

Multiple Cyclic Queuing and Forwarding

new-finn-multiple-CQF-0921-v02

Norman Finn

Huawei Technologies Co. Ltd

October 26, 2021

Abstract

Multi-CQF, an expansion of Cyclic Queuing and Forwarding (CQF, IEEE Std 802.1Q-2018 Annex T) is presented. Multi-CQF requires that all nodes run at the same frequency, but does not require that all output ports' cycles be in phase. By using three buffers or more, multi-CQF supports long links, as well as intentional fixed delays for short/long path recombination. Running multiple instances of multi-CQF on one port at different cycle rates gives good latency and bandwidth utilization for a mix of streams with varied bandwidth requirements. This expansion requires little alteration of IEEE 802.1Q, and like the original 2-buffer synchronized CQF, requires no per-hop per-stream dynamic state. Multi-CQF is presented for inclusion into an amendment to IEEE Std 802.Q. The determinism of multi-CQF is demonstrated, and a set of parameters is suggested for network management purposes. The paternoster algorithm of Mick Seaman is offered as a useful addition to the multi-CQF algorithm.

Table of contents

Abstract	1
1 Introduction.....	4
1.1 Deterministic Quality of Service	4
1.2 Continuous streams.....	4
1.3 Implementation requirements	5
1.3.1 Store-and-forward	5
1.3.2 Frequency synchronization.....	5
1.3.3 Continuous transmission	5
1.4 Multi-CQF vs. CQF.....	5
2 CQF timing model.....	6
2.1 Output timeline 1	8
2.2 Receive timeline 2	10
2.3 Storing frames timeline 3	10
2.4 Calculation of T_B	11
2.5 Transmitting frames timeline 4	12
2.6 More than 3 output buffers.....	12
3 Multiple CQF classes of service	12
3.1 Multiple T_C model.....	12
3.2 Preemption and interference	15
3.3 T_C computation	15
3.4 Why integer multiples for T_C ?.....	16
4 Deterministic behavior of Multi-CQF.....	16
4.1 Basic requirement for determinism	17
4.2 Admission control for multiple T_C values.....	18
4.3 Computing the actual end-to-end latency for Multi-CQF.....	19
5 Parameterization of Multi-CQF	19
5.1 Per-stream and per-port static state	19
5.2 Cycle wander	20
5.3 Link delay variation.....	21
5.4 Calculating the number of buffers required.....	21
5.5 Initial buffer phase	22

5.6	Externally-visible Multi-CQF managed objects and protocol items.....	23
5.6.1	Cycle and priority structure managed objects.....	23
5.6.2	Cycle phase managed objects.....	23
5.6.3	Cycle variation information	24
5.6.4	Dead time / bandwidth balance information	24
6	Other issues.....	25
6.1	Frame size problem	25
6.2	Tailored bandwidth offerings	26
6.3	Overprovisioning is not always bad.....	26
6.4	CQF and credit-based shaper	27
6.5	Fundamental CQF pros and cons.....	27
7	Paternoster and multi-CQF.....	27
7.1	Ingress conditioning	27
7.2	Changing cycle times	28
7.3	Aggregation and dis-aggregation.....	29

1 Introduction

NOTE: This document is a revision of [1-21-0059-00-ICne-multiple-cyclic-queuing-and-forwarding](#).

The remainder of section 1 defines the domain of interest of this paper. Section 2 provides a detailed timing model for Multiple Cyclic Queuing and Forwarding (based on CQF, IEEE Std 802.1Q-2018 Annex T). It shows how two, three, or more buffers can be used to manage the allocable bandwidth and per-hop delay, and why the systems' CQF cycles do not have to operate in phase. Section 3 shows how multiple instances of CQF can be run on the same output port, with different cycle times, in order to efficiently serve streams with a wide range of bandwidth and latency requirements. Section 4 analyzes the determinism of CQF, and Section 5 suggests a set of control parameters. Section 6 addresses miscellaneous issues related to multi-CQF. Section 7, discusses how Mick Seaman's paternoster algorithm can greatly improve the usefulness of a multi-CQF network.

This paper assumes the reader is reasonably familiar with CQF as described in IEEE Std 802.1Q Annex T. The frame timestamps used in this paper are described in IEEE Std 802.3-2018 clause 90. Preemption, or interspersed express traffic, is described in IEEE Std 802.3-2018 clause 99 and IEEE Std 802.1Q-2018 clause 6.7.2.

1.1 Deterministic Quality of Service

TSN (and IETF DetNet) supply a Quality of Service (QoS) to a critical data flow, or "stream". This QoS is:

- a. An absolute upper bound on the end-to-end latency to frames belonging to the stream. (Bounded Latency); and
- b. A guarantee that no frames of the stream will be discarded due to a buffer being full when the frame arrives at an intermediate hop (Zero Congestion Loss).

This QoS, which we will call the Deterministic QoS, is made possible by a promise, made by the source of a stream, to not exceed a contracted bandwidth and maximum frame size. This guarantee allows the network to run a resource reservation procedure that dedicates resources to a particular stream (or sometimes, to a class of similar streams) at every hop through the network, before the first frame of the stream can be transmitted.

1.2 Continuous streams

We can divide the streams that can make use of the Deterministic QoS into two classes:

- Continuous streams can be usefully characterized by a maximum frame size, and a maximum bandwidth.

- Scheduled streams transmit on a regular, repeating schedule.

Note that these two categories do not encompass all possible data flows. Bursty, irregular flows are (by definition) not streams, in the sense that it is difficult, in the presence of multiple of these flows, to guarantee Deterministic QoS except by overprovisioning and an extensive analysis of worst-case inter-stream interference scenarios.

Scheduled streams can be handled by using IEEE Std 802.1Qbv (now IEEE Std 802.1Q-2018 clause 8.6.8.4, Enhancements for Scheduled Traffic) to schedule traffic class transmission in detail. They are of no interest to this paper. This paper is concerned only with continuous streams.

1.3 Implementation requirements

1.3.1 Store-and-forward

By definition, CQF is a store-and-forward technology; no buffer entry has frames both being stored and being transmitted at the same time. So cut-through forwarding is inapplicable to CQF. That does not mean that CQF and cut-through forwarding cannot be used on the same bridge, or even on the same port—only that we do not address the issue, here.

1.3.2 Frequency synchronization

The IEEE Std 802.1Q Annex T CQF assumes that every bridge has a system clock that is synchronized with the other bridges' system clocks in the network, so that configuring a transmission time in one bridge has meaning in another bridge. Multi-CQF does not require synchronization of the system clocks, but does require frequency lock. That is, there is a maximum difference in the elapsed time between two events as measured by two bridges' system clocks, no matter the length of that elapsed time.

1.3.3 Continuous transmission

Delivering the deterministic QoS using multi-CQF depends upon a transmitting port being able to select the correct frame to transmit according to strict priority among the CQF priority levels, and initiate all transmissions in that order, without introducing extra inter-frame gap time. Since, with CQF, no buffer has frames both arriving and being transmitting at the same instant, this should pose no insurmountable problems for implementors.

1.4 Multi-CQF vs. CQF

The differences between CQF, as defined in Annex T of IEEE Std 802.1Q-2018, and multi-CQF, as defined here, are:

- a) For a given value of the cycle time T_C on a given output port, Annex T provides two buffers, each implemented as a separate class of service queue. Scheduled output gates (8.6.8.4 of IEEE Std 802.1Q-2018) are configured to enable the two queues to output alternately, each given time T_C to drain.

