

Draft Standard for Local and Metropolitan Area Networks— Bridges and Bridged Networks— Amendment: Congestion Isolation

Sponsor

LAN/MAN Standards Committee
of the
IEEE Computer Society

Unapproved draft

Prepared by Paul Congdon

All participants in IEEE standards development have responsibilities under the IEEE patent policy and should familiarize themselves with that policy, see <http://standards.ieee.org/about/sasb/patcom/materials.html>

As part of our IEEE 802® process, the text of the PAR (Project Authorization Request) and CSD (Criteria for Standards Development) is reviewed regularly to ensure their continued validity. A vote of “Approve” on this draft is also an affirmation that the PAR is still valid. It is included in these cover pages.

The text proper of this draft begins with the title page (1). The cover pages (a), (b), (c) etc. are for 802.1 WG information, and will be removed prior to Sponsor Ballot.

Editors' Foreword

This draft standard is the subject of an amendment project. The scope of changes to the initial standard as modified by its standardized amendments is thus strictly limited, as detailed in the [PAR](#).

Information on participation in this project, and in the IEEE 802.1 Working Group can be found [here](#).

This second draft is an individual contribution issued for Task Group discussion.

Participation in 802.1 standards development

Comments on this draft are encouraged. **NOTE: All issues related to IEEE standards presentation style, formatting, spelling, etc. are routinely handled between the 802.1 Editor and the IEEE Staff Editors prior to publication, after balloting and the process of achieving agreement on the technical content of the standard is complete.** Readers are urged to devote their valuable time and energy only to comments that materially affect either the technical content of the document or the clarity of that technical content. Comments should not simply state what is wrong, but also what might be done to fix the problem.

Full participation in the work of IEEE 802.1 requires attendance at IEEE 802 meetings. Information on 802.1 activities, working papers, and email distribution lists etc. can be found on the 802.1 Website:

<http://ieee802.org/1/>

Use of the email distribution list is not presently restricted to 802.1 members, and the working group has a policy of considering ballot comments from all who are interested and willing to contribute to the development of the draft. Individuals not attending meetings have helped to identify sources of misunderstanding and ambiguity in past projects. The email lists exist primarily to allow the members of the working group to develop standards, and are not a general forum. All contributors to the work of 802.1 should familiarize themselves with the IEEE patent policy and anyone using the mail distribution will be assumed to have done so. Information can be found at <http://standards.ieee.org/db/patents/>

Comments on this document may be sent to the 802.1 email exploder, to the Editor of this individual contribution, or to the Chairs of the 802.1 Working Group and Time Sensitive Networking Task Group.

Paul Congdon

Editor

paul.congdon@tallac.com

János Farkas

Chair, 802.1 TSN Task Group

janos.farkas@ericsson.com

Glenn Parsons

Chair, 802.1 Working Group

glenn.parsons@ericsson.com

NOTE: Comments whose distribution is restricted in any way cannot be considered, and may not be acknowledged.

All participants in IEEE standards development have responsibilities under the IEEE patent policy and should familiarize themselves with that policy, see

<http://standards.ieee.org/about/sasb/patcom/materials.html>

As part of our IEEE 802 process, the text of the PAR and CSD (Criteria for Standards Development, formerly referred to as the 5 Criteria or 5C's) is reviewed on a regular basis in order to ensure their continued validity. A vote of "Approve" on this draft is also an affirmation by the balloter that the PAR is still valid.

Project Authorization Request, Scope, Purpose, and Criteria for Standards Development (CSD)

The complete PAR, as approved by IEEE NesCom 28th September 2017, can be found at:

<https://development.standards.ieee.org/get-file/P802.1Qcz.pdf?t=98411700003>

The CSD including the "5 Criteria" that were approved by 802.1 and the 802 EC at PAR submission can be found here:

<http://www.ieee802.org/1/files/public/docs2018/cz-CSD-0718-v01.pdf>

As part of our IEEE 802 process, the text of the PAR and the CSD should be reviewed on a regular basis in order to ensure their continued validity. A vote of "Approve" on this draft is assumed also to be an affirmation by the balloter that the text of the PAR and CSD are still valid

Scope:

The scope of this standard specifies Bridges that interconnect individual LANs, each supporting the IEEE 802 MAC Service using a different or identical media access control method, to provide Bridged Networks and VLANs.

Purpose:

Bridges, as specified by this standard, allow the compatible interconnection of information technology equipment attached to separate individual LANs.

Need for the Project:

There is significant customer interest and market opportunity for large scale, low-latency, lossless Ethernet data centers to support high-performance computing and distributed storage applications. Congestion is the primary cause of loss and delay. These environments currently use higher layer end-to-end congestion control coupled with priority-based flow control at Layer 2 to avoid performance degradation from packet loss due to congestion. As the Ethernet data center network scales in size, speed and number of concurrent flows, the current environment creates head-of-line blocking for flows sharing the same traffic class. Isolating flows that cause congestion reduces latency for flows not causing congestion and improves the scale and performance of the Ethernet data center network. This amendment will support the identification and isolation of the higher layer protocol flows that are creating congestion. The amendment will interoperate with existing congestion management. Use of a consolidated Ethernet data center network will realize operational and equipment cost benefits.

