirrors the dynamic nature of the network by getting updated as customers get connected & disconnected, optical nodes get segmented, network equipment gets commissioned & decommissioned, as well as the happening of any other impactful network change. Having all the relational information in one source, and combining the physical & logical elements in a single view allows analyzing the access network at any level of network topology (e.g., service group, fiber node, amplifier, tap, drop) on a use case basis. The graph database technology also allows enrichment of the data with ease by overlaying device telemetry, Cable Modem Termination System (CMTS) telemetry, and maintenance data on top in order to implement algorithms for business intelligence and troubleshooting (e.g., root cause analysis).

Building the graph database required combining and reconciling data across many different sources of the organization without identified primary keys (ids) and creating algorithms to automate inference of connections. In this paper, we share Comcast's journey into this process that is currently scaled to cover ~20% of our footprint. We present the very first use case of utilizing the network graph to study the effects of the amplifier cascade length on radio frequency (RF) performance in the upstream (US) and downstream (DS). The interest in this investigation falls within a broader question on the operational effort required to maintain nodes with large cascade of amplifiers (both in terms of depth and breadth). Our findings reveal that, predictably, longer amplifier cascade lengths exhibit degraded RF signal-to-noise ratio (SNR) -- yet with no significant impact on quality of service, likely due to the mitigating impact of the profile management application (PMA) system currently deployed in Comcast.

We acknowledge contributions to the project from former team members Doga Kerestecioglu, Matt Lord, and Athanasios Tsiaras.

2. Constructing the Graph

The Comcast access network data platform is an enterprise-scale graph data platform that maps the entire access network from the CMTS to the Customer Premise Equipment (CPE) while incorporating all of the physical and logical elements that form part of the network. At a very high level, the graph platform brings together site/headend topology, computer aided design (CAD) physical plans, telemetry, and billing systems together to construct the access network graph as one connected entity (see illustration in Figure 1).

While the individual data sources that form the building blocks of this graph data platform are highly mature with rich offerings, they have varied goals and are managed by various teams. There are no common keys that connect the boundaries between these data sources to form one connected access network graph. Hence, the need for such a graph data platform assumes significant importance to drive various data science applications and serve as a source of truth for deriving relationships between various access network components.

The platform supports a rich set of use cases from anomaly/network deterioration detection, triangulation, capacity planning [1] and various aggregation use cases. In addition, the platform uses several statistical

and data science techniques to bridge boundaries between individual data sources by inferring relationships and provides a single unified view of the access network.

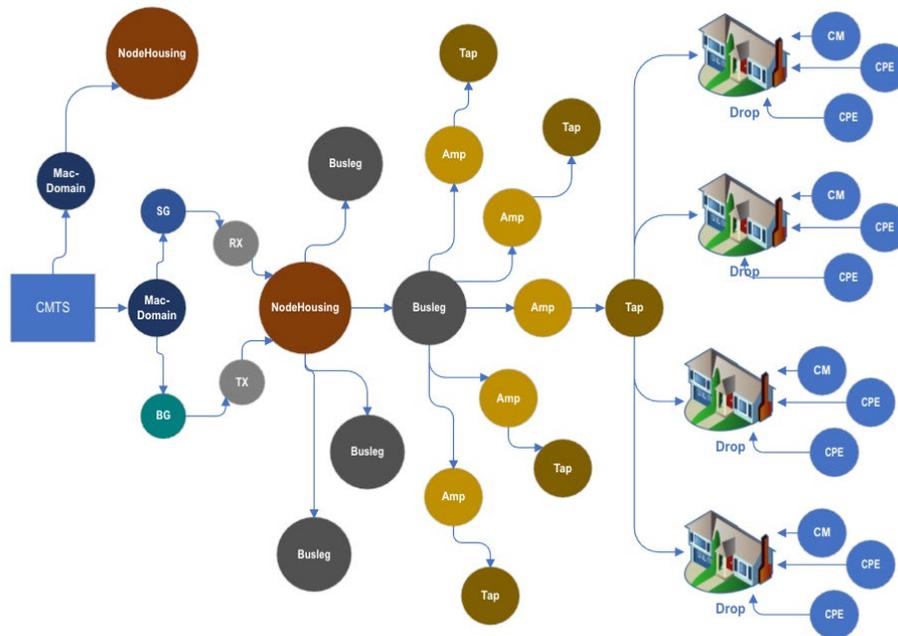


Figure 1 - Schematic representation of some of the main entities captured in the graph database.

By the virtue of being a scalable graph platform that is being purpose-built within the cloud, we're able to build multiple data pipelines to ingest various data sources and rapidly experiment and iterate through various algorithms, data processing, and refreshing techniques that feed into the end-product. The graph defines vertices and corresponding attributes for CMTS, nodes, RF equipment such as bus legs (port on a node), amplifiers, passives, taps, drops, and customer entities such as household, devices entitled, and so on. Each of these vertices are connected by edges, and in the case of equipment, these edges define the attributes of RF cables that form the connection (e.g., cable length). This graph platform allows us to query and traverse the network in either traffic direction and start at the most granular level (i.e., the CPE) and go all the way up to the CMTS or start at the CMTS and terminate at the CPE.

One of the many functions that the graph platform realizes is to diagrammatically connect physical address drops that are defined in CAD documents, to the appropriate households, which are logical elements defined in billing systems. Since there are no common keys connecting these 2 elements, the platform performs address standardization, further employs various Natural Language Processing (NLP) techniques and coordinate-based proximity to identify the right households and match them to the appropriate drops. The telemetry data set that consists of online cable modem inventory and CMTS configurations helps cross-validate these inferred relationships using reported MAC (Media Access Control) domain, bonding group, service groups, and so on.

