

Ethernet Congestion Manager

Davide Bergamasco

Cisco Systems, Inc.

Abstract

This document describes a congestion management mechanism aiming at controlling congestion in short-range, high-speed Ethernet networks such as Data Center Networks. Such mechanism, called Ethernet Congestion Manager (ECM), includes three components: (1) congestion detectors associated with bridge transmission queues, (2) rate limiters associated with NICs transmission queues to control the traffic injection rate, and (3) a signaling protocol to convey congestion control information from detectors to rate limiters. When congestion arises in some queue in the network, congestion control signals are sent from such queue to the NICs originating the flows causing congestion. Such signals will cause rate limiters to slow down the offending flows to a rate compatible with the transmission rate from the congested queue, effectively bringing down the queue depth to a desired level.

Acknowledgemtns: Many thanks to Bruce Kwan (Broadcom), Mitch Gusat (IBM), and Manoj Wadekar (Intel) for their kind contribution to this document.

Modification History

Revision	Date	Originator	Comments
0.1	2/5/07	Davide Bergamasco	Initial Draft

Table of Contents

1		
2	1	Introduction 3
3	2	ECM Overview 3
4	2.1	Terminology 4
5	3	ECM Components and Operations..... 5
6	3.1	Signaling..... 5
7	3.2	Detection..... 8
8	3.2.1	Congestion Point Pseudo-code 10
9	3.3	Reaction..... 11
10	3.3.1	Rate Control Algorithm 13
11	3.3.2	Exceptions and Non-linear Rate Adjustments 15
12	3.3.3	Reaction Point Pseudo-code 16
13	4	Conclusions 17
14	5	References 18
15	6	Glossary..... 18
16		

List of Figures

17		
18	Figure 1	– Example of a congested Data Center Network and ECM messaging..... 4
19	Figure 2	– Types of ECM messages 5
20	Figure 3	– ECM Frame Format..... 6
21	Figure 4	– Format of the CM-Tag..... 7
22	Figure 5	- Congestion Detection Process and Message Generation at a Congestion Point..... 8
23	Figure 6	– Sampling Process 8
24	Figure 7	– Data-path structure of a Reaction Point..... 12
25	Figure 8	– Example of timeout and random restart..... 15
26		
27		

1 Introduction

Data Center networks are a peculiar environment because of their high speed (at least 1 Gbps) and low latency (a few tens of microseconds of round trip time). In certain cases, such networks may make use of 802.3X link-level flow control to deliver near-zero packet loss to applications. Such requirements make congestion management in data center networks quite challenging because:

- High speed coupled with small buffers (required to provide low latency) causes such buffers to fill up extremely quickly when congestion arises;
- If link-level flow-control is being used, congestion spreads over a saturation tree almost instantly causing severe head of the line blocking, possibly live-lock, or – in the worst of scenarios – even dead-lock.

Traditional congestion control techniques such as RED [1] and ECN [2] have been shown not to work well with small buffers because of the extremely compressed dynamics exhibited by such buffers. In fact, under congestion conditions a buffer in a typical data center network may fill up in a few hundreds of microseconds, forcing RED and ECN to work in the region of maximum drop/mark probability. This, in turn, causes the traffic flows to back off too much, and – consequently – substantial loss of throughput. Also, RED and ECN are layer 3 and 4 congestion management mechanisms which work only in presence of a cooperating transport protocols such as TCP. Since in data center network there is a substantial presence of non-TCP traffic RED and ECN are ineffective at controlling congestion caused by such traffic.

To overcome the above limitations, we propose *Ethernet Congestion Manager* (ECM), a layer 2 congestion control mechanism conceived to operate in networks limited in scope as per the IEEE P802.1Qau Project Authorization Request (PAR) [3]. The founding principles of ECM are:

- Pushing congestion from the core of the network towards the edge, where there is less traffic aggregation and more resources to deal with it more effectively;
- Using rate-limiters at the edge to control the rate of traffic injection for flows causing congestion;
- Tuning rate-limiter parameters based on continuous feedback coming from the congestion points.

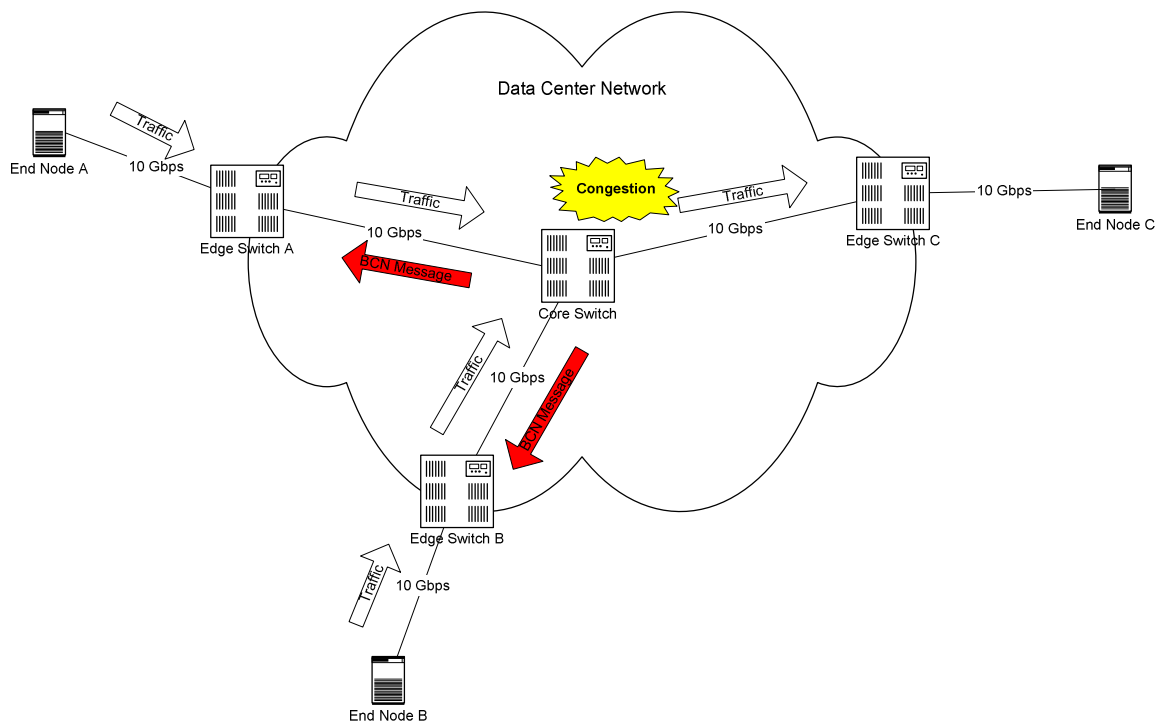
The rest of this document is organized as follows: section 2 provides an overview of the ECM mechanism; section 3 describes ECM in detail, particularly section 3.1 discusses the signaling component and the underlying protocol, section 3.2 addresses the congestion detection component, and section 3.3 examines the reaction component. Finally, conclusions are presented in section 4.