Multi-CQF allows two or more buffers per output port per cycle time. It is not practical to dedicate one class of service queue for each of these buffers; there are only eight available. We assume here that alterations to IEEE Std 802.1Q will be made to enable one class of service queue to provide any number of CQF buffers for a single cycle time (see [new-finn-pulsed-queuing-0821-v03](#)).

- b) Annex T assumes that one value of T_C is sufficient for any given output port.

Multi-CQF allows multiple values of T_C , one for each Multi-CQF class of service. The output cycles are constrained, on any given port, so that an integral, and never a fractional, number of shorter cycles are contained within any given longer cycle. We will assume that one class of service queue serves a single value of T_C .

- c) Annex T uses synchronized time so that every output port in a TSN network switches buffers simultaneously.

Multi-CQF allows each bridge to perform its buffer switching at different times, subject to the above constraint b).

- d) Annex T uses the same cycle time and phasing for the input gates as for the output gates. The input gates select which of the two output buffers on a given port stores a received frame.

Multi-CQF runs the same cycle time T_C for input and output gates, but adjusts the phase of the input gates on a port to match the phase of the frames arriving from the output gates of the bridge transmitting to that receiving port.

2 CQF timing model

We have two nodes, A and B. Both are running an instance of Cyclic Queuing and Forwarding on each of multiple ports, more-or-less as described in IEEE Std 802.1Q-2018 Annex T. We will assume that nodes A and B are frequency locked (see 1.4). We do not assume that the output buffers switch in synchrony; they can be out of phase.

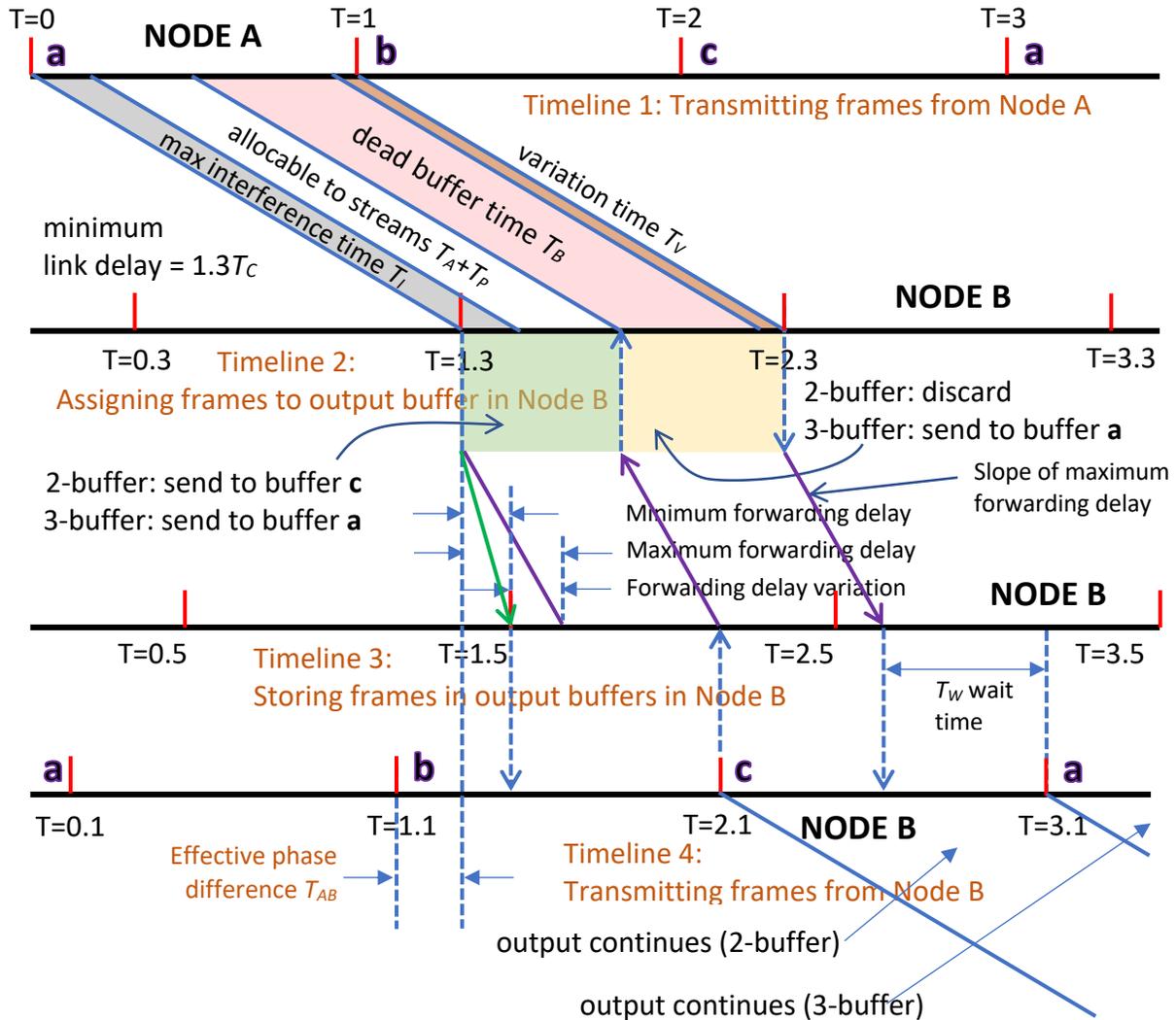
After a gate open/close event on a particular port, node A transmits all of the frames in one cyclic buffer towards receiving node B, not necessarily in a single burst. After some gap following the transmission of the last frame in the buffer, another gate open/close event is performed. At this point, it starts transmitting the frames from the next cyclic buffer. The gate open/close events in both nodes happen regularly, with the same period T_C . At the next hop, node B must be able to assign each received frame to a transmit buffer such that 1) frames that were in the same buffer in node A, and are transmitted on the same port from node B, are placed into the same buffer in node B; and 2) frames in different buffers in node A are placed in different buffers in node B.

Figure 1 shows an example of Cyclic Queuing and Forwarding. Node A and Node B are transmitting at the same frequency, but are offset by $0.1T_C$, as shown by timelines 1 and 4. In Figure 1, we use the following notation for time intervals:

T_C	nominal (intended) period of the buffer-swapping cycle
T_I	maximum interference from lower-priority queues, one frame or one fragment
T_V	sum of the variation in output delay, link delay, clock accuracy, and timestamp accuracy
T_A	the part of the cycle allocable to (reservable by) streams
T_P	worst-case time taken by additional bytes if this traffic class is preemptable
T_B	end-of-cycle buffer dead time optionally imposed on node A by node B
T_W	wait time during which buffer is neither receiving nor transmitting frames
T_{AB}	effective phase difference between cycle start times for input from A and output from B

Following the definitions of output gates in IEEE Std 802.1Q-2018, the red ticks in timelines 1 and 4 in Figure 1 represent the earliest possible moment at which the first bit of the destination address of the first frame of the cycle can be transmitted. These ticks are ultimately driven by the system clock. They are the basis for all cyclic buffer transmissions. If Enhancements for Scheduled Traffic (ETS, IEEE Std 802.1Q-2018 8.6.8.4) are used for controlling the output buffers, the ticks are the points in time when the output gate of one queue is closed, and the next queue's gate is opened. These are the points in time as programmed into the managed objects that control ETS. An implementation may need to schedule events in anticipation of the time specified in the managed objects in order to maximize throughput. Note that the preamble of an IEEE 802.3 Ethernet frame can be transmitted before gate open event.