1 Contents

2	1.	Overview.....	h
3	1.3	Introduction.....	h
4	2.	Definitions	i
5	4.	Abbreviations and acronyms	j
6	5.	Conformance.....	k
7	6.	Support of the MAC Service	l
8	8.	Principals of bridge operation.....	m
9	12.	Bridge management	p
10	12.1	Congestion Isolation managed objects	p
11	98.	Congestion Isolation	r
12	98.1	Congestion isolation requirements and objectives.....	s
13	98.2	Identifying congested flows.....	t
14	98.3	Registering congested flows with the 802.1CB stream identification function	t
15	98.4	Congestion Isolation Entity Operation	t
16	98.5	Congestion Isolation Protocol.....	t
17	Annex D	v
18	D.1	Requirements of the IEEE 802.1 Organizationally Specific TLV sets.....	v
19	D.2	Organizationally Specific TLV definitions.....	v
20	D.3	IEEE 802.1 Organizationally Specific TLV management.....	w
21	D.4	PICS proforma for IEEE 802.1 Organizationally Specific TLV extensions	x
22	Annex X	y
23	X.1	Queue Markers for Order Preservation.....	z

1 Figures

2 Figure 98-1 Congestion Isolation Model r
3 Figure 98-2 Congestion Isolation reference diagram u
4 Figure X-1 Out-of-order packet example y
5 Figure X-2 Using queue markers and counters to preserve order z

1 Tables

2	Table 12.1	Congestion Isolation component managed object	p
3	Table 12.2	Congestion Isolation Point component managed object.....	q
4	Table D-1	IEEE 802.1 Organizationally Specific TLVs	v

1

2 Draft Standard for
3 Local and metropolitan area networks—

4 **Bridges and Bridged Networks—**

5

6 **Amendment: Congestion Isolation**

7 [This amendment is based on IEEE Std 802.1Q™-2018]

8 NOTE—The editing instructions contained in this amendment define how to merge the material contained therein into
9 the existing base standard and its amendments to form the comprehensive standard.

10 The editing instructions are shown in *bold italic*. Four editing instructions are used: change, delete, insert, and replace.
11 *Change* is used to make corrections in existing text or tables. The editing instruction specifies the location of the change
12 and describes what is being changed by using ~~strikethrough~~ (to remove old material) and underscore (to add new
13 material). *Delete* removes existing material. *Insert* adds new material without disturbing the existing material. Deletions
14 and insertions may require renumbering. If so, renumbering instructions are given in the editing instruction. *Replace* is
15 used to make changes in figures or equations by removing the existing figure or equation and replacing it with a new
16 one. Editing instructions, change markings, and this NOTE will not be carried over into future editions because the
17 changes will be incorporated into the base standard.

18 **1. Overview**

19 **1.3 Introduction**

20 *Insert the following text after the eleventh paragraph of 1.3 and renumber accordingly*

21 This standard specifies protocols, procedures and managed objects that support the isolation of congested
22 data flows within data center environments. This is achieved by enabling systems to individually identify
23 flows creating congestion, adjust transmission selection for packets of those flows, and signal to neighbors.
24 This mechanism reduces head-of-line blocking for uncongested flows sharing a traffic class in lossless
25 networks. Congestion Isolation is intended to be used with higher layer protocols that utilize end-to-end
26 congestion control in order to reduce packet loss and latency. This amendment also addresses errors and
27 omissions in the description of existing functionality. To this purpose it:

- 28 bd) Defines a means for VLAN-aware Bridges that support congestion isolation for identifying flows
29 that are creating congestion.
- 30 be) Defines a means for adjusting transmission selection for packets of congested flows
- 31 bf) Provides for a means for discovering peer VLAN-aware Bridges and stations that support
32 congestion isolation
- 33 bg) Defines a means for signaling congestion isolation to supporting peer Bridges and stations.

34

1 2. Definitions

2 *Insert the following definitions in the appropriate collating sequence, re-numbering as*
3 *appropriate:*

4 **Congestion Isolation Aware System:** A bridge component conforming to the congestion isolation
5 provisions of this standard.

6 **Congested Flow:** A sequence of frames the end-to-end congestion controlled higher-layer protocol treats as
7 belonging to a single flow that is experiencing congestion within a Congestion Isolation Aware System.

8 **Congestion Isolation Message (CIM):** A message transmitted by a Congestion Isolation Aware System,
9 conveying Congested Flow information used by the upstream peer Congestion Isolation Aware System.

10 **Congestion Isolation Point (CIP):** A Congestion Isolation Aware System that monitors a set of queues for
11 Congested Flows and can generate Congestion Isolation Messages

12

1 **4. Abbreviations and acronyms**

2 *Insert the following acronym(s) and abbreviation(s), in the appropriate collating*
3 *sequence:*

4 CF Congested Flow

5 CI Congestion Isolation

6 CIM Congestion Isolation Message

7 CIP Congestion Isolation Point

8

1 5. Conformance

2 *Insert the following subclause after Clause 5.4.6:*

3 5.4.7 VLAN Bridge requirements for congestion isolation

4 A VLAN-aware Bridge implementation that conforms to the provisions of this standard for congestion
5 isolation (XX) shall:

- 6 a) Support, on one or more Ports, the creation of at least one Congestion Isolation Point (xx.x.x);
- 7 b) Support, at each Congestion Isolation Point, the generation of Congestion Isolation Messages (xx.x);
- 8 c) Support the ability to configure the variables controlling the operation of each Congestion Isolation
9 Point (xx.x.x);
- 10 d) Conform to the required capabilities of the LLDP of 5.2 of IEEE Std 802.1AB-2009;
- 11 e) Support the use of the Congestion Isolation TLV in LLDP (xx.x.x)