2 ECM Overview

Figure 1 shows a sample data center network composed of a core switch, a number of edge switches and some end nodes. End Nodes A and B are simultaneously sending traffic at line rate (10 Gbps) to end node C. Since the aggregate traffic rate exceeds the capacity of the link connecting the Core Switch to Edge Switch C, this link is subject to congestion and the queue(s) associated with it start filling up. The *detection* component of ECM associated with that queue samples arriving frames with a certain probability. Based on a configurable threshold that define the ideal queue depth, if the current queue length exceeds the threshold, ECM *signals* “slow-

1 down” by sending messages destined to end nodes A and B. Such messages are processed at end
2 nodes A and B. The *reaction* to a “slow-down” message is the instantiation of a rate limiter (or a
3 further slow down if one is already instantiated) at the processing point. The purpose of the rate
4 limiter is to slow down a congesting traffic flow to mitigate congestion at the core switch.

5 Frames that have been rate limited carry this information with them. ECM continues to monitor
6 the queue length and, as the congestion begins to abate, signals the nodes which have rate limited
7 their traffic “speed up”. This is to avoid under-utilizing the bandwidth at the congestion point.
8



9
10 **Figure 1 – Example of a congested Data Center Network and ECM messaging**

11 2.1 Terminology

12 This section provides the definition of a number of terms that will be used throughout this
13 document.

- 14 • **Congestion Management Domain:** The *contiguous* set of Layer 2 devices that support
15 ECM. In a given Congestion Management Domain (CMD), ECM may be enabled or
16 disabled on each of the eight IEEE 802.1Q priorities independently.
- 17 • **Congestion Point:** Place where an uncontrolled accumulation of data frames occurs
18 because of the mismatch in the arrival rate and the departure rate. An example of
19 congestion point (CP) is one of the output queues of a switch.
- 20 • **Detection:** Component of ECM residing at a CP which detects a congestion condition and
21 generates “slow-down” (or “speed-up”) signals to bring such condition under control.
- 22 • **Reaction Point:** Place where signals generated by the ECM detection component are
23 processed and terminated. Reaction points (RPs) are usually located in Network Interface

1 Cards (NICs) and consist of a set of queues and an equal number of rate limiters
 2 associated with them.

- 3 • **Reaction:** Component of ECM residing in an RP which implements congestion mitigation
 4 actions. The reaction component process ECM signals and accordingly control the rate of
 5 injection of traffic by adjusting the current rate of the rate limiters.
- 6 • **Signaling:** Component of ECM used to carry congestion control messages from the CPs to
 7 the RPs.

8 The next section describes in detail the three components of the ECM mechanism, namely
 9 signaling, detection, and reaction.

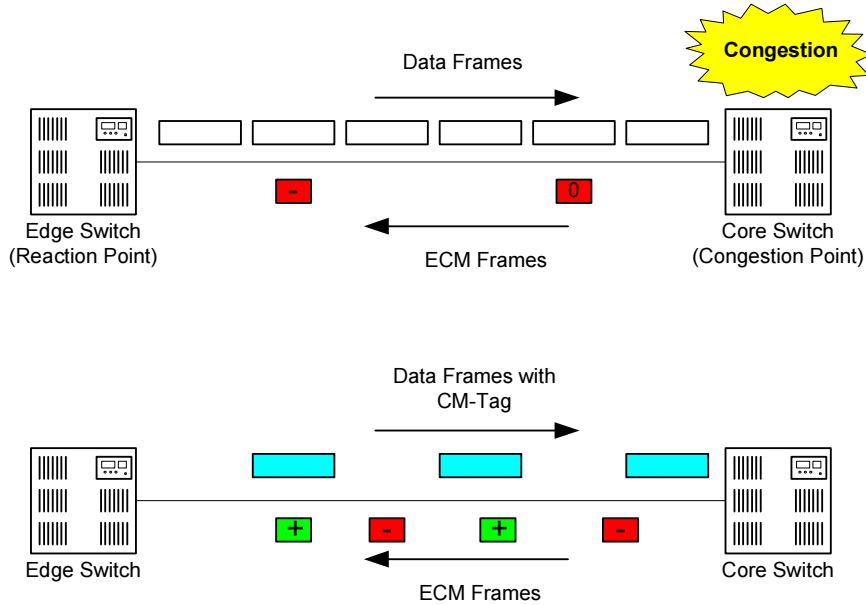
10 3 ECM Components and Operations

11 3.1 Signaling

12 Figure 2 shows the exchange of messages between a CP and a RP. As soon as a CP detects
 13 congestion, as described later in section 3.2, it starts sending explicit feedback messages to the
 14 RPs associated with the traffic flows causing such congestion. The feedback message is an
 15 Ethernet frame known as the *ECM Frame*. A possible format for the ECM Frame is shown in
 16 Figure 3.

17 An ECM Frame is generated by a CP by sampling incoming frames, as described in section 3.2.
 18 The ECM Frame has *Destination Address (DA)* equal to the Source Address of the sampled
 19 frame, and a *Source Address (SA)* equal to a MAC address associated with the CP (usually the
 20 MAC address of the Management Entity of the switch where the CP is located).