A significant engineering challenge that was solved during the development of this platform is to account for graph refreshes. The access network consists of components that independently refresh at different rates. Hence, we needed a decoupled solution that independently refreshes various sub-graphs, even within the footprint of a single CMTS. We use property-graph architecture to implement our data model and use Apache Gremlin Tinkerpop to perform traversals, look ups, or aggregation queries. We often

notice that the depth of traversals for our queries is more than 40 layers deep, indicating the extent of cascading and density of our access network. These traversals would never be effective for any relational or NoSQL database. In contrast, through the use of large-scale data processing platforms, we're able to complete the required traversals or queries on the graph platform in a matter of minutes, for the entire footprint.

Figure 2 is a screen capture from a graphical interactive tool used to explore the graph database. It allows retrieval of attributes for vertices and edges by clicking through the graphical interface.

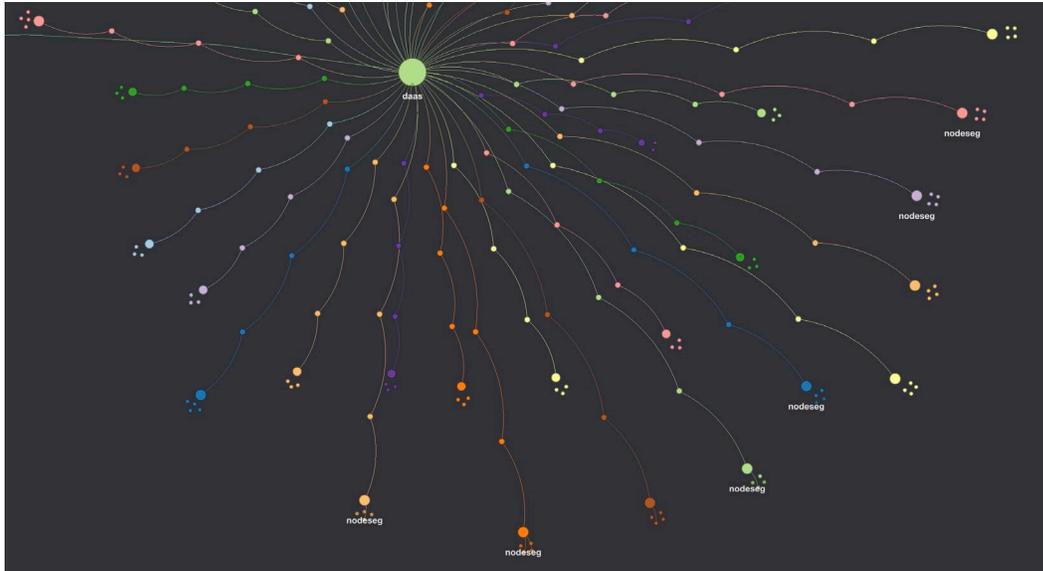


Figure 2 - Screen capture from Graphistry—an interactive graph map that allows browsing attributes of the different vertices and edges by hovering over the element with the mouse tip.

3. Visualizing the Graph

The ability to visualize the network in the form of a graph, which encompasses vertices and edges, serves a multitude of purposes. These include validating the graph construction algorithms by thorough visual inspection of the outcome, uncovering errors & inconsistencies where they may appear, and supporting use cases related to troubleshooting network issues by layering key telemetry data on top of the basic topological view. There exists a host of libraries for the purpose of graph visualization—both in the form of stand-alone proprietary software as well as open-source packages that integrate with Data Science programming language such as R and Python. For this analysis, the R packages *tidygraph* and *ggraph* were adopted for creating graph visuals. They both follow the data principles established by the popular *tidyverse* family of R packages used for data wrangling and visualization [2], thus, allowing enrichment of the graph data with ease.

The graph database at Comcast establishes relationships that span thousands of CMTSs, hundreds of thousands of optical nodes, and millions of households. For visualization purposes, the task is limited to producing a visual for a very small subset of the full network at a time and on demand. Typically, the interest is in plotting the tree-like structure that connects devices to the same fiber optics node (also referred to as a busleg/port on a node clamshell). The pathway from the device to the node includes elements such as amplifiers, passives, taps, and drops. This view is very relevant to troubleshooting of access network problems because issues originating in elements under a node may impact multiple

customers (this is especially true in the upstream direction in which noise is known to funnel). In contrast, customers connected to different nodes are usually isolated from each other. Figure 3 shows an example of a graph visual for a single node. The power of the graph is manifested in its ability to combine physical elements and logical elements in the same view. For example, the customer's billing address (shown as black square) is included in the graph and hierarchically positioned between the drop and the cable modem (tree leaf). Notice that some customers have multiple cable modems under the same billing address. These include the household internet gateway and one or multiple DOCSIS video set-top boxes.

