# Adapting D4.0 Scheduler to Achieve the Optimal Latency

A Technical Paper prepared for SCTE by

**Hongbiao Zhang**
Architect Wireline Solutions
Casa Systems
100 Old River Rd
Andover, MA01886
978 688 6706
hongbiao.zhang@casa-systems.com

**Chain Lee**
Principal Software Engineer, R&D SW Cable
Casa Systems
100 Old River Rd
Andover, MA01886
978 688 6706
clee@casa-systems.com

**Vishnuvinod Sambasivan**
Principal SQA Engineer, R&D SW Cable
Casa Systems
100 Old River Rd
Andover, MA01886
978 688 6706
vishnuvinod.sambasivan@casa-systems.com

# Table of Contents

## List of Figures

## List of Tables

# 1. Introduction

The internet today has seen a blast of cloud-based services that has an increasing demand for fast access and short response time, a.k.a. low latency requirements. These services include online gaming, cloud computing, online trading and monetary transactions, video conferencing, VR/AR, etc. To cable networks, they appear as Over-the-Top (OTT) applications that are not managed by operators. Therefore, they cannot be serviced with higher priority like the MSO-offered voice or video. This implies that the cable systems need to provide low latency not only to MSO-managed QoS services, but also to all best effort services just the same.

In the past, DOCSIS cable is known for large latency as compared with competition technologies such as FTTP, especially the media access latency in the upstream direction. In D3.1, Low Latency DOCSIS (LLD) was proposed to alleviate the problem. This includes Proactive Grant Service (PGS) that potentially reduces request-grant cycle by 2/3. Despite the advantages, to this day the authors have not seen large deployment of this mechanism, possibly and understandably due to the fact that the current sub-split and mid-split plants still have limited upstream capacity, and thus favor bandwidth utilization over low latency.

The introduction of high-split and ultra-split DOCSIS pushes the equilibrium to a new level: with the advent of additional upstream bandwidth that is several times larger than the current, it is affordable to pre-allocate some grants prior to any requests so as achieve a smaller access latency. Similarly, it is also affordable to allocate some more grants on top of the requested bytes. In other words, it does not require CMTS to use accurate grant counts to match the requests.

On the other hand, ultra-split DOCSIS as defined in D4.0 requires Minimum Grant Bandwidth (MGB) to be enforced for any transmitting CM at any time of its transmission across the entire extended spectrum, whenever the extended spectrum exceeds 192 MHz. This requirement imposes new challenges to upstream scheduling, including that of PGS.

The paper assumes D4.0 Frequency Division Duplex (FDD) is adopted, yet most of the discussions and solutions apply to Full Duplex (FDX) DOCSIS as well. In this paper we show that both latency and data speed of user applications' might be significantly impacted due to the requirements of MGB. We further discuss various methods in D4.0 upstream scheduling that tackle the MGB issue, and how PGS could be realized using these methods. Thus, it is expected to achieve an optimal latency as well as the desired data speed in operators' D4.0 networks.

# 2. Proactive Granting

In the past, DOCSIS cable is known for large latency as compared with competition technologies such as FTTP, especially the media access latency in the upstream direction. The lengthy media access time is due to the request-and-grant based mechanism, which requires a request being sent by a CM and a subsequent grant being assigned by the CMTS before any upstream data can be transmitted by the CM. One of the reasons for adopting this mechanism since day one is possibly because early stage DOCSIS systems have a relatively low data rate in the upstream direction. Meanwhile there are a great number of users sharing the same upstream channel, with varying distances from the CMTS. These attributes suggest that upstream scheduling design would favor bandwidth utilization over latency, and the CM and CMTS have to use accurate byte counts in calculation.

In D3.1, Low Latency DOCSIS (LLD) was proposed to address the latency issues. Among a set of LLD features, Proactive Grant Service (PGS) was introduced and targeted at reducing the request-grant cycle

([4]). In a nutshell, PGS allows certain amount of traffic to take advantage of proactive grants without going through the request-grant cycle, thus reducing media acquisition time. This improvement in latency comes at a cost of bandwidth utilization, as the pre-allocated grants will be wasted if the targeted service flow has nothing to transmit at the designated grant time. However as we've shown in our experiments ([3]), we can achieve great mean latency and jitter even when the instantaneous traffic rate goes several times over the proactive grant rate. In other words, in PGS operations, latency is insensitive to the fluctuation of the offered traffic rate, as long as the traffic rate doesn't exceed the Maximum Sustained Rate (MSR) limit. Therefore, it is possible to pre-allocate only a small amount of bandwidth regardless of the actual traffic rate, yet achieving an ideal balance between low latency and bandwidth utilization.