21



22

23

24

Figure 2 –ECM signaling

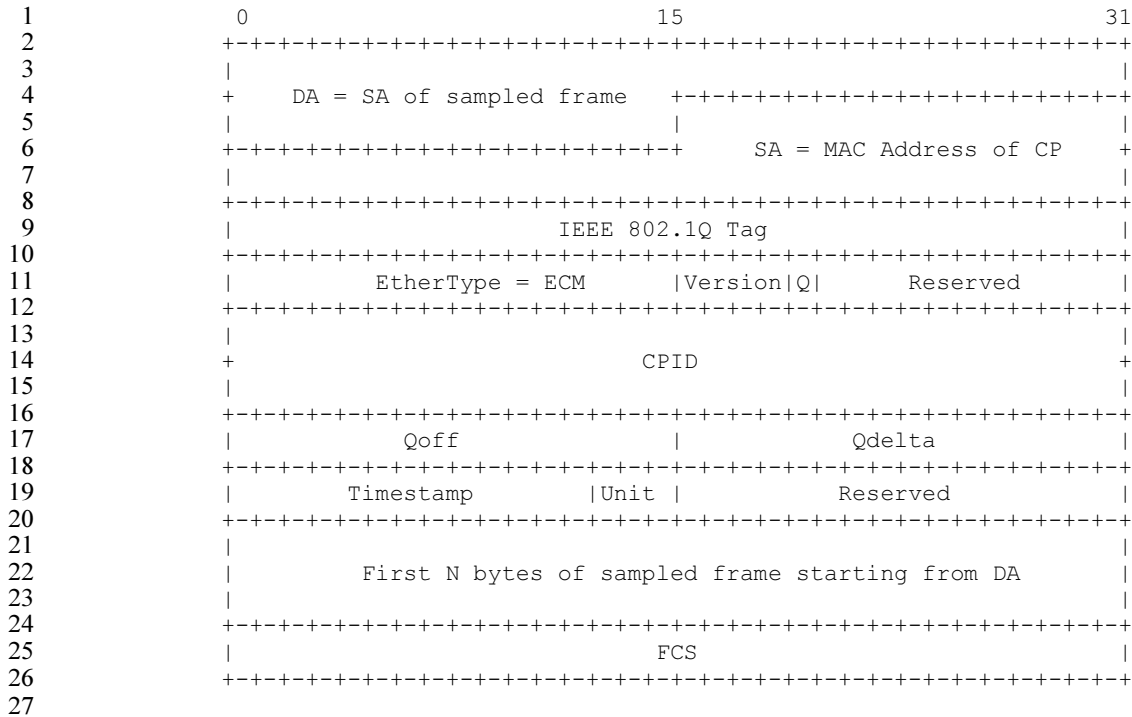


Figure 3 – ECM Frame Format

This ensures that the ECM Frame is forwarded back to the source of the traffic causing congestion with a valid source address.

The *IEEE 802.1Q Tag* is copied from the sampled frame. The Priority field of the ECM Frame 802.1Q Tag is set either to the priority of the sampled frame or to a configurable priority. It is preferable to use the highest priority in order to minimize the latency experienced by ECM Frames.

The *EtherType* of the ECM Frame is set to 0xXXXX, identifying the frame as being an ECM Feedback message.

The *Version* field indicates the version of the ECM protocol. Currently only version 0 (zero) is defined. Subsequent versions of the ECM protocol must be backward compatible. If a ECM implementation version X receives a ECM frame with version Y and $Y > X$, such an implementation must use only the fields defined for version X.

The *Q* bit indicates that the *Qdelta* field has saturated, i.e., its value is either equal to $-2Qeq$ or $2Qeq$. The meaning of such bits will be clarified in section 3.3.1.

The bits in the *Reserved* fields are currently not used. They must be set to 0 (zero) on transmission and ignored on reception. Future versions of the BCN protocol may redefine all or some of the reserved bits.

The *CPID* field is the *Congestion Point Identifier* and its purpose is to univocally identify a congested entity – usually a queue – within CMD. This information has to be propagated to the RP in order to create a bi-univocal association between the CP and RP(s). The CPID field must be unique across the network but, since it is an opaque object, its format is only relevant to the CP

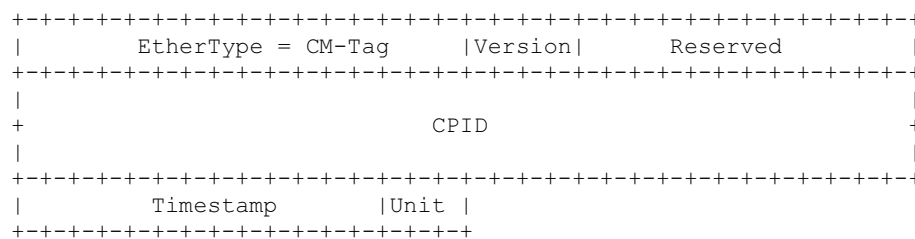
1 that assigns it. The CPID should at least include a MAC address associated with the switch where
 2 the CP resides to ensure global uniqueness, plus a local identifier to ensure local uniqueness.

3 The *Qoff* and *Qdelta* fields contain the actual feedback information conveyed by the Congestion
 4 Point to the Reaction Point. The use of such fields will be described in the next two sections.

5 The *Timestamp* and *Unit* fields are copied from the corresponding fields from CM-Tag (see below
 6 for its definition) of the sampled frame. If the sampled frame does not carry such a tag, the
 7 *Timestamp* and *Unit* fields are set to 0 (zero).

8 The payload of a BCN Frame consists of the first *N* bytes (where *N* is a configurable parameter)
 9 of the sampled frame starting from the DA. The minimum value of *N* is 24, to ensure that a BCN
 10 Frame is at least 64 bytes long. The purpose of such payload is to convey to the RP enough
 11 information to exert the finest congestion mitigation action possible (it should at least include DA,
 12 SA, and 802.1Q Tag).

13 When a RP receives an ECM Frame from a CP, and such message causes a congestion mitigation
 14 action to be performed on a particular traffic flow (usually the activation of a rate limiter or the
 15 adjustment of an existing one), the CPID field from the ECM Frame is stored in a local register
 16 associated with the corresponding rate limiter. All the frames belonging to that flow subsequently
 17 injected by the RP in the network will carry a *Congestion Management Tag* (CM-Tag) containing
 18 the CPID from the register (see section 3.3).



29 **Figure 4 – Format of the CM-Tag**

30
 31 The CM-Tag is identified by the value 0xXXXX in the *EtherType* field and it should be located
 32 after the 802.1Q Tag. Its main purpose is to complete the bi-univocal association between a CP
 33 and a RP. The purpose of this association is to prevent a RP from receiving positive feedback
 34 from multiple CPs for the same flow. In fact, a CP is supposed to generate positive ECM
 35 Feedback messages only for frames that carry a CM-Tag with a CPID matching its own ID.

36 The *Version* field indicates the version of the ECM protocol. Currently only version 0 (zero) is
 37 defined. Subsequent version of the ECM protocol must be backward compatible. If a ECM
 38 implementation version X receives a ECM frame with version Y and Y > X, such an
 39 implementation must use only the fields defined for version X.

40 The *Timestamp* field is used by a RP to estimate the round trip time from the CP it is associated
 41 with. Every time a RP inserts a CM-Tag in the frame it is going to transmit, the current value of a
 42 local free running timer is copied into the Timestamp field. The *Unit* field indicates the time unit
 43 used by the free running timer according to the following equation:

44
 45
$$\text{Time_unit} = 2^{\text{Unit}} \cdot 1 \mu\text{s}$$

3.2 Detection

Figure 5 shows how the detection process works and how ECM Frames are generated at a Congestion Point.