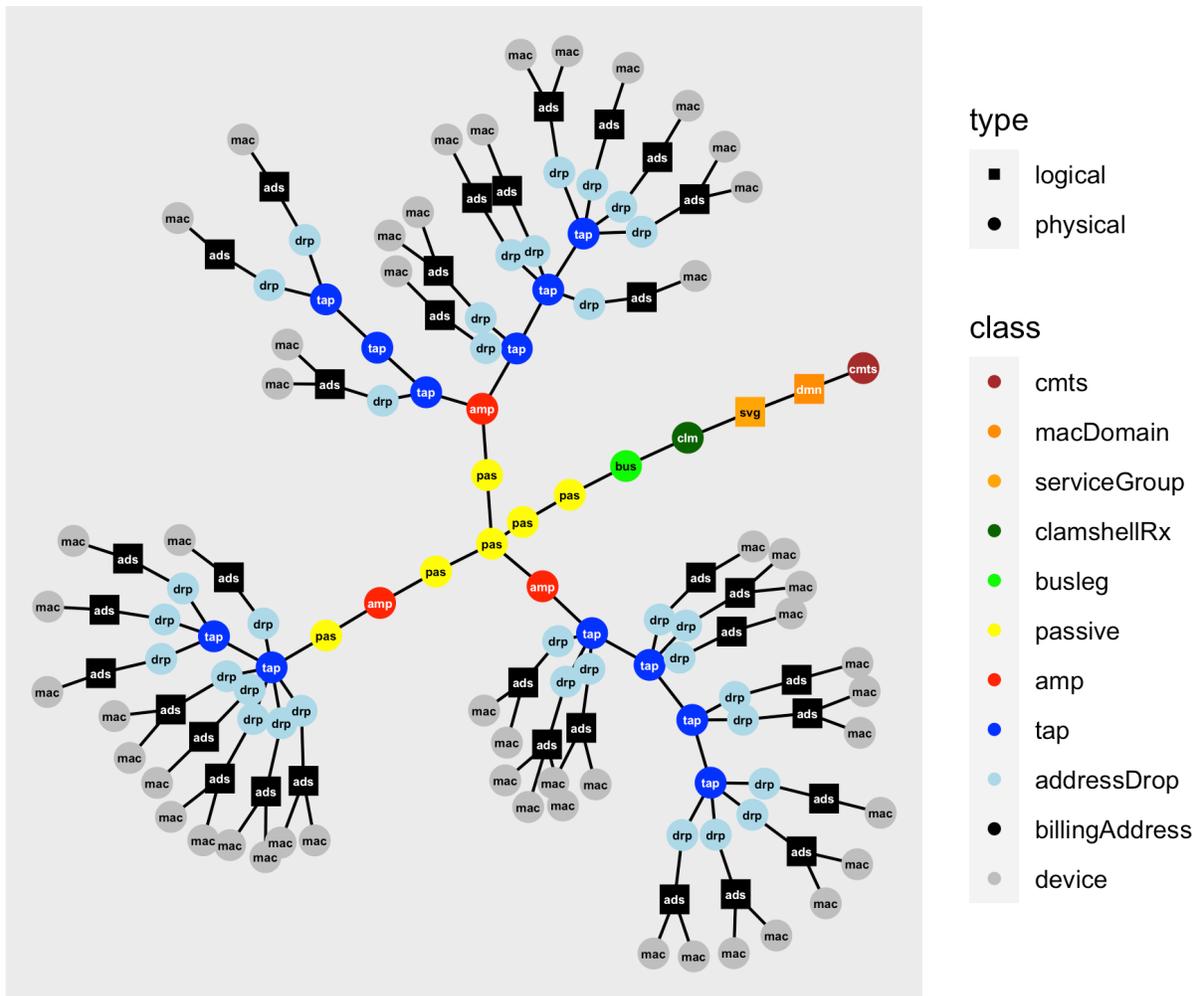


Figure 3 - Example of a basic graph visual for a relatively small size node. The tree trunk (brown circle) corresponds to the CMTS, and the leaves (gray circles) correspond to the cable modems/IP devices. The pathway connecting the two traverses multiple physical and logical elements as shown here.

Figure 3 represents a basic view of the graph that can be enriched in many ways. A few illustrative examples are considered henceforth. In the first example, shown in Figure 4, the thickness of the edge is utilized to designate the traffic volume as measured by the number of devices that transmit/receive traffic through that link. In the same figure, the number of devices is also annotated on the graph vertices. Alternatively, the weights can be easily adjusted to correspond to actual traffic volume rather than number of connected devices.