Figure 1 Reference timelines for time-based CQF



2.1 Output timeline 1

Figure 1 shows an interference delay T_i (the gray area) between the gate event (the red ticks in Figure 1) and the transmission of the first bit of the first stream frame's destination MAC address. The interference is from frames transmitted from lower-priority queues. It is equal to the time required for one maximum-length transmission unit over all lower-priority queues. That maximum transmission unit is either a maximum-length fragment, for preemptable lower-priority queues, or the maximum-length frame, for non-preemptable queues. The value of T_i depends upon the configuration of lower-priority queues.

It is possible that the class of service illustrated in Figure 1 is, itself, a preemptable class. In that case, a higher-priority class of service can preempt transmission of frames in this class.

Preempting a frame adds additional bytes to the resultant fragments, which must be accounted for when allocating bandwidth to a class of service. T_P represents the worst-case additional time required to transmit these extra bytes caused by preempting frames belonging to a CQF stream. This value is always bounded. See below, section 3.2.

There can be some variation in the time from the selection of a frame for output in node A to the timestamp moment, when the first bit of the destination MAC address is transmitted (see IEEE Std 802.3-2018 clause 90). This is called output delay variation. The total time between the transmission of the first bit of the frame and the reception of that first bit at the next hop is called the link delay. Depending on the medium and the length of the link, there can be variations in link delay. The worst-case variation between the two node's clocks caused by accumulated frequency variations, asymmetrical links, etc., causes uncertainty between the transmitting and receiving nodes' clocks, and in the determination of the link delay. The inaccuracy in converting between IEEE 802.3 transmit and receive timestamps and the local clock that drives the gate open/close events also contributes to cycle accuracy. The worst-case combination of these four items, output delay variation, link delay variation, clock/frequency uncertainty, and timestamp conversion inaccuracies, is labeled, T_V .

All of the contributions to T_V are lumped together at the end of the cycle, even though contributions to T_V are made throughout the cycle.

As described below (section 2.4) the next hop can impose a buffer dead time T_B on this hop. This is a time at the end of the cycle, during which no frames can be transmitted from the cyclic buffer, so that the last frame of the cycle can be received earlier than the end of the cycle.

It is necessary, in order to know how much data can be transmitted in one cycle, that an implementation be able to transmit all of the frames in a cyclic output buffer together, at line rate, with no interference from lower-priority queues on the same output port. (Interference from higher-priority queues is described in section 3.2.) Given that is true, then the total time per cycle that can be used for transmitting streams is:

$$T_A = T_C - T_I - T_P - T_B - T_V.$$

This T_A is a maximum, local to a particular output port on a node. It guarantees that the last frame of cycle (plus a possible preamble of the first frame of the next cycle) will be on the wire before the output gate closes. All of the components of T_A can be calculated by an implementation from its configuration and from knowledge of the implementation, except for T_B and parts of T_V . T_B is supplied by configuration, or by the node to which the output port is connected. T_V can be supplied either by the time sync implementation, by configuration, by summing the contributions of node A and node B, or by the specification of a maximum allowed value by a standard or an equipment purchaser.

Note that T_A , as defined here, includes the entire transmission time of stream data, including one 12-byte inter-frame gap and one 8-byte preamble for every frame. The preamble of the

first frame of a cycle is counted in the previous cycle due to the way in which the output gates are defined in IEEE Std 802.1Q.

2.2 Receive timeline 2

The timeline at the receiving port is timeline 2 in Figure 1. The red ticks represent the earliest possible moment that the first bit of the destination MAC address of the first frame of a cycle can be received. In terms of IEEE 802.1Qci (IEEE Std 802.1Q-2018 clause 8.6.5.1 Per-stream Filtering and Policing), a timed input gate must open no later than this point.

On timeline 2, each frame is assigned to a buffer on an output port based on the timestamp (IEEE Std 802.3-2018 clause 90) on the frame.

A critical aspect of timeline 2 is its offset from timeline 4, the output timeline. This offset is shown as T_{AB} in Figure 1. It is clear from the figure that T_{AB} must be known in order to compute T_B and T_W . T_{AB} can be computed by 1) synchronizing the clocks of nodes A and B, and 2) measuring the link delay from node A to node B using PTP. Other methods are also possible, e.g. that described in [1-21-0056-00-ICne-input-synchronization-for-cyclic-queueing-and-forwarding](#).

Once T_{AB} is known, all of the timing relationships shown in Figure 1 can be computed. The phasing of the nodes' output buffer cycles certainly does affect the end-to-end latency of any stream, so that phasing must be known when the latency is computed. The end-to-end latency is no longer an integer multiple of the cycle time. It is even possible to adjust the phasing to favor certain paths through the network.

For a node B that is connected to and receiving cyclic frames from n other nodes, we have n assignment problems to solve, one for each input port on node B.

If a frame (belonging to a stream) is received that straddles a cycle (first bit in one cycle on timeline 2 of Figure 1, and end frame plus inter-frame gap plus a preamble time occurs in the next cycle), then either 1) some part of that frame was transmitted from node A outside the cycle window T_C , or 2) one or more of the constants, measurements, or calculations above is incorrect. Either way, the frame must be discarded or marked down to best-effort service, or else it can cause disruption of delivery guarantees farther along in the network.

2.3 Storing frames timeline 3

The timeline at the point where frames are stored into an output buffer is timeline 3 in Figure 1. The red ticks on timeline 3 mark the earliest point at which the first frame transmitted from a particular buffer could reach the output buffers (neglecting transmission time on the input medium). These ticks are offset from timeline 2 by the minimum forwarding delay, required to forward the frame from the input port to the output queue. The maximum forwarding delay is

also shown. The forwarding delays shown in Figure 1 include the time to install the frame in the output buffer and for its presence to filter through to the point that it can be selected for output.

There are two possible buffer assignment methods shown in Figure 1: the two-buffer method, in which the frames received from node **A** buffer **a** are assigned to buffer **c** in node **B**, and the three-buffer method, where those same frames are assigned to buffer **a** in node **B**. The slope of the maximum forwarding delay allows us to compute the latest moment at which frames received from buffer **a** on node **A** can be stored into buffer **c** on node **B**. The shaded areas just below timeline 2 in Figure 1 show the time windows for buffer assignment. If two output buffers are used, then frames received from buffer **a** on node **A** can be assigned on input (timeline 2) to buffer **c** only as long as they are assured of being placed into buffer **c** before node **B** starts transmitting buffer **c**. As shown, frames from buffer **a** can be assigned to buffer **a** (three-buffer mode) during the entire length of the cycle on timeline 2. Time T_W in Figure 1 is the time during which, in three-buffer mode, buffer **c** is holding frames, neither filling nor emptying. In 3-buffer mode, the dead buffer time T_B is 0, and T_B , the allocable transmission time, encompasses both the T_B (white) and T_B (red) regions in Figure 1.

Note that an implementation may require a minimum offset between timeline 3 and timeline 4. That is, a time lag may be required between the last opportunity to store a frame in a buffer, and the earliest time at which the first bit of a frame from that buffer can appear on the link. Some time could, for example, be necessary in order to schedule the transmission of frames across multiple queues in order to ensure that the requirements of strict frame priority and back-to-back frame transmission (1.3.3) can be met.

