

Draft Standard for Local and metropolitan area networks—

Bridges and Bridged Networks

Amendment: Enhancements to Cyclic Queuing and Forwarding

Developed by the

LAN/MAN Standards Committee

of the

IEEE Computer Society

Unapproved draft

Prepared by the Security Task Group of IEEE 802.1

This and the following cover pages are not part of the draft. They provide revision and other information for IEEE 802.1 Working Group members and participants in the IEEE Standards Association ballot process, and will be updated as convenient. New participants: Please read these cover pages, they contain information that should help you contribute effectively to this standards development project.

The text proper of this draft begins with the [Title page](#).

Important Notice

This document is an unapproved draft of a proposed IEEE Standard. IEEE hereby grants the named IEEE SA Working Group or Standards Committee Chair permission to distribute this document to participants in the receiving IEEE SA Working Group or Standards Committee, for purposes of review for IEEE standardization activities. No further use, reproduction, or distribution of this document is permitted without the express written permission of IEEE Standards Association (IEEE SA). Prior to any review or use of this draft standard, in part or in whole, by another standards development organization, permission must first be obtained from IEEE SA (stds-copyright@ieee.org). This page is included as the cover of this draft, and shall not be modified or deleted.

IEEE Standards Association
445 Hoes Lane
Piscataway, NJ 08854, USA

1 Editors' Foreword

2 Throughout this document any notes with a Cyan background (as in this paragraph) are temporary, inserted
3 by the Editors for a variety of purposes. They will be removed prior to SA Ballot and are not part of the
4 normative text. To avoid gratuitous changes to page numbering between final WG ballot and initial SA Ballot,
5 the number of cover pages should remain unchanged, with pages intentionally left blank marked as such. The
6 SA ballot cover pages may contain information for the stages of that ballot (including, if appropriate,
7 cross-references to text changed in the course of SA balloting). The records of participants in the
8 development of the standard will be added pre-publication.

9 This draft is a proposed amendment to an approved standard. All that it has to show are the proposed
10 changes (including additions) to the standard that it amends. However experience has shown that the
11 development of an amendment that includes the minimum amount of text needed to meet this goal is
12 undesirable. First, such a minimal amendment hands the task of combining the amended text with the base
13 standard not just to an editor rolling up the base text and outstanding amendments into a new edition, but also
14 to everyone who wants to use the standard before that rolled up edition is available, which might be ten years
15 in the future. Second, few if any reviewers have the time to mentally undertake that roll-up process when
16 reviewing each successive draft. Much of the base text can remain out of sight and out of mind, with the
17 consequence that a developed amendment may add material that does not take advantage of material
18 already in the approved, duplicate that material, or even contradict it. If the changes consist of many small
19 fragments, the result may prove barely readable when the merge is done. Accordingly this amendment may
20 contain more of the base text than may appear strictly necessary. The eventual aim is to include sufficient text
21 to make the context of the additions clear without repeated reference to the base text, thus making the
22 intended use of the amendment easier. In early drafts more material can be included, with the aim of making
23 sure that all the text that needs to be reviewed or appreciated when contributing to draft development is
24 readily available to reviewers. There is a known drawback to including this additional text. Commenters tend
25 to assume that any text shown can be amended. Only new text introduced by an *Insert* editing instruction can
26 be freely changed. Where base text is included as part of a *Change* editing instruction, changes are restricted
27 to those that are within the Scope of the project (refer to the PAR).

28 Participation in 802.1 standards development

29 All participants in the standardization activities of IEEE 802.1 should be aware of the Working Group Policies
30 and Procedures, and the fact that they have obligations under the IEEE Patent Policy, the IEEE Standards
31 Association (SA) Copyright Policy, and the IEEE SA Participation Policy. For information on these policies see
32 1.ieee802.org/rules/ and the slides presented at the beginning of each of our Working Group and Task Group
33 meeting.

34 As part of our IEEE 802® process, the text of the PAR (Project Authorization Request) and CSD (Criteria for
35 Standards Development) of each project is reviewed regularly to ensure their continued validity. The PAR is
36 summarized in these cover pages and a links are provided to the full text of both PAR and CSD. A vote of
37 "Approve" on this draft is also an affirmation that the PAR and CSD for this project are still valid.

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

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

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

47 Use of the email distribution list is not presently restricted to 802.1 members, and the working group has a
48 policy of considering comments from all who are interested and willing to contribute to the development of the
49 draft. Individuals not attending meetings have helped to identify sources of misunderstanding and ambiguity
50 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 email distribution list will be assumed to have done so. Information can be found at <http://standards.ieee.org/db/patents/>. Comments on this draft may be sent to the 802.1 email exploder, to the Editor, or to the Chairs of the 802.1 Working Group and TSNTask Group.

Norm Finn
Editor, P802.1Qdv
Email: nfinn@nfinnconsulting.com

Janos Farkas
Chair, 802.1 TSN Task Group
Email: Janos.Farkas@ericsson.com

Glenn Parsons
Chair, 802.1 Working Group
+1 514-379-9037
Email: 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.

Draft development

During the early stages of draft development, 802.1 editors have a responsibility to attempt to craft technically coherent drafts from the resolutions of ballot comments and from the other discussions that take place in the working group meetings. Preparation of drafts often exposes inconsistencies in editor's instructions or exposes the need to make choices between approaches that were not fully apparent in the meeting. Choices and requests by the editors' for contributions on specific issues will be found in the editors' [Introduction to the current draft](#) and at appropriate points in the draft.

The ballot comments received on each draft, and the editors' proposed and final disposition of comments on working group drafts, are part of the audit trail of the development of the standard and are available, along with all the revisions of the draft on the 802.1 website (for address see above).

During the early stages of draft development the proposed text can be moved around a great deal, and even minor rearrangement can lead to a lot of 'change', not all of which is noteworthy from the point of the reviewer, so the use of automatic change bars is not very effective. In early drafts change bars may be omitted or applied manually, with a view to drawing the readers attention to the most significant areas of change. Readers interested in viewing every change are encouraged to use Adobe Acrobat to compare the document with their selected prior draft. Note that the FrameMaker change bar feature is useless when it comes to indicating changes to Figures.

1 Project Authorization Request, Scope, Purpose, and Criteria for Standards 2 Development (CSD)

3 The complete PAR, as approved by IEEE NesCom 21st September 2022, can be found at:

4 <https://development.standards.ieee.org/myproject-web/public/view.html#pardetail/10027>

5 and the CSD (Criteria for Standards Development) at:

6 <https://mentor.ieee.org/802-ec/dcn/22/ec-22-0083-00-ACSD-p802-1qdt.pdf>

7 extracts of relevant material from the PAR and CSD follow.

8 PAR Scope, Purpose, and Need

9 The Scope of the standard (IEEE Std 802.1Q) as amended by this project remains unchanged and is shown
10 below. The Purpose (clause 1.3) of IEEE Sd 802.1Q is not changed by this project.

11 Scope:

12 This standard specifies Bridges that interconnect individual LANs, each supporting the IEEE 802 MAC
13 Service using a different or identical media access control method, to provide Bridged Networks and
14 VLANs.

15 Scope of the Project:

16 This amendment specifies procedures, protocols and managed objects to enhance Cyclic Queuing and
17 Forwarding, comprising: a transmission selection procedure that organizes frames in a traffic class output
18 queue into logical bins that are output in strict rotation at a constant frequency; a procedure for storing
19 received frames into bins based on the time of reception of the frame; a procedure for storing received
20 frames into bins based on per-flow octet counters; a protocol for determining the phase relationship between
21 a transmitter's and a receiver's bin boundaries in time; managed objects, Management Information Base
22 (MIB), and YANG modules for controlling these procedures; and an informative annex to provide guidance
23 for applying these procedures. This amendment also addresses errors and omissions in the description of
24 existing IEEE Std 802.1Q functionality.

25 Purpose:

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

28 Need for the Project:

29 The existing Cyclic Queuing and Forwarding (CQF) functionality in IEEE Std 802.1Q provides bounded
30 end-to-end delays, allows simple delay analysis methods, and does not depend on per-flow state. These
31 properties are critical for scaling up Time-Sensitive Networking to large networks with a high number of
32 simultaneous flows, such as service provider networks. This amendment extends the existing CQF
33 functionality to support long physical links with high delay, processing delay variations in bridges,
34 non-time-synchronized ingress traffic, and/or flows with a wider range of latency requirements. These
35 properties enhance the suitability of CQF for large networks.

36 CSD broad market potential

37 The features of this standard broaden the applicability of Time-Sensitive Networking (TSN) to networks
38 with simpler bridges than are possible with the existing, deployed TSN features, and to service provider
39 networks, a large market so far untapped by TSN.

40 The interest expressed by vendors and users in IEEE 802.1 indicates that sufficient interest will exist outside
41 IEEE 802.1 for this standard to succeed.CSD compatability

1 CSD technical feasibility

2 The existing Asynchronous Traffic Shaping and Cyclic Queuing and Forwarding provisions of IEEE Std
3 802.1Q bound, on either side, the complexity of this standard. Both are deployed, indicating the feasibility
4 of this standard.

¹ **Introduction to the current draft**

² This is an initial draft of P802.1Qdv.

Draft Standard for Local and metropolitan area networks—

Bridges and Bridged Networks

Amendment: Enhancements to Cyclic Queuing and Forwarding

Unapproved draft, prepared by the
Time-Sensitive Networking (TSN) Task Group of IEEE 802.1

Sponsored by the
LAN/MAN Standards Committee
of the
IEEE Computer Society

Copyright ©2023 by the IEEE.
3 Park Avenue
New York, NY 10016-5997
USA

All rights reserved.

This document is an unapproved draft of a proposed IEEE Standard. As such, this document is subject to change. USE AT YOUR OWN RISK! IEEE copyright statements SHALL NOT BE REMOVED from draft or approved IEEE standards, or modified in any way. Because this is an unapproved draft, this document must not be utilized for any conformance/compliance purposes. Permission is hereby granted for officers from each IEEE Standards Working Group or Committee to reproduce the draft document developed by that Working Group for purposes of international standardization consideration. IEEE Standards Department must be informed of the submission for consideration prior to any reproduction for international standardization consideration (stds.ipr@ieee.org). Prior to adoption of this document, in whole or in part, by another standards development organization, permission must first be obtained from the IEEE Standards Department (stds.ipr@ieee.org). When requesting permission, IEEE Standards Department will require a copy of the standard development organization's document highlighting the use of IEEE content. Other entities seeking permission to reproduce this document, in whole or in part, must also obtain permission from the IEEE Standards Department.

IEEE Standards Activities Department
445 Hoes Lane
Piscataway, NJ 08854, USA

1 **Abstract:** This amendment enhances Cyclic Queuing and Forwarding. It specifies a transmission
2 selection procedure that organizes frames in a traffic class output queue into logical bins that are
3 output in strict rotation at a constant frequency; a procedure for storing received frames into bins
4 based on the time of reception of the frame; a procedure for storing received frames into bins based
5 on per-flow octet counters; and protocol for determining the phase relationship between a
6 transmitter's and a receiver's bin boundaries.

7 **Keywords:** CQF, Cyclic Queuing and Forwarding, IEEE 802.1Q™, LAN, local area network, Time-
8 Sensitive Networking, TSN, Virtual Bridged Network, virtual LAN, VLAN Bridge

The Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA

Copyright © 2023 by the Institute of Electrical and Electronics Engineers, Inc.
All rights reserved. Published dd month year. Printed in the United States of America.

IEEE and 802 are registered trademarks in the U.S. Patent & Trademark Office, owned by the Institute of Electrical and Electronics Engineers, Incorporated.

Print: ISBN 978-X-XXX-XXX-X STDXXXXX
PDF: ISBN 978-X-XXX-XXX-X STDPDXXXXX

IEEE prohibits discrimination, harassment, and bullying.

For more information, visit <http://www.ieee.org/web/aboutus/whatis/policies/p9-26.html>.

No part of this publication may be reproduced in any form, in an electronic retrieval system or otherwise, without the prior written permission of the publisher.

1 Important Notices and Disclaimers Concerning IEEE Standards Documents

2 IEEE Standards documents are made available for use subject to important notices and legal disclaimers.
3 These notices and disclaimers, or a reference to this page (<https://standards.ieee.org/ipr/disclaimers.html>),
4 appear in all standards and may be found under the heading “Important Notices and Disclaimers Concerning
5 IEEE Standards Documents.”

6 Notice and Disclaimer of Liability Concerning the Use of IEEE Standards 7 Documents

8 IEEE Standards documents are developed within the IEEE Societies and the Standards Coordinating
9 Committees of the IEEE Standards Association (IEEE SA) Standards Board. IEEE develops its standards
10 through an accredited consensus development process, which brings together volunteers representing varied
11 viewpoints and interests to achieve the final product. IEEE Standards are documents developed by
12 volunteers with scientific, academic, and industry-based expertise in technical working groups. Volunteers
13 are not necessarily members of IEEE or IEEE SA, and participate without compensation from IEEE. While
14 IEEE administers the process and establishes rules to promote fairness in the consensus development
15 process, IEEE does not independently evaluate, test, or verify the accuracy of any of the information or the
16 soundness of any judgments contained in its standards.