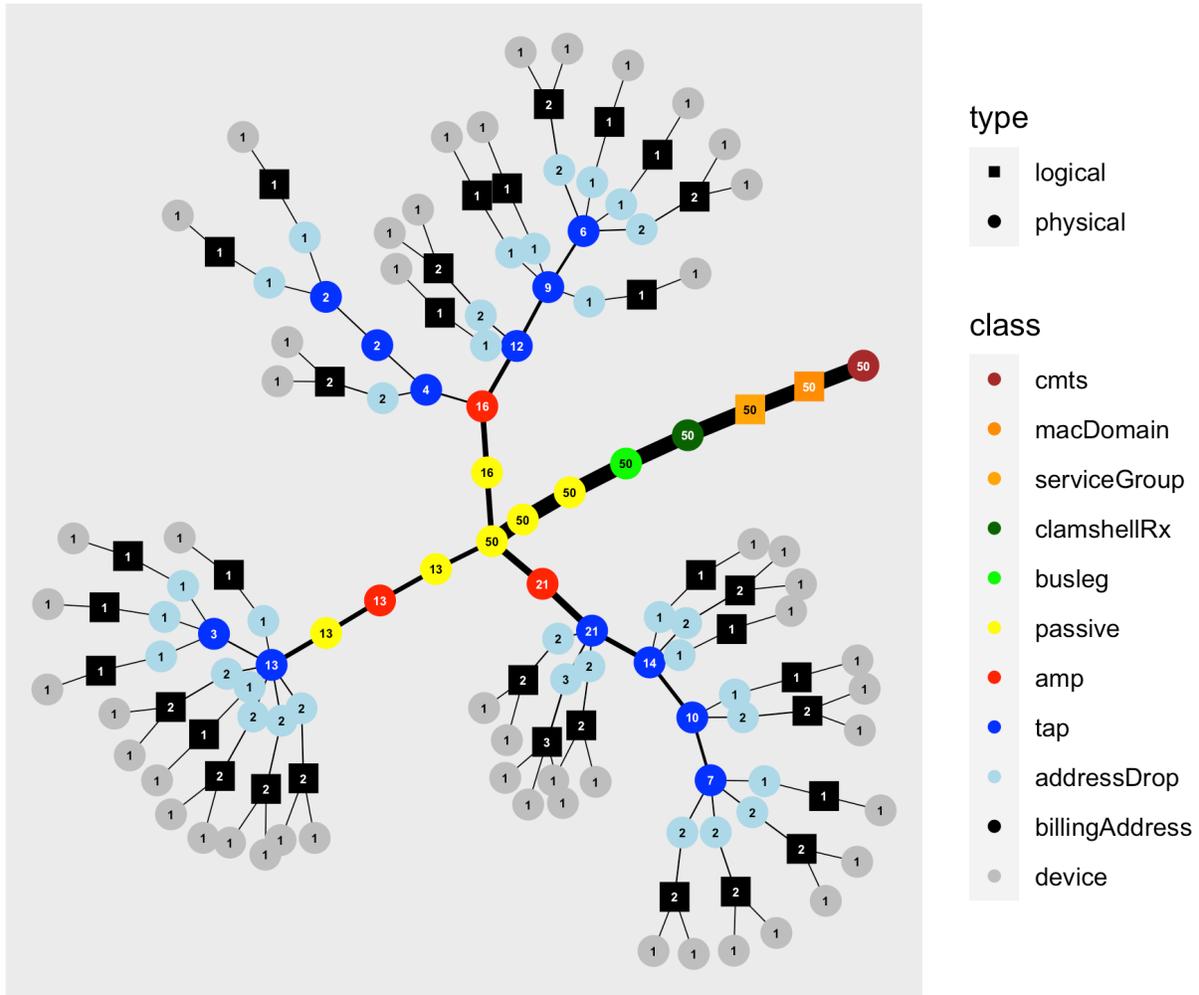


Figure 4 - Example of graph visual in which the thickness of the link (edge) is weighted by number of cable modems connected through the link. The label on each vertex represents the total number of cable modems that are hierarchically located underneath the vertex.

The example shown in Figure 5 overlays telemetry data onto the topological view. In this example, the US device-level SNR is coded as the color of the device symbol. Such view is useful for visual identification of “hot spots”. These could be a cluster of devices suffering from the same impairment and in which the graph visual provides a clue to the problem root cause.

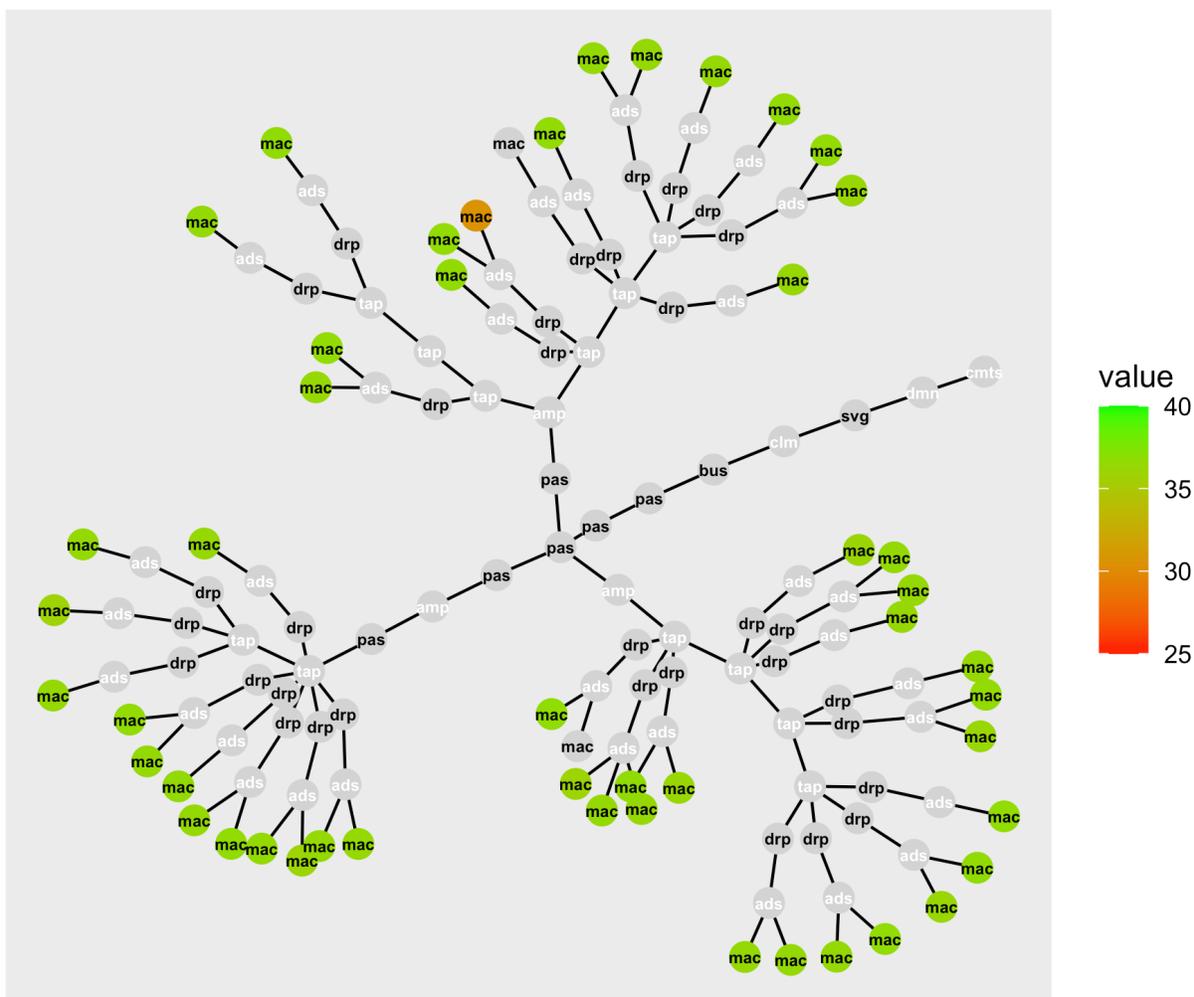


Figure 5 - Example of a graph visual in which the color of the cable modem symbol indicates the upstream SNR on the 25-to-40 dB scale shown in the adjacent color bar. The one cable modem colored gray had missing telemetry data for that polling time sample.