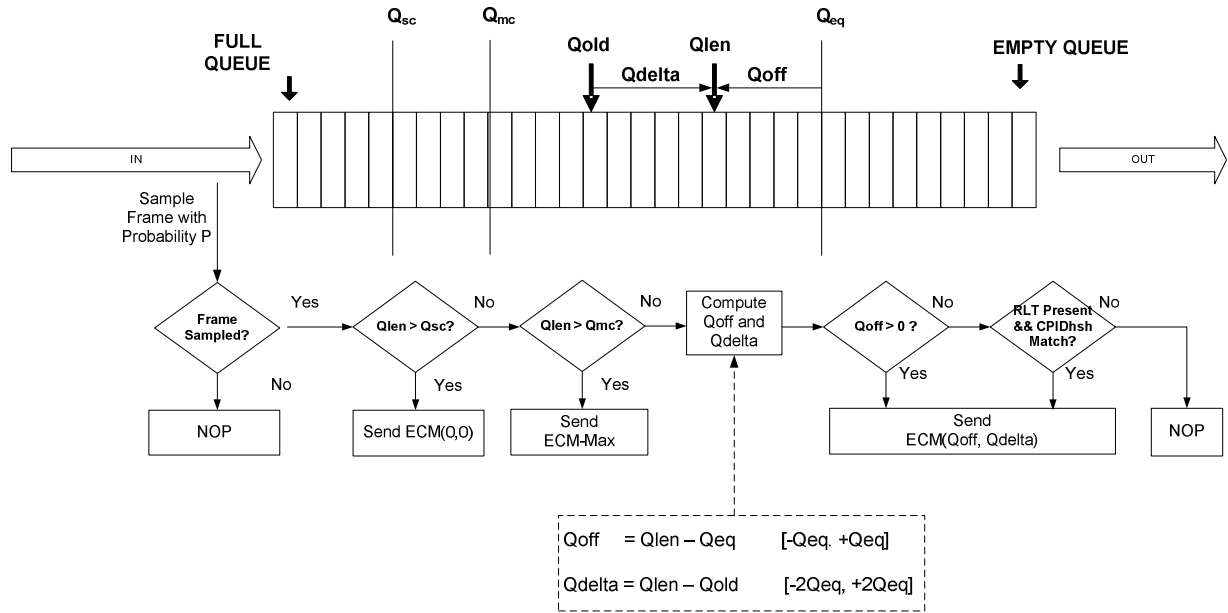


Figure 5 - Congestion Detection Process and Message Generation at a Congestion Point

An *equilibrium* threshold Q_{eq} defines the operating point of a queue under congestion conditions. In other words, Q_{eq} is the target level around which the queue length should oscillate under normal congestion conditions. A severe congestion threshold Q_{sc} defines the level at which the queue is subject to extreme congestion conditions.

Incoming frames are sampled with a certain probability P , e.g., 0.01. Sampling is performed on a byte arrival basis. That is, if the average frame length is $E[L]$, then a frame is sampled on average every $E[L]/P$ byte received. If we assume an average frame length of 1000 bytes, then the average sampling rate is going to be one frame every 100 KB of data received. This can be easily implemented using a *Fixed Interval* followed by a *Random Interval* as shown in Figure 6.

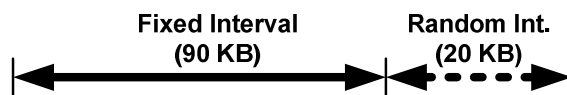


Figure 6 – Sampling Process

Initially and after every sample, the effective sampling interval S is calculated by adding the two intervals:

$$S = S_f + S_r$$

1

2 The length of every frame arriving at the queue is accumulated in L . The frame that makes $L \geq S$
3 is sampled. A new random interval Sr is picked and L is set to zero. The fixed interval Sf should
4 be configurable in the range $[0, 256]$ KB with 1 byte increments. The random interval Sr should
5 be generated in the range $[0, 64]$ KB with 1 byte increments.

6 When the queue length is zero (0) or it is above the *mild congestion threshold* Qmc , the sampling
7 probability is increased by a factor $Sscale$, i.e., S is divided by $Sscale$ ¹. This is called *over-*
8 *sampling* and it's purpose is to speedup the response to transient conditions (i.e., queue going full
9 or empty) by generating more ECM frames when such conditions occur.

10 When a frame is sampled, the current queue length $Qlen$ is compared with the Qsc threshold. If
11 $Qlen$ is greater than Qsc , the queue is under severe congestion conditions, and a special ECM
12 message, i.e., ECM(0,0), is generated. As discussed later in section 3.3, such message causes a
13 rate limiter to temporarily drop its rate to zero.

14 If $Qlen$ is below Qsc but above Qmc , the ECM message corresponding to the maximum negative
15 feedback, i.e., ECM(Qeq, 2*Qeq) is generated. This message, also known as *ECM-Max* will
16 cause the maximum rate decrement to occur at a rate limiter receiving it. ECM-Max, along with
17 over-sampling, allows for a faster response to sudden and quick positive changes in the queue
18 length.

19 If the queue is not operating under severe or mild congestion conditions, the two components of
20 the ECM feedback, $Qoff$ and $Qdelta$, are computed. As shown in Figure 5, $Qoff$ is the offset of the
21 current queue length with respect to the equilibrium threshold Qeq . $Qoff$ must be saturated at
22 $+Qeq$ and $-Qeq$. $Qdelta$ is the change in length of the queue since the last sampled frame, and it
23 must be saturated at $+2Qeq$ and $-2Qeq$. When $Qdelta$ saturates, the Q bit in the ECM Frame is set.
24 The unit of Qeq , $Qoff$, and $Qdelta$ is multiples of 64 bytes.

25 If $Qoff$ is positive, i.e., the queue is above the equilibrium threshold, a ECM message containing
26 $Qoff$ and $Qdelta$ is generated. If this is not the case, a ECM message has to be generated only if
27 the CM-Tag of the sampled frame contains an RL option and the CPIDhsh field matches the
28 CPIDhsh associated with the queue. The rationale behind this is the following: If $Qoff$ is positive,
29 the queue is above the equilibrium and therefore an ECM message has to be generated anyway. In
30 the other cases, the queue is either emptying out, or it is filling up but it has not yet reached the
31 equilibrium threshold. In such cases, an ECM message has to be generated only on those flows
32 that are currently rate limited and associated with this particular CP. This check is necessary to
33 reduce as much as possible the generation of "false positive" ECM messages, i.e., positive ECM
34 messages for non rate-limited flows, or for rate limited flows associated with other CPs.

35 In certain networks Qeq may be set differently for different CPs. In order to generate ECM
36 messages carrying a normalized feedback, a scaling factor $Qscale$ is used to multiply the values of
37 $Qoff$ and $Qdelta$ copied into an ECM frame. In other words, $Qoff$ and $Qdelta$ are calculated as
38 described above and the actual ECM frame will contain $Qscale \cdot Qoff$ and $Qscale \cdot Qdelta$ ². For
39 example, a CP with a smaller buffer than other CPs may have a lower Qeq . This will result in a
40 smaller range for the $Qoff$ and $Qdelta$ values generated by such CP compared with other CPs. To
41 compensate for that, such CP will use a $Qscale$ value greater than one.

¹ From the implementation perspective such division can be safely approximated with a shift operation as, most likely, $Sscale$ is a power of 2.

² From the implementation perspective such multiplications can be safely approximated with shift operations as, most likely, $Qscale$ is a power of 2.

1 When Q_{off} and Q_{delta} are both zero, no ECM message is generated. As described earlier, the
2 message ECM(0,0) has a special meaning.
3 Every ECM frame generated by a CP will carry in the payload the CM-Tag copied from the
4 sampled frame. This information is used by the RP as described in the next section.