12 A VLAN Bridge implementation that conforms to the provisions of this standard for congestion isolation
13 may:

- 14 a) Support the creation of up to four CIPs on a Bridge Port (xx.x.x)
- 15 a) Support Congested Traffic Enhanced Traffic Selection (yy.y.y)
- 16 a) Support the Congestion Isolation YANG model (xx.x.x))

17

1 **6. Support of the MAC Service**

2 **6.10.1 Data Indication**

3 *Insert the following sentence after the first sentence:*

4 If the PIP is congestion isolation aware (5.4.7) and the initial octets of the mac_service_data_unit contain a
5 valid CIM encapsulation, the received frame is processed according to (xx.xx)

6

1 8. Principals of bridge operation

2 8.6.5 Flow classification and metering

3 *Replace the last sentence of the third paragraph with to following:*

4 Item e), specifying a `connection_identifier`, is only applicable to bridges that support PSFP and/or
5 congestion isolation.

6 *Insert the following clause after 8.6.5.1:*

7 8.6.5.2 Congestion Isolation flow classification

8 A Bridge or an end station may support Congestion Isolation that allows traffic class modification of
9 congested flows and subsequent frame queuing decisions (8.6.6.2) to be made on a per-stream basis for
10 received frames.

11 Support of Congestion Isolation requires implementation of the Stream identification function specified in
12 Clause 6 of IEEE Std 802.1CB-2017 [B14], as the `stream_handle` provided by this function is used to
13 identify frames from congested flows and queuing decisions taken by congestion isolation.

14 8.6.5.2.1 Congestion isolation use of the stream filter instance table

15 Congestion isolation uses a single stream filter instance from the PSFP stream filter instance table (8.6.5.1.1)
16 to select a stream gate instance that will modify the priority of congested flows. The stream filter instance
17 for congestion isolation contains the following elements:

- 18 a) A *stream filter instance identifier*. An integer value that uniquely identifies the filter instance. Since
19 there is a single filter instance for each monitored non-congested traffic class, the ordinal value of
20 this identifier is insignificant.
- 21 b) A *stream_handle specification*. A common `stream_handle` value, as specified in IEEE Std 802.1CB.
- 22 c) A priority specification. A single priority value that specifies the monitored non-congested traffic
23 class.
- 24 d) A *stream gate instance identifier*. Identifies a stream gate instance that is configured specifically for
25 congestion isolation. The stream gate allows congestion isolation to specify the internal priority
26 value of the congested traffic class for the matching stream filter instance. The state of the stream
27 gate will always be open for congestion isolation.
- 28 e) There will be no *filter specifications* used by congestion isolation.
- 29 f) A count of frames matching both the `stream_handle` and the priority specification. The other frame
30 counters specified by the stream filter instance table (8.6.5.1.1) are not used by congestion isolation.
- 31 g) The *StreamBlockedDueToOversizeFrameEnable* parameter that is set FALSE to disable this
32 function not used by congestion isolation.

33 The `stream_handle` and priority parameters associated with a received frame select the stream filter instance
34 of congestion isolation for a particular monitored non-congested traffic class. The purpose of the stream
35 filter is to select the stream gate instance that will modify the priority of congested flow frames.

36 8.6.5.2.2 Congestion isolation use of the stream gate instance table

37 Congestion isolation uses a single stream gate instance from the PSFP stream gate instance table (8.6.5.1.2)
38 to modify the priority of congested flows. The single stream gate instance for congestion isolation contains
39 the following elements:

- 1 a) A *stream gate instance identifier*. An integer value identifying the stream gate instance
- 2 b) An operational and an administrative *stream gate state* (8.6.10.4, 8.6.10.5, 12.31.3) that is set to
3 Open for congestion isolation.
- 4 c) An operational and an administrative *internal priority value* specification (IPV, 8.6.10.6, 8.6.10.7,
5 12.31.3). The IPV is used in place of the priority value associated with the frame to determine the
6 frame's traffic class, using the Traffic Class Table as specified in 8.6.6.
- 7 d) The *GateClosedDueToInvalidRxEnable* parameter that is set to FALSE to disable this function not
8 used by congestion isolation.
- 9 e) The *GateClosedDueToOctetsExceededEnable* parameter that is set to FALSE to disable this
10 function not used by congestion isolation.
- 11 f) A null *stream gate control list* since this feature is not used by congestion isolation.

12 8.6.6 Queuing Frames

13 *Replace the last paragraph with to following paragraph:*

14 In a congestion aware Bridge (Clause 30) or a congestion isolation aware Bridge (Clause XX), the act of
15 queuing a frame for transmission on a Bridge Port can result in the Forwarding Process generating a CNM or
16 a CIM. The CNM is and CIM are injected back into the Forwarding Process (8.6.1) as if it had been received
17 on that Bridge Port.

18 *Rename clause 8.6.6.1 as follows:*

19 8.6.6.1 PSFP and Congestion Isolation queuing

20 *Replace the first paragraph of 8.6.6.1 with to following paragraph:*

21 If PSFP (8.6.5.1) or Congestion Isolation (8.6.5.2) are supported, and the IPV associated with the stream
22 filter that passed the frame is anything other than the null value, then that IPV is used to determine the traffic
23 class of the frame, in place of the frame's priority, via the Traffic Class Table specified in 8.6.6. In all other
24 respects, the queuing actions specified in 8.6.6 are unchanged.