4. Amplifier Cascade Length Analysis

A question that arises frequently in discussions around network design best practices is the impact on the customer experience of a large cascade of amplifiers. From an RF design perspective, it is accepted that amplifiers introduce distortions to the signal—there is no such thing as an “ideal amplifier”. Albeit there is an implicit assumption that distortions are tolerated if they are deemed to be within an acceptable design range. However, the reality is that network growth is dictated by customer demand and geography, and often deviates from original plans, causing limitations on node size to be exceeded. In this context, no hard written rule exists on the maximum allowable amplifier cascade length. Furthermore, there is no straightforward approach to quantifying the impact of the amplifier cascade length on customer experience. In fact, examining this question was the very first use case of the network graph database.

The analysis dataset contains the ~20% of our footprint’s CMTSs that are fully captured in the graph database at the time of writing this paper (summer of 2021.) The basic idea behind the analysis is to correlate the node size with key telemetry data and explore this relationship in depth. The very first task within the analysis requires defining what is meant by “node size” in relation to the RF amplifiers. It was decided to explore the following two features:

- **Amplifier cascade length:** This is a device-level feature that represents the total number of amplifiers traversed in the path between a cable modem and the node (see Figure 6).
- **Total number of amplifiers:** This is a node-level feature that represents the count of all amplifiers within the node.

The rationale behind the choice of cascade length and total amplifier count was to accommodate the different noise accumulation behaviors between downstream and upstream paths. In the downstream, noise does not funnel between devices. Therefore, what matters from a distortion perspective is the number of amplifiers in a device’s path. Whereas in the upstream, noise funneling causes all amplifiers to potentially contribute to the distortion for any given device in the node.

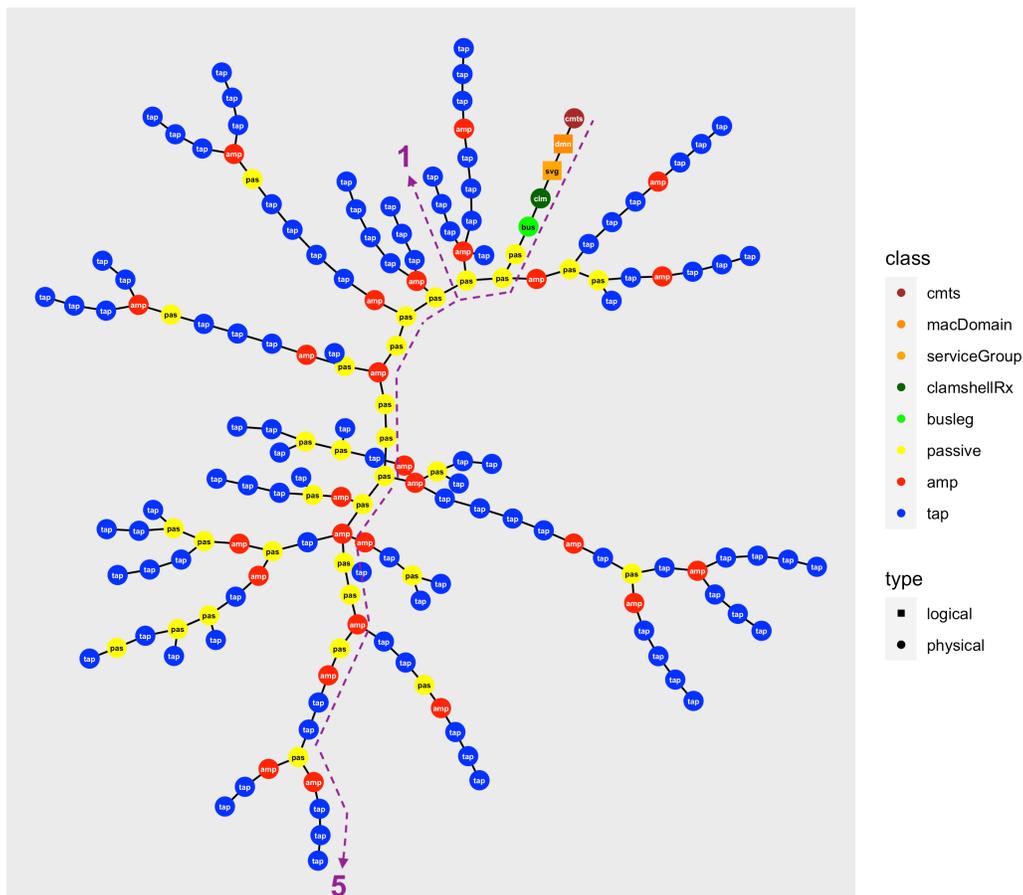


Figure 6 - Visual example of a large node in which devices, billing addresses, and drops were removed for clarity. Two paths from CMTS to tap are highlighted in purple. One path traverses a single amplifier and the other traverses 5 amplifiers. Devices attached to the latter are expected to experience larger amplifier distortions in the downstream path.

Figure 7 shows the distribution of the two features across the ~20% of CMTSs mapped in the graph database. The amplifier cascade length distribution is unimodal with a peak at 2, meaning that the majority of devices connect to the node via a cascade of 2 amplifiers. The distribution in the figure was intentionally cut off at 10, even though the data contains outliers with cascade lengths that exceed this value. The total amplifier distribution has a peak at 1 and gradually diminishes at ~40 amplifiers. It may be surprising to see ~3000 nodes served by a single amplifier. These are small size nodes in terms of either the number of customers or the geographical extent of the node (or both).