It should be noted that the usage of PGS is not limited to special applications, instead, it could be used for all best effort traffic. Even with service flows configured with Best Effort (BE) scheduling type, it is allowed for CMTS to treat them as PGS services and allocate proactive grants without CMs noticing it.

## 3. MGB Requirements and Implications

In D4.0, Minimum Grant Bandwidth refers to the smallest grant that a CMTS is allowed to assign to a CM in the FDX or FDD upstream channels ([1], [2]). It is needed to ensure fidelity of a D4.0 CM operating in FDX or FDD spectrum. For example, for a total of 2 FDD channels, or 192 MHz in spectrum, the MGB is 10.4 MHz, or 26 minislots across the two channels. For a total of 3 FDD channels, or 288MHz in spectrum, the MGB becomes 16.0 MHz, or 40 minislots across the 3 channels. For a total of 4 FDD channels, or 384 MHz in spectrum, the MGB is 21.2 MHz, or 53 minislots across the 4 channels. The MGB has to be enforced at any time when the corresponding CM is transmitting. The numerical values are shown in the table below.
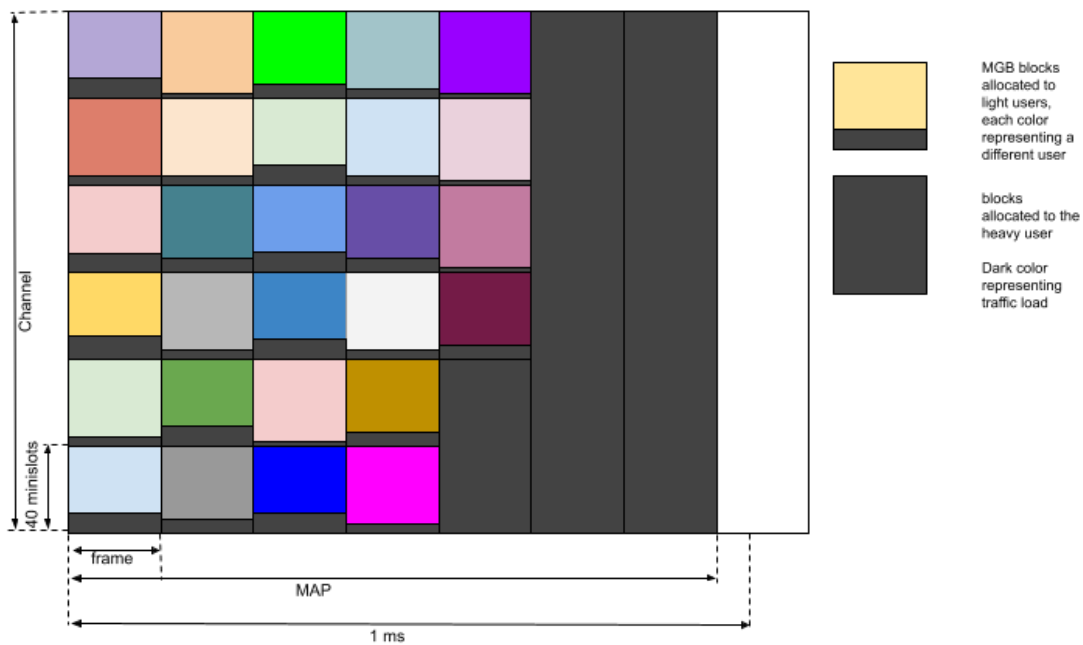
**Table 1 - FDD Minimum Grant Bandwidth**

| Split Name | FDD Upstream Spectrum (MHz) | Minimum Grant Bandwidth (MHz) | Equivalent Minislots |
|---|---|---|---|
| UHS-300 | 192 | 10.4 | 26 |
| UHS-396 | 288 | 16.0 | 40 |
| UHS-492 | 384 | 21.2 | 53 |

As compared with a non-FDD channel where a grant allocated for a transmitting CM could be as small as 1 minislot, enforcing a minimum of 26, 40, or 53 minislots for FDD channels makes the scheduling choices much more restrictive. We can see that with MGB imposed for each transmitting CM, there could only be a maximum of 18 transmitting CMs at any instance across the entire extended spectrum. With 2K FFT and a Cyclic Prefix (CP) of 2.5us, there could be between 2.5 - 7.4 OFDMA frames per millisecond depending on the frame size (K), which means there could be at most 44 - 133 CMs on average being scheduled in 1 millisecond time. This is barely enough to support proactive granting for just a couple of hundreds of users alone if the Guaranteed Grant Interval is chosen as 1ms.