2.4 Calculation of T_B

Timeline 3 in Figure 1 shows the calculation of T_B , which applies only to two-buffer mode. The starting point of is the moment that the output cycle starts (the tick on timeline 4), backed up by the worst-case forwarding delay. This is the last moment on timeline 3 that a frame can be assigned to buffer **c** in the example in Figure 1. The end of T_B is the end of the cycle T_C , less the variation time T_V . In three-buffer mode, T_B is zero.

T_B is only known to node **B**. Its effect on the allocable bandwidth T_A must be taken into account when admitting new streams. If a network uses a peer-to-peer control structure using, e.g. IEEE Std 802.1Q-2018 MSRP, then the value of T_B must be made available to the previous node **A** so that node **A** does not exceed the reduced T_A .

There are many ways to deal with this issue. Here are three:

1. The value of T_B can be propagated backwards to the previous node, either via management or via an extension of the reservation protocol.

2. A node can compute the value of T_B and decide whether to employ 2-buffer or 3-buffer mode, depending on how much bandwidth has been allocated, so far. This, of course, can change previously-computed stream's end-to-end latency.
3. All nodes in a network can be configured with a reasonable maximum value for T_B . If a particular input/output port pair on a particular node computes a value for T_B that exceeds this maximum, then 3-buffer operation is required.

2.5 Transmitting frames timeline 4

Depending on whether two-buffer or three-buffer mode is used, one can trade off reduced total available bandwidth against per-hop delay. Timeline 4 in Figure 1 shows the two options for the choice of which output cycle in node **B** is used to transmit frames that were transmitted from buffer **a** in node **A**.

2.6 More than 3 output buffers

The discussion over Figure 1, so far, assumes that the variation in forwarding delay is small, relative to T_C . If this is not the case, node **B** can use more than 3 output buffers, and assign received frames to buffers whose output is scheduled far enough ahead in time to ensure that, in the worst case, they will arrive in the buffer before the buffer begins transmitting. This works only because the buffer assignment decision is made based on time-of-arrival of the frame at the input port, not the time-of-arrival of the frame at the output port.

In certain situations, e.g. when stream is split and traverses two paths of different lengths using IEEE Std 802.1CB Frame Replication and Elimination for Reliability (FRER), it can be desirable to purposely delay a stream's frames in order to match the total delay for the stream along the two paths. In this case, more than 3 output buffers can be allocated, and used to impose a delay of an arbitrary number of cycle times T_C on every frame.

It will be shown in section 5.1 that it is not difficult to implement CQF so that each output port in a node, and each output port along the path of a stream, can have a different number of buffers, whether 2, 3, or 50. Not only that, but one flow can use 3 buffers on an output port, while another flow, which needs a path-matching delay, can use 12 buffers on the same port.

3 Multiple CQF classes of service

3.1 Multiple T_C model

With CQF as it is described in IEEE Std 802.1Q-2018 Annex T, we are limited to a single class of service (a single value of T_C) and to 2-buffer operation, only. We have already discussed 3-buffer (or more) operation. We will now discuss the simultaneous use of more than one value of T_C on the same output port.

It can be difficult to pick a single value of T_C for a network. If the chosen value is small, then only a few streams can be accommodated on any one port, because all frames for all streams sharing a port must fit into a single T_C period. If the value chosen for T_C is large, then more streams can be accommodated, with a wide variation in allocated bandwidth, but the larger T_C increases the per-hop latency. In the ideal case, of course, every stream would have a T_C value chosen so that exactly one frame of a stream is transmitted on each cycle T_C .

While this is not always possible, we can apply multiple values of T_C to a single output port, as shown in Figure 2.

Figure 2 Multiple T_C values

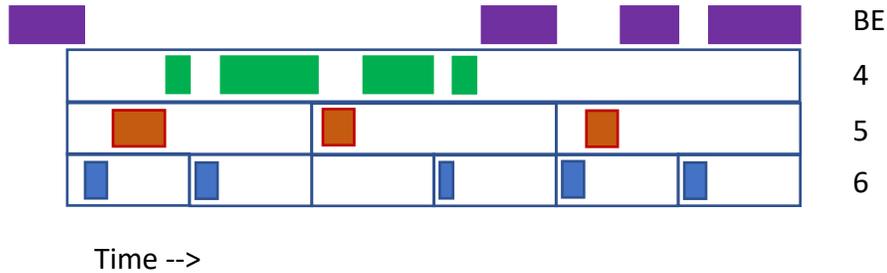
h															3												
f					g					f					4												
d				e				d				e				d				e				5			
a	b	c	a	b	c	a	b	c	a	b	c	a	b	c	a	b	c	a	b	c	a	b	c	a	b	c	6

In Figure 2, we have a schematic timeline. We are running four values of T_C simultaneously. The fastest (call it, " T_{C6} ") runs at the highest priority (6). T_{C5} is slower by a factor of 4 from T_{C6} in this example, and its buffers run at priority 5 (less important than priority 6). T_{C4} is slower by a factor of 2 from T_{C5} , and by a factor of 8 from T_{C6} . T_{C3} is 24 times slower than T_{C6} . The letters in Figure 2 label which buffer is output during the cycle. There are 9 buffers a through i (buffer i is not shown). In this example, priority 6 uses three buffers, because the timing is tight; the others use two each.

We assume here that the receiver of a frame can identify the particular CQF instance (T_C value) to which the frame belongs by inspecting the frame. A TSN bridge could use the L2 priority of the field, for example. An IP router could use the DSCP. IEEE Std 802.1CB and IEEE Std 802.1Q provide for the use of other fields in the frame, e.g. IP 5-tuple.

Since the total bandwidth of the link is not oversubscribed by streams, each cycle, fast high-priority and slow low-priority, is guaranteed to be able to transmit all of its frames within the duration of its cycle. For example: If 50% of T_{C5} is reserved, and 30% of T_{C3} is reserved, then 80% of the total bandwidth has been reserved, leaving only 20% for other streams, best effort traffic, and dead time. This is shown in Figure 3, where we illustrate the timing of transmission of frames from three levels of CQF and the best-effort (BE) level. Note that CQF traffic can be delayed within its window by interference from both higher priorities (e.g. the first priority 4 frame) and lower priorities (e.g. the first priority 6 frame), but that it will always get out before the window closes, assuming that the bandwidth is not oversubscribed.

Figure 3 Transmission timing



A given stream is allocated a fixed number of bits that it can transmit per cycle T_{Cn} . A scheduler would typically assign each stream to the highest-numbered (fastest) CQF instance such that, at the stream's bandwidth and frame size, the stream occupies some space in every buffer at that level. Then, CQF will maintain one or two frames in its buffers per stream, the best possible latency is given that stream, and the buffer space is not wasted in unused cycles.

Of course, it is the "best possible" latency only to a certain extent. The potential mismatch between the stream's frame rate and frame size to the available values of T_{Cn} requires some overprovisioning.

Streams are allocated to, and thus use up the bandwidth available to, each cycle separately. Any cycle can allocate up to 100% of the bandwidth of that cycle's T_A , but the percentages allocated to all of the cycles must, of course, add up to less than 100%. The total amount of buffer space required depends on the allocation of streams to priority values. If all streams are slow and are allocated to T_{C4} up to a total of 100%, then full-sized buffers must be used for buffers h and i. If all streams are fast and are allocated to T_{C6} , then only three small buffers are used—buffers a, b, and c are rapidly re-used.

NOTE—There are many ways to allocate buffer space to individual frames. Running CQF at 5 levels does not increase the buffer memory requirements beyond that of 1-level CQF. Allocating bandwidth to slow cycle times uses more buffer space, of course, because frames dwell for a longer time.