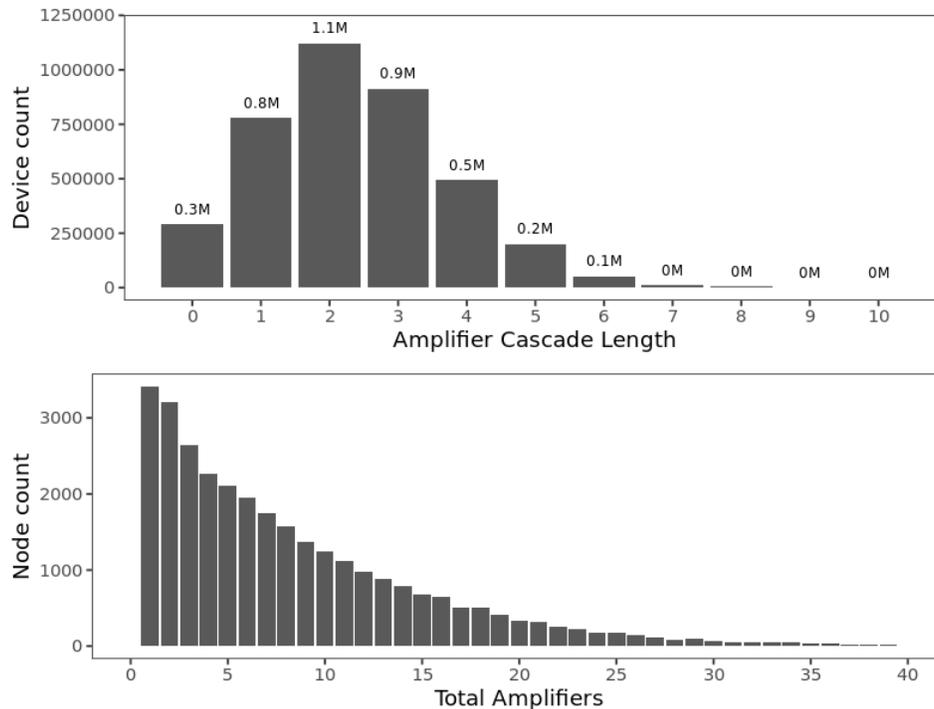


Figure 7 - Distribution of amplifier cascade length across devices (top panel) and distribution of total number of amplifiers across nodes (bottom panel). This data was generated by querying the graph database.

With the features derived from the graph database extracted, the next step involved identifying and examining telemetry data in relation to the graph features. The following variables were considered in the analysis:

- **Upstream Signal-to-Noise Ratio (US SNR):** 10th percentile of a device’s upstream SNR samples (i.e., data is collected from multiple time polls and aggregated)
- **Upstream Forward Error Correction (US FEC):** The percentage of time a device polled when transmitting upstream experiences an uncorrectable codeword error rate > 0
- **Upstream Partial Service:** The percentage of time a polled device goes into partial service with respect to an US channel
- **Upstream Power:** mean device power level
- **Downstream Receive Modulation Error Ratio (RxMER):** A device’s OFDM channel RxMER samples (data resolution by 50 KHz subcarrier)

The correlations between graph features and the variables above were explored and are shown in Table 1. One glaring outcome was the lack of a correlation between graph features and the FEC metric. The

absence of correlation with FEC is attributed to the mitigating effect of Comcast’s Profile Management Application (PMA). PMA was deployed for both DS D3.1 (OFDM) and US D3.0 since early 2020 [3,4]. Under PMA, channels that exhibit degraded spectrum get assigned a configuration that ensures that devices continue to use the spectrum without experiencing unacceptable levels of uncorrectable errors. While the algorithms are quite different between PMA for DS D3.1 vs. US D3.0, the result is the same: for those degraded channels, capacity is traded off for robustness. Hence, the lack of correlation is an assuring sign that PMA is doing its intended job.

Table 1 - Correlations between the node features and device telemetry.

	US SNR	US Rx Power	US FEC	US Partial Service	DS OFDM MER
Amplifier Cascade Length	-0.04	0.03	0.01	0.00	-0.12
Total Amplifiers	-0.05	0.04	0.02	0.01	-0.07

The strongest correlation exists between amplifier cascade length and OFDM RxMER. In Figure 8, two views supporting this trend are presented. The first shows different aggregating RxMER percentiles vs. amplifier cascade length. In all of these, there exists a trend of decreasing RxMER with increasing amplifier cascade length. Trend lines are included in the plot, which one can use to estimate an effect of ~2 dB for every 10 amplifiers added to the path (a more proper calculation, based on regression model, is presented below). The second view shows the distribution of modulation assigned to each subcarrier based on the standard D3.1 RxMER-to-modulation mapping for DS [5]. In this view, no aggregation is done as every subcarrier contributes to the distribution. Once again, there’s a clear effect that is most visible in the diminishing proportion of 4096-QAM as the amplifier cascade length increases.

A question that remains is to quantify the effect of this relationship on customer experience. Given that PMA is mitigating the impact of low RxMER by dynamically managing OFDM profiles, we restrict the definition of “customer experience” to the “capacity” dimension. Below, a model is introduced that quantifies the impact of the amplifier cascade length on available capacity. In the first step, a linear regression model is fitted to that RxMER data as presented in Figure 9. The linear model yields a statistically significant relationship (with a p -value $< 2 \times 10^{-16}$) between RxMER and the amplifier cascade length. The relationship is outlined in the following equation:

$$\text{RxMER} = 37.7 \text{ dB} - 0.37N_{\text{Amps}} , \quad (1)$$

in which N_{Amps} is the length of the cascade. In other words, every additional 10 amplifiers reduce RxMER, on average, by 3.7 dB.

An interesting feature of the overall distribution of RxMER is that the 38 dB line falls through the middle of the distribution (notice that the linear fit intercept is 37.7 dB). This level happens to be the PMA threshold for assigning a modulation of 4096-QAM, the highest possible under the current vendor implementation of D3.1. The threshold was deliberately set to be 3 dB lower than the recommended value under the D3.1 specification [5] (i.e., it is more aggressive). This means that for more than half the population (at and below this level), reduction in SNR due to increasing amplifier cascade length may

cause demotion to lower modulations. Indeed, this explanation agrees with the trend shown in the lower panel of Figure 8.