5 3.2.1 Congestion Point Pseudo-code

```

6
7
8 initialize()
9 {
10     queue_len = 0;           // Queue length. Incremented on packet arrival by packet length (in
11                             // pages) and decremented on packet departure.
12     samp_byte_acc = 0;      // Sampling process byte arrival accumulator.
13     queue_old = 0;         // Queue length at previous sample.
14
15     ecm_sampling_interval = ECM_FIXED_SAMPLING_INT +
16                             (ECM_RANDOM_SAMPLING_INT * (urand() * 2 - 1)); // urand(): random number
17                                     // uniformly distributed in [0,1)
18 }
19
20
21 foreach ( IncomingFrame = frame_arrival() )
22 {
23     if ( Enable_ecm_generation &&
24         IncomingFrame.Ethertype != ECM ) // Sample only frames subject to ECM
25     {
26         samp_byte_acc += len( IncomingFrame );
27         if ( samp_byte_acc > ( ecm_sampling_interval >>
28             ((queue_len > ECM_Q_MC || queue_len == 0) ? ECM_S_SCALE : 0) ))
29         {
30             /* Frame has been sampled */
31             need_gen_ecm_frame = 1; // Assume an ECM frame has to be generated
32
33             /* Setup next sampling interval */
34             samp_byte_acc = 0;
35             ecm_sampling_interval = ECM_FIXED_SAMPLING_INT +
36                                     (ECM_RANDOM_SAMPLING_INT * (urand() * 2 - 1));
37
38             if ( queue_len > ECM_Q_SC ) // ECM(0,0)
39             {
40                 ECMFrame.Qoff = 0;
41                 ECMFrame.Qdelta = 0;
42             }
43             else if ( queue_len > ECM_Q_MC ) // ECM-Max
44             {
45                 ECMFrame.Qoff = ECM_Q_EQ;
46                 ECMFrame.Qdelta = 2 * ECM_Q_EQ;
47             }
48             else // Regular ECM Frame
49             {
50                 qoff = queue_len - ECM_Q_EQ;
51                 if (qoff < -ECM_Q_EQ) qoff = -ECM_Q_EQ;
52                 ECMFrame.Qoff = qoff * ECM_Q_SCALE;
53
54                 qdelta = queue_len - queue_old;
55                 if ( qdelta > 2 * ECM_Q_EQ )
56                 {
57                     qdelta = 2 * ECM_Q_EQ;
58                     if ( ECM_qsat_enable ) ECMFrame.Q = 1; // NB: experimental feature,
59                                                         // enable needed
60                 }
61                 if ( qdelta < -2 * ECM_Q_EQ )
62                 {
63                     qdelta = -2 * ECM_Q_EQ;
64                     if ( ECM_qsat_enable ) ECMFrame.Q = 1;
65                 }
66                 ECMFrame.Qdelta = qdelta * ECM_Q_SCALE;

```

```

1
2      /* Filter out spurious feedback */
3      if ((qoff == 0 && qdelta == 0) || // No rate change
4          ((qoff < 0) && // Positive Fb and ...
5              (!has_cmtag(SampledFrame) || // ... no CM-Tag or CPID mismatch
6              (has_cmtag(SampledFrame) && SampledFrame.CMTag.CPID != CPID))))
7          {
8              Need_gen_ecm_frame = 0;
9          }
10     }
11     queue_old = queue_len;
12
13     if ( NeedGenECMFrame )
14     {
15         ECMFrame.DA = IncomingFrame.SA;
16         ECMFrame.SA = SWITCH_MAC_ADDRESS; // MAC address of the switch generating
17                                             // ECM frame
18         ECMFrame.8021QTag.Priority = HIGH_PRI; // May be priority of sampled frame
19         ECMFrame.CPID = CPID;
20         ECMFrame.Timestamp = SampledFrame.CMTag.Timestamp;
21         ECMFrame.Unit = SampledFrame.CMTag.Unit;
22         forward( ECMFrame );
23     }
24     } // if ( samp_byte_acc >= ECM_sampling_interval )
25 } // if ( Enable_ecm_generation ...
26 } // foreach()

```

27 3.3 Reaction

28 Figure 7 shows the structure of an RP data-path which may be implemented in the egress port of
29 NIC. A set of filters, F_1 through F_n , divert the traffic that matches a particular filtering criterion
30 (e.g., {VLAN, DA, SA, Priority}, {VLAN, Priority}, {Priority}, etc.) from the regular data path
31 (“No Match” in the figure) to a set of corresponding queues. Traffic is drained from such queues
32 by a set of rate limiters, R_1 through R_n , whose rate is controlled by the ECM Frames coming
33 from CP. Note that, in order to avoid out-of-order frames, the “No Match” path must not queue
34 traffic, i.e., it has absolute priority with respect to the rate-limited paths. Besides controlling the
35 rate of traffic, the rate limiters also add the CM-Tag to all the transmitted frames in order to elicit
36 feedback from the CPs they are currently associated with.