Given the ideal allocation described, each stream is allocated one frame in each cycle of one row. It thus gets the optimal latency for its allocated bandwidth, which may be somewhat oversubscribed. If the end-to-end latency requirements of the streams permit, a stream can be assigned to a slower (lower-numbered) cycle. This will reduce the overprovision factor, since the overprovision factor depends on the number of frames per cycle. It also increases buffer usage, of course.

Any such overprovision can equally be thought of as an increased latency for that same stream. That is, if that oversubscribed stream was the only stream, then the T_C cycle time could be

shortened to exactly the point of 1 frame per cycle, with 0 overprovision, and thus give a faster latency. Overprovision = lower latency, in this case.

The maximum reserved bandwidth is supported by allocating a stream multiple frames per cycle, as allowed by the stream's required end-to-end latency, thus minimizing overprovision.

3.2 Preemption and interference

Frame preemption is described in IEEE Std 802.3-2018 clause 99 and IEEE Std 802.1Q-2018 clause 6.7.2. Not all of the bandwidth in a cycle T_C can be allocated. The smaller the cycle time, the greater the impact of the interference time (T_I in section 2 and Figure 1) on the allocable bandwidth.

T_I is equal to the worst-case transmit time for a single transmission from a lower-priority queue. This interference can occur only at the beginning of a cycle. Since this value must obviously be bound, it places a requirement, that must be enforced, on all lower-priority queues that they either have a maximum frame size or that frame preemption is applied to the lower-priority queues. If preemption is used, the maximum interference is the maximum fragment size (about 150 bytes, see IEEE Std 802.3). The interference time is shown as a gray parallelogram attached to timeline 1 in Figure 1.

The other time is the preemption time T_P , which applies only to streams that are preemptable. This case is not typical, but is possible if a large fraction of the available bandwidth is to be assigned to one or a few high-bandwidth streams, and lower-priority streams use larger frames. T_P is the product of (the maximum number of highest-priority transmission windows that can open during a single window for the level being computed) * (the per-preemption penalty). Thus, in Figure 2, if priority 4 is preemptable, then there are 8 level 6 windows that can open. This means that there can be 8 preemption events during one level 4 window, so the total preemption time T_P is 8 times the preemption penalty. (It doesn't matter which specific frames are preempted; only how many such events occur during the cycle.) The preemption penalty is the number of bytes added when a frame is preempted, which is 4 (CRC on preempted fragment) + 20 (inter-frame gap) + 8 (preamble for continuation fragment) = 32 bytes.

3.3 T_C computation

If the time per cycle that is allocable to streams is T_A , then we can now state the computation for T_C , given T_A , or for T_A , given T_C , at each level in Figure 2:

$$T_C = T_A + T_P + T_I + T_B + T_V$$

The sum of all stream's bits-per-cycle allocation must be less than or equal to T_A .

4.1 Basic requirement for determinism

Multi-CQF guarantees the Deterministic QoS by the following argument.

We assume that the Talker uses Multi-CQF. Non-CQF Talkers are discussed in 7.1.

We consider only one value of T_c along the path of a given stream from Talker to Listener. Changing the TC value requires re-conditioning, as described in 7.3.

The contract between the Talker and the network is in terms of 1) a maximum frame size, and 2) a maximum number of bit-times on the medium per cycle time. For Ethernet, the number of bit times for a given frame is equal to (the frame size from destination MAC address through Frame Check Sequence, plus 20 bytes for preamble and inter-frame gap) times 8 bits per byte.

A number of considerations reduce the fraction of the total time T_c that can actually be used to transmit data. See section 2 for details. For example, the maximum frame size of each stream allows us to determine the worst-case interference that a given stream can have on higher-priority streams. All of these considerations are bounded; if an implementation cannot bound one or more of these considerations, then it cannot guarantee the Deterministic QoS in a CQF network.

In a detailed timing analysis, we will assume that the primary rule of IEEE Std 802.1Q Scheduled Transmissions is adhered to: the first bit of the MAC address of a frame is never transmitted before the start of the window time (according to the local time in the transmitter) and the last bit of the interframe gap (always) and the preamble of the next frame (if any) are transmitted before the end of the window.

In order to obtain Deterministic QoS for each stream, we must ensure that no buffer is ever asked to hold more data than it can transmit during one cycle time T_c . Since the amount of data supplied by any given stream in one cycle is set by contract, we can accomplish this as follows:

- 1) The Talker contract is enforced when a Talker's frames are first placed into a CQF output buffer after entry to the network. That is, the frames from a given stream do not exceed the Talker contract in the first CQF output buffer in the network.

Ingress conditioning and/or policing is discussed in 7.1.

- 2) Frames belonging to the same stream that are in the same CQF output buffer in one bridge in the network are placed in the same CQF output buffer in all subsequent bridges along a shared path.

Section 2, and particularly Figure 1, show the details of how this is accomplished. The key is to get the input gates synchronized with the output gates of the transmitting

system, offset by the link delay. Frames received during one input cycle are always placed in the same buffer. If the input cycle is synchronized with the previous hop's output cycle, then cycle integrity is maintained. (Of course, this only works for point-to-point links.)

- 3) There is no fan-in for a particular stream.

We assume that the path of a stream reservation through the network is known and does not change. A given stream enters a bridge through one port only, although it may be a multicast stream, and thus be enqueued and transmitted on more than one port.

- 4) Admission control ensures that, on any given output port and cycle time T_c , the total bits times for all streams passing through that port and T_c value does not exceed the available transmission time on that port. (This assumes that no bridge has a limitation on available *receive* time on an input port that is smaller than the attached output port's available transmit time. The implications of such a limitation are obvious.)

4.2 Admission control for multiple T_c values

Section 3 describes the operation of Multi-CQF with multiple T_c values operating simultaneously on one output port. Figure 3 shows an example of a sequence of transmissions. We observe that the shortest cycle times operate at the highest priority, and the longest at the lowest priority. Because different CQF priority levels may have different maximum frame sizes, and because some may enable preemption, different priority levels may have different amounts of time during one cycle that cannot be allocated to stream transmission. Clearly, allocating time for any CQF priority level reduces the time allocable to other priority levels; there is only one physical link.

An administrator may wish to restrict allocation of CQF transmission times to leave room for transmitting non-CQF frames, either best-effort traffic or other, lower-priority TSN traffic.

For a new stream to be admitted, it must be true that the available transmission times over all of the CQF levels on all of the output ports through which the stream travels have not been exhausted. At any given CQF priority level x , one can add the bits allocated to all streams in one cycle at CQF priority level x , plus the sum over all more-important CQF priority levels y (faster cycles), of the product of the number of bits per cycle allocated at that level times the number of cycles at that level contained within one cycle at level x . At every level, the total must not exceed the maximum number of allocable bits at that level.

(This calculation is simpler if, at every CQF priority level, there is the same percentage of dead time and slop for inaccuracies, but this is not necessarily the case.)

4.3 Computing the actual end-to-end latency for Multi-CQF

After adjusting to get the receiving window aligned with the previous-hop transmitting window, a bridge knows the “effective phase difference T_{AB} ” described in section 2. Referring to Figure 1, this allows the bridge to compute the difference, in time, between the start of an input window for the stream, and the start of the output window in which a frame received in that input window will be transmitted. This is the dwell time for the frame in this bridge. Adding this to the one-way link delay gives the per-hop delay for frames in the stream. At egress from the CQF network, there is a margin of one cycle time less one frame transmission time for delivery of the frame, as the frame can be transmitted at any point during the cycle, but must both start and finish its transmission within the cycle. The delay at ingress is somewhat more complicated to measure, as it depends upon the method used by the Talker and the ingress bridge to shape its transmissions.