We assume that a D4.0 ultra-split network offers more than enough upstream bandwidth to meet users' overall demands of traffic, had it not needed to fulfill MGB. Under this condition, we will demonstrate how enforcing MGB might impact the performance of user applications. These applications appear as OTT to the cable system and are treated with best effort and with the same priority. Therefore round robin is used among the service flows regardless of the offered rate. We will present a concrete example with numerical data in a later section, but here are some intuitions. First, we exam a scenario when there are

numerous SFs and traffic load is balanced among them, in which case the scheduler will ensure all of them get serviced in a timely fashion. Even if the total number of SFs exceed what can be scheduled in a scheduling cycle (due to MGB and smaller number of scheduling opportunities), the ones not serviced in the current scheduling cycle may be deferred to the next. As a result, their queues may accumulate quickly, but then get drained quickly as well when they get the scheduling opportunity (again due to MGB), in a round robin fashion. This may cause additional latency but only by a few scheduling cycles.

Now we exam a different scenario when traffic load among the various SFs is not balanced. For example, if there are a group of SFs with light but frequent traffic, and one or more additional SFs with heavy traffic. In that case the latter SFs may experience substantial traffic loss and unbounded latency, as the former SFs occupy much more bandwidth than necessary due to the effect of MGB, resulting in the latter being underserved. An example with MGB of 40 minislots is illustrated in the figure below.



**Figure 1 – Effect of MGB**

The second scenario may occur when a group of upstream users are transmitting light but frequent traffic, whereas one additional user starts transmitting at a peak rate. This additional user will suffer both traffic loss and large latency, including any delay sensitive traffic therein. Even if the corresponding SF is configured with a Minimum Reserved Traffic Rate and so that portion can be treated with a higher priority, it wouldn't help. This is because the delay sensitive traffic is not distinguishable from others within the same SF, therefore there is no guarantee that it can be served with a priority.

For simplicity of illustration, we assumed in the above context that each cable modem contains one upstream service flow only, therefore the limit of total CMs in a scheduling cycle becomes the total SFs permitted. In the next section we will discuss the case with multiple service flows per cable modem.

## 4. D4.0 Upstream Scheduling Enhancements

In this section, we will present enhancement methods in D4.0 upstream scheduling algorithms that target at mitigating the effect of MGB. A CMTS may adopt one or more of these methods per its design choices.

### 4.1. Filling up Extra Space

The D4.0 specification suggests that the "minimum grant bandwidth can be met through any combination of probe, ranging, OUDP testing SID, and data grant allocations across any of the Extended Upstream Channels in the CM's Extended Transmit Channel Set" ([1]). However probing signals may span the entire channel for only a subset of symbols. They can't be combined with other bursts to fulfil MGB. Besides, all of probing, ranging and OUDP testing are supposed to be allocated per CMTS's own schedules, and it serves no purpose if they're offered too often or at random. Instead, it makes sense to fill up the remaining space of a MGB block using data grant allocations. These extra grants help to reduce access latency for newly arrived packets, as the bandwidth would otherwise be wasted anyway. Per the specification, a CM is required to fill up the extra space of a grant with paddings ([1]).

### 4.2. Selecting OFDMA Frame Size

As mentioned above, there could only be a maximum of 18 simultaneous transmitting CMs at any time using the FDD channels. To accommodate more users, it is desirable to adopt a smaller OFDMA frame size, e.g., K = 6, which gives more granularity to grant allocations. With a 2K FFT, i.e., symbol time is 20us, with K = 6, and CP = 2.5 us, the resulted OFDMA frame time is 6 x (20 + 2.5) = 135 microseconds. Therefore, the maximum number of CMs that can be accommodated in a 1 millisecond scheduling cycle is 18 x 1000/135 = 133.2 CMs on average.