37

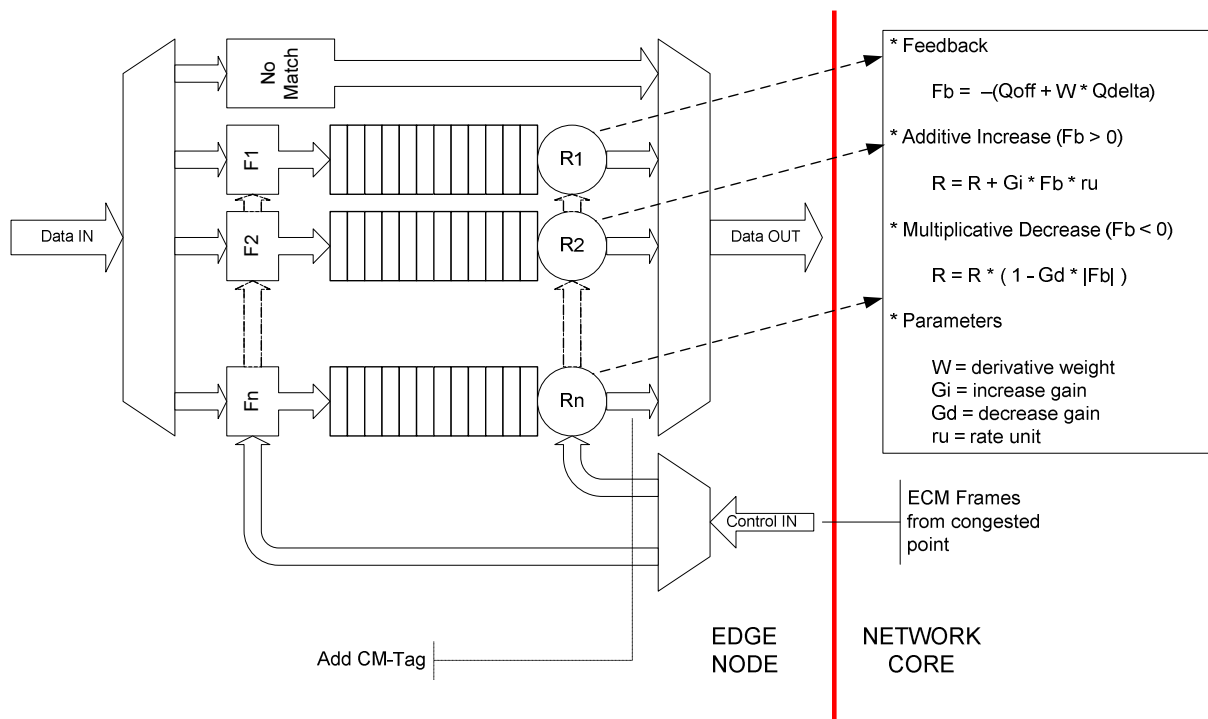


Figure 7 – Data-path structure of a Reaction Point

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Every time an ECM frame is received by an RP, the information necessary to identify a traffic flow (e.g., DA, SA, VLAN, Priority) is extracted from the header of the sampled frame carried inside the ECM frame. This information is combined in some way (for example using a hash function) to obtain a *compressed flow identifier (FlowID)*. The FlowID is compared with the FlowIDs stored in the currently active filters. If there is no match, this flow is not currently being rate limited by the RP. A filter/rate-limiter pair is instantiated and the FlowID and CPID received in the ECM frame are stored in registers associated with this pair. Also, the rate of such rate limiter is set to a configurable initial rate Ri^3 .

If there is a match, however, this particular flow is already being rate limited by the RP. The value of the feedback from the CP is then calculated as described later in this section. If the feedback is negative, the rate-limiter rate is adjusted as per feedback, and CPID from the ECM frame are stored with the filter/rate-limiter pair. If the feedback is positive, instead, the rate is adjusted if and only if the CPID of the ECM frames matches the CPID currently stored with the rate-limiter. In other words, all ECM frames carrying a negative feedback are honored, while the ECM frames carrying a positive feedback are processed only if they have been generated by the CP currently associated with the RP.

Since FlowID is a compressed flow identifier, multiple flows may be identified the same FlowID and end up in the same rate-limited queue. Different implementations may choose the degree of flow aggregation by using a different number of rate-limited queues (n) and a different hash function to obtain FlowID.

³ Usually a fraction of the link capacity C such as $\frac{1}{2}$ or $\frac{1}{4}$.

1 To ensure that (potentially) positive feedback is generated only by the CP currently associated
2 with a filter/rate-limiter, the RP adds the CM-Tag to the frames it transmits containing the CPID
3 currently stored in the rate-limiter register.

4 An active filter may change its association with a CP over time. As mentioned above, the
5 association can be changed only when a ECM frame conveying negative feedback is received
6 from a CP different from the one currently associated with the filter. For example, if a traffic flow
7 is subject to congestion at CP1 – and, therefore, is rate controlled by CP1 – starts experiencing
8 congestion at CP2, the latter will generate negative ECM frames for that flow, causing its filter to
9 change association from CP1 to CP2. After some time, the negative feedback generated by one of
10 the two Congestion Points will prevail over the other and the filter will settle its association with
11 the prevailing one.

12 Every time an RP receives an ECM frame, the *Round Trip Time* (RTT) between the RP and the
13 CP (defined as the difference between the current time and the Timestamp field carried by the
14 ECM Frame) is calculated. The RTT is then filtered through an *exponential weighted moving*
15 *average* (EWMA) as follows:

$$16 \quad \text{RTTavg} \leftarrow (1 - 2^{-W_{\text{rtt}}}) \text{RTTavg} + 2^{-W_{\text{rtt}}} * \text{RTT}$$

17
18
19 and made available to control software through a read-only register. This measure may be used to
20 dynamically adjust the value of some of the control loop parameters. The algorithm to carry out
21 such an adjustment is still being investigated.

22 At every RP there are only a limited number (n) of rate limiters available. Thus, it may happen
23 that, at a certain moment, all the rate limiters are in use and an ECM frame arrives that would
24 cause the instantiation of a new rate limiter. In such a case, a per-priority rate limiter (a.k.a.
25 *coarse rate limiter*) gets instantiated, and all the rate limiters (a.k.a. *fine rate limiters*) currently
26 associated with such a priority are set to operate at the maximum rate to force a rapid transition to
27 the per-priority rate-control. When all the queues of the individual rate-limiters are empty, they
28 may be released.

29 Once a rate limiter has been instantiated, it may be reclaimed once two conditions are satisfied:
30 (1) the queue of the rate limiter is empty, and (2) its rate is at or above line-rate. These two
31 conditions are necessary to avoid out of order packet delivery.

32 **3.3.1 Rate Control Algorithm**

33 The rate control algorithm employed by ECM works according to an *Additive Increase*
34 *Multiplicative Decrease* (AIMD) scheme loosely derived from the one employed by TCP. TCP
35 increases its rate linearly over time in absence of congestion and halves its rate every time it
36 receives negative feedback, either explicit (i.e., ECN), or implicit (i.e., packet drop). The
37 granularity of this AIMD scheme is quite coarse and it has been shown that in many cases it may
38 lead to link underutilization. In contrast, ECM employs an AIMD scheme with a much finer
39 granularity. Every time a ECM Frames arrives at a rate limiter, a Feedback signal is calculated
40 according to the following equation:

$$41 \quad Fb = -(Q_{\text{off}} + w \cdot Q_{\text{delta}})$$

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

where w is a parameter used to weight the delta component more or less with respect to the offset component. w should be configurable in the interval $[1/8, 8]$ with $1/8$ increments. Based on the sign of the Feedback signal Fb , the rate is increased or decreased as follows:

- If $Fb > 0$ $R \leftarrow R + Gi \cdot Fb \cdot Ru$
- If $Fb < 0$ $R \leftarrow R \cdot (1 - Gd \cdot |Fb|)$

where Gi and Gd are the *Increase Gain* and *Decrease Gain* respectively, and Ru is the *Rate Unit* (i.e., the granularity of the rate adjustment) employed by the rate limiters. Ru should be configurable in the range $[1, 100]$ Mbps with increments of 1 Mbps. Both Gi and Gd are fractional values. Since their granularity is unknown, they should be represented in fixed point notation with the largest number of bits possible.

Since Gd and Qeq may be chosen independently, to limit the maximum negative rate adjustment to a fraction $\alpha < 1$, the product $Gd \cdot |Fb|$ must be saturated to α . Likewise, since Gi and Qeq may be independently chosen, the maximum positive rate adjustment should be limited to a certain fraction $\beta < 1$ of the link capacity C^4 . Therefore the product $Gi \cdot Fb \cdot Ru$ must be saturated to $\beta \cdot C$. Given such constraints, the previous equations can be rewritten as:

- If $Fb > 0$ $R \leftarrow R + \min(Gi \cdot Fb \cdot Ru, \beta \cdot C)$
- If $Fb < 0$ $R \leftarrow R \cdot (1 - \min(Gd \cdot |Fb|, \alpha))$

Just like Gi and Gd , α and β are fractions and they should be represented in fixed point notation with as many bits as possible.