If we look again at Figure 1, we can see that the difference between using two and three buffers for a given input-output port pair is really a matter of rounding up the link delay to an integral number of cycle times. If the sum of link delay and phase delay between output cycles is negligible, or happens to be very nearly an integer multiple of the cycle time, then the yellow “discard” area is small, and two buffers can be used. If sum is larger, then one necessarily chooses between a smaller allocation (large discard area) and increased delay.

5 Parameterization of Multi-CQF

5.1 Per-stream and per-port static state

In the general case, the minimum number of buffers required (usually 2 or 3) depends on the relative phase of the input and the output cycle start times. But different input ports generally will have different phases. Thus, the number of buffers used by any given output port will vary with the input port; an output port can have three buffers, for example, but for some input ports, there are never frames from that port in more than two buffers.

For the present time, we will assume the following method for receiving a frame and assigning it to a buffer. It is important to stress that there are many ways to accomplish the same task.

Let B_o be the number of physical output buffers on port o . We compute N , the least common multiple over all B_o in the system. Each input port i assigns each received frame a buffer selector S , which is an integer in the range 0 through $N-1$, and which increments (modulo N) each input cycle. Thus, frames transmitted from the same buffer are assigned the same S value at the receiving end of the link.

At the output port o , each of the B_o buffers is identified by a buffer number in the range 0 through B_o-1 . A variable X_o indicates which buffer is currently transmitting. X_o increments once modulo B_o each output cycle.

When a frame arrives at an output port, it is assigned to a buffer b using the formula:

$$b = (S + P_{io}) \bmod B_n$$

Where P_{io} is the cycle phase offset from input port i to output port o . See 5.5 for the determination of P_{io} . Note that in the extreme case of all output ports using two buffers, all synchronized, and all input cycles in phase with the output cycles, the table P_{io} reduces to a single value, 0 or 1, and we have CQF from Annex T.

It is desirable in some cases to deliberately use more buffers than are required for insurance against congestion loss in order to match the end-to-end delay of a stream across different paths through the network. If such delay matching is performed per-stream, instead of per-input port, then per-stream P_{io} values are required for buffer selection.

P_{io} is not dynamic, though its values may change when the relative phasing between an input port cycle and the transmitter feeding it change suddenly. Such a change will always disrupt the CQF service guarantees.

Let us go through the exercise of initializing an input/output port pair for Multi-CQF. In the process, we will collect a set of parameters that can be used with protocols and/or network management to monitor and control the operation of Multi-CQF.

5.2 Cycle wander

Adjacent bridges must be frequency locked as described in 1.3.2. For any given port, there is a worst-case system clock difference, $sysClockVar$, between this bridge's system clock and the neighbor system attached to the port. Its units are a time difference. We will assume that this parameter is configured by management, based on network design parameters and system data sheets. It is possible that this parameter can be adjusted during network operation. A bridge could have more than one system clock, and be connected to another system by multiple links, but there is only one value for $sysClockVar$ for any given port, because we assume point-to-point links. We will assume that the variation can be in either direction, this-end-late or this-end-early.

IEEE Std 802.1Q scheduled transmissions and scheduled input gates are assumed to operate under control of a clock that is local to a port. In 802.1Q, the management controls that configure the schedule are defined in terms of the system clock. The bridge aligns the port clock(s) with the system clock either periodically or continuously. There is thus a worst-case excursion of the actual start of a cycle from the time configured in terms of the system clock. We parameterize this with four parameters for each port, $cycleInMaxEarly$, $cycleInMaxLate$, $cycleOutMaxEarly$, and $cycleOutMaxLate$, all measures of time. $cycleIn^{***}$ is for the input gate error, $cycleOut^{***}$ for the transmission error. $^{***}Early$ is the worst-case for starting the cycle

before the system clock time, and ***Late the worst-case for starting after the system clock time. Having both early and late parameters allows for different implementation methods for aligning the port schedule to the system clock. This parameter is computed by the bridge from knowledge of the implementation, and is constant over the lifetime of a physical connection.

5.3 Link delay variation

The time taken for a frame to travel from the transmitter to the receiver can vary for two reasons: the actual delay can change, due for example to temperature variations in a multi-kilometer link, and the measurement of the link delay can vary due to various clock inaccuracies. We will deal only with actual variations, not measurement variations. [1-21-0056-00-ICne-input-synchronization-for-cyclic-queueing-and-forwarding](#) explains why this is possible.

5.4 Calculating the number of buffers required

The procedure to calculate the number of buffers needed on an output port to support one particular input port is as follows:

- 1) Establish a Nominal Input Cycle Start time (NICS) for the input port, and a Nominal Output Cycle Start time (NOCS) for the output port. The NICS and NOCS each repeat every T_c seconds, according to the system clock. We will assume that the offset between them is a constant (i.e., they are both driven by the same system clock).
- 2) Compute the earliest time, relative to the NICS, at which the first frame of a cycle can receive its IEEE Std 802.3 clause 90 timestamp. This frame is assumed to be a minimum-length frame (64 bytes plus overhead).
- 3) Compute the earliest time, relative to the NICS, at which a buffer on the output port must be eligible to receive the frame. This is equal to the timestamp time in bullet (2) plus the minimum time required to move the frame through the bridge to the output buffer.
- 4) Compute the latest time, relative to the NICS, at which the last frame of a cycle can receive its timestamp. This frame is assumed to be a minimum-length frame.
- 5) If the difference between the earliest timestamp and the latest timestamp is greater than or equal to the cycle time T_c , then dead time must be imposed on the transmitter, at the end of the cycle, to reduce the difference.
- 6) Compute the latest time, relative to the NICS, at which the last frame of a cycle can be stored into an output buffer and be ready for selection for transmission, given the worst-case forwarding delay through the bridge.
- 7) Convert these earliest (2) and latest (4) arrival times to times relative to the NOCS of the output port.
- 8) Arbitrarily label an input port NICS event $NICS_0$. Determine the latest subsequent NOCS event, which we will label $NOCS_0$, during which the earliest-arriving frame of $NICS_0$ must be stored in the output queue.

- 9) Determine the earliest subsequent NOCS event, which we will label NOCS_n , before which the latest-arriving frame from NICS_0 can be stored in the queue, and still be available for transmission at the start of cycle NOCS_n .
- 10) The number of cycles NOCS_0 through NOCS_n , inclusive, is the number of buffers required for the input/output port pair, B_{io} .

The number of buffers required can sometimes be reduced by:

- a) Imposing a larger dead time on the transmitter feeding the input port, at the end of every cycle;
- b) Altering the phase of the output port's cycle; and/or
- c) Imposing implementation-specific limitations on the flows, e.g. reducing fan-in to an output port, or restricting bridging/routing features to reduce forwarding delay variation.

Finally, let us observe that large link delay variations can be accommodated by varying the above calculation. Assuming that the variations take place slowly, and that changes in relative phase between transmitter and receiver are detected using a protocol (e.g. that in [new-finn-CQF-sync-method-09-21-v1](#)), the difference between the maximum and minimum link delay can be added to the difference between the earliest- and latest- arriving frames to increase the number of buffers allocated. The phase of the input gate can be altered by small increments as the protocol detects the phase differences, without gaining or losing cycles in the transfer. Of course, the maximum adjustment made per phase adjustment event must be removed from the allocable bandwidth.