### 4.3. Scheduling Across Multiple Channels

Prior to D4.0, upstream SC-QAM channels and OFDMA channels are typically configured separately and operating independently. In D4.0 however, with the MGB requirements, it is desirable to concatenate all the FDD channels in order to make a more efficient usage of the extended spectrum. For example, if MGB is 40 minislots, a remaining 30 minislots on the first channel can be allocated to the next requesting CM, with an additional 10 or more minislots from a second channel being allocated to the same CM. In this way the MGB requirement is satisfied for the CM and the remaining 30 minislots on the first channel will not be wasted.

The implication is that all the FDD channels have to operate synchronously so that their symbols are aligned. Additionally, the scheduling cycle has to be identical for all of them, so that they can be scheduled as a unified resource pool. We require that these channels are configured with the same frame size K and the same Cyclic Prefix, and that the CMTS software forces the same OFDMA frame boundary across all the channels.

Note it is not required to synchronize between FDD channels and non-FDD channels.

### 4.4. Aggregating Service Flows

Prior to D4.0, each service flow is scheduled independently, with no regard to other SFs from the same CM. As MGB is enforced on a per CM basis, it makes sense to consider requests from all SFs of the same CM concurrently and allocate grants to multiple of them in the same frame, expecting the summation of all the grants satisfying MGB, or at least close to it. The aggregation results in a more efficient usage of channels' resources.

A question comes up as how we should handle the situation if multiple SFs of a CM have different scheduling priorities. When we aggregate these SFs, we may have ignored their priorities as compared with other SFs from different CMs. In other words, we may favor some lower priority SFs from one CM over a higher priority SF from another CM.

One possible method is to extend the concept of Aggregate Service Flow (ASF) and hierarchical QoS so that there could be e.g. 1 – 4 ASFs per cable modem, and there are multiple individual service flows in each ASF, with each SF having its own QoS parameters. For example, if there are 2 ASFs for each CM, ASF 1 may contain all SFs used for management, signaling protocols and voice, whereas ASF 2 may contain all other SFs. In each scheduling cycle, strict priority may be used between ASF 1 and ASF 2 regardless of the CMs. The total requests of an ASF may be scheduled on either FDD channels or non-FDD channels (see the next subsection). In the former case, MGB will be satisfied for grants allocated to the current ASF, with the remaining space possibly filled up by another ASF with a lower priority and from the same CM. Among the total allocated grants for the ASF, the CMTS will satisfy individual SFs one by one using round robin, weighted round robin, or strict priority, per QoS parameters of the ASF and individual SFs.

In a way, an ASF is similar to the concept of "T-CONT" in GPON, whereas an individual SF is similar to a GEM port.

As an extreme case, each SF is an ASF, meaning that each SF is scheduled independently regardless of the CMs, thus reverting to the original upstream scheduling scenario. As another extreme case, all SFs from the same CM belong to one ASF, meaning all CMs of the same ASF type have the same relative priority, regardless of individual SFs therein. By carefully grouping SFs into ASFs, we are able to achieve a balance between prioritization and aggregation.

## 4.5. Skipping Small Bursts

In the past, a scheduler may prefer channels in a higher spectrum over those in a lower spectrum, i.e., it will schedule grants on a higher spectrum channel until the current map space is filled, at which point it will schedule grants on a lower spectrum channel. In D4.0 however, to further mitigate the issue with MGB, we may deliberately leave small data grants to non-FDD channels, and only schedule large data grants on FDD channels. That is, when scheduling a FDD channel, if the total bytes of all aggregated requests from a CM or an ASF are way smaller than what is required to satisfy MGB, we may skip such requests in the current channel, so they will be picked up by non-FDD channels. Similarly, when scheduling a non-FDD channel, if the total bytes of all aggregated requests from a CM or an ASF are close to or exceed MGB, we may skip such requests in the current channel, so they will be picked up by FDD channels. The threshold setting of what size to pick up or skip may depend on relative utilization levels between the FDD channels and the non-FDD channels.

Note in the context above and thereafter, a "request" may refer to either an explicit or piggyback request from a CM, or an artificial request that is inserted by the CMTS itself to satisfy the Guaranteed Grant Rate of a PGS SF.