Besides the changes driven by feedback from CPs, the current rate of a rate limiter is also subject to a periodical *self-increase*. Every time interval T_d (e.g., 1 ms), the rate is increased by a small amount R_d (e.g., 1 Mbps). Such self-increase is useful for a number of reasons:

1. Speeds up convergence to fairness, as small flows receive substantially larger relative increments compared with large flows;
2. Allows for the reclamation of stale rate limiters. In fact, a rate limiter may stop receiving ECM frames because two main reasons: (1) the traffic stream that such rate limiter was controlling has suddenly ended, and (2) routing issues in the network prevent ECM Frames from reaching the rate limiter. When this happens, the rate limiter will remain stuck at the current rate forever. Instead, the self-increase will bring the rate-limiter rate back to line-rate and will cause its decommissioning;
3. Improves the recovery after a ECM(0,0) message is received (see next subsection).

T_d should be programmable in the range $[1 \mu s, 10 ms]$ with $1 \mu s$ increments, while R_d should be programmable in the range $[1 Mbps, 100 Mbps]$ with 1 Mbps increments.

Different (and more complex) self-increase strategies may be employed by a reaction point (see [7]). The one described in this document has been chosen for its simplicity.

⁴ C is the capacity of the link draining the rate limiter.

3.3.2 Exceptions and Non-linear Rate Adjustments

When a CP is subject to severe congestion, it may send the special ECM frame ECM(0,0) (i.e., a ECM message with $Q_{off}=0$ and $Q_{delta}=0$). When a rate limiter receives such a message, as shown in Figure 8, it sets its current rate R to 0 and starts a random timer whose range is determined by a parameter T_{max} (e.g., 10 μ s). When the timer started by the ECM(0,0) frame expires, the rate limiter is set to operate at a minimum rate R_{min} (e.g., 1/100 of line rate). This should restart the traffic flow towards the congestion point and trigger – hopefully positive – feedback. During the random timeout period the automatic self-increase of the rate is suspended, and it is resumed only after the timer expiration. Also, all ECM messages, including ECM(0,0) must be ignored during a timeout period.

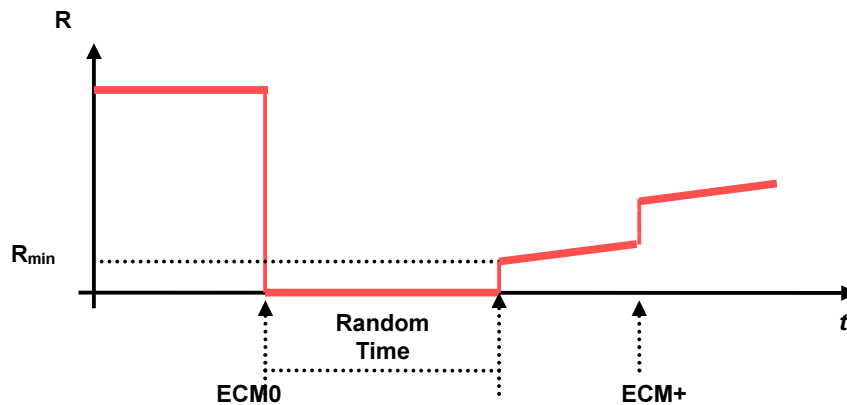


Figure 8 – Example of timeout and random restart

After the timer expiration, T_{max} is doubled and R_{min} is halved, so that the next ECM(0,0) (if any) will cause the random timer to have a longer duration and the rate limiter to restart from a slower rate, effectively realizing an exponential back-off. The initial values of T_{max} and R_{min} are restored upon the reception of the first positive feedback.

T_{max} should be configurable in the range [1 μ s, 1s] with 1 μ s increments, while R_{min} should be configurable in the range [1 Mbps, 1 Gbps] with 1 Mbps increments. The initial values of T_{max} and R_{min} are 10 μ s and 100 Mbps respectively.

The timeout with random restart has been introduced with the goal to mimic two behaviors which have been proven very successful when dealing with severe congestion caused by multiple traffic sources:

1. TCP retransmission timeout, which causes the silencing of most of the sources contributing to congestion, and
2. Ethernet CSMA/CD algorithm, which helps desynchronize traffic sources restarting after a timeout.

Special handling of ECM messages is also required when the Q is set in the ECM Frame, signaling that the Q_{delta} feedback component is saturated at $2Q_{eq}$ or $-2Q_{eq}$. When this happens, a stronger rate adjustment must be performed because the system is working outside of the linear region. The following rate adjustment is performed based on the sign of Q_{delta} :

- 1
- 2 • If $Q_{delta} < 0$ $R \leftarrow R + 2 \cdot \beta \cdot C$
- 3 • If $Q_{delta} > 0$ $R \leftarrow R \cdot (1 - \min(2 \cdot \alpha, 1))$
- 4

5 In other words, when Q_{delta} saturates a rate adjustment twice as big as the maximum rate
6 adjustment in either direction is performed. Since in the case of the decrease $2 \cdot \alpha$ may be larger
7 than 1, the resulting rate may be negative. To avoid this, the product $2 \cdot \alpha$ is saturated at 1.

8 3.3.3 Reaction Point Pseudo-code

9

```

10 initialize()
11 {
12     RL[*].state = INACTIVE;           // * = all rate limiters
13     RL[*].flowid = -1;                // no flow id
14     RL[*].rate = C;
15     RL[*].CPID = 0;
16     RL[*].RTTavg = 0;
17     RL[*].Rmin = Rmin;
18     RL[*].Tmax = Tmax;
19 }
20
21
22 processECMFrame ( ECMFrame )
23 {
24     FlowId = some_hash( ECMFrame.Payload );    // ECM frame payload contains the header of
25                                                // sampled frame
26
27     rliidx = getRateLimiterIndx( FlowID );     // Returns the index of the RL associated with
28                                                // FlowID, or the index of the next available
29                                                // RL if no FlowID match. Note that FlowID may
30                                                // be used directly as the index in the RL table
31
32     if ( ( ECMFrame.Qoff == 0 && ECMFrame.Qdelta == 0 ) &&    // ECM(0,0)
33         RL[rliidx].state != TIMEOUT )
34     {
35         RL[rliidx].state = TIMEOUT;
36         RL[rliidx].rate = 0;
37         RL[rliidx].CPID = ECMFrame.CPID;
38         Tmax_timer_set( rliidx, RL[rliidx].Tmax * urand() );
39     }
40     else
41     {
42         Fb = -(ECMFrame.Qoff + W * ECMFrame.Qdelta);
43
44         if ( Fb < 0 )
45         {
46             if ( RL[rliidx].state = INACTIVE )
47             {
48                 RL[rliidx].state = ACTIVE;
49                 RL[rliidx].flowid = FlowId;
50                 RL[rliidx].rate = Ri;
51                 RL[rliidx].CPID = ECMFrame.CPID;
52                 RL[rliidx].RTTavg = 0;
53             }
54             else
55             {
56                 RL[rliidx].rate *= 1 - ( ECMFrame.Q == 1 ? min(2·α, 1) : min(-Fb·Gd, α) );
57                 if ( RL[rliidx].rate == 0 ) RL[rliidx].rate = RL[rliidx].Rmin; // Saturate to Rmin
58                 if ( RL[rliidx].CPID == ECMFrame.CPID )
59                 {
60                     RL[rliidx].RTTavg = calc_movavg( RL[rliidx].RTTavg,           // old value
61                                                       (now() - ECMFrame.Timestamp), // new value
62                                                       2 ^ -Wrtt );                    // weight
63                 }
64             }
65         }
66     }
67 }