17 IEEE makes no warranties or representations concerning its standards, and expressly disclaims all
18 warranties, express or implied, concerning this standard, including but not limited to the warranties of
19 merchantability, fitness for a particular purpose and non-infringement. In addition, IEEE does not warrant
20 or represent that the use of the material contained in its standards is free from patent infringement. IEEE
21 standards documents are supplied “AS IS” and “WITH ALL FAULTS.”

22 Use of an IEEE standard is wholly voluntary. The existence of an IEEE Standard does not imply that there
23 are no other ways to produce, test, measure, purchase, market, or provide other goods and services related to
24 the scope of the IEEE standard. Furthermore, the viewpoint expressed at the time a standard is approved and
25 issued is subject to change brought about through developments in the state of the art and comments
26 received from users of the standard.

27 In publishing and making its standards available, IEEE is not suggesting or rendering professional or other
28 services for, or on behalf of, any person or entity, nor is IEEE undertaking to perform any duty owed by any
29 other person or entity to another. Any person utilizing any IEEE Standards document, should rely upon his
30 or her own independent judgment in the exercise of reasonable care in any given circumstances or, as
31 appropriate, seek the advice of a competent professional in determining the appropriateness of a given IEEE
32 standard.

33 IN NO EVENT SHALL IEEE BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
34 EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO: THE
35 NEED TO PROCURE SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR
36 BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY,
37 WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR
38 OTHERWISE) ARISING IN ANY WAY OUT OF THE PUBLICATION, USE OF, OR RELIANCE UPON
39 ANY STANDARD, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE AND
40 REGARDLESS OF WHETHER SUCH DAMAGE WAS FORESEEABLE.

1 Translations

2 The IEEE consensus development process involves the review of documents in English only. In the event
3 that an IEEE standard is translated, only the English version published by IEEE is the approved IEEE
4 standard.

5 Official statements

6 A statement, written or oral, that is not processed in accordance with the IEEE SA Standards Board
7 Operations Manual shall not be considered or inferred to be the official position of IEEE or any of its
8 committees and shall not be considered to be, nor be relied upon as, a formal position of IEEE. At lectures,
9 symposia, seminars, or educational courses, an individual presenting information on IEEE standards shall
10 make it clear that the presenter's views should be considered the personal views of that individual rather than
11 the formal position of IEEE, IEEE SA, the Standards Committee, or the Working Group.

12 Comments on standards

13 Comments for revision of IEEE Standards documents are welcome from any interested party, regardless of
14 membership affiliation with IEEE or IEEE SA. However, **IEEE does not provide interpretations,**
15 **consulting information, or advice pertaining to IEEE Standards documents.**

16 Suggestions for changes in documents should be in the form of a proposed change of text, together with
17 appropriate supporting comments. Since IEEE standards represent a consensus of concerned interests, it is
18 important that any responses to comments and questions also receive the concurrence of a balance of
19 interests. For this reason, IEEE and the members of its Societies and Standards Coordinating Committees
20 are not able to provide an instant response to comments, or questions except in those cases where the matter
21 has previously been addressed. For the same reason, IEEE does not respond to interpretation requests. Any
22 person who would like to participate in evaluating comments or in revisions to an IEEE standard is welcome
23 to join the relevant IEEE working group. You can indicate interest in a working group using the Interests tab
24 in the Manage Profile & Interests area of the [IEEE SA myProject system](#).¹ An IEEE Account is needed to
25 access the application.

26 Comments on standards should be submitted using the [Contact Us](#) form.²

27 Laws and regulations

28 Users of IEEE Standards documents should consult all applicable laws and regulations. Compliance with the
29 provisions of any IEEE Standards document does not imply compliance to any applicable regulatory
30 requirements. Implementers of the standard are responsible for observing or referring to the applicable
31 regulatory requirements. IEEE does not, by the publication of its standards, intend to urge action that is not
32 in compliance with applicable laws, and these documents may not be construed as doing so.

33 Data privacy

34 Users of IEEE Standards documents should evaluate the standards for considerations of data privacy and
35 data ownership in the context of assessing and using the standards in compliance with applicable laws and
36 regulations.

1. Available at: <https://development.standards.ieee.org/myproject-web/public/view.html#landing>.

2. Available at: <https://standards.ieee.org/content/ieee-standards/en/about/contact/index.html>.

1 Copyrights

2 IEEE draft and approved standards are copyrighted by IEEE under U.S. and international copyright laws.
3 They are made available by IEEE and are adopted for a wide variety of both public and private uses. These
4 include both use, by reference, in laws and regulations, and use in private self-regulation, standardization,
5 and the promotion of engineering practices and methods. By making these documents available for use and
6 adoption by public authorities and private users, IEEE does not waive any rights in copyright to the
7 documents.

8 Photocopies

9 Subject to payment of the appropriate fee, IEEE will grant users a limited, non-exclusive license to
10 photocopy portions of any individual standard for company or organizational internal use or individual, non-
11 commercial use only. To arrange for payment of licensing fees, please contact Copyright Clearance Center,
12 Customer Service, 222 Rosewood Drive, Danvers, MA 01923 USA; +1 978 750 8400. Permission to
13 photocopy portions of any individual standard for educational classroom use can also be obtained through
14 the Copyright Clearance Center.

15 Updating of IEEE Standards documents

16 Users of IEEE Standards documents should be aware that these documents may be superseded at any time
17 by the issuance of new editions or may be amended from time to time through the issuance of amendments,
18 corrigenda, or errata. An official IEEE document at any point in time consists of the current edition of the
19 document together with any amendments, corrigenda, or errata then in effect.

20 Every IEEE standard is subjected to review at least every ten years. When a document is more than ten years
21 old and has not undergone a revision process, it is reasonable to conclude that its contents, although still of
22 some value, do not wholly reflect the present state of the art. Users are cautioned to check to determine that
23 they have the latest edition of any IEEE standard.

24 In order to determine whether a given document is the current edition and whether it has been amended
25 through the issuance of amendments, corrigenda, or errata, visit [IEEE Xplore](#) or [contact IEEE](#).³ For more
26 information about the IEEE SA or IEEE's standards development process, visit the IEEE SA Website.

27 Errata

28 Errata, if any, for all IEEE standards can be accessed on the [IEEE SA Website](#).⁴ Search for standard number
29 and year of approval to access the web page of the published standard. Errata links are located under the
30 Additional Resources Details section. Errata are also available in [IEEE Xplore](#). Users are encouraged to
31 periodically check for errata.

32 Patents

33 IEEE Standards are developed in compliance with the [IEEE SA Patent Policy](#).⁵

34 Attention is called to the possibility that implementation of this standard may require use of subject matter
35 covered by patent rights. By publication of this standard, no position is taken by the IEEE with respect to the
36 existence or validity of any patent rights in connection therewith. If a patent holder or patent applicant has
37 filed a statement of assurance via an Accepted Letter of Assurance, then the statement is listed on the

3. Available at: <https://ieeexplore.ieee.org/browse/standards/collection/ieee>.

4. Available at: <https://standards.ieee.org/standard/index.html>.

5. Available at: <https://standards.ieee.org/about/sasb/patcom/materials.html>.

1 IEEE SA Website at <http://standards.ieee.org/about/sasb/patcom/patents.html>. Letters of Assurance may
2 indicate whether the Submitter is willing or unwilling to grant licenses under patent rights without
3 compensation or under reasonable rates, with reasonable terms and conditions that are demonstrably free of
4 any unfair discrimination to applicants desiring to obtain such licenses.

5 Essential Patent Claims may exist for which a Letter of Assurance has not been received. The IEEE is not
6 responsible for identifying Essential Patent Claims for which a license may be required, for conducting
7 inquiries into the legal validity or scope of Patents Claims, or determining whether any licensing terms or
8 conditions provided in connection with submission of a Letter of Assurance, if any, or in any licensing
9 agreements are reasonable or non-discriminatory. Users of this standard are expressly advised that
10 determination of the validity of any patent rights, and the risk of infringement of such rights, is entirely their
11 own responsibility. Further information may be obtained from the IEEE Standards Association.

12 **IMPORTANT NOTICE**

13 IEEE Standards do not guarantee or ensure safety, security, health, or environmental protection, or ensure
14 against interference with or from other devices or networks. IEEE Standards development activities consider
15 research and information presented to the standards development group in developing any safety
16 recommendations. Other information about safety practices, changes in technology or technology
17 implementation, or impact by peripheral systems also may be pertinent to safety considerations during
18 implementation of the standard. Implementers and users of IEEE Standards documents are responsible for
19 determining and complying with all appropriate safety, security, environmental, health, and interference
20 protection practices and all applicable laws and regulations.

1 Participants

2 <<The following lists will be updated in the usual way prior to publication>>

3 At the time this standard was completed, the IEEE 802.1 working group had the following membership:

4 **Glenn Parsons, *Chair***
5 **Jessy Royer, *Vice Chair***
6 **Janos Farkas, *TSN Task Group Chair***
7 **Lily Lv, *Editor***
8

9 The following members of the individual balloting committee voted on this standard. Balloters may have
10 voted for approval, disapproval, or abstention.

A.N. Other

11 <<The above lists will be updated in the usual way prior to publication>>

12

1

2 When the IEEE-SA Standards Board approved this standard on <dd> <month> <year>, it had the following
3 membership:

4

Chair

5

Vice-Chair

6

Past Chair

7

Secretary

8

*Member Emeritus

9

<<The above lists will be updated in the usual way prior to publication>>

10

1 Introduction

2

This introduction is not part of IEEE Std 802.1Qdv-20XX, IEEE Standard for Local and metropolitan area networks—Bridges and Bridged Networks—Amendment: Enhancements to Cyclic Queuing and Forwarding

3 This standard amends IEEE Std 802.1Q™-2022 as previously amended by IEEE Std 802.1Qcz™-2022. In
4 particular it enhances capabilities introduced by IEEE Std 802.1Qch™-2017.

5 This standard contains state-of-the-art material. The area covered by this standard is undergoing evolution.
6 Revisions are anticipated within the next few years to clarify existing material, to correct possible errors, and
7 to incorporate new related material. Information on the current revision state of this and other IEEE 802
8 standards may be obtained from

9 Secretary, IEEE-SA Standards Board
10 445 Hoes Lane
11 Piscataway, NJ 08854-4141
12 USA

Contents

1	1.	Overview.....	21
3	1.3	Introduction.....	21
4	2.	Normative references	22
5	3.	Definitions	23
6	4.	Abbreviations.....	24
7	5.	Conformance.....	25
8	5.4	VLAN Bridge component requirements.....	25
9	5.4.1	VLAN Bridge component options	25
10	5.13	MAC Bridge component requirements.....	25
11	5.13.1	MAC Bridge component options	25
12	5.14	End station requirements for count-based CQF.....	26
13	8.	Principles of Bridge operation	27
14	8.6	The Forwarding Process	27
15	8.6.5	Flow classification and metering	27
16	8.6.8	Transmission selection	27
17	99.	CQF Phase Alignment Protocol.....	28
18	99.1	Overview of CPAP	28
19	99.2	CPAP procedures.....	29
20	99.3	CPAP message timing	30
21	99.4	CPAP message frame formats	30
22	99.5	CPAP managed objects.....	31
23	100.	Cyclic queuing and forwarding.....	32
24	100.1	CQF managed objects.....	32
25	100.1.1	Cycle and priority structure managed objects	32
26	100.1.2	Cycle phase managed objects	32
27	100.1.3	Cycle variation information	32
28	100.1.4	Dead time	33
29	100.1.5	CQF forwarding delays	33
30	100.2	CQF LLDP TLVs	34
31	Annex A (normative)	PICS proforma—Bridge implementations	35
32	Annex T (informative)	Cyclic queuing and forwarding	36
33	T.1	Principles of CQF	36
34	T.1.1	Overview	36
35	T.1.2	CQF transmission selection	36
36	T.1.3	Bin selection	37

1	T.2	CQF in multiple queues on one output port.....	37
2	T.2.1	Multiple TC model	37
3	T.2.2	Integer multiples for TC	39
4	T.2.3	Admission control for multiple TC values	39
5	T.2.4	Implementation requirements	40
6	T.3	Time-based CQF.....	40
7	T.3.1	Frequency lock requirement	40
8	T.3.2	Timeline for time-based bin assignment	40
9	T.3.3	Preemption and interference	44
10	T.3.4	TC computation	44
11	T.3.5	Calculation of dead time TD	45
12	T.3.6	More than 3 output bins	45
13	T.3.7	Deterministic behavior of time-based CQF	46
14	T.3.8	Changing TC values along the path of Stream	47
15	T.3.9	Computing the actual end-to-end latency for time-based CQF	47
16	T.3.10	Output bin selection	47
17	T.3.11	Parameterization of time-based CQF	48
18	T.4	Count-based CQF	50
19	T.4.1	Calculating allocable time TA	51
20	T.4.2	Dead time TD	51
21	T.4.3	Number of output bins	51
22	T.4.4	Using both count-based and time-based CQF	51
23	T.5	Stream Aggregation.....	52
24	T.5.1	CQF Stream aggregation	53
25	T.5.2	CQF Stream disaggregation	53
26	T.5.3	CQF delay variation and disaggregation buffers	54
27	T.5.4	CQF Stream mixing	54
28	T.6	Additional considerations	54
29	T.6.1	Computing the CQF reservation	54
30	T.6.2	Frame size problem	55
31	T.6.3	Tailored bandwidth offerings	55
32	T.6.4	Overprovisioning to improve latency	55
33	T.6.5	CQF and credit-based shaper	56
34	T.6.6	Interactions among CQF, ATS, and control traffic	56
35	Annex ZY (informative)	Bibliography.....	57
36	Annex ZZ (informative)	Commentary.....	58