## 4.6. Allocating Proactive Grants

With the mechanism introduced in subsection 4.5, we can now adopt PGS to further reduce access latency, especially that of initial bursts. As mentioned in section 2, it is sufficient and necessary to pre-allocate a small amount of bandwidth for a SF, therefore a PGS grant typically contains a small number of bytes as compared with MGB. When scheduling a FDD channel, these PGS bytes will be aggregated with any other requests from the same or different SFs belonging to the same CM or ASF, if there are any. If the total bytes are large enough, they will be scheduled on the current channel and will conform to MGB. If however there're no other requests to aggregate, or if the aggregated total bytes are too small as compared with MGB, they will be skipped on the current channel and left to non-FDD channels.

## 5. Experiments and Analysis

In this section, we use a concrete example to illustrate how MGB may impact performance of an SF. We further demonstrate by adopting some of the methods introduced in section 4, it is possible to mitigate the effect of MGB and achieve an optimal performance.

As a full implementation of D4.0 is not ready yet for either CMTS or CM vendors, we perform our experiments with simulation on a high-split system. We also design an analytical model and validate it with simulation. Using that, we may derive further results for scenarios that cannot be experimented.

Here again for simplicity, we assume that each SF belongs to a different CM or ASF, so each SF will be allocated with grants of at least MGB. Also we assume the traffic demands from all users not exceeding the overall channels' capacity if we ignore the effect of MGB.

### 5.1. The Test Cases

Consider a Casa CMTS system configured with 2 OFDMA channels and one ATDMA channel on an upstream port, with their specifications as follows:

- Channel 1: ATDMA, 5.8 – 12.2 MHz, 64-QAM
- Channel 2: OFDMA, 70 – 101 MHz, K = 6, 1024-QAM
- Channel 3: OFDMA, 108 – 204 MHz, K = 6, 1024-QAM
- MAP interval: 1ms

Channel 3 is used to simulate an FDD channel in a UHS-300 plant, so that the grants allocated to any SF in any frame will span at least 26 minislots in the frequency domain. Consequently, there can be at most 9 SFs allocated in an OFDMA frame, and on average, at most 66 of them can be accommodated in one scheduling cycle.

A high-split CM is attached to the upstream port, with channel bonding that spans all 3 channels. A test SF is created on this CM, with QoS parameters as follows:

- Scheduling type: PGS
- Guaranteed Grant Rate: 5 mbps
- Guaranteed Grant Interval: 1ms

This test SF has an offered traffic rate of 700 mbps and a packet size of 1024 bytes from a Spirent test equipment. We measure the throughput and latency of it using the same equipment.

Consider a group of N synthetical SFs with a PGS scheduling type, a Guaranteed Grant Rate (GGR) of 1.5 mbps, and a Guaranteed Grant Interval (GGI) of 1ms. These SFs are simulated at the scheduler by inserting data grants in each scheduling cycle that can accommodate the specified parameters. That is, if using Channel 3 each of them will be assigned with 26 minislots in each scheduling cycle, and if using Channel 2 each of them will be assigned with 4 minislots. Among the N synthetic SFs, we denote M as the number of SFs that are granted on Channel 2.

In the first test case (denoted as test case A), we increase N from 0 to 100, and allocate grants for all of the synthetic SFs on Channel 3, i.e., M = 0. We observe how performance of the test SF degrades. Then in the next test case (denoted as test case B), we fix N and gradually move the synthetic SFs from Channel 3 to Channel 2, i.e., increase M from 0 to N, and observe how performance of the test SF is recovered.

### 5.2. Test A: Performance Degradation with MGB

In this test we increase N from 0 to 100 and allocate grants for all of the N synthetic SFs on Channel 3. As can be seen from the figure below, initially the test SF is the only one serviced by all channels, so it has 100% throughput and a minimum latency of 1.408 ms. When we increase N up to 20, some traffic of the test SF is pushed to Channel 1 and 2, but the total can still be satisfied by the remaining bandwidth, so the impact to its performance is negligible. When N continues to increase, traffic of the test SF can no longer be satisfied by the remaining bandwidth, and its performance degrades quickly. Buffer limit is reached and the latency increases as the throughput decreases. When N reaches 70 and above, traffic of the test SF is primarily serviced by Channel 1 and 2, therefore the performance metrics stay almost flat.