25 *Add the following new clause after 8.6.8.4:*

26 8.6.8.5 Enhancements for congestion isolation

27 A Bridge or an end station may support enhancements to isolate the frames of congested flows to a
28 designated congested traffic class. In order to meet the ordering requirements of 8.6.6, when a frame's
29 priority is modified by the congestion isolation stream gate instance (8.6.5.2.2), a transmission gate is
30 associated with each monitored non-congested queue and each congested queue. The state of the
31 transmission gate determines whether or not queued frames can be selected for transmission. For a given
32 queue, the transmission gate can be either *Open* or *Closed* as described in 8.6.8.4.

33 Congestion isolation specifies that congested and non-congested queues use the same transmission selection
34 algorithm. Additionally, the congested queues have lower priority than the monitored non-congested queues.
35 When the transmission selection algorithm is strict priority, the state of the transmission gate is permanently
36 *open*. When the transmission selection algorithm is anything other than strict priority the state of the
37 transmission gate is controlled by the state machines specified in (8.6.11)

38 *Add the following new clause after 8.6.10:*

1 **8.6.11 Congestion Isolation transmission gate state machines**

2 <<Editor's note: lots of new text, or simply punt on this detail and refer to the example in Annex X? The
3 preference would be to simply describe the transmission gate control variables, but not provide a detailed
4 state-machine that sets them.>>

5 **8.6.11.1 Isolate state machine**

6 **8.6.11.2 De-isolate state machine**

7

1 12. Bridge management

2 *Insert Congestion Isolation objects into the existing list of managed objects in 12.1.1*

- 3 m) The ability to create and delete the functional elements of congestion isolation and to control their
4 operation.

5

6 *Insert Congestion Isolation in the list of VLAN Bridge Objects of 12.2*

- 7 s) The congestion isolation entities (12.xx)

8

9 *Insert a new clause 12.xx Congestion Isolation managed objects*

10 12.1 Congestion Isolation managed objects

11 Several variables control the operation of Congestion Isolation in a congestion isolation aware bridge. The
12 managed objects are as follows:

- 13 a) CI component managed object (12.1.1)
14 b) Congestion Isolation Point (CIP) component managed object (12.1.2)

15 12.1.1 CI component managed object

16 A single instance of the CI component managed object shall be implemented by a Bridge component or end
17 station that is congestion isolation aware. It comprises all the variables included in the CI component
18 variables (XX.X) as illustrated in Table 12-1.

19

Table 12.1—Congestion Isolation component managed object

Name	Data type	Operations supported	Conformance	References
ciMasterEnable	Boolean	RW		XX.XX
ciCimTransmitPriority	Unsigned integer [0..7]	R		XX.XX

20 12.1.1 Congestion Isolation Point component managed object

21 There is one congestion isolation point (CIP) managed object for each CIP in a Bridge component or end
22 station that is congestion isolation aware. The CIP managed object comprises the some of the variables
23 included in the CI variables (XX.X) as illustrated in table 12-2.

24

25

Table 12.2—Congestion Isolation Point component managed object

Name	Data type	Operations supported	Conformance	References
cipMonitoredQueues	Unsigned integer [0..255]	RW		XX.XX
cipCongestedQueue	Unsigned integer [0..7]	RW		XX.XX
cipMinHeaderOctets	Unsigned integer	RW		XX.XX
cipTransmittedCims	counter	R		XX.XX
cipReceivedCims	counter	R		XX.XX

1 98. Congestion Isolation

2 In current data center networks, traffic can be a mix of various multi-tenant TCP and UDP flows across both
3 the physical underlay and virtual overlay network. Intermittent congestion within the network can be caused
4 by the unfortunate mix of flows across the fabric. A small number of long duration ‘elephant’ flows can
5 align in such a way to create queuing delays for the larger number of short but critical ‘mice’ flows. The
6 queuing delays deter the end-to-end congestion control loop and, in a lossless environment, cannot prevent
7 priority based flow control (PFC) from being invoked. When buffers fill and eventual flow-control kicks in,
8 mice flows can be blocked by the unfortunate burst alignment of elephant flows. If PFC is not being used,
9 packet loss on short mice flows can result in full retransmission timeouts, significantly penalizing the
10 latency of mice flows used for application control and synchronization.

11 Congestion Isolation (CI) avoids head-of-line blocking caused by the frequent-use of PFC while supporting
12 lossless behavior. CI identifies the flows that are causing congestion, isolates them to a separate lower
13 priority traffic class and signals to the upstream neighbor to do the same. CI effectively moves the congested
14 flows out of the way, temporarily delaying the delivery of congested frames, while the higher-layer
15 congestion feedback control loop has time to take effect.