1 **Figures**

2	Figure 99-1	Aligning the transmitter and receiver CQF cycle start times.....	28
3	Figure 99-2	CQF Phase Alignment Protocol sequence	29
4	Figure T-1	Multiple TC values on multiple queues on one CQF port	37
5	Figure T-2	Transmission timing.....	38
6	Figure T-3	Variable TC.....	39
7	Figure T-4	Example of timelines for time-based CQF	41
8	Figure T-5	CQF Stream Aggregation example	53

¹ Tables

1

2

3 Draft Standard for 4 Local and Metropolitan Networks —

5 Bridges and Bridged Networks

6 Amendment: Enhancements to Cyclic 7 Queuing and Forwarding

8 (Amendment to IEEE Std 802.1Q™–2022 as amended by IEEE Std 802.1Qcz™–2022)

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

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

20 The contents of this initial Framemaker ‘kit’ clause are taken from an SA ballot copy of P802.1Qcz. This
21 provides an example of how to modify 1.1 Scope and 1.3 Introduction. See the EDITOR-PLEASE-READ file
22 and the 802-1Qxx-conditional-tags.fm file in this book for the various conditional tag views of this clause. This
23 paragraph is in paragraph style ‘Text’, character style ‘Editor’, with conditional tag ‘Editor comment’.

1. Overview

1.3 Introduction

Insert the following text at the end of 1.3 and reletter accordingly:

This amendment specifies procedures, protocols and managed objects for Cyclic Queuing and Forwarding (CQF). To this end, it

- a) Specifies a transmission selection procedure that organizes frames in a traffic class output queue into logical bins that are output in strict rotation at a constant frequency;
- b) Specifies a procedure for storing received frames into bins based on the time of reception of the frame.
- c) Specifies a procedure for storing received frames into bins based on per-flow octet counters.
- d) Specifies a protocol for determining the phase relationship between a transmitter's and a receiver's bin boundaries in time.
- e) Provides managed objects, Management Information Base (MIB), and YANG modules for controlling these procedures.
- f) Provides an informative annex to provide guidance for applying these procedures.

2. Normative references

The contents of this initial Framemaker 'kit' clause are taken from an SA ballot copy of P802.1Qcz. This provides an example of how to add to the list of references. See the EDITOR-PLEASE-READ file and the 802-1Qxx-conditional-tags.fm file in this book for the various conditional tag views of this clause. The Final text view may be useful in identifying missing references, failure to update existing references, and avoiding duplicates. New references should be added to the base text in collating order.

Insert the following items into the list of Normative References:

8

3. Definitions

The contents of this initial Framemaker 'kit' clause are taken from an SA ballot copy of P802.1Qcz. This provides an example of how to add to the list of definitions. See the EDITOR-PLEASE-READ file and the 802-1Qxx-conditional-tags.fm file in this book for the various conditional tag views of this clause. The Final text view may be useful in identifying missing definitions, failure to update existing definitions, and avoiding duplicates. New definitions should be added to the base text in collating order. Before adding new definitions, or modifying existing definitions search for all (if any) text conditionally tagged 'Delete' or 'Change remove' and delete it (use Edit>Find/Change...Conditional Text, and in the 'Find Conditional Text' pop up select the tag to be found, and then Change 'To Text' leaving the change to field blank and selecting 'Change All'). Then convert (using the Conditional Tags panel) any Change add or Insert tagged text to 'Base hide' (re remove tags entirely from it, if it is to be shown).

Insert the following definitions in the appropriate collating sequence, renumbering accordingly:

13

4. Abbreviations

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

The contents of this initial Framemaker 'kit' clause are taken from an SA ballot copy of P802.1Qcz. This provides an example of how to modify this clause. See Clause 3 for preliminary steps before editing for your amendment.

CPAP	Cyclic queuing and forwarding (CQF) Phase Alignment Protocol
------	--

1 5. Conformance

2 5.4 VLAN Bridge component requirements

3 5.4.1 VLAN Bridge component options

4 *Insert the following three sections at the end of 5.4.1 and renumber accordingly:*

5 5.4.1.12 Cyclic Queuing and Forwarding (CQF) requirements

6 A VLAN Bridge component implementation that conforms to the provisions of this standard for CQF (see
7 Annex T) shall

- 8 a) Support the ATS transmission selection algorithm as specified in 8.6.8.5.
- 9 b) Support the ATS scheduler state machines as specified in 8.6.11.
- 10 c) Support the requirements for Per-Stream Filtering and Policing (PSFP) as stated in 5.4.1.8.
- 11 d) Support the management entities for CQF as specified in §12.TBD.
- 12 e) Support the management entities for PSFP as specified in 12.31.

13 5.4.1.13 Time-based CQF requirements

14 A VLAN Bridge component implementation that conforms to the provisions of this standard for time-based
15 CQF (see Annex T) shall

- 16 a) Support CQF as stated in 5.4.1.12.
- 17 b) Support time-based CQF bin selection as specified in §8.TBD.
- 18 c) Support the management entities for time-based CQF as specified in §12.TBD.

19 5.4.1.14 Count-based CQF requirements

20 A VLAN Bridge component implementation that conforms to the provisions of this standard for count-based
21 CQF (see Annex T) shall

- 22 a) Support CQF as stated in 5.4.1.12.
- 23 b) Support count-based CQF bin selection as specified in §8.TBD.
- 24 c) Support the management entities for count-based CQF as specified in §12.TBD.

25 5.13 MAC Bridge component requirements

26 5.13.1 MAC Bridge component options

27 *Insert the following three sections at the end of 5.13.1 and renumber accordingly:*

28 5.13.1.4 Cyclic Queuing and Forwarding (CQF) requirements

29 A MAC Bridge component implementation that conforms to the provisions of this standard for CQF (see
30 Annex T) shall

- 31 a) Support the ATS transmission selection algorithm as specified in 8.6.8.5.
- 32 b) Support the ATS scheduler state machines as specified in 8.6.11.
- 33 c) Support the requirements for Per-Stream Filtering and Policing (PSFP) as stated in 5.4.1.8.
- 34 d) Support the management entities for CQF as specified in §12.TBD.
- 35 e) Support the management entities for PSFP as specified in 12.31.

1 5.13.1.5 Time-based CQF requirements

2 A MAC Bridge component implementation that conforms to the provisions of this standard for time-based
3 CQF (see Annex T) shall

- 4 a) Support CQF as stated in 5.13.1.4.
- 5 b) Support time-based CQF bin selection as specified in §8.TBD.
- 6 c) Support the management entities for time-based CQF as specified in §12.TBD.

7 5.13.1.6 Count-based CQF requirements

8 A MAC Bridge component implementation that conforms to the provisions of this standard for count-based
9 CQF (see Annex T) shall

- 10 a) Support CQF as stated in 5.4.1.12.
- 11 b) Support count-based CQD bin selection as specified in §8.TBD.
- 12 c) Support the management entities for count-based CQF as specified in §12.TBD.

13 *Insert the following section at the end of Clause 5 and renumber accordingly:*

14 5.14 End station requirements for count-based CQF

15 An end station implementation that conforms to the provisions of this standard for count-based CQF (see
16 Annex T) shall

- 17 a) Support the ATS transmission selection algorithm as specified in 8.6.8.5.
- 18 b) Support the ATS scheduler state machines as specified in 8.6.11.
- 19 c) Support the requirements for Per-Stream Filtering and Policing (PSFP) as stated in 5.4.1.8.
- 20 d) Support the management entities for CQF as specified in §12.TBD.
- 21 e) Support the management entities for PSFP as specified in 12.31.
- 22 f) Support count-based CQF bin selection as specified in §8.TBD.
- 23 g) Support the management entities for count-based CQF as specified in §12.TBD.

24

1 8. Principles of Bridge operation

2 8.6 The Forwarding Process

3 8.6.5 Flow classification and metering

4 *Change 8.6.5 as follows:*

5 << Editor's note: There are a great many references to ATS in 8.6.5, because of the split between "General
6 flow classification and metering" 8.6.5.1 and "per-Stream classification and metering" 8.6.5.2, which was
7 introduced by ATS. The editor believes that CQF will want to use 8.6.5.2 per-Stream C&M. Perhaps we can
8 present CQF as a special case of ATS to minimize changes. Perhaps we change lots of references in 8.6.5
9 from "ATS" to "ATS or CQF". Perhaps we invent a new term that encompasses both ATS and CQF for use in
10 subclauses 8.6.5.2 through 8.6.5.4.. Comments are solicited.>>

11 << Editor's note: It is not settled whether, in Clause 8, CQF should use bin number or eligibility time when
12 queuing frames. Using bin number simplifies the CQF bin assignment clauses somewhat, but likely means
13 less reuse of the existing ATS descriptions, especially in 8.6.8 Transmission selection. Comments are
14 solicited. >>

15 *Insert the following at the end of 8.6.5:*

16 8.6.5.7 CQF time-based bin assignment

17 << Editor's note: It is not settled whether, in Clause 8, CQF should use bin number or eligibility time for storing
18 frames. The title of this subclause may well be, "CQF time-based eligibility time assignment". >>

19 *Insert the following at the end of 8.6.5:*

20 8.6.5.8 CQF count-based bin assignment

21 << Editor's note: It is not settled whether, in Clause 8, CQF should use bin number or eligibility time for storing
22 frames. The title of this subclause may well be, "CQF count-based eligibility time assignment". >>

23 8.6.8 Transmission selection

24 *Insert the following at the end of 8.6.8:*

25 8.6.8.6 CQF transmission selection algorithm

26 << Editor's note: It is not settled whether, in Clause 8, CQF should use bin number or eligibility time for storing
27 frames. If we use eligibility time, there may be no need for this subclause. >>

28

1 *Insert the following Clause:*

2 99. CQF Phase Alignment Protocol

3 99.1 Overview of CPAP

4 See Annex T for an explanation of Cyclic Queuing and Forwarding (CQF), to which the CQF Phase
5 Alignment Protocol (CPAP) applies.

6 Figure 99-1 is an abridgment of Figure T-4. It shows an example of CQF. Bridge A is transmitting from a
7 CQF queue to Bridge B. Whether it is using time-based CQF (T.3) count-based CQF (T.4) or some other
8 method to fill its bins is irrelevant. Bridge B is using time-based CQF to assign frames to its output bins (not
9 shown). The time ticks on each timeline in Figure 99-1 indicate the start/end of a cycle of duration T_C . These
10 ticks are defined in terms of the transmit timestamp values of transmitted frames (IEEE Std 802.3 clause 90).
11 As described in T.3.2, Bridge B needs to assign each received Stream frame to an output queue bin, based
12 solely on the time of arrival of the frame at Bridge B's input port (IEEE Std 802.3 clause 90). In order to
13 assign frames to bin based only on time, Bridge B runs its output cycles with exactly the same period T_C as
14 Bridge A (see T.3.1), though not necessarily in phase (synchronized). The problem to be solved by CPAP is
15 described in T.3.2.2. Bridge B needs to establish the frame arrival time at its input port that corresponds to a
16 transmit cycle boundary in Bridge A. That is, Bridge B wants to know what input timestamp it would expect
17 to see on a frame Bridge A transmitted at exactly the start of a cycle in Bridge A.