Note in reality, the lower channels are most likely occupied by legacy CMs, meaning the test SF will have to compete with them for resources, therefore its performance will be even worse.



**Figure 2 – Test A: Throughput and Latency vs N**

### 5.3. Test B: Performance Enhancements

In this test we fix N as 70, meaning that the system has to accommodate 70 synthetic SFs. Now increase M from 0 to 70, meaning we gradually move these SFs from Channel 3 to Channel 2. This corresponds to the method introduced in subsection 4.5.

As we can see from the figure below, performance of the test SF starts to improve as M increases. When M reaches 60, its performance has recovered almost completely, with a full throughput and a minimum latency of 1.586 ms.

Note in reality, the lower channels are most likely occupied by the legacy CMs, meaning the M SFs will have to compete resources with them. However this is much better than the case where light users occupy large chucks of resources on FDD channels, and push heavy users to the lower channels.

**Figure 3 – Test B: Throughput and Latency vs M with N = 70**

### 5.4. Further Analytical Results

We've designed an analytical model to calculate the performance metrics of various scenarios. Before using that, let's first validate the model with simulation results measured in Test A and Test B. Figure 4 below depicts a comparison of theoretical and simulation results for both throughput and latency of the test SF in Test A, and Figure 5 depicts the same in Test B. As we can see from the figures, the theoretical model matches well with the simulation.



**Figure 4 – Test A Theoretical vs Simulation Results**

**Figure 5 – Test B Theoretical vs Simulation Results**

We can then use the analytical model to derive performance metrics under various conditions, as follows.

### 5.4.1. UHS-300

The system under study now is equipped with the following:

- Channel 1: Not present
- Channel 2: OFDMA, 8 – 101 MHz, K = 6, 1024-QAM
- Channel 3: OFDMA, 108 – 204 MHz, K = 6, 1024-QAM
- Channel 4: OFDMA, 204 – 300 MHz, K = 6, 1024-QAM

The system contains 2 FDD channels and an upstream upper edge of 300 MHz, i.e., UHS-300. N represents all light but frequent users whose traffic can be accommodated with PGS-like service with about 1.5 Mbps grant rate and 1ms grant interval. M represents how many of the N users are served using the non-FDD channel. Under this condition we evaluate the total bandwidth available to a heavy user under test. The table below enumerates the throughput of it with different N and M combinations.

**Table 2 – Throughput vs N and M (UHS-300)**

| N / M | 0 | 20 | 40 | 60 | 80 | 100 | 120 | 140 | 160 | 180 | 200 | 220 | 240 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1933.9 | 1743.0 | 1552.0 | 1361.1 | 1170.1 | 979.2 | 788.3 | 637.6 | 636.4 | 635.5 | 634.8 | 634.2 | 633.7 |
| 20 | 0.0 | 1904.5 | 1713.6 | 1522.7 | 1331.7 | 1140.8 | 949.8 | 758.9 | 608.2 | 607.1 | 606.2 | 605.4 | 604.9 |
| 40 | 0.0 | 0.0 | 1875.2 | 1684.2 | 1493.3 | 1302.3 | 1111.4 | 920.4 | 729.5 | 578.8 | 577.7 | 576.8 | 576.1 |
| 60 | 0.0 | 0.0 | 0.0 | 1845.8 | 1654.8 | 1463.9 | 1273.0 | 1082.0 | 891.1 | 700.1 | 549.5 | 548.3 | 547.4 |
| 80 | 0.0 | 0.0 | 0.0 | 0.0 | 1816.4 | 1625.5 | 1434.5 | 1243.6 | 1052.6 | 861.7 | 670.8 | 520.1 | 518.9 |
| 100 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1787.0 | 1596.1 | 1405.2 | 1214.2 | 1023.3 | 832.3 | 641.4 | 490.7 |
| 120 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1757.7 | 1566.7 | 1375.8 | 1184.8 | 993.9 | 802.9 | 612.0 |
| 140 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1728.3 | 1537.3 | 1346.4 | 1155.5 | 964.5 | 773.6 |
| 160 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1698.9 | 1508.0 | 1317.0 | 1126.1 | 935.1 |
| 180 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1669.5 | 1478.6 | 1287.6 | 1096.7 |
| 200 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1640.2 | 1449.2 | 1258.3 |
| 220 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1610.8 | 1419.8 |
| 240 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1581.4 |

The MGB is 26 minislots. As shown in the table, the heavy user may achieve a peak rate of ~1.9 Gbps, but the rate quickly drops to ~630 Mbps when the number of light users N reaches 133. The lowest rate of ~630 Mbps is what the lower channel has to offer. When the light users are moved to the lower channel, the data rate is gradually recovered, up till ~1.58 Gbps.