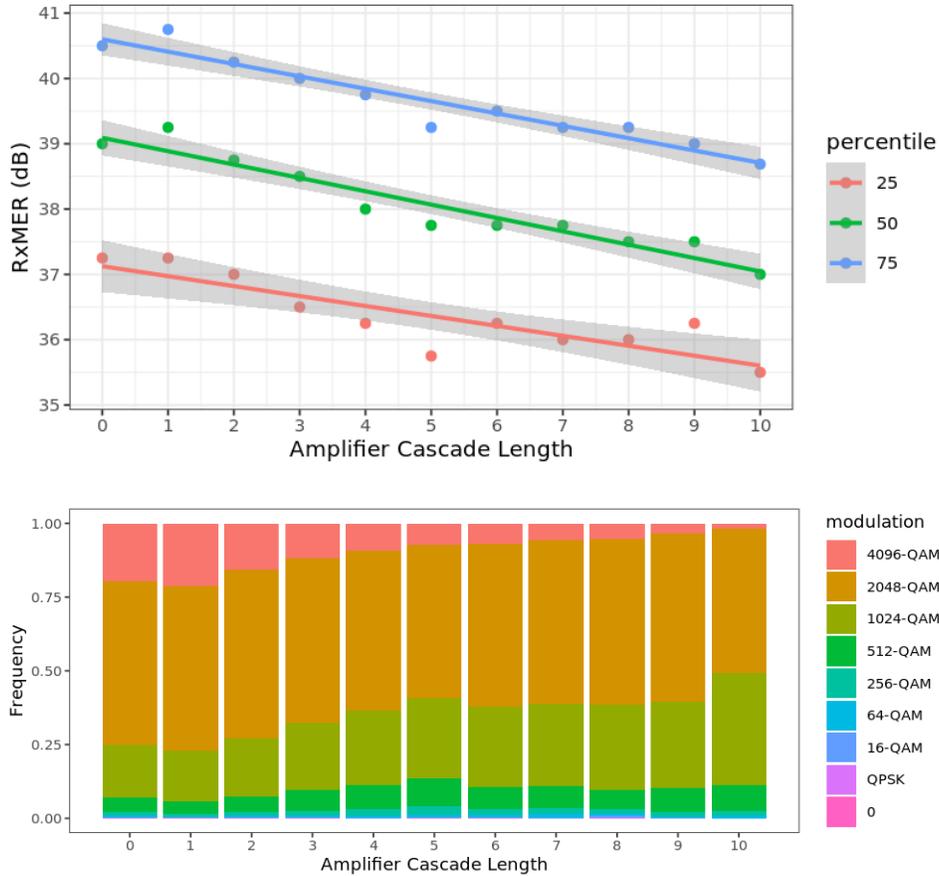


Figure 8 - The top panel shows the relationship between RxMER and amplifier cascade length for different aggregating percentiles (25th, 50th, and 75th). There’s a clear downward trend, as highlighted by the linear best-of-fit lines. The bottom panel shows the corresponding modulation distribution. Once again, the trend is visible by the decreasing ratio of 4096-QAM as the amplifier cascade length increases.

To translate the results from the linear model into a capacity impact figure-of-merit, we turn to the Shannon capacity theorem [7]. In a previous SCTE contribution [3,4], we argued that D3.1, with its low-density parity check (LDPC) error correction algorithm, operates close to the Shannon limit. The Shannon theorem can be approximated in the large SNR regime as follows:

$$C \approx 0.332 \cdot B \cdot \text{SNR}(\text{dB}), \quad (2)$$

where C is capacity, B is bandwidth, and SNR is measured in dB. The derivative of the equation above yields:

$$\Delta C \approx 0.332 \cdot B \cdot \Delta \text{SNR}(\text{dB}), \quad (3)$$

which provides a recipe for translating small changes in SNR to changes in capacity. Based on the above, for a standard 96 MHz wide OFDM channel, the results from the linear model translate into ~118 Mbps reduction in capacity for every 10 amplifiers traversed in the node-to-device DS path.

As demonstrated in this analysis, the effect is measurable and statistically significant. However, it is not impactful in the context of DS capacity for several reasons. First, the 10 amplifiers represent the upper bound of the distribution, i.e., this is an extreme scenario, as the bulk of the population falls below a cascade length of 5 amplifiers. Second, the spread in RxMER values within each “bucket” is much wider than the effect of increasing cascade length. This implies that there are other pressing issues one can solve before turning attention to cascade lengths. Third, even with loss of ~100 Mbps of capacity, D3.1-capable devices have access to sufficient D3.0 and D3.1 spectrum to support our highest speed tiers.

The same analysis was conducted for the US, exploring the relationship between US SNR and the total number of amplifiers in the node. Once again, the linear regression model shown in Figure 10 reveals a statistically significant relationship. However, the effect of increasing the number of amplifiers is even less impactful on customers compared to the DS, because of the large safety net that is intrinsic to a D3.0 US: while the highest possible modulation under D3.0 US is 64-QAM, the distribution of SNR lies ~10 dB above that level.