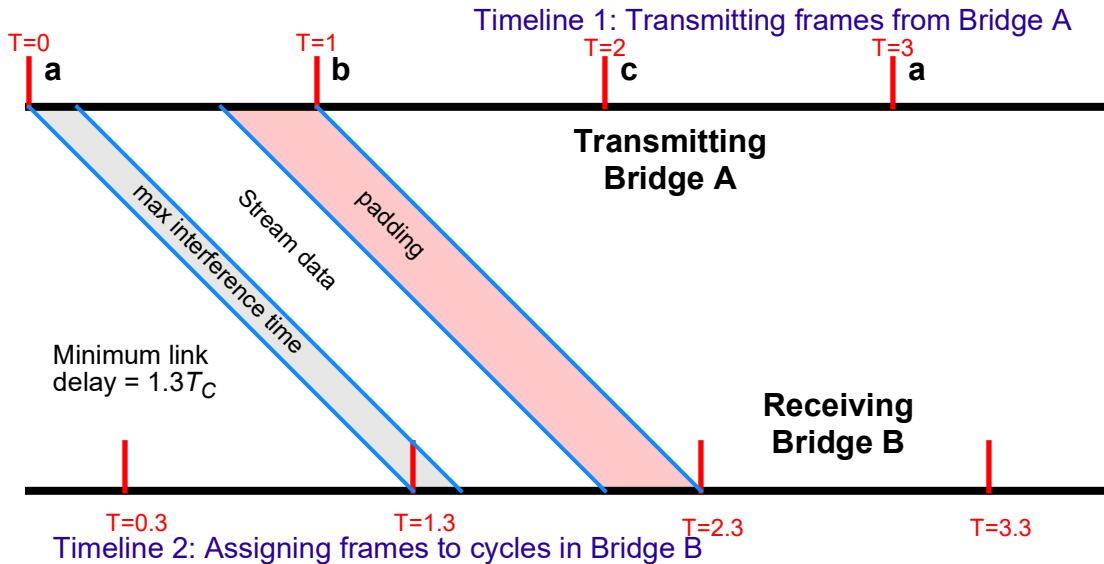


Figure 99-1—Aligning the transmitter and receiver CQF cycle start times

18 The transmit timestamp in Bridge A is a local matter; its format and frequency are not visible outside
19 Bridge A, and the timestamp is not a part of the transmitted frame. Similarly, the input timestamp in
20 Bridge B is local to that Bridge and bears no relationship to the transmit timestamp in Bridge A, except that,
21 when translated to seconds of elapsed time, both Bridges' timestamps advance at a rate close to their
22 respective CQF cycles. (Such variations are included in the definition of T_f in T.3.2.)

99.2 CPAP procedures

Figure 99-2 illustrates the operation of CPAP. Two systems are involved, a CPAP transmitter and a CPAP receiver. The CPAP transmitter is presumed to also be transmitting, or preparing to transmit, data frames from one or more CQF-enabled queues. The transmission of two CPAP messages are required, both transmitted from the CPAP transmitter: a CPAP Time Marker Frame, and a CPAP Phase Offset Message. See 99.4 for the formats of these messages.

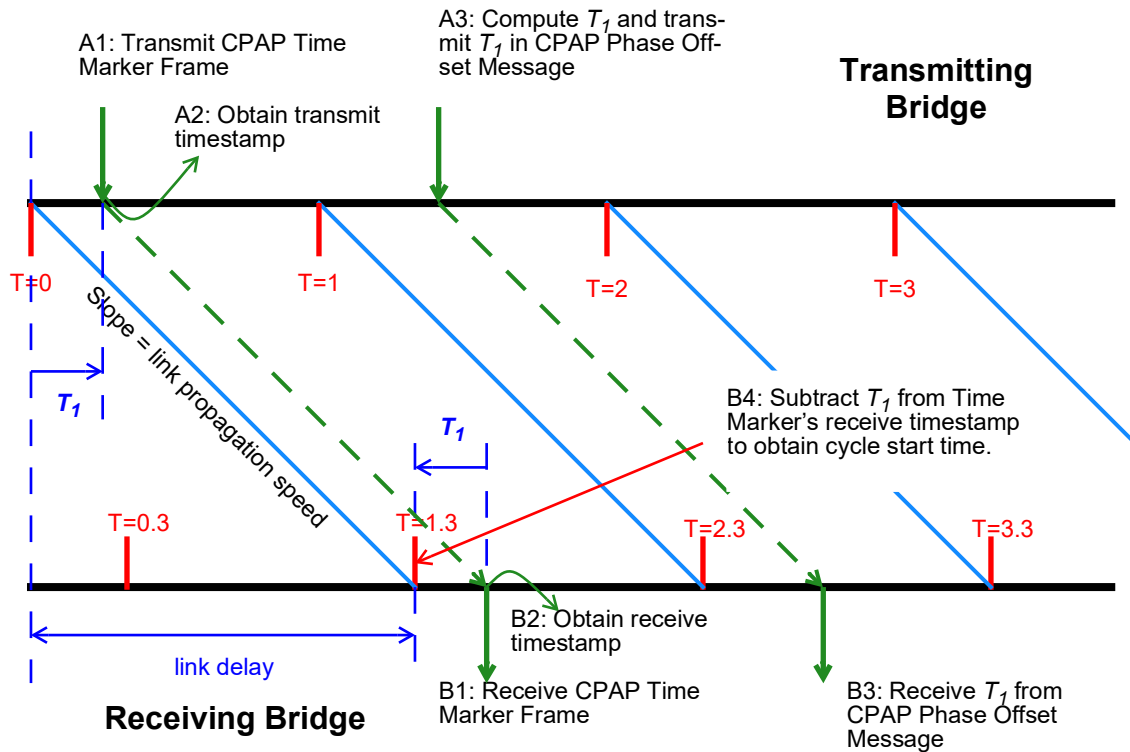


Figure 99-2—CQF Phase Alignment Protocol sequence

There is no implied relationship among the relative timing of cycle start times and CPAP messages beyond the requirements in 99.3. The two message need not be transmitted either in the same or in different cycles. It is a system administrator's choice as to what priority CPAP message are sent, and whether a Stream reservation is established for them. The default is priority 0, which is presumed to be a best-effort priority.

One CPAP sequence consists of the transmission of a CPAP Time Marker Frame followed by the transmission of a CPAP Phase Offset Message. The combination of the time between the two messages of a CPAP sequence, and the frequency of the CPAP sequences, contribute to the accuracy of the resultant. See 99.4.

If the CQF transmitter is operating more than one class of service queue with CQF enabled, the CPAP sequence is applied only to the slowest cycle (largest value of T_C).

1 The sequence of events in the CPAP transmitter, for one CPAP sequence, is as follows:

- 2 a) (Event A1 in Figure 99-2.) The CPAP transmitter transmits a CPAP Time Marker Frame.
- 3 b) (Event A2 in Figure 99-2.) The CPAP transmitter obtains a value, in locally-meaningful units, for
4 the time at which the first bit of the CPAP Time Marker Frame was placed on the medium
5 connecting the two CPAP systems (see IEEE Std 802.3 Clause 90).
- 6 c) (Event A3 in Figure 99-2.) The CPAP transmitter calculates time interval T_I , which is the time at
7 which the CPAP Time Marker frame was sent minus the start time of a CQF cycle. While, in
8 principle, any cycle start time in the past or future could be chosen, the CPAP transmitter shall
9 choose a cycle such that $-T_C \leq T_I \leq T_C$, where T_C is the CQF cycle time. This time can be positive or
10 negative, depending on whether it is compared to a past or a future CQF cycle.

11 The sequence of events in the CPAP receiver, for that same CPAP sequence, is as follows:

- 12 d) (Event B1 in Figure 99-2.) The CPAP receiver receives a CPAP Time Marker Frame.
- 13 e) (Event B2 in Figure 99-2.) The CPAP obtains a value, in locally-meaningful units, for the time at
14 which the first bit of the CPAP Time Marker Frame was received.
- 15 f) (Event B3 in Figure 99-2.) The CPAP receiver receives a CPAP Phase Offset Message containing
16 the time offset T_I .
- 17 g) (Event B4 in Figure 99-2.) The CPAP receiver uses the time of receipt from step e) and the time
18 offset T_I to compute the start time of a receive CQF cycle that is aligned with the CPAP transmitter.

19 << Editor's note: Unless the editor receives convincing opinions in the form of ballot comments to the
20 contrary, he has no intention of creating state machines or C code for CPAP. His claim is that this would add
21 complexity to the document, time to its development, and additional work for the reader, without increasing
22 the likelihood of interoperability. >>

23 99.3 CPAP message timing

24 << Editor's note. Input is encouraged on this subject. The editor does not believe that we should discuss the
25 matter of accuracy to a depth approaching that in IEEE Std 802.1AS. Exactly what should be said, here, is
26 therefore problematical, at present. >>

27 << Editor's note. In the editor's opinion, this document should NOT take an approach, which is possible, of
28 using the CPAP protocol to control the timing of the receiver's *output* cycles. That would be re-inventing IEEE
29 Std 1588,. >>

30 99.4 CPAP message frame formats

31 << Editor's note: Input is required for the editor to finish this section. At least the following possibilities can be
32 identified:

- 33 a) We could define an EtherType that would serve for both frames. This is a problem, as EtherTypes are
34 becoming an endangered species.
- 35 b) No data is carried by the CPAP Time Marker Frame except its identity as such. We could define a TLV
36 to be carried, for example, in one of the IEEE 802.1AS and/or IEEE 1588 Precision Time Protocol
37 (PTP) frame types. This has the advantage that one or the other of these protocols are often used
38 with TSN, and that they are often used with the IEEE 802.3 timestamp function, but of course this is a
39 disadvantage if neither is used. If PTP is not running, but we use PTP to carry CPAP messages, we
40 must ensure that we do not accidentally start up PTP or trigger other errors.
- 41 c) We could piggyback CPAP on LLDP.

42 >>

43 << Editor's note: It seems that a CPAP Phase Offset Message should somehow be tied to a particular
44 previous CPAP Time Marker Frame. We could assume that it applies to the previously-received one, if any.
45 We could add a serial number to the CPAP Time Marker Frame and include it in the CPAP Phase Offset
46 Message. We could include a time interval in the CPAP Phase Offset Message such that the message applies

¹ only to the one CPAP Time Marker Frame received within that interval before receipt of the CPAP Phase
² Offset Message.

³ Readers are solicited for their opinion(s) on these possibilities, or for a suggestion for others. >>

⁴ **99.5 CPAP managed objects**

⁵ << Editor's note: To Be Done when the other issues are settled. >>

1 Insert the following Clause:

2 100. Cyclic queuing and forwarding

3 100.1 CQF managed objects

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

10 100.1.1 Cycle and priority structure managed objects

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

- 12 a) The cycle time of the slowest CQF priority value (as a rational number of nanoseconds).
- 13 b) The priority value of the slowest CQF cycle.

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

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

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

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

23 100.1.2 Cycle phase managed objects

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

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

26 This variable establishes the phase of the input or output cycle. Typically, this variable would be manageable
27 the network administrator for output ports. For time-synchronized systems, it can be administered for input
28 ports, as well, in order to adjust for link delay. Alternatively, the input phase can be determined dynamically
29 (Clause 99), and be read-only for the network administrator.

30 100.1.3 Cycle variation information

31 For each output port only, we have:

- 32 a) The largest offset from the nominal (system clock) Nominal Output Cycle Start time (NOCS,
33 T.3.11.3) event to the actual cycle start time, in the negative (actual earlier than NOCS) direction.
- 34 b) The largest offset from the nominal (system clock) NOCS event to the actual cycle start time, in the
35 positive (actual later than NOCS) direction.

36 There are other ways to express the information in these two items. These values must be known to the
37 connected input port in order for that system to compute its buffer space and dead time requirements. This

1 information transfer could be accomplished by means of a protocol, managed objects, or by restrictions on
2 implementations.

3 **100.1.4 Dead time**

4 Given the context of dead time determination described in T.3.11.5, the following items are required by
5 CQF:

- 6 a) Per input port, per priority level, the total dead time that must be provided by the adjacent
7 transmitter at the end of each transmit cycle.
8 There is a component of this dead time computed according to T.3.11.5, as well as one computed in
9 item e) of T.3.11.3. The sum of these must be known to the adjacent transmitting port.
- 10 b) Per output port, per priority level, the total dead time that is to be provided at the end of each
11 transmit cycle.
12 This can be configured, obtained from the adjacent input system, or be a maximum of these values.
- 13 c) The allocable bandwidth for this input port and priority level.
14 This has three components, the minimum of the allocable bandwidth over all output ports reachable
15 from this input port (in the input port's own system), any limitations imposed by the input port
16 implementation, and any maximum imposed by management. Whether this is computed by, received
17 by, or even known by the output port, or whether allocable bandwidth is the concern only of the
18 admission control system, is an open question.
- 19 d) The allocable bandwidth for this output port and priority level.
20 This can be configured, computed from the adjacent input system's requirements, or be a minimum
21 of these values. Whether this is computed by, received by, or even known by the output port, or
22 whether allocable bandwidth is the concern only of the admission control system, is an open
23 question.

24 **100.1.5 CQF forwarding delays**

25 **100.1.5.1 Minimum CQF delay**

26 Read-only per Stream. See T.3.2 and Figure T-4. This is the CQF-caused delay. The starting point of this
27 delay is the IEEE Std 802.3 Clause 90 receive timestamp moment for a minimum-length frame transmitted
28 at the earliest possible moment in a cycle by the adjacent transmitting Bridge. The ending point of the
29 minimum CQF delay is the earliest time when the bin, into which the frame is stored, could be enabled for
30 output, assuming that the extra delay (100.1.5.3) is 0.

31 Thus, this time does not include the time required to empty a bin, the link delay, or any extra imposed delay.

32 **100.1.5.2 Maximum CQF delay**

33 Read-only per Stream. See T.3.2 and Figure T-4. This is the CQF-caused delay. The starting point of this
34 delay is the IEEE Std 802.3 Clause 90 receive timestamp moment for a minimum-length frame transmitted
35 at the latest possible moment in a cycle by the adjacent transmitting Bridge. The ending point of the
36 minimum CQF delay is the latest possible time when the bin, into which the frame is stored, could be
37 enabled for output, assuming that the extra delay (100.1.5.3) is 0.