### 5.4.2. UHS-396

Here we consider a system with 3 FDD channels and an upstream upper edge of 396 MHz, i.e., UHS-396, as follows:

- Channel 1: Not present
- Channel 2: OFDMA, 8 – 101 MHz, K = 6, 1024-QAM
- Channel 3: OFDMA, 108 – 204 MHz, K = 6, 1024-QAM
- Channel 4: OFDMA, 204 – 300 MHz, K = 6, 1024-QAM
- Channel 5: OFDMA, 300 – 396 MHz, K = 6, 1024-QAM

The MGB becomes 40 minislots in this case. The throughput with any combination of N and M is enumerated in the following table. Explanation is omitted here.

**Table 3 – Throughput vs N and M (UHS-396)**

| N / M | 0 | 20 | 40 | 60 | 80 | 100 | 120 | 140 | 160 | 180 | 200 | 220 | 240 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2586.72 | 2292.96 | 1999.2 | 1705.44 | 1411.68 | 1117.92 | 824.16 | 642.2094 | 640.484 | 639.1399 | 638.0633 | 637.1815 | 636.4461 |
| 20 | 0 | 2557.344 | 2263.584 | 1969.824 | 1676.064 | 1382.304 | 1088.544 | 794.784 | 612.8334 | 611.108 | 609.7639 | 608.6873 | 607.8055 |
| 40 | 0 | 0 | 2527.968 | 2234.208 | 1940.448 | 1646.688 | 1352.928 | 1059.168 | 765.408 | 583.4574 | 581.732 | 580.3879 | 579.3113 |
| 60 | 0 | 0 | 0 | 2498.592 | 2204.832 | 1911.072 | 1617.312 | 1323.552 | 1029.792 | 736.032 | 554.0814 | 552.356 | 551.0119 |
| 80 | 0 | 0 | 0 | 0 | 2469.216 | 2175.456 | 1881.696 | 1587.936 | 1294.176 | 1000.416 | 706.656 | 524.7054 | 522.98 |
| 100 | 0 | 0 | 0 | 0 | 0 | 2439.84 | 2146.08 | 1852.32 | 1558.56 | 1264.8 | 971.04 | 677.28 | 495.3294 |
| 120 | 0 | 0 | 0 | 0 | 0 | 0 | 2410.464 | 2116.704 | 1822.944 | 1529.184 | 1235.424 | 941.664 | 647.904 |
| 140 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2381.088 | 2087.328 | 1793.568 | 1499.808 | 1206.048 | 912.288 |
| 160 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2351.712 | 2057.952 | 1764.192 | 1470.432 | 1176.672 |
| 180 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2322.336 | 2028.576 | 1734.816 | 1441.056 |
| 200 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2292.96 | 1999.2 | 1705.44 |
| 220 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2263.584 | 1969.824 |
| 240 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2234.208 |

### 5.4.3. UHS-492

We now consider a system with 4 FDD channels and an upstream upper edge of 492 MHz, i.e., UHS-492, as follows:

- Channel 1: Not present
- Channel 2: OFDMA, 8 – 101 MHz, K = 6, 1024-QAM
- Channel 3: OFDMA, 108 – 204 MHz, K = 6, 1024-QAM
- Channel 4: OFDMA, 204 – 300 MHz, K = 6, 1024-QAM
- Channel 5: OFDMA, 300 – 396 MHz, K = 6, 1024-QAM
- Channel 6: OFDMA, 396 – 492 MHz, K = 6, 1024-QAM

The MGB is 53 minislots in this case. The throughput with different combinations of N and M is listed as follows. Explanation is omitted.