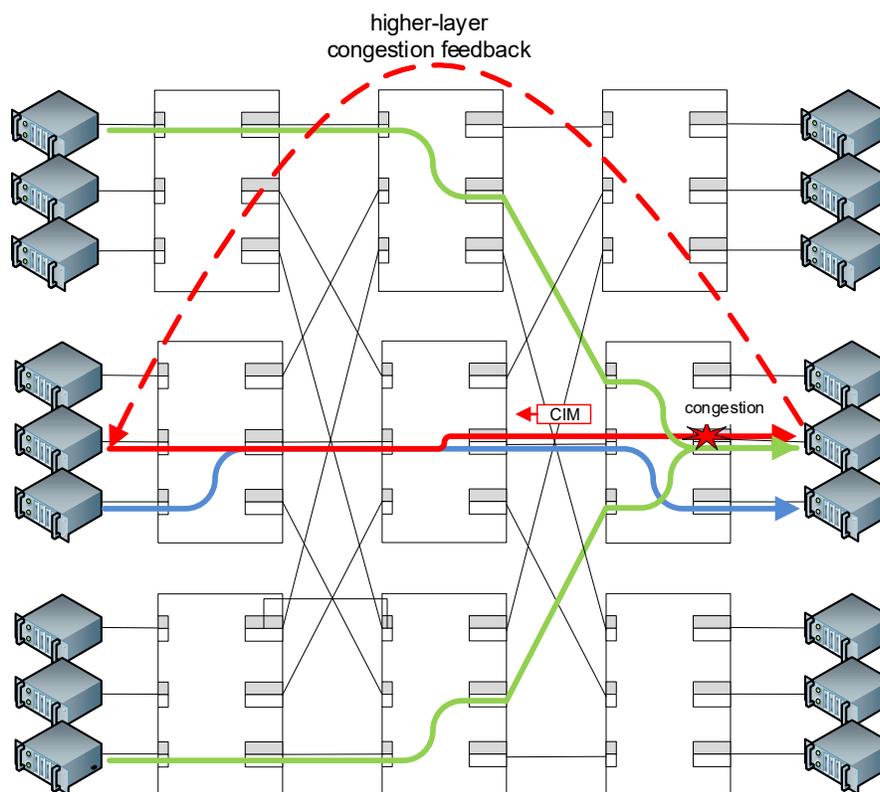


Figure 98-1—Congestion Isolation Model

16 Figure 98-1 shows the operation of CI. In the figure, server to server traffic is flowing from left to right
17 across the CLOS fabric. When flows unfortunately collide at the egress port of a bridge, congestion is
18 detected, and the offending flows are identified. Higher-layer congestion feedback is provided by the
19 destination system, but takes time to impact the rate of injection by the source system. Subsequent packets
20 from the offending flows are routed through a dedicated congested flow queue (i.e. they are effectively
21 moved out of the way). Once the congested flow queue reaches a threshold, the CI functionality signals to
22 the upstream bridge using a Congestion Isolation Message (CIM) containing flow description information

1 necessary for the upstream bridge to identify the same congested flow. The upstream bridge also isolates the
2 same flow by placing subsequent packets into the lower priority congested queue and begins to monitor the
3 depth of the congested flow queue. The packets in the congested flow queue are drained at a lower priority
4 than other non-congested queues, so when congestion persists, the congested flow queue may fill. When the
5 congested flow queue fills, the ingress port feeding that queue can issue PFC to avoid packet loss. Flow
6 control is only blocking the congested flow queues and other well-behaved mice and elephant flows are free
7 to traverse the fabric via non-congested queues

8 This clause introduces the concepts and protocols essential to congestion isolation as follows:

- 9 a) The requirements and objectives for congestion isolation (98.1)
- 10 b) Methods for identifying congested flows
- 11 c) Registering congested flows with the 802.1CB stream identification function
- 12 d) Process for modifying the priority of congested flows
- 13 e) Signaling congestion isolation messages to peers
- 14 f) Relationship and comparison with Congestion Notification

15 **98.1 Congestion isolation requirements and objectives**

16 The operation, procedures and protocols of congestion isolation are designed to meet the following
17 objectives by category:

18 — Functionality

- 19 1) With high probability, identify the flows that are causing congestion
- 20 2) Quickly adjust the traffic class of congested flows
- 21 3) Avoid head-of-line blocking of non-congested flows by signaling to upstream peers the
22 information needed to isolate the same congested flows
- 23 4) Reduce the frequency of invoking PFC in a lossless environment

24 — Compatibility

- 25 1) Work in legacy environments by automatically detecting legacy peers and automatically
26 disabling functionality
- 27 2) Work in existing lossless environments using Priority-based flow control without requiring
28 additional traffic classes
- 29 3) Work in conjunction with higher-layer end-to-end congestion control protocols (e.g Explicit
30 Congestion Notification, RDMA over Converged Ethernet, QCN)
- 31 4) Coexist with existing traffic scheduling paradigms on other traffic classes

32 — Performance

- 33 1) Reduce average flow completion times across the data center network
- 34 2) Reduce the amount of pause time when PFC is enabled
- 35 3) Reduce overall frame loss when PFC is not enabled
- 36 4) Reduce head-of-line blocking of victim flows at upstream peers from PFC
- 37 5) Reduce overall congestion control signaling
- 38 6) Increase link utilization

39 — Scale

- 40 1) Work in arbitrary data center network topologies with a mix of link speeds
- 41 2) Limit messaging overhead by restricting message propagation to hop-by-hop
- 42 3) Reduce stream identification table requirements by only requiring the registration of congested
43 flows and providing facilities to rapidly remove flows that are no longer congested.

- 1 — Implementation complexity
- 2 1) Limit the impact of existing traffic selection algorithms
- 3 2) Achieve the benefits of congestion isolation without additional buffer requirements
- 4 3) Support implementations of existing traffic classes
- 5 4) Leverage existing standard functionality for congested flow identification and stream
- 6 identification
- 7 — Manageability
- 8 1) Only require a small set of configuration parameters which are consistent across multiple
- 9 bridge deployments
- 10 2) Eliminate the ability to configure an inoperable environment
- 11 3) Provide auto discovery of peer capabilities using existing LLDP messages and without
- 12 creating additional hello and auto-configuration protocols

13 **98.2 Identifying congested flows**

14 An essential step in the process of congestion isolation is the identification of congested flows by an Active
15 Queue Management (AQM) scheme. There are many potential methods of identifying congested flows and
16 interoperable implementations can exist using different approaches. IEEE Std 802.1Q defines the CP
17 algorithm (30.2.1) for detecting congested controlled flows in congestion aware bridges. This approach may
18 be used to detect congested flows in a congestion isolation aware system. A number of other possible
19 approaches, including those that support the end-to-end ECN congestion control, are discussed in IETF RFC
20 7567.