5.5 Initial buffer phase

The number of buffers required on an output port is the maximum required over all input ports. This may be further increased by intentional delays (2.6). When initializing an input port, a correspondence must be made between the input and output ports, so that a frame received on the input port will be stored in a particular buffer in the output port, the one that will become the transmitting buffer in the appropriate number of output cycles in the future.

The phasing between input and output ports' cycles, and thus the number of buffers in port o used by port i , is determined by the P_{io} table defined in 5.1. We compute P_{io} when initializing CQF, or when the relative phase of the input and output ports change significantly, by selecting a time T that coincides with the start of an input cycle on input port i and computing:

$$P_{io} = (X_o - S_i - B_{io} + 1) \bmod N$$

Where X_o is the identity of the transmitting buffer on output port o at time T , B_{io} is the total number of buffers required of output port o by input port i (including the transmit buffer), S_i is the value of buffer ID S assigned by port i during the input cycle starting at time T , and N is the range of S_i , the least common multiple of the number of physical buffers over all output ports.

5.6 Externally-visible Multi-CQF managed objects and protocol items

The following list includes both objects of interest to a network manager, and information elements that might be exchanged using a link-local protocol. Most items could be carried in a protocol as a check on proper configuration of adjacent ports, with varying degrees of utility for different items. Some items can only be computed by one system, and must also be known to the adjacent system. It is for further study what protocols would be used for such information transfers, or and/or whether the transfers are best accomplished using network management.

5.6.1 Cycle and priority structure managed objects

For each output port and each input port, separately, we have:

- a) The cycle time of the slowest CQF priority value.
- b) The priority value of the slowest CQF cycle.

For each priority level running Multi-CQF on an input port or an output port (separately), we have:

- c) The layer 2 priority value
- d) The number of cycles at this priority level contained within one next-lower priority value cycle.

There are other, equivalent, ways to formulate this same information. We can divorce layer 2 priority code point from importance, for example.

These parameters are not expected to change over the lifetime of a data stream. A system would not be expected to obtain this configuration information from a neighbor through a CQF-specific protocol, though exchanging this information could be done to discover of configuration errors.

5.6.2 Cycle phase managed objects

For each output port and input port, separately, we have:

- a) The start time of the slowest Multi-CQF cycle, in terms of the system clock.

This variable establishes the phase of the input or output cycle (NICS and NOCS, see 5.4). Typically, this variable would be managed the network administrator for output ports. For time-synchronized systems, it can be administered for input ports, as well. Alternatively, the input phase can be determined dynamically, and be read by the network administrator. Because it is in terms of the system clock, it is of no interest to

neighboring systems except, perhaps, as a configuration error check for time-synchronized networks.

5.6.3 Cycle variation information

For each output port only, we have:

- a) The largest offset from the nominal (system clock) NOCS event to the actual cycle start time, in the negative (actual earlier than NOCS) direction.
- b) The largest offset from the nominal (system clock) NOCS event to the actual cycle start time, in the positive (actual later than NOCS) direction.

There are other ways to express the information in these two items. These values must be known to the connected input port in order for that system to compute its buffer and dead time requirements (`cycleOutMaxEarly`, `cycleOutMaxLate` in 5.2). This information transfer could be accomplished by means of a protocol, managed objects, or by restrictions on implementations.

5.6.4 Dead time / bandwidth balance information

There remains the balancing of conflicting goals between dead the percentage of a cycle that is available to transmit critical data streams, and the number of buffers required on the output port. Increasing the dead time can reduce the number of buffers required, and thus the end-to-end latency of a data stream, as described in 5.4. There are, at the very least, the following ways to make this decision:

- 1) Configure the output cycle phase and number of buffers to use for all bridges, in order to establish a constant per-hop delay in a network with short links. Let each system compute the dead time on each input port required to make this work, and the bandwidth available for allocation. Convey the required dead time either by protocol or by management to the transmitters, and the available bandwidth to the admission control system.
- 2) Configure the output cycle phase on all bridges. Configure minimum and maximum allocable bandwidth values for each CQF priority level. Let each system compute the minimum number of buffers required to meet the minimum bandwidth value, taking advantage of the maximum bandwidth value to compute a dead time value that minimizes the number of buffers required. This would be useful in a network with very long links. Convey the resultant dead time to the transmitter via protocol, and the resultant allocable bandwidth to the admission control system.
- 3) Using data sheet information, configure all parameters via network management. Adjust the output port cycle phasing to optimize the delay for certain specific streams.

Given that context, the following items are required by Multi-CQF:

- a) Per input port, per priority level, the total dead time that must be provided by the adjacent transmitter at the end of each transmit cycle.

There is a component of this dead time computed in this section, as well as one computed in item (5) of 5.4. The sum of these must be known to the adjacent transmitting port.

- b) Per output port, per priority level, the total dead time that is to be provided at the end of each transmit cycle.

This can be configured, obtained from the adjacent input node, or be a maximum of these values.

- c) The allocable bandwidth for this input port and priority level.

This has three components, the minimum of the allocable bandwidth over all output ports reachable from this input port (in the input port's own system), any limitations imposed by the input port implementation, and any maximum imposed by management. Whether this is computed by, received by, or even known by the output port, or whether allocable bandwidth is the concern only of the admission control system, is an open question.

- d) The allocable bandwidth for this output port and priority level.

This can be configured, computed from the adjacent input node's requirements, or be a minimum of these values. . Whether this is computed by, received by, or even known by the output port, or whether allocable bandwidth is the concern only of the admission control system, is an open question.

6 Other issues

6.1 Frame size problem

The above discussion has largely assumed that each stream consists of frames of a uniform size, equal to the stream's maximum frame size. Of course, this is not always true.

The advantage of uniform frame size is that, in the ideal case, one can allocate a stream one frame per cycle, and choose the cycle time and/or the stream's bandwidth reservation so that there is no wasted bandwidth. Similarly, if we imagine that a stream alternates frames of 4000 bit times and 800 bit times, we can allocate 4800 bit times per T_c and still get perfect results.

But, in a service provider situation where we are allocating a certain bandwidth per customer, but the frame sizes are essentially random, things are not so simple. Let us suppose that the maximum frame for a stream is 13000 bit times, which is approximately equal to a maximum-length Ethernet frame, and that the cycle time $T_C = 100\mu\text{s}$. $13000/100\mu\text{s} = 130$ Mbits/s. But, allocating a bandwidth of $13000 \text{ bits}/T_C$ will not give the stream 130 Mb/s. In the worst case, one 13000 bit frame followed by one minimum-length frame = 672 bits, the stream gets $(13000+672)/(200 \mu\text{s}) = 68.36$ Mb/s.

We could overprovision the stream by a factor of almost 2, keep the same T_C , and get minimal latency. However, we could also assign the stream to a longer T_C . In the worst case, there are $(13000-8)$ wasted bits in each cycle. Therefore, we can guarantee 130 Mb/s using a cycle time of $500\mu\text{s}$ by provisioning $(5*13000 + 13000 - 8)/(5*100\mu\text{s})$, or 156 Mb/s, which is a 20% overprovisioning, rather than a 90% overprovisioning, at the cost of five times the per-hop latency.

This overprovisioning/latency tradeoff is only needed for streams that have variable frame sizes, such as service provider streams. But, for those streams, the lengths of the links may be a larger source of latency than the queuing delays, so the situation may not be so bad. Also, any unused bandwidth is available to non-TSN data, so overprovisioning may not be a serious concern.