**Table 4 – Throughput vs N and M (UHS-492)**

| N / M | 0 | 20 | 40 | 60 | 80 | 100 | 120 | 140 | 160 | 180 | 200 | 220 | 240 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3239.52 | 2850.288 | 2461.056 | 2071.824 | 1682.592 | 1293.36 | 904.128 | 646.8391 | 644.5386 | 642.7465 | 641.311 | 640.1354 | 639.1549 |
| 20 | 0 | 3210.144 | 2820.912 | 2431.68 | 2042.448 | 1653.216 | 1263.984 | 874.752 | 617.4631 | 615.1626 | 613.3705 | 611.935 | 610.7594 |
| 40 | 0 | 0 | 3180.768 | 2791.536 | 2402.304 | 2013.072 | 1623.84 | 1234.608 | 845.376 | 588.0871 | 585.7866 | 583.9945 | 582.559 |
| 60 | 0 | 0 | 0 | 3151.392 | 2762.16 | 2372.928 | 1983.696 | 1594.464 | 1205.232 | 816 | 558.7111 | 556.4106 | 554.6185 |
| 80 | 0 | 0 | 0 | 0 | 3122.016 | 2732.784 | 2343.552 | 1954.32 | 1565.088 | 1175.856 | 786.624 | 529.3351 | 527.0346 |
| 100 | 0 | 0 | 0 | 0 | 0 | 3092.64 | 2703.408 | 2314.176 | 1924.944 | 1535.712 | 1146.48 | 757.248 | 499.9591 |
| 120 | 0 | 0 | 0 | 0 | 0 | 0 | 3063.264 | 2674.032 | 2284.8 | 1895.568 | 1506.336 | 1117.104 | 727.872 |
| 140 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3033.888 | 2644.656 | 2255.424 | 1866.192 | 1476.96 | 1087.728 |
| 160 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3004.512 | 2615.28 | 2226.048 | 1836.816 | 1447.584 |
| 180 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2975.136 | 2585.904 | 2196.672 | 1807.44 |
| 200 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2945.76 | 2556.528 | 2167.296 |
| 220 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2916.384 | 2527.152 |
| 240 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2887.008 |

## 6. Conclusion

Proactive granting is an important tool that significantly reduces access latency for any best effort traffic, especially that of initial bursts. However, proactive granting in D4.0 may result in inefficient channel utilization on FDD channels due to the MGB restrictions. Our study demonstrates that by adopting various enhancement methods with upstream scheduling, it is possible to reduce the effect of MGB to the minimum, thus achieving the best utilization of channels' bandwidth. It is therefore possible to obtain an optimal latency for all traffic.

Our study also brings forth new questions, as there are other complexities in D4.0 that might require operators to reconsider their service offerings and SLAs. For example, the choice between channel utilization and prioritization becomes non-trial. Besides, an accurate realization of weighted round robin scheduling based on numerical values e.g., MSR may become difficult. These topics are left for future studies.

# Abbreviations

| | |
|---|---|
| ASF | Aggregate Service Flow |
| ATDMA | Advanced Time Division Time Access |
| BE | Best Effort |
| CM | Cable Modem |
| CMTS | Cable Modem Termination System |
| CP | Cyclic Prefix |
| FDD | Frequency Division Duplex |
| FDX | Full Duplex |
| GGI | Guaranteed Grant Interval |
| GGR | Guaranteed Grant Rate |
| LLD | Low Latency DOCSIS |
| MGB | Minimum Grant Bandwidth |
| MSR | Maximum Sustained traffic Rate |
| OFDMA | Orthogonal Frequency Division Multiple Access |
| OTT | Over-the-Top |
| PGS | Proactive Grant Service |
| SF | Service Flow |
| SLA | Service Level Agreement |
| UHS | Ultra-high Split |
| WRR | Weighted Round Robin |

# Bibliography & References

1. *Data-Over-Cable Service Interface Specifications DOCSIS® 4.0, MAC and Upper layer Protocols Interface Specification*, CM-SP-MULPIv4.0-I07-230503
2. *Data-Over-Cable Service Interface Specifications DOCSIS® 4.0, Physical layer Specification,* CM-SP-PHYv4.0-I06-221019
3. *H. Zhang P. Wolff, V. Sambasivan, Obtaining Low Latency in Upstream DOCSIS Transmissions,* SCTE Technical Journal, Vol 2, Num 1, Mar 2022.
4. *Data-Over-Cable Service Interface Specifications DOCSIS® 3.1, MAC and Upper layer Protocols Interface Specification*, CM-SP-MULPIv3.1-I21-201020