21 Once a congested flow has been identified, it is necessary for the implementation to assert the
22 ciCongestedFlow variable and provide the initial 64 bytes of the frame immediately following the MAC
23 header. The ciCongestedFlow variable along with the contents of the frame are used to invoke the CIM send
24 statemachine (x.x.x) to cause the transmission of a CIM message to the upstream congestion isolation peer.

25 **98.3 Registering congested flows with the 802.1CB stream identification function**

26

27 **98.4 Congestion Isolation Entity Operation**

28 << Editor's Note: The text should cover the following:

- 29 a) Congestion Isolation aware Bridge Forwarding Process diagram (see next slide)
- 30 b) Congestion Isolation Point (CIP)
- 31 c) Congestion Isolation Input Multiplexer - how CIMs are decoded and received.
- 32 d) Congested Flow Identification and Table

33 >>

34

35 **98.5 Congestion Isolation Protocol**

36 << Editor's Note - This section needs to include:

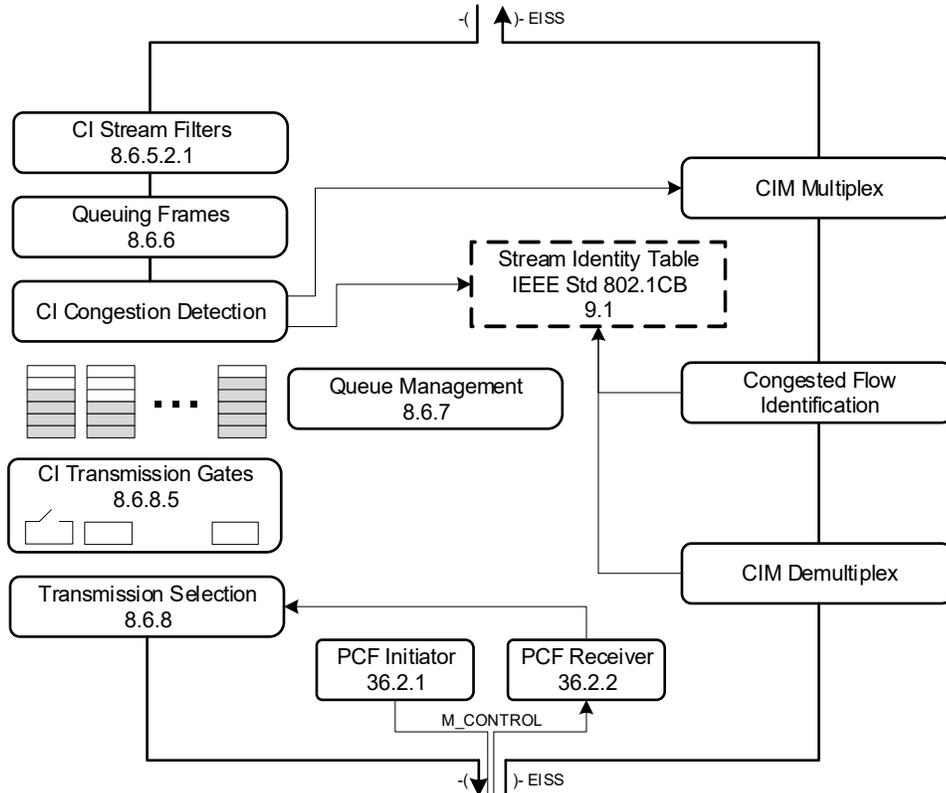


Figure 98-2—Congestion Isolation reference diagram

- 1 a) Variables controlling operation
- 2 b) State machine and associated variables and procedures that control the CIM and congested flow table
- 3
- 4 c) The Congestion Isolation Protocol
- 5 1)Generating CIMs - State Machines
- 6 2)Creating and Deleting entries in the Congested Flow Table
- 7 3)Queuing frames via EM_UNITDATA.request
- 8 4)Processing received CIMs
- 9 e) Encoding of CIM PDUs
- 10 f) Congestion Isolation LLDP TLV definition
- 11 g) UML and YANG model

12 >>

13 **98.5.1 Congestion Isolation LLLDP TLV**

14

1 Annex D

2 (normative)

3 IEEE 802.1 Organizationally Specific TLVs

4 D.1 Requirements of the IEEE 802.1 Organizationally Specific TLV sets

5 *Add the following row to Table D-1*

Table D-1—IEEE 802.1 Organizationally Specific TLVs

IEEE 802.1 subtype	TLV name	TLV set name	TLV reference	Feature clause reference
xx	Congestion Isolation	ciSet	D.2.16	Clause 98

6 << Editor's Note: the value xx will be assigned at sponsor ballot >>

7 D.2 Organizationally Specific TLV definitions

8 *Add the following clauses to the end of D.2*

9 D.2.16 Congestion Isolation TLV

10 The Congestion Isolation TLV is an optional TLV that allows an IEEE 802.1Q-compliant bridge and an
11 IEEE 802.1Q-compatible IEEE 802 LAN station to discover each other and exchange configuration
12 information.