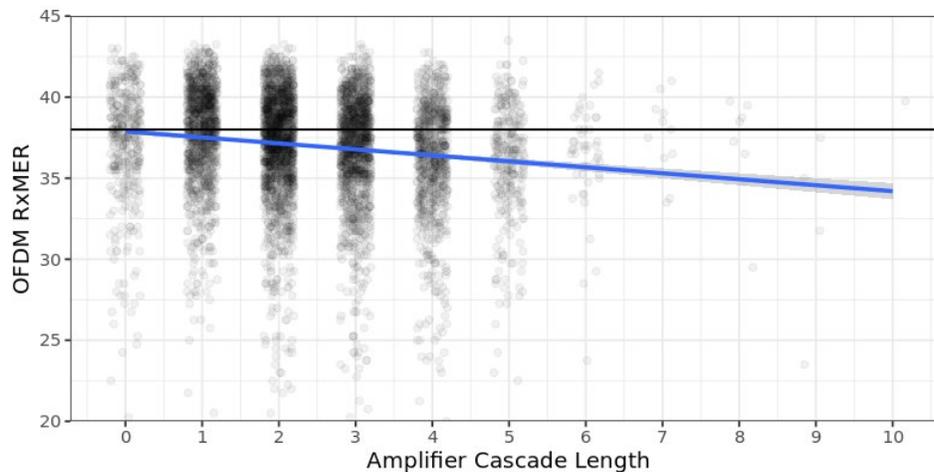


Figure 9 - Linear Fit to the RxMER vs. amplifier cascade length data (blue line with gray confidence band). The horizontal black line, at 38 dB, is the current threshold in PMA for assigning a 4096-QAM modulation.

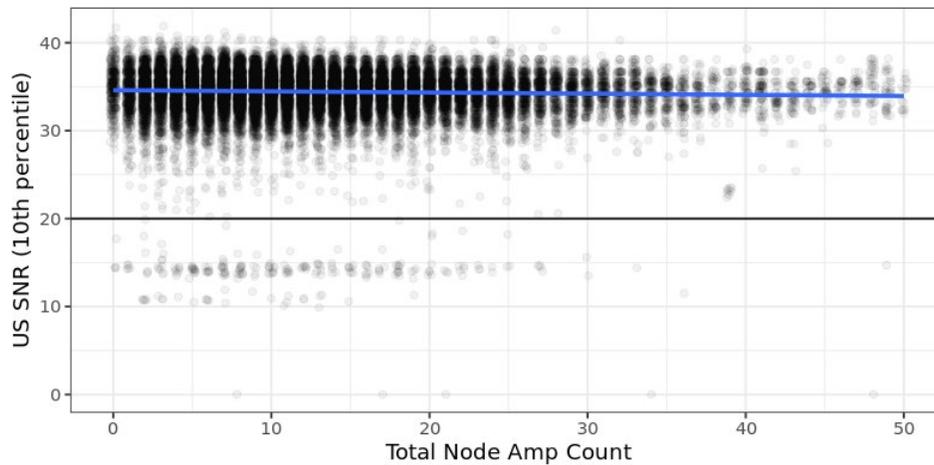


Figure 10 - Linear Fit to the US SNR vs. total number of amplifiers (blue line with gray confidence band). The horizontal black line at ~20 dB is the current threshold in PMA for assigning a 64-QAM modulation.

5. Conclusion

This paper highlights the power of the graph database in getting a better understanding of our network, and especially to track all physical and logical elements within the network. Using the relevant information, we were able to analyze and measure the impact of the amplifier cascade length on RF performance. The outcome was in the form of guidance on the effect of increased cascade length on capacity. While the impact on customers is assessed to be minimal today, future evolution of the network will bring service offerings that push the physical bandwidth to its limit (e.g., symmetrical Gbps service). These developments will make it critical to have a thorough understanding of the impact of network topology on customer experience—beyond the particular “amplifier cascade length” feature.

The graph database is being scaled up to cover our entire service footprint. In addition to the use case presented in this paper, there is much excitement about utilizing the graph database to build algorithms for root cause analysis and noise triangulation. These examples constitute a “holy grail” for us and the industry, as they promise to significantly cut down the labor-intensive troubleshooting processes involved in locating sources of noise/ingress. Given that machine learning techniques that utilize graph data representation exist today and are mature enough to support the intended use cases, the remaining bottleneck to conquer is to complete the construction and maintenance of the high-quality graph database presented in this paper.

Abbreviations

API	Application Programming Interface
CMTS	Cable Modem Termination System
CPE	Customer Premise Equipment
D3.0	DOCSIS 3.0
D3.1	DOCSIS 3.1
DOCSIS	Data Over Cable Service Interface Specification
DS	Downstream
MAC	Medium Access Control

NLP	Natural Language Processing
OFDM	Orthogonal Frequency Division Multiplexing
PMA	Profile Management Application
RF	Radio Frequency
RxMER	Receive Modulation Error Ratio
US	Upstream
SNR	Signal to Noise Ratio

Bibliography & References

1. *“Access Capacity Planning: Staying Well Ahead Of Customer Demand Helped Ensure Stability During COVID-19”*, B. Baker, C. Bou Abboud, E. Neeld. NCTA technical paper, 2020.
2. *“R for Data Science”*, H. Wickam, G. Grolemond, O-Reilly, 2017.
3. *“A Machine Learning Pipeline for D3.1 Profile Management”*, M. Harb, J. Ferreira, D. Rice, B. Santangelo, and R. Spanbauer, NCTA technical paper, 2019.
4. *“Full Scale Deployment of PMA”*, M. Harb, B. Santangelo, D. Rice, J. Ferreira, NCTA technical paper, 2020.
5. *“Data-Over-Cable Service Interface Specifications DOCSIS 3.1, PHY Layer Specificaiton, CM-SP-PHYv3.1-I18-210125”*, Cable Labs, <https://community.cablelabs.com/wiki/plugins/servlet/cablelabs/alfresco/download?id=f00df402-7367-4f86-a35c-5c22a2bfbaed>
6. *“Practical Lessons from D3.1 Deployments and a Profile Management Application (PMA)”*, NCTA technical paper, 2019.
7. *“A Mathematical Theory of Communication”*, C.E.Shannon, The BellSystem Technical Journal, vol27, pp. 379-423, 623-656 July, October 1948