38 Thus, this time does not include the time required to empty a bin, the link delay, or any extra imposed delay.

1 **100.1.5.3 Extra delay**

2 Configurable for each priority level, input/output port pair, and stream_handle, a read-write object is
3 required to specify the number of extra bins, beyond that computed/configured by other means, that frames
4 are to be stored in, in order to increase their delivery delay in this Bridge.

5 **100.2 CQF LLDP TLVs**

6 << Editor's note: In general, it would be good for two CQF devices to exchange information to allow them to
7 verify that they are both configured with the same priority levels, T_C values, etc. LLDP seems a reasonable
8 choice for this. Comments/suggestions are welcome. >>

9

¹ **Annex A**

² (normative)

³ **PICS proforma—Bridge implementations¹**

⁴ << Editor's note: To Be Done. >>

⁵

⁶

¹

Annex T

(informative)

Cyclic queuing and forwarding

Replace the contents of Annex T with the following:

T.1 Principles of CQF

T.1.1 Overview

Cyclic queuing and forwarding (CQF) is a method of transmission selection that can deliver deterministic, and easily calculated, latency for time-sensitive traffic streams. It is based on the following principles:

- a) A Bridge output queue using CQF (a “CQF queue”) is notionally divided into bins. The bins are enabled for output serially, at a fixed interval T_C , which same (or nearly the same) value is used for some number of Bridges along the path of a Stream, said path constituting a CQF segment of a network. At any given instant in time, a particular output bin can be available for accepting frames for later transmission, or enabled for transmitting frames to the associated medium, or neither, but never both. See T.1.2.
- b) Each Stream utilizing a CQF segment is allocated a certain number of bit times per transmission interval T_C . Steps are taken to ensure that no bin contains frames for any Stream that will take, in total, longer than that Stream’s allocated bit times to transmit. Resource reservation ensures that the total bit times allocated over all Streams passing through a CQF queue do not exceed T_C , even including possible interference from other queues on the port. See T.1.3.
- c) Frames assigned to the same bin at ingress to a CQF Segment remain together in the same bin at each hop along the CQF segment. Two methods are provided to accomplish this, time-based bin assignment and count-based bin assignment.

Taken together, these principles mean that no frames conforming to a Stream’s bit time allocation are dropped due to congestion, and that the end-to-end delivery delay varies by little more than $\pm T_C$. End-to-end delay calculation largely reduces to a hop count (T.3.9). These properties have significant consequences in larger networks, because they support the aggregation of Streams (T.5), which can reduce end-to-end delivery times and/or reduce network resource requirements. Different queues on a single port can operate at different T_C values (T.2) to provide CQF facilities for different levels of latency and bandwidth requirements. T.2

T.1.2 CQF transmission selection

A CQF queue is described in this standard as being divided into bins, because this simplifies the procedures described for assigning frames to bins. In this formulation, each received frame is assigned an integer output bin number b when queued for output on a port. This assignment can just as well be described in terms used for Asynchronous Transmission Selection (ATS) in §Clause 8. In ATS, each received frame is assigned a transmission time. If some number of frames are all assigned the same transmission time, selected from a range of future times separated by integral multiples of time T_C , this integer multiple is equivalent to the bin number b .

In previous versions of this standard, through IEEE Std 802.1Q-2022, each of the bins in the present standard were implemented using an entire class of service queue, and transmission gates were used to swap between queues, thus rotating the bins. There was also a requirement that all of the Bridges in a network

1 synchronize their transmission gates, and rotate the output bins (queues) at the same time. This is a perfectly
2 valid method for implementing CQF. The present standard describes CQF as a one or more individual class
3 of service queues, each with multiple bins. This formulation offers a wider range of services.

4 CQF class of service queues can be utilized on the same port with other transmission selection methods;
5 strict priority determines which queue is selected for transmission.

6 T.1.3 Bin selection

7 When a Stream frame is received and forwarded to a class of service output queue that is enabled for CQF,
8 the frame is assigned to a particular bin in that queue. There are two methods for assigning a frame to a bin,
9 time-based (T.3.1) and count-based (T.4). The same frame can be assigned to bins on two different output
10 ports using two different methods; the same bin can have frames from different Streams assigned to it using
11 different methods. All of the frames in the same Stream received from the same port and transmitted on the
12 same port use the same method. A Bridge can be configured for the bin selection method to use for all
13 frames received from a given port, regardless of the output port. It can be configured on an input-output port
14 pair basis. The selection method can be configured for specific Streams.

15 T.2 CQF in multiple queues on one output port

16 T.2.1 Multiple T_C model

17 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
18 Streams can be accommodated on any one port, because all frames for all Streams sharing a port must fit
19 into a single T_C period. If the value chosen for T_C is large, then more Streams can be accommodated, with a
20 wide variation in allocated bandwidth, but the larger T_C increases the per-hop latency. In the ideal case, of
21 course, every Stream would have a T_C value chosen so that exactly one frame of a Stream is transmitted on
22 each cycle T_C .

23 Instead of picking a single value for T_C that is sub-optimal for most Streams, we can apply multiple values
24 of T_C to a single output port, as shown in Figure T-1.

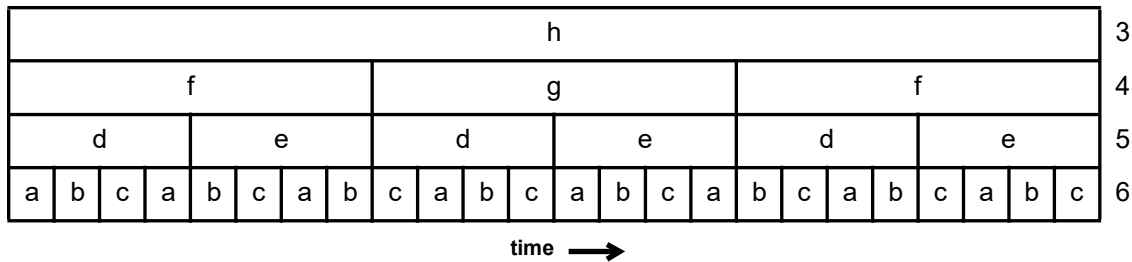


Figure T-1—Multiple T_C values on multiple queues on one CQF port

25 In Figure T-1, we have a schematic timeline. Four class of service queues have been configured for CQF,
26 each with a different value of T_C . The fastest (call it, " T_{C6} ") runs at the highest priority (6). T_{C5} is slower by
27 a factor of 4 from T_{C6} in this example, and its bins run at priority 5 (less important than priority 6). T_{C4} is
28 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
29 in Figure T-1 label which bin is output during the cycle. There are 9 bins a through i. Bin i, the second bin at
30 priority 3, is not shown. In this example, priority 6 uses three bins, because the timing is tight; the others use
31 two each.

1 We assume here that the receiver of a frame can identify the particular CQF instance (T_C value) to which the
2 frame belongs by inspecting the frame. A TSN Bridge could use the priority field of a VLAN tag, or it could
3 use the DSCP field of an IP packet. IEEE Std 802.1CB provides for the use of other fields in the frame, e.g.
4 IP 5-tuple.

5 Since the total bandwidth of the link is never oversubscribed by Streams, each cycle, fast high-priority and
6 slow low-priority, is guaranteed to be able to transmit all of its frames within the duration of its cycle. For
7 example: If 50% of T_{C5} is reserved, and 30% of T_{C3} is reserved, then 80% of the total bandwidth has been
8 reserved, leaving only 20% for other Streams, best effort traffic, and dead time. This is shown in Figure T-2,
9 where we illustrate the timing of transmission of frames from three levels of CQF and the best-effort (BE)
10 level. Note that CQF traffic can be delayed within its window by interference from both higher priorities
11 (e.g. the first priority 4 frame) and lower priorities (e.g. the first priority 6 frame), but that it will always get
12 out before the window closes, assuming that the bandwidth is not oversubscribed.

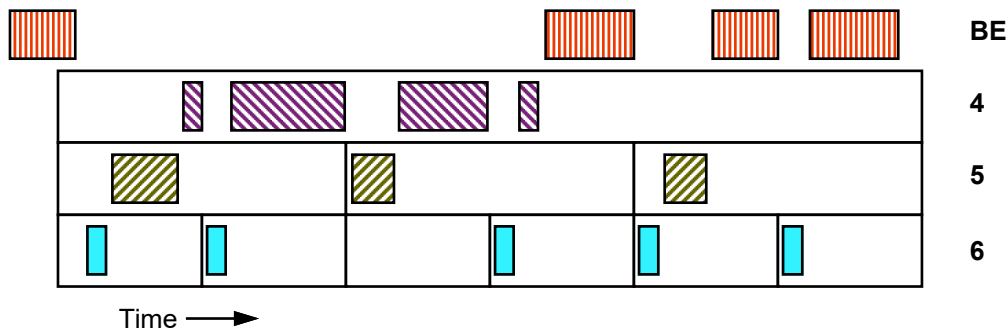


Figure T-2—Transmission timing

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

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

20 Streams are allocated to, and thus use up the bandwidth available to, each cycle separately. Any cycle can
21 allocate up to 100% of the bandwidth of that cycle's T_A , but the percentages allocated to all of the cycles
22 must, of course, add up to less than 100%. The total amount of buffer space required depends on the
23 allocation of Streams to priority values. If all Streams are slow and are allocated to T_{C4} up to a total of
24 100%, then full-sized bins must be used for bins h and i. If all Streams are fast and are allocated to T_{C6} , then
25 only three small bins are used—bins a, b, and c are rapidly re-used.

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

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

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

1 point of one frame per cycle, with no overprovisioning, and thus give a faster latency. Overprovision =
2 higher latency, in this case.

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

5 T.2.2 Integer multiples for T_C

6 The ideal would for each Stream S to have its own T_{CS} that requires no overprovisioning. But, that winds up
7 being equivalent to a per-Stream-shaper solution such as Asynchronous Traffic Shaping or IntServ. The
8 reason can be seen in Figure T-3.

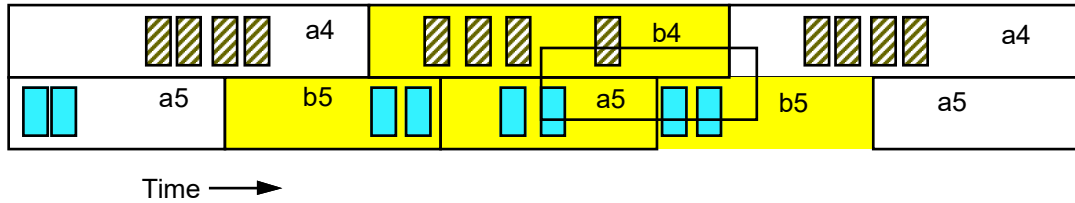


Figure T-3—Variable T_C

9 In Figure T-3, we have allocated 40% of the link bandwidth to the Streams using priority 5, and 50% of the
10 link bandwidth to the Streams using priority 4. The cycles do not line up with an integral number of faster
11 cycles in each period of slower cycle. Since we cannot predict exactly where, during a cycle, frames can be
12 emitted (see T.6.5), we can get the situation shown, in the shaded bins. Bins b5, a5, and then again, b5 emit
13 their frames (at high priority) at the indicated times. Even though the priority 5 Streams take up only 40% of
14 each level-2 cycle, they can output 6 frames over the course of cycle b4, thus taking up 60% of the
15 bandwidth during that period. There is, therefore, 110% of the bandwidth that must be output during the
16 period that b4 is transmitting. b4 cannot output all of its data. Some of it must be somehow delayed, but
17 there is no place to put that data. Deterministic QoS is not obtained.

18 Having an integral number of cycles at each layer fitting exactly into one cycle at the next-slower layer
19 ensures that the lower-priority, slower cycle, will always have sufficient time to output all of its frame,
20 because the problem in Figure T-3 is avoided. Integral multiples fitting exactly means that, at the moment a
21 cycle starts and ends at one priority level, a cycle starts and ends at each higher priority level, as illustrated in
22 Figure T-1. This scheme also bounds the number of preemption events that can steal bandwidth from a given
23 priority level (see T.3.3).

24 T.2.3 Admission control for multiple T_C values

25 T.2.1 describes the operation of CQF with multiple T_C values operating simultaneously on one output port.
26 Figure T-2 shows an example of a sequence of transmissions. We observe that the shortest cycle times
27 operate at the highest priority, and the longest at the lowest priority. Because different CQF priority levels
28 may have different maximum frame sizes, and because some may enable preemption, different priority
29 levels may have different amounts of time during one cycle that cannot be allocated to Stream transmission.
30 Clearly, allocating time for any CQF priority level reduces the time allocable to other priority levels; there is
31 only one physical link.

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

34 For a new Stream to be admitted, it must be true that the available transmission times over all of the CQF
35 levels on all of the output ports through which the Stream travels have not been exhausted. At any given

1 CQF priority level x , one can add the bits allocated to all Streams in one cycle at CQF priority level x , plus
2 the sum over all more-important CQF priority levels y (faster cycles), of the product of the number of bits
3 per cycle allocated at that level times the number of cycles at that level contained within one cycle at level x .
4 At every level, the total must not exceed the maximum number of allocable bits at that level.