In some use cases, it would be useful to “bundle” frames. That is, combine and or split frames to form constant-sized transmission units perfectly suited to particular multi-CQF cycle time. We will not expand upon this idea, here.

6.2 Tailored bandwidth offerings

In a service provider environment, overprovisioning can also be improved by offering the customer only a specific set of choices for a bandwidth contract, corresponding to the values of T_C implemented in the provider’s network. This way, the overprovisioning required for meeting an arbitrary distribution of requirements using a small set of T_C values is eliminated. (Or, at least, shifted to the customer’s shoulders.)

6.3 Overprovisioning is not always bad

Overprovisioning the bandwidth (allocating more of T_A than is necessary) is not always a bad thing:

- a. Allocating a stream to a higher priority (smaller T_C) than it needs reduces its worst-case latency. This may be necessary to meet a stream’s end-to-end latency requirement. That is, one can overprovision the frame rate in order to obtain a reduced latency. The unused bandwidth is still available for best-effort traffic. Not all TSN transmission selection schemes have this feature.

- b. If the total bandwidth required by critical streams is relatively low, using faster-than-necessary T_c values will both improve latency and reduce buffer requirements in the network. The allocated-but-unused bandwidth is still available to best-effort traffic, and thus may be of no consequence.

6.4 CQF and credit-based shaper

Looking at Figure 3, we see that, once the major cycle at priority level 4 begins transmitting, the best-effort traffic is interrupted until all of the CQF level 4 data is transmitted. At some point, as the amount of traffic in a very slow CQF cycle increases, the burstiness of the best-effort transmission opportunities could, in theory, become a problem. This can be mitigated by applying a credit-based shaper function to the slowest multi-CQF cycle(s). However, the parameters of this shaper must be adjusted as the load on the slow CQF cycle(s) changes, because a bridge must *always* finish transmitting all of the data in a CQF buffer. Thus, adding a credit-based shaper would detract from the most-significant advantage of multi-CQF—its freedom from requiring reconfiguring a bridge each time a flow is added.

6.5 Fundamental CQF pros and cons

The obvious downside of CQF is that it requires clock frequency synchronization and per-port time-based gating. On the other hand, CQF requires no per-stream per-hop active state machines. A new stream can be provisioned by a network controller without any interaction between the network controller and any of the network's relay systems, except for configuring one system for ingress policing/conditioning. Furthermore, calculation of the worst-case end-to-end latency is trivial, and the calculation made for one allocated stream is never affected by any other allocations or deallocations.

7 Paternoster and multi-CQF

There are several cases where Mick Seaman's Paternoster algorithm greatly aids multi-CQF operations. This algorithm is described in [cr-seaman-paternoster-policing-scheduling-0519-v04](#). We can characterize it, for the purposes of the present document, as providing a counter state machine for each stream that allows a stream to store no more than its contracted amount of data per cycle into any given CQF buffer. Frames above that limit are stored in subsequent buffers, up to the maximum amount of buffer space allowed that stream, whereupon excess data is discarded.

7.1 Ingress conditioning

A reserved stream entering a bridge from a correctly-behaving bridge or end station that runs CQF at the same cycle time as the receiving bridge needs no examination beyond being forwarding to the correct port(s). However, a multi-CQF bridge could receive input from a

bridge, a talker, a router, or any other device that uses some deterministic algorithm(s) to condition its critical streams, but uses an algorithm other than CQF. We assume that reservations (contracts) for these streams can be translated into CQF terms, with perhaps some overprovisioning required. Long term, the data adheres to the contract. But, for any given input window, an excessive number of bits can be received for a particular stream. We would like to accommodate such input.

The paternoster algorithm makes this possible. A bridge uses the same buffer structure and output methods described in the preceding sections, but instead of obtaining the buffer selector S from the time of receipt of the input frame, as described in 5.1, it uses a state machine dedicated to each non-CQF stream, running the paternoster algorithm, to determine the buffer selector.

More specifically, a stream's frames are stored in the buffer that is next to be output until either a) that buffer is switched to become the output buffer, or b) storing a frame would exceed the number of transmission bit times allocated to that stream. In the latter case, the data is stored and counted in the next available buffer. If the number of buffers required exceeds a limit specific to the stream, the frame is discarded and an error condition is noted. A more detailed description of the algorithm is given in Seamans paper, and will not be given, here.

7.2 Paternoster Class of Service vs CQF Class of Service

A given output buffer can accept input from both CQF and paternoster streams, as long as they share the same cycle time; separate paternoster and CQF buffers or queues are not necessary. In addition, a paranoid network administrator could very well configure a paternoster shaper on every stream in a purely-CQF network, in order to guard against misbehaving bridges or talkers. That is, while paternoster can be thought of as separate algorithm from CQF, it can also be thought of as a protection mechanism for CQF that can be employed as need, and when employed everywhere, removes the restriction that bridges operate at exactly (or even, approximately) the same frequency.

This author would suggest that, in many networks, paternoster and CQF should be the same priority level. The choice between paternoster and CQF can be made on a bridge-by-bridge basis, and not be visible to the talker, the listener, or the user.

7.3 Changing cycle times

If a stream enters a bridge using a cycle time T_C , and is being transmitted on an output port with cycle time $n * T_C$, then n successive input cycles can be deposited in the same output buffer with no problem, as long as the larger cycle time's dead time requirements are met. (This is not a trivial exception, as the larger cycle's dead time occurs at the end of the large cycle, and thus may take up much or even all of one small cycle.) Equivalently, the input port can be

configured with the slower cycle time to match the output port in the same system. Of course, when making the reservation for that stream, the adjustment of its contract must be made; it is allocated n times the number of bits in the slower cycle than in the faster cycle.

In all other cases, when a stream changes cycle times, the stream must pass through a paternoster shaper to ensure that the stream never exceeds its contract in the new cycle time.

7.4 Aggregation and dis-aggregation

A particular form of stream aggregation is useful for both reducing the number of streams tracked by a bridge, and for improving the per-hop latency of a stream. In this type of aggregation, a number of streams are treated as a single stream, with a single reservation, traversing a single path, for some portion of their journey through the network. It does not matter whether the frames are actually encapsulated in some common wrapper, or whether they are simply treated identically (e.g. given the same IEEE Std 802.1CB stream_identifier).

What does matter for multi-CQF is that the aggregate stream, which has a bandwidth equal to the sum of its component streams, can be assigned a CQF level with a faster T_C than its components could make use of. Less buffer space is used for the aggregate than for the separate streams, because only one frame of the aggregate need fit in a cycle, instead of one frame per component stream. The faster T_C lowers the latency and reduces the buffer requirements for all its component streams.

In general, this requires that the aggregate stream pass through a paternoster state machine when it is formed from its components, and that each component pass through a paternoster state machine if and when it is again separated as an individual stream and passes, presumably, to a slower T_C value.

7.5 Ingress Paternoster

Seaman describes paternoster as if it operates on output buffers. In certain situations, as for example in a bridge that has only one of its many ports that connects to a non-CQF transmitter, it may be more convenient to attach the paternoster state machine to the input port, rather than to the output port. In this scenario, the assignment of a buffer selector S at the input gate, based on time as described in 5.1, is replaced by an assignment based on a paternoster state machine. This, of course, requires that a given stream enter the bridge on exactly one port (or that the ports' paternoster state machines be coordinated somehow). However, it may be an economical alternative to providing the paternoster capability on all output ports, when only one or a few input ports requires conditioning.