```



```

1         }
2     else
3     {
4         RL[rldix].CPID = ECMFrame.CPID;
5         RL[rldix].RTTavg = 0;
6     }
7 }
8 }
9 else if ( RL[rldix].CPID == ECMFrame.CPID &&
10         RL[rldix].state == ACTIVE )
11 {
12     RL[rldix].rate += ECMFrame.Q == 1 ? 2·β·C : min(Gi·Fb·Ru, β·C);
13     RL[rldix].Rmin = Rmin;
14     RL[rldix].Tmax = Tmax;
15     RL[rldix].RTTavg = calc_movavg( RL[rldix].RTTavg,           // old value
16                                     (now() - ECMFrame.Timestamp), // new value
17                                     2 ^ -WrTT );                // weight
18 }
19 }
20 }
21
22 /* Timers */
23
24 foreach ( rldix = Tmax_timeout() )
25 {
26     RL[rldix].state = ACTIVE;
27     RL[rldix].rate = RL[rldix].Rmin;
28     RL[rldix].Rmin /= 2;
29     RL[rldix].Tmax *= 2;
30 }
31
32
33 foreach ( Td_timeout() )
34 {
35     if ( RL[*].state == ACTIVE )           // * = all Rate Limiters
36     {
37         RL[*].rate += Rd;
38         if ( RL[*].rate > C ) RL[*].rate = C; // saturate @ C
39     }
40     Td_timer_set( Td );
41 }
42
43
44 /* Frame departure from RL queue */
45
46 foreach ( rldix = frame_departure_from_rl() )
47 {
48     insertCMTag( OutgoingFrame );
49     OutgoingFrame.CMTag.CPID = RL[rldix].CPID;
50     OutgoingFrame.CMTag.Timestamp = now();
51
52     if ( RL[rldix].queue_len == 0 && RL[rldix].rate == C )
53     {
54         RL[rldix].state = INACTIVE;
55         RL[rldix].flowid = -1;
56     }
57 }
58

```

59 4 Conclusions

60 This document describes ECM, a backward notification-based mechanism for congestion
61 management in data center networks. Such networks are peculiar because of their high speed, low
62 latency, and, in certain cases, zero traffic loss. In such environments traditional congestion
63 management mechanism such as RED [1] and ECN [2] have been shown not to work particularly
64 well.

1 Simulation evidence [4] [5] shows that ECM works substantially better than the above mentioned
2 alternatives, especially in data center networks where TCP and non-TCP traffic share the same
3 infrastructure. This is because traditional congestion management schemes work only when the
4 vast majority of traffic is TCP, i.e., they assume a congestion-responsive transport layer. Since
5 ECM does not make any assumption on the transport layer, it can deal even with non-responsive
6 protocols.

7 ECM has also been studied from the analytical standpoint using techniques commonly used in
8 Control Theory. It has been analytically shown that the ECM control loop is stable in a wide
9 region as determined by its parameters [6].

10 ECM has been originally presented to 802.1 in May 2005 for review and to gather feedback from
11 the standards community⁵ [4][5][6]. A Project Authorization Request (PAR) [1], along with a
12 tutorial on ECM, was presented to the IEEE 802 Plenary in July 2006. The PAR was approved
13 and a Task Force named 802.1Qau has been formed with the charter to develop a congestion
14 management framework for Ethernet. ECM is currently one of the proposals being considered by
15 the task force for such framework.

16 5 References

- 17 [1] Floyd, S., and Jacobson, V., “[Random Early Detection gateways for Congestion](#)
18 [Avoidance](#)”, V.1 N.4, August 1993, p. 397-413
- 19 [2] Ramakrishnan, K.K., Floyd, S., and Black, D.,” [The Addition of Explicit Congestion](#)
20 [Notification \(ECN\) to IP](#)”, RFC 3168
- 21 [3] [P802.1Qau](#) (C/LM) Standard for Local and Metropolitan Area Networks – Virtual
22 Bridged Local Area Networks - Amendment 10: Congestion Notification
- 23 [4] D. Bergamasco, “Presentation of Backward Congestion Notification to IEEE 802.1”,
24 [http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-backward-congestion-
26 notification-0505.pdf](http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-backward-congestion-
25 notification-0505.pdf)
- 26 [5] D. Bergamasco, “Updates on BCN For the IEEE 802 July 2005 Plenary Meeting”,
27 [http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-bcn-july-plenary-
29 0705.ppt](http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-bcn-july-plenary-
28 0705.ppt)
- 29 [6] D. Bergamasco, R. Pan, “Presentation of BCN V2.0 to IEEE 802.1 Sep 2005 Interim
30 Meeting”, [http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-bcn-
32 september-interim-rev-final-0905.ppt](http://www.ieee802.org/1/files/public/docs2005/new-bergamasco-bcn-
31 september-interim-rev-final-0905.ppt)
- 32 [7] Y. Lu, R. Pang, B. Prabhakar, D. Bergamasco, V. Alaria, A. Baldini, “[Congestion Control](#)
33 [in Networks with No Congestion Drops \(I\)](#)”, 44th Annual Allerton Conference, September
34 2006

35 6 Glossary

36 The following list describes acronyms and definitions for terms used throughout this document:

37 **AIMD**: Additive Increase Multiplicative Decrease

⁵ At that time, ECM was known as BCN, or *Backward Congestion Notification*. The name has been recently changed into ECM to avoid confusion with the generic concept of sending congestion notifications in the opposite direction of the traffic causing congestion.

- 1 **AQM:** Active Queue Management
- 2 **CMD:** Congestion Management Domain
- 3 **CM-Tag:** Congestion Management Tag
- 4 **CP:** Congestion Point
- 5 **CPID:** Congestion Point IDentifier
- 6 **CSMA/CD:** Carries Sense Multiple Access with Collision Detection
- 7 **DA:** Destination Address
- 8 **ECM:** Ethernet Congestion Management
- 9 **ECN:** Explicit Congestion Notification
- 10 **FlowID:** Flow Identifier
- 11 **NIC:** Network Interface Card
- 12 **RED:** Random Early Detection
- 13 **RL:** Rate Limiter
- 14 **RL Option:** Rate Limited Option
- 15 **RP:** Reaction Point
- 16 **RTT:** Round Trip Time
- 17 **SA:** Source Address
- 18 **TCP:** Transmission Control Protocol
- 19 **VLAN:** Virtual Local Area Network