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

7 **T.2.4 Implementation requirements**

8 The admission control calculations presented here depend upon the transmitting port being able to select the
9 correct frame to transmit according to strict priority among the CQF priority levels, and initiate all
10 transmissions in that order, at line rate, without introducing extra inter-frame gap time. Since, with CQF, no
11 bin has frames both arriving and being transmitting at the same instant, this should pose no insurmountable
12 problems for implementors.

13 **T.3 Time-based CQF**

14 **T.3.1 Frequency lock requirement**

15 CQF does not require synchronization of the system clocks, but does require frequency lock. That is, the
16 number of CQF cycles in two Bridges that are frequency locked must be the same, over an arbitrarily long
17 interval of time.

18 **T.3.2 Timeline for time-based bin assignment**

19 We have two Bridges, A and B. Both are running time-based CQF on each of multiple ports.

20 When a CQF cycle starts on a particular port, Bridge A transmits all of the frames in one bin towards
21 receiving Bridge B, not necessarily in a single burst. After some gap following the transmission of the last
22 frame in the bin, and at time T_C after the cycle started, another cycle starts. At this point, it starts transmitting
23 the frames from the next bin. The cycle in both Bridges happen regularly, with the same period T_C . At the
24 next hop, Bridge B must be able to assign each received frame to a transmit bin such that 1) frames that were
25 in the same bin in Bridge A, and are transmitted on the same port from Bridge B, are placed into the same
26 bin in Bridge B; and 2) frames in different bins in Bridge A are placed in different bins in Bridge B.

27 Figure T-4 shows an example of CQF. Bridge A and Bridge B are transmitting at the same frequency, but are
28 offset by $0.1T_C$, as shown by timelines 1 and 4. In Figure T-4, we use the following notation for time
29 intervals:

- 30 T_C nominal (intended) period of the CQF cycle
- 32 T_I maximum interference from lower-priority queues, either one frame or one preemption fragment
- 34 T_V sum of the variation in output delay, link delay, clock accuracy, and timestamp accuracy
- 36 T_A the part of the cycle allocable to (reservable by) Streams
- 38 T_P worst-case time taken by additional bytes added to Stream data if this traffic class is preemptable
- 40 T_D end-of-cycle dead time optionally imposed on Bridge A by Bridge B
- 42 T_W wait time during which the bin is neither receiving nor transmitting frames
- 44 T_{AB} effective phase difference between cycle start times for input from A and output from B

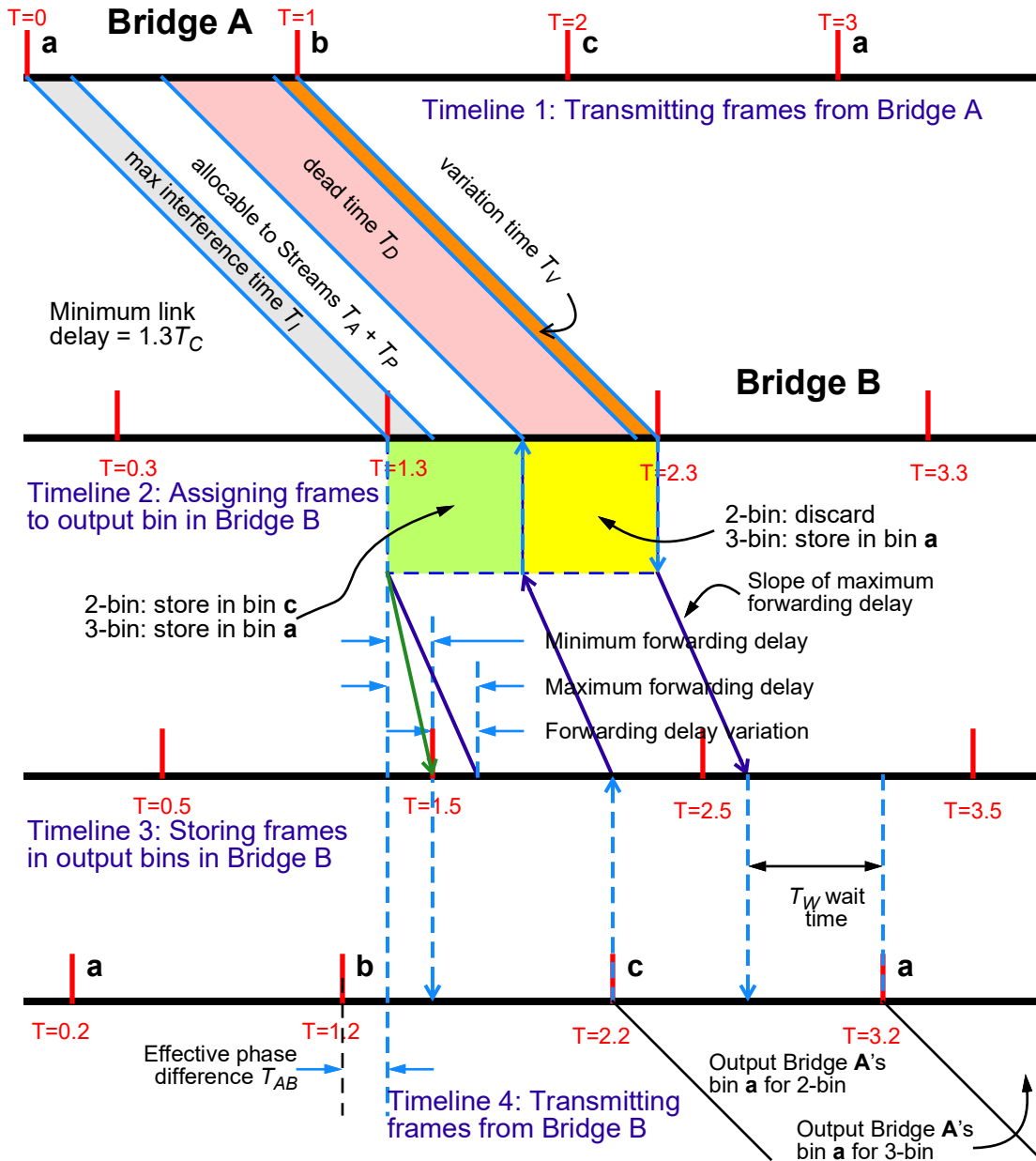


Figure T-4—Example of timelines for time-based CQF

1 For time-based CQF, T_{AB} must remain constant; that is, any variation in T_{AB} is included in T_V . Bounding this
2 variation is another way of saying that all Bridges' T_C values are exactly equal.

3 Following the definitions of transmission gates in §8.6.8.4, the red ticks in timelines 1 and 4 in Figure T-4
4 represent the earliest possible moment at which the first bit of the destination address of the first frame of the
5 cycle can be transmitted. These ticks are ultimately driven by the frequency-locked clock. They are the basis
6 for all bin transmissions. If Enhancements for Scheduled Traffic (ETS, §8.6.8.4) are used for controlling the
7 output bins, the ticks are the points in time when the transmission gate of one queue is closed, and the next
8 queue's transmission gate is opened. These are the points in time as programmed into the managed objects
9 that control ETS. An implementation may need to schedule cycle start times in anticipation of the time

1 specified in the managed objects in order to maximize throughput. Note that the preamble of an IEEE Std
2 802.3 Ethernet frame can be transmitted before the start of a cycle.

3 T.3.2.1 Output timeline 1

4 Figure T-4 shows an interference delay T_I (the gray area) between the start of Bridge A's cycle (the red ticks
5 in Figure T-4) and the transmission of the first bit of the first Stream frame's destination MAC address. The
6 interference is from frames transmitted from lower-priority queues. It is equal to the time required for one
7 maximum-length transmission unit over all lower-priority queues. That maximum transmission unit is either
8 a maximum-length fragment, for preemptable lower-priority queues, or the maximum-length frame, for non-
9 preemptable queues. The value of T_I depends upon the configuration of lower-priority queues.

10 It is possible that the class of service illustrated in Figure T-4 is, itself, a preemptable class. In that case, a
11 higher-priority class of service can preempt transmission of frames in this class. Preempting a frame adds
12 additional bytes to the resultant fragments, which must be accounted for when allocating bandwidth to a
13 class of service. T_P represents the worst-case additional time required to transmit these extra bytes caused by
14 preempting frames belonging to an CQF Stream. This value is always bounded. See T.3.3.

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

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

28 As described in T.3.5, the next hop can impose a dead time T_D on this hop. This is a time at the end of the
29 cycle, during which no frames can be transmitted from the bin, so that the last frame of the cycle can be
30 received earlier than the end of the cycle.

31 The total time per cycle that can be used for transmitting Streams is, then:

$$32 T_A = T_C - T_I - T_P - T_D - T_V.$$

33 This T_A is a maximum, local to a particular class of service and output port on a Bridge. It guarantees that the
34 last frame of cycle (plus a possible preamble of the first frame of the next cycle) will be on the wire before
35 the start of the dead time. All of the components of T_A can be calculated by an implementation from its
36 configuration and from knowledge of the implementation, except for T_D and parts of T_V . T_D is supplied by
37 configuration, or by the Bridge to which the output port is connected. T_V can be supplied either by the time
38 sync implementation, by configuration, by summing the contributions of Bridge A and Bridge B, or by the
39 specification of a maximum allowed value by a standard or an equipment purchaser.

40 Note that T_A , as defined here, includes the entire transmission time of Stream data, including one 12-byte
41 inter-frame gap and one 8-byte preamble for every frame. The preamble of the first frame of a cycle is
42 counted in the previous cycle due to the way in which the transmission gates are defined in §8.6.8.4.

1 The last frame transmitted from the bin has to complete transmission within the period marked T_A in
2 Figure T-4. (But, see T.6.6.)

3 T.3.2.2 Receive timeline 2

4 The timeline at the receiving port is timeline 2 in Figure T-4. The red ticks represent the earliest possible
5 moment that the first bit of the destination MAC address of the first frame of a cycle can be received.

6 On timeline 2, this standard assumes that each frame is assigned to a bin on an output port based on the
7 timestamp (Clause 90 of IEEE Std 802.3-2018) on the frame. Other means of assigning an arrival time to a
8 frame can be used.

9 A critical aspect of timeline 2 is its offset from timeline 4, the output timeline. This offset is shown as T_{AB} in
10 Figure T-4. It is clear from the figure that T_{AB} must be known in order to compute T_D and T_W . T_{AB} can be
11 computed by 1) synchronizing the clocks of Bridges A and B, and 2) measuring the link delay from
12 Bridge A to Bridge B using PTP. Other methods are also possible, e.g. that described in Clause 99.

13 Once T_{AB} is known, all of the timing relationships shown in Figure T-4 can be computed. The phasing of the
14 Bridges' output bin cycles affects the end-to-end latency of any Stream, so that phasing must be known
15 when the end-to-end latency is computed. However, the end-to-end latency is not necessarily an integer
16 multiple of the cycle time, because cycle start times are not necessarily synchronized among the Bridges in a
17 network. One could even adjust the phasing (by adjusting the phase of timeline 4) to favor certain paths
18 through the network.

19 For time-based CQF, if a frame (belonging to a Stream) is received that straddles a cycle (first bit in one
20 cycle on timeline 2 of Figure T-4, and end frame plus inter-frame gap plus a preamble time occurs in the
21 next cycle), then either 1) some part of that frame was transmitted from Bridge A outside the cycle window
22 T_C , or 2) one or more of the constants, measurements, or calculations above is incorrect. Either way, unless
23 the frame is discarded or marked down to best-effort service, it can cause disruption of delivery guarantees
24 farther along in the network.

25 T.3.2.3 Storing frames timeline 3

26 The timeline at the point where frames are stored into an output bin is timeline 3 in Figure T-4. The red ticks
27 on timeline 3 mark the earliest point at which the first frame transmitted from a particular bin could reach
28 the output bins (neglecting transmission time on the input medium). These ticks are offset from timeline 2
29 by the minimum forwarding delay, required to forward the frame from the input port to the output queue.
30 The maximum forwarding delay is also shown. The forwarding delays shown in Figure T-4 include the time
31 to install the frame in the output bin and for its presence to filter through to the point that it can be selected
32 for output.

33 For a Bridge B that is connected to and receiving CQF frames from n other Bridges, we have n bin
34 assignment problems to solve, one for each input port on Bridge B. The problem for each is to determine
35 how many bins are needed, and to which bin each frame is to be assigned.

36 There are two bin assignment methods shown in Figure T-4: the 2-bin method, in which the frames received
37 from Bridge A bin a are assigned to bin c in Bridge B, and the 3-bin method, where those same frames are
38 assigned to bin a in Bridge B. The slope of the maximum forwarding delay allows us to compute the latest
39 moment at which frames received from bin a on Bridge A can be stored into bin c on Bridge B. The shaded
40 areas just below timeline 2 in Figure T-4 show the time windows for bin assignment. If two output bins are
41 used, then frames received from bin a on Bridge A can be assigned on input (timeline 2) to bin c only as
42 long as they are assured of being placed into bin c before Bridge B starts transmitting bin c . As shown,
43 frames from bin a can be assigned to bin a (3-bin mode) during the entire length of the cycle on timeline 2.