13 Figure D-16 shows the VLAN Name TLV format.

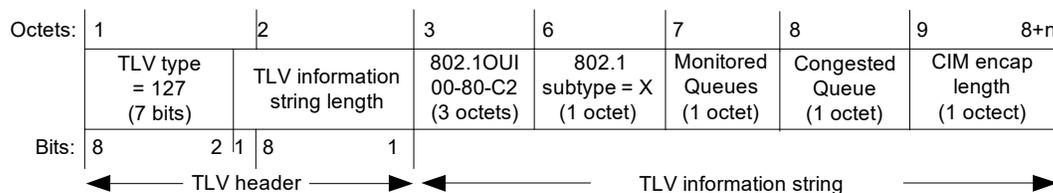


Figure D-16—VLAN Name TLV format

14 << Editor's Note: the subtype value of X will be assigned at sponsor ballot >>

15 D.2.16.1 TLV type

16 A 7-bit integer value occupying the most-significant bits of the first octet of the TLV. Always contains the
17 value 127

1 **D.2.16.2 TLV information string length**

2 The TLV information string length field of the Congestion Isolation TLV is fixed and shall contain the value
3 7.

4 **D.2.16.3 Monitored queues**

5 A bit vector, one per priority value, containing all eight of the traffic classes support by the Bridge or end
6 station. The LSB of the octet carries priority 0, and the MSB is that of priority 7.

7 **D.2.16.4 Congested queue**

8 The numerical value in the range of 0 to 7 that represents the traffic class that will act as the congested
9 queue. The congested queue is lower priority than the monitored queues.

10 **D.2.16.5 CIM encap length**

11 A single octet unsigned integer representing the requested length in octets of the data from the frame of a
12 congested flow to be encapsulated into a CIM message by a peer. The default value is 64.

13 **D.2.16.6 Congestion Isolation TLV usage rules**

14 The priority of the congested queue shall be lower than the priority of all monitored queues.

15 **D.3 IEEE 802.1 Organizationally Specific TLV management**

16 **D.3.2 IEEE 802.1 managed objects—TLV variables**

17 *Add the following clauses to the end of D.3.2*

18 **D.3.2.11 Congestion Isolation TLV managed objects**

- 19 a) **monitored queues:** see D.2.16.3.
20 b) **congested queue:** see D.2.16.4.
21 c) **CIM encap length:** see D.2.16.5.

1 D.4 PICS proforma for IEEE 802.1 Organizationally Specific TLV extensions

2 D.4.3 Major capabilities and options

3 Add the following to the end of D.4.3

Item	Feature	Status	References	Support
ciSet	Is the IEEE 802.1 Organizationally Specific TLV ciSet implemented?	O.1	D.1, Table D-1	Yes [] No []
ciQueuePri	Are the monitored queues higher priority than the congested queue?	ciSet:M	D.2.16.6	Yes []

1 Annex X

2 (informative)

3 Maintaining Packet Order with Congestion Isolation

4 The process of congestion isolation involves identifying the packets of a congested flow and subsequently
5 modifying the egress traffic class of those packets based on the level of congestion in the monitored and
6 congested queues. During this process, it is possible that packets for the same flow can exist in multiple
7 queues at the same time, resulting in the possibility of an out-of-order packet delivery. As an illustration,
8 consider the following example depicted in Figure X-1.

9

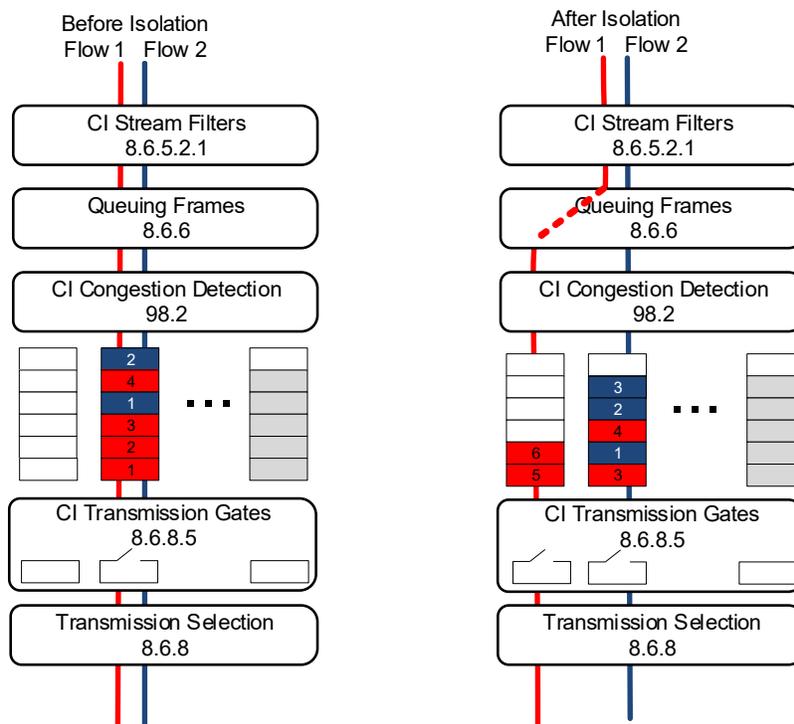


Figure X-1—Out-of-order packet example

10 In the example in Figure X-1, the packets of two flows, red and blue, are intermixed while traversing a
11 common monitored queue. As the monitored queue fills, the red flow is determined to be a congested flow
12 and subsequent packets of that flow will be reclassified and queued in the congested queue. Previously
13 received packets for that flow, numbered 1 through 6, may reside in the monitored queue. Since the
14 congested queue is empty, the subsequent packets, numbered 7 and 8, are placed at the head of the congested
15 queue. Depending upon the traffic selection algorithm, it may be possible for packets 7 and 8 to be selected
16 for transmission before some of packets 1 through 6.