1 Time T_W in Figure T-4 is the time during which, in 3-bin mode, bin C is holding frames, neither filling nor
2 emptying. In 3-bin mode, the dead time T_D is 0, and T_A , the allocable transmission time, encompasses both
3 the T_A (white) and T_D (red) regions in Figure T-4.

4 Unlike timeline 1, timeline 2, or timeline 3, the red ticks on timeline 3 are not hard boundaries. The
5 forwarding delay variation shown in Figure T-4 could, in theory, be longer than one cycle time T_C . See T.3.6.

6 Note that an implementation may require a minimum offset between timeline 3 and timeline 4. That is, a
7 time lag may be required between the last opportunity to store a frame in a bin, and the earliest time at which
8 the first bit of a frame from that bin can appear on the link. Some time could, for example, be necessary in
9 order schedule the transmission of frames across multiple queues in order to ensure that the requirements of
10 strict frame priority and back-to-back frame transmission (T.2.4) can be met.

11 T.3.2.4 Transmitting frames timeline 4

12 Depending on whether 2-bin or 3-bin mode is used, one can trade off reduced total available bandwidth
13 against per-hop delay. Timeline 4 in Figure T-4 shows the two options for the choice of which output cycle
14 in Bridge B is used to transmit frames that were transmitted from bin a in Bridge A.

15 T.3.3 Preemption and interference

16 Not all of the bandwidth in a cycle T_C can be allocated. The smaller the cycle time, the greater the impact of
17 the interference time (T_I in T.3.1 and Figure T-4) on the allocable bandwidth. Frame preemption is described
18 in §6.7.2 and in Clause 99 of IEEE Std 802.3-2018. Preemption can reduce the interference time.

19 T_I is equal to the worst-case transmit time for a single transmission from a lower-priority queue. This
20 interference can occur only at the beginning of a cycle. Since this value must be bound, it places a
21 requirement, that must be enforced, on all lower-priority queues that they either have a maximum frame size
22 or that frame preemption is applied to the lower-priority queues. If preemption is used, the maximum
23 interference is the maximum fragment size (about 150 bytes, see IEEE Std 802.3). The interference time is
24 shown as a gray parallelogram attached to timeline 1 in Figure T-4.

25 The other time is the preemption time T_P , which applies only to Streams that are preemptable. This case is
26 not typical, but is possible if a large fraction of the available bandwidth is to be assigned to one or a few
27 high-bandwidth Streams, and lower-priority Streams use larger frames. T_P is the product of (the maximum
28 number of highest-priority transmission windows that can open during a single window for the level being
29 computed) * (the per-preemption penalty). Thus, in Figure T-1, if priority 4 is preemptable, then there are 8
30 level 6 windows that can open. This means that there can be 8 preemption events during one level 4 window,
31 so the total preemption time T_P is 8 times the preemption penalty. (It doesn't matter which specific frames
32 are preempted; only how many such events occur during the cycle.) The preemption penalty is the number
33 of bytes added when a frame is preempted, which is 4 (CRC on preempted fragment) + 20 (inter-frame gap)
34 + 8 (preamble for continuation fragment) = 32 bytes.

35 T.3.4 T_C computation

36 We can also compute a suitable value for T_C , given a desired value for T_A :

$$37 T_C = T_A + T_P + T_I + T_D + T_V$$

1 Annex T CQF assumes the 2-bin scheme, and so assumes that T_D and T_V are small enough and T_C large
2 enough to leave a useful T_A . Assuming that one's goal is the smallest possible T_C :

- 3 a) T_D can be eliminated by using the 3-bin scheme of CQF.
- 4 b) Implementation steps can be taken to reduce T_V . This may include steps to reduce the variability of
5 the forwarding delay, the delay between selection-for-output and first-bit-on-the-wire at the previous
6 hop, or increased accuracy of the synchronized clock.
- 7 c) T_I can be reduced by restricting the maximum frame size of lower-priority Streams, or by enabling
8 frame preemption.

9 T.3.5 Calculation of dead time T_D

10 Timeline 3 in Figure T-4 shows the calculation of T_D , which applies only to 2-bin mode. The starting point
11 of T_D is the moment that the output cycle starts (the tick on timeline 4), moved backward by the worst-case
12 forwarding delay. This is the last moment on timeline 3 that a frame can be assigned to bin C in the example
13 in Figure T-4. The end of T_D is the end of the cycle T_C , less the variation time T_V . In 3-bin mode, T_D is zero.

14 T_D can only be computed by Bridge B. Its effect on the allocable bandwidth T_A must be taken into account
15 when admitting new Streams. If a network uses a peer-to-peer control structure using, e.g. MSRP (Clause
16 §35), then the value of T_D must be made available to the previous Bridge A so that Bridge A does not exceed
17 the reduced T_A .

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

- 19 a) The value of T_D can be propagated backwards to the previous Bridge, either via management or via
20 an extension of the reservation protocol.
- 21 << Editor's note: No such mechanism exists, at this point. >>
- 22 b) A Bridge can compute the value of T_D and decide whether to employ 2-bin or 3-bin mode,
23 depending on how much bandwidth has been allocated, so far. This, of course, can change a
24 previously-computed Stream's end-to-end latency.
- 25 c) All Bridges in a network can be configured with a reasonable maximum value for T_D . If a particular
26 input/output port pair on a particular Bridge computes a value for T_D that exceeds this maximum,
27 then 3-bin operation is required.

28 T.3.6 More than 3 output bins

29 So far, the discussion of Figure T-4 assumes that the variation in forwarding delay is small, relative to T_C . If
30 this is not the case, Bridge B can use more than 3 output bins, and assign received frames to bins whose
31 output is scheduled far enough ahead in time to ensure that, in the worst case, they will arrive in the proper
32 bin before the bin begins transmitting. This works only because the bin assignment decision is made based
33 on time-of-arrival of the frame at the input port, not the time-of-arrival of the frame at the output port.

34 In certain situations, e.g. when a Stream is replicated and traverses two paths of different lengths using IEEE
35 Std 802.1CB Frame Replication and Elimination for Reliability (FRER), it can be desirable to purposely
36 delay a Stream's frames in order to match the total delay for the Stream along the two paths (see C.9 of IEEE
37 Std 802.1CB-2020). In this case, extra output bins can be allocated, and used to impose a delay of an
38 arbitrary number of cycle times T_C on every frame.

39 Each output port in a Bridge, and each output port along the path of a Stream, can have a different number of
40 bins, whether 2, 3, or 50. Furthermore, one Stream can use (e.g.) 3 bins on an output port, while another
41 Stream, which needs a path-matching delay, can use 12 bins on the same port. (Of course, this requires per-
42 Stream configuration.)

1 **T.3.7 Deterministic behavior of time-based CQF**

2 CQF guarantees the Deterministic QoS by the following argument.

3 We assume that the Talker uses CQF. Non-CQF inputs to a Bridge are discussed in T.4.4.

4 We consider only one value of T_C along the path of a given Stream from Talker to Listener. T.2.3 and T.3.8
5 deal with exceptions to this assumption.

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

10 A number of considerations reduce the fraction of the total time T_C that can actually be used to transmit data.
11 See T.3.1 for details. For example, the maximum frame size of each Stream allows us to determine the
12 worst-case interference that a given Stream can have on higher-priority Streams. All of these considerations
13 are bounded; if an implementation cannot bound one or more of these considerations, then it cannot
14 guarantee the Deterministic QoS in a CQF network.

15 In a detailed timing analysis, we will note that the first bit of the MAC address of a frame is never
16 transmitted before the start of the window time (according to the local time in the transmitter) and the last bit
17 of the interframe gap (always) and the preamble of the next frame (if any) are transmitted before the end
18 of the window.

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

- 22 a) The Talker contract is enforced when a Talker's frames are first placed into a CQF output bin after
23 entry to the network. That is, the frames from a given Stream do not exceed the Talker contract in the
24 first CQF output bin in the network.
25 Ingress conditioning and/or policing is discussed in T.4.4.
- 26 b) Frames belonging to the same Stream that are in the same CQF output bin in one Bridge in the
27 network are placed in the same CQF output bin in all subsequent Bridges along a shared path.
28 T.3.1, and particularly Figure T-4, show the details of how this is accomplished. The key is to get the
29 Stream gates synchronized with the transmission gates of the transmitting system, offset by the link
30 delay. Frames received during one input cycle are always placed in the same bin. If the input cycle is
31 synchronized with the previous hop's output cycle, then cycle integrity is maintained. (Of course,
32 this only works for point-to-point links.)
- 33 c) There is no fan-in for a single Stream.
34 We assume that the path of a Stream reservation through the network is known and does not change.
35 A given Stream enters a Bridge through one port only, although it may be a multicast Stream, and
36 thus be enqueued and transmitted on more than one port.
- 37 d) Admission control ensures that, on any given output port and cycle time T_C , the total bits times for
38 all Streams passing through that port and T_C value does not exceed the available transmission time
39 on that port. (This assumes that no Bridge has a limitation on available receive time on an input port
40 that is smaller than the attached output port's available transmit time. The implications of such a
41 limitation are obvious.)

1 T.3.8 Changing T_C values along the path of Stream

2 If a Stream enters a Bridge using a cycle time T_C , and is being transmitted on an output port with cycle time
3 $n \cdot T_C$, then n successive input cycles can be deposited in the same output bin with no problem, as long as the
4 larger cycle time's dead time requirements are met. (This is not a trivial exception, as the larger cycle's dead
5 time occurs at the end of the large cycle, and thus may take up much or even all of one small cycle.)
6 Equivalently, the input port can be configured with the slower cycle time to match the output port in the
7 same system. Of course, when making the reservation for that Stream, the adjustment of its contract must be
8 made; it is allocated n times the number of bits in the slower cycle than in the faster cycle.

9 In all other cases, when a Stream changes cycle times, the Stream must pass through a conditioning step,
10 such as a count-based CQF step (see T.4.4), to ensure that the Stream never exceeds its contract in the new
11 cycle time.

12 T.3.9 Computing the actual end-to-end latency for time-based CQF

13 After adjusting to get the receiving window aligned with the previous-hop transmitting window, a Bridge
14 knows the "effective phase difference T_{AB} " described in T.3.1. Referring to Figure T-4, this allows the
15 Bridge to compute the difference, in time, between the start of an input window for the Stream, and the start
16 of the output window in which a frame received in that input window will be transmitted. This is the dwell
17 time for the frame in this Bridge. A maximum and minimum time for this delay is given in 100.1.5.

18 Link delay is relevant to the computation of end-to-end delay, but it can be hidden by using time-based CQF
19 in time-synchronized Bridges, and using dead time, so that the link delay is accounted for within the CQF
20 cycle time T_C . If the link delay does need to be added to the delay, it is the one-way link delay that is added.
21 Typically, this is measured using the PTP. At egress from the network, there is a margin of one cycle time
22 less one frame transmission time for delivery of the frame, as the frame can be transmitted at any point
23 during the cycle, but must both start and finish its transmission within the cycle. The delay at ingress is
24 somewhat more complicated to measure, as it depends upon the method used by the Talker and the ingress
25 Bridge to shape its transmissions.

26 If we look again at Figure T-4, we can see that the difference between using two and three bins for a given
27 input-output port pair is really a matter of rounding up the link delay to an integral number of cycle times. If
28 the sum of link delay and phase delay between output cycles is negligible, or happens to be very nearly an
29 integer multiple of the cycle time, then the yellow "discard" area is small, and two bins can be used. If sum
30 is larger, then one necessarily chooses between a smaller allocation (large discard area) and increased delay.

31 T.3.10 Output bin selection

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

36 We describe here one method for receiving a frame and assigning it to a bin. There are many ways to
37 accomplish the same task.

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

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

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

$$4 \ b = (S + P_{io}) \bmod B_n$$

5 Where P_{io} is the cycle phase offset from input port i to output port o and B_n is the number of output bins on
6 the port. See T.3.11.4 for the determination of P_{io} . Note that in the extreme case of all output ports using two
7 bins, all synchronized, and all input cycles in phase with the output cycles, the table P_{io} reduces to a single
8 value, 0 or 1.

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

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

15 **T.3.11 Parameterization of time-based CQF**

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

19 **T.3.11.1 Cycle wander**

20 Adjacent Bridges must be frequency locked as described in T.3.1. For any given port, there is a worst-case
21 system clock difference between this Bridge's system clock and the neighbor system attached to the port. Its
22 units are a time difference. We will assume that this parameter is configured by management, based on
23 network design parameters and system data sheets. It is possible that this parameter can be adjusted during
24 network operation. A Bridge could have more than one system clock, and be connected to another system by
25 multiple links, but there is only one value for the difference for any given port, because we assume point-to-
26 point links. We will assume that the variation can be in either direction, this-end-late or this-end-early.