17 The priority of the congested queue is intended to be lower than the priority of the monitored queue. The
18 strict priority traffic selection algorithm (8.6.8.1) will assure that no out-of-order packet delivery occurs,
19 however there is a risk of starvation for congested flows and alternative traffic selection algorithms may be
20 desired. The enhanced transmission selection algorithm (8.6.8.3) or other vendor specific algorithms may
21 not assure out-of-order delivery on their own.

1 Congestion isolation defines a transmission gate for the monitored and congested queues that make those
 2 queues available to the transmission selection algorithm. When the transmission selection algorithm is strict
 3 priority, the transmission gate is permanently open. The transmission gate is controlled by the ciGateControl
 4 variable for other transmission selection algorithms that can not assure out-of-order delivery. The
 5 management of ciGateControl is implementation dependent, but must be asserted in a way to assure the
 6 externally visible behavior of the bridge supporting congestion isolation is to maintain packet order.

7 The following informative description provides an example mechanism to preserve packet order for
 8 transmission selection algorithms other than strict priority.

9 X.1 Queue Markers for Order Preservation

10 The mechanism below provides control of the ciGateControl variable for a congested queue in order to
 11 preserve the order of packets for a congested flow. It involves a queue position marker and marker counter
 12 for both the monitored and congested queues. The mechanism is described using the example shown in
 13 Figure X-2.

14

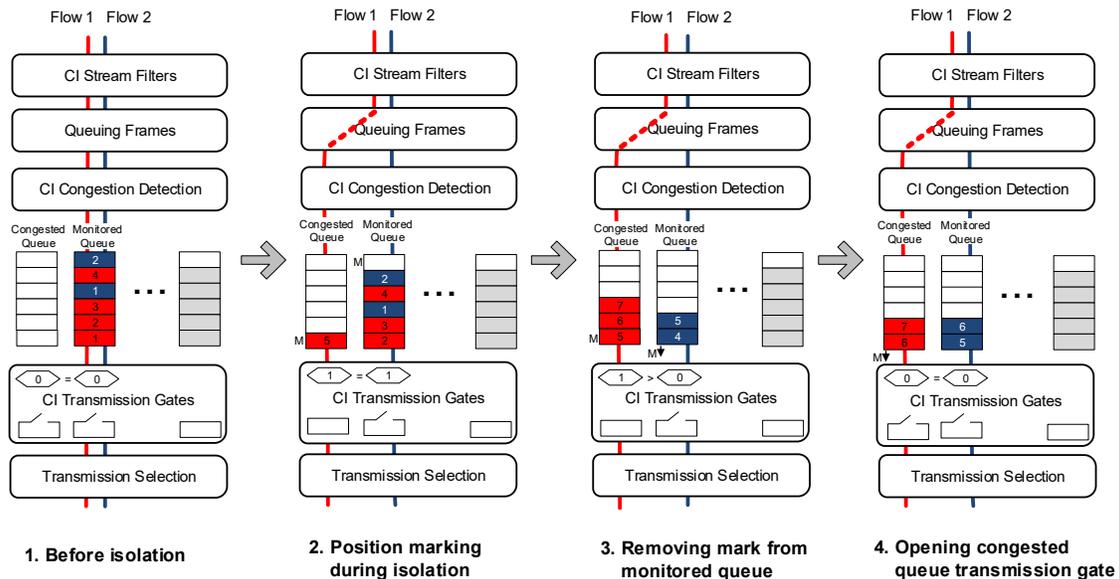


Figure X-2—Using queue markers and counters to preserve order

15 The example in Figure X-2 shows the state of the position markers and marker counters during four different
 16 phases of congestion isolation operation; before a flow is isolated, position marking during the isolation of a
 17 flow, the closing of the congested queue transmission gate and the opening of the congested queue
 18 transmission gate. When a queue transmission gate is open, that queue is available to the traffic selection
 19 algorithm. When it is closed, the queue is not available to the traffic selection algorithm.

20 In the example, the packets of two flows are intermixed in the monitored queue as they traverse the switch.
 21 Since no flows have been isolated yet, the position marker counters of both the congested queue and the
 22 monitored queue are set to 0. Once a flow has been isolated and subsequent packets of that flow are placed
 23 in the congested queue, a marker is placed in both the congested queue and the monitored queue. The marker
 24 counter is incremented for both queues once the position is marked. The congested queue is empty at the
 25 time the first flow is isolated, so the marker will be at the head of the queue. When a marker is at the head of
 26 the congested queue and the marker counters are equal, the ciGateControl variable for the congested queue
 27 is set to closed. The monitored queue continues to drain and eventually the position marker will reach the

1 head of the monitored queue. The marker counter for the monitored queue will be decremented when the
2 packet associated with the position marker is scheduled for transmission and exits the queue. When the
3 value of the congested queue marker counter is greater than the monitored queue position counter, it is
4 possible to set the ciGateControl variable to open and begin to schedule the congested queue. When ever a
5 packet that aligns with a position marker is scheduled for transmission, the associated marker counter is
6 decremented.

7

8