27 We assume that bin rotation operate under control of a clock that is local to a port. The management controls
28 that configure the rotation are defined in terms of a system clock. The Bridge can align the port clock(s) with
29 the system clock either periodically or continuously. There is thus a worst-case excursion of the actual start
30 of a cycle from the time configured in terms of the system clock. This feeds into the calculation of T_A in
31 T.3.2.1.

32 **T.3.11.2 Link delay variation**

33 The time taken for a frame to travel from the transmitter to the receiver can vary for two reasons: the actual
34 delay can change, due for example to temperature variations in a multi-kilometer link, and the measurement
35 of the link delay can vary due to various clock inaccuracies. We will deal only with actual variations, not
36 measurement variations.

1 T.3.11.3 Calculating the number of bins required

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

- 4 a) Establish a Nominal Input Cycle Start time (NICS) for the input port, and a Nominal Output Cycle
5 Start time (NOCS) for the output port. The NICS and NOCS each repeat every T_C seconds,
6 according to the system clock. We will assume that the offset between them is a constant (i.e., they
7 are both driven by the same system clock).
- 8 b) Compute the earliest time, relative to the NICS, at which the first frame of a cycle can receive its
9 IEEE Std 802.3 clause 90 timestamp. This frame is assumed to be a minimum-length frame (64
10 bytes plus overhead).
- 11 c) Compute the earliest time, relative to the NICS, at which a bin on the output port must be eligible to
12 receive the frame. This is equal to the timestamp time in bullet b) plus the minimum time required to
13 move the frame through the Bridge to the output bin.
- 14 d) Compute the latest time, relative to the NICS, at which the last frame of a cycle can receive its
15 timestamp. This frame is assumed to be a minimum-length frame.
- 16 e) If the difference between the earliest timestamp and the latest timestamp is greater than or equal to
17 the cycle time T_C , then dead time must be imposed on the transmitter, at the end of the cycle, to
18 reduce the difference.
- 19 f) Compute the latest time, relative to the NICS, at which the last frame of a cycle can be stored into an
20 output bin and be ready for selection for transmission, given the worst-case forwarding delay
21 through the Bridge.
- 22 g) Convert these earliest b) and latest d) arrival times to times relative to the NOCS of the output port.
- 23 h) Arbitrarily label an input port NICS event NICS0. Determine the latest subsequent NOCS event,
24 which we will label NOCS0, during which the earliest-arriving frame of NICS0 must be stored in
25 the output queue.
- 26 i) Determine the earliest subsequent NOCS event, which we will label NOCSn, before which the
27 latest-arriving frame from NICS0 can be stored in the queue, and still be available for transmission
28 at the start of cycle NOCSn.
- 29 j) The number of cycles NOCS0 through NOCSn, inclusive, is the number of bins required for the
30 input/output port pair, B_{io} .

31 The number of bins required can sometimes be reduced by:

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

36 Finally, let us observe that large link delay variations can be accommodated by varying the above
37 calculation. Assuming that the variations take place slowly, and that changes in relative phase between
38 transmitter and receiver are detected using a protocol (e.g. that in Clause 99), the difference between the
39 maximum and minimum link delay can be added to the difference between the earliest- and latest- arriving
40 frames to increase the number of bins allocated. The phase of the Stream gate can be altered by small
41 increments as the protocol detects the phase differences, without gaining or losing cycles in the transfer. Of
42 course, the maximum adjustment made per phase adjustment event must be removed from the allocable
43 bandwidth.

44 T.3.11.4 Initial bin phase

45 The number of bins required on an output port is the maximum required over all input ports. This may be
46 further increased by intentional delays (T.3.6). When initializing an input port, a correspondence must be

1 made between the input and output ports, so that a frame received on the input port will be stored in a
2 particular bin in the output port, the one that will become the transmitting bin in the appropriate number of
3 output cycles in the future.

4 The phasing between input and output ports' cycles, and thus the number of bins in port o used by port i , is
5 determined by the P_{io} table defined in T.3.10. We compute P_{io} when initializing CQF, or when the relative
6 phase of the input and output ports change significantly, by selecting a time T that coincides with the start of
7 an input cycle on input port i and computing:

$$8 P_{io} = (X_o - S_i - B_{io} + 1) \bmod N$$

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

13 T.3.11.5 Dead time / bandwidth balance calculation

14 There remains the balancing of conflicting goals between dead the percentage of a cycle that is available to
15 transmit critical data Streams, and the number of bins required on the output port. Increasing the dead time
16 can reduce the number of bins required, and thus the end-to-end latency of a data Stream, as described in
17 T.3.11.3. There are, at the very least, the following ways to make this decision:

- 18 — Configure the output cycle phase and number of bins to use for all Bridges, in order to establish a
19 constant per-hop delay in a network with short links. Let each system compute the dead time on each
20 input port required to make this work, and the bandwidth available for allocation. Convey the
21 required dead time either by protocol or by management to the transmitters, and the available
22 bandwidth to the admission control system.
- 23 — Configure the output cycle phase on all Bridges. Configure minimum and maximum allocable
24 bandwidth values for each CQF priority level. Let each system compute the minimum number of
25 bins required to meet the minimum bandwidth value, taking advantage of the maximum bandwidth
26 value to compute a dead time value that minimizes the number of bins required. This would be
27 useful in a network with very long links. Convey the resultant dead time to the transmitter via
28 protocol, and the resultant allocable bandwidth to the admission control system.
- 29 — Using data sheet information, configure all parameters via network management. Adjust the output
30 port cycle phasing to optimize the delay for certain specific Streams.

31 T.4 Count-based CQF

32 As described in T.3.1, time-based bin assignment assigns frames to output bins based on the time of arrival
33 of the frame, and requires that the output queues of successive hops along an CQF path run at exactly the
34 same frequency, in order to ensure that no bin's capacity can be exceeded. This requirement can be relaxed,
35 at the cost of implementing a state machine for each Stream passing through each output port. Then, the
36 output queues along the path can run nearly the same frequency, and their relative phases (T_{AB} in Figure T-4)
37 can diverge.

38 Count-based bin assignment (count-based CQF) is an alternative description of the paternoster algorithm
39 defined in [B1]. It provides a counter state machine for each Stream that allows that Stream to store no more
40 than its contracted amount of data per cycle into any given CQF bin. Frames above that limit are stored in
41 subsequent bins, up to the maximum amount of buffer space allowed that Stream, whereupon excess data is
42 discarded.

1 T.4.1 Calculating allocable time T_A

2 Count-based CQF computes the length of the portion of a cycle that can be allocated to Stream data, T_A , in a
3 manner similar to that used for time-based CQF in T.3.2.1 and Figure T-4:

$$4 T_A = T_C - T_I - T_P - T_D - T_X.$$

5 Again, T_C is the nominal cycle length, T_I is the interference from lower layers (T.3.2.1, T.3.3), T_P is the
6 penalty incurred if this priority level is preemptable (T.3.3), is the T_D is the dead time imposed by Bridge to
7 which this Bridge is transmitting (T.3.2.1). However, the time-based calculation uses T_V for the last time, a
8 catch-all for discrepancies including output delay, link delay, clock accuracy, and timestamp accuracy. The
9 count-based calculation uses:

10 T_X worst-case difference between the receiver's actual T_C values and the T_C value by which the
12 Talker's reservation is defined.

13 All of the items included in T_V in T.3.2.1 are irrelevant to count-based CQF:

- 14 a) Output delay, link delay, and clock accuracy affect only the phase relationship between the
15 transmitter's cycle the receiver's output ports' cycles. (In time-based CQF, there is no long-term
16 clock error.)
- 17 b) Timestamps are not relevant to count-based CQF.

18 Persistent differences in T_C , however, are important. If the transmitter's cycle time is shorter than the a
19 receiver cycle time, and if the Talker is generated data over the long term that keeps every transmitter cycle
20 full to the limit of a Stream's reservation, then the receiver would eventually have to drop frames. The term
21 T_X ensures that each count-based CQF hop can serve the Streams allocated to it. If a transmitting port is
22 faster than the receiving port, and thus builds up a extra frame in the receiver's bin(s), then in the long term,
23 even if the Talker runs continuously, the transmitter will eventually run ahead of the Talker, and have a less-
24 than-full cycle. This gives the net-hop receiver a chance to catch up.

25 T.4.2 Dead time T_D

26 Count-based CQF in a receiving Bridge cannot impose dead time (T.3.5) on the transmitting Bridge; it has
27 no need to. However, it may have dead time imposed upon it if it transmits to Bridge using time-based CQF.

28 T.4.3 Number of output bins

29 In an ideal world, only two bins are required per output queue for count-based CQF, one filling and one
30 transmitting. The fact that successive Bridges employing count-based CQF have slightly different actual
31 values for T_C makes a third bin necessary, because all of the frames destined for one bin (the one that is
32 filling) do not necessarily all arrive during the time when the filling bin is open. If the forwarding delay
33 variation shown in Figure T-4 is non-0, as it is in most implementations, at least one more bin, the fourth, is
34 necessary. Additional bins can be added to accommodate input Streams from devices that are not
35 transmitting using CQF. Assuming that such input Streams can be characterized by a committed burst size
36 (§8.6.5.5), this committed burst size can be added to the basic three bins to calculate the total number of bins
37 required.

38 T.4.4 Using both count-based and time-based CQF

39 A Stream entering a Bridge from a correctly-configured Bridge or end station that runs CQF, and that has
40 reserved bin space allocated for it, will not disrupt the deterministic behavior of CQF. However, an CQF
41 Bridge could receive input from a Bridge, a Talker, a router, or any other device that uses some deterministic

1 algorithm(s) to condition its Streams, but uses an algorithm other than CQF. We assume that reservations
2 (contracts) for these Streams can be translated into CQF terms, with perhaps some overprovisioning
3 required. In the long term, the Talker adheres to the contract, and will not disrupt determinism. But, since the
4 transmitter is not using CQF, we cannot use just a frame's arrival time to assign it to a bin, except by
5 overprovisioning CQF sufficiently to accommodate the worst-case burst behavior of the algorithm
6 employed by the sender. We would like to accommodate such input.

7 The count-based bin assignment makes this possible. A Bridge uses the same bin structure and output
8 methods described for time-based bin assignment in T.3.1, but instead of obtaining the bin selector S from
9 the time of receipt of the input frame, as described in T.3.10, it uses a state machine dedicated to each CQF
10 Stream using count-based bin assignment, to determine the bin selector. Extra bins are provided to accept
11 such bursts. At the next hop, time-based CQF can be used.

12 A given output bin can accept input from both time-based and count-based Streams, as long as they share the
13 same cycle time; separate count-based and time-based bins or queues are not necessary. In addition, a
14 paranoid network administrator could very well configure count-based bin assignment on every Stream in a
15 frequency-locked network, in order to guard against misbehaving Bridges or Talkers. That is, while count-
16 based bin assignment can be thought of as separate algorithm from time-based bin assignment, it can also be
17 thought of as a protection mechanism for time-based assignment that can be employed as need, and when
18 employed everywhere, removes the restriction that Bridges operate at exactly the same frequency.

19 In many networks, count-based and time-based bin assignment can be used at the same priority level in one
20 network. The choice between count-based and time-based bin assignment can be made on a Bridge-by-
21 Bridge basis, and not be visible to the Talker, the Listener, or the user.

22 T.5 Stream Aggregation

23 Stream aggregation is useful for both scaling up the number of Streams that a network can support, and for
24 decreasing the end-to-end latency of Streams that are aggregated. In this type of aggregation, a number of
25 Streams are treated as a single Stream, with a single reservation, traversing a single path, for some portion of
26 their journey through the network. Stream aggregation can work whether the frames are encapsulated in
27 some common wrapper, or whether they are simply treated identically (e.g. all given the same IEEE Std
28 802.1CB stream_identifier).

29 This standard does not specify any protocol for encapsulating aggregated Streams.

30 The aggregated Stream has a single reservation that is the union its component Streams' reservations (T.6.1).
31 Ideally, this higher-bandwidth Stream can be assigned a CQF priority level with a faster T_C than its
32 components can use. For example, instead forwarding 10 Streams, each with one frame, in a 1 millisecond
33 T_C bin, a Bridge could be forwarding one aggregated Stream that has one frame in a 100 microsecond bin,
34 Buffer space requirements are cut by 90%, per-hop delay by up to 90%, and state machines, e.g. count-based
35 bin assignment machines, by 90%.

36 In general, this requires that the aggregate Stream pass through a count-based bin assignment state machine
37 when it is formed from its components, and that each component pass through a count-based bin assignment
38 state machine if and when it is again separated as an individual Stream and passes, presumably, to a slower
39 T_C value.

40 Figure T-5 illustrates the value and the limitations of Stream aggregation. In this figure, there are three
41 Streams, all entering Bridge A, and all three traversing the same path at least as far as Bridge E. Bridge A
42 operates on three separate Streams in the usual manner for CQF, placing one frame from each stream into
43 each output bin. In Bridge B, the three Streams are aggregated into a single aggregated Stream. This Stream

1 is forwarded through Bridges C and D. Bridge E dissolves the aggregated Stream, distributing the
2 component Streams' frames on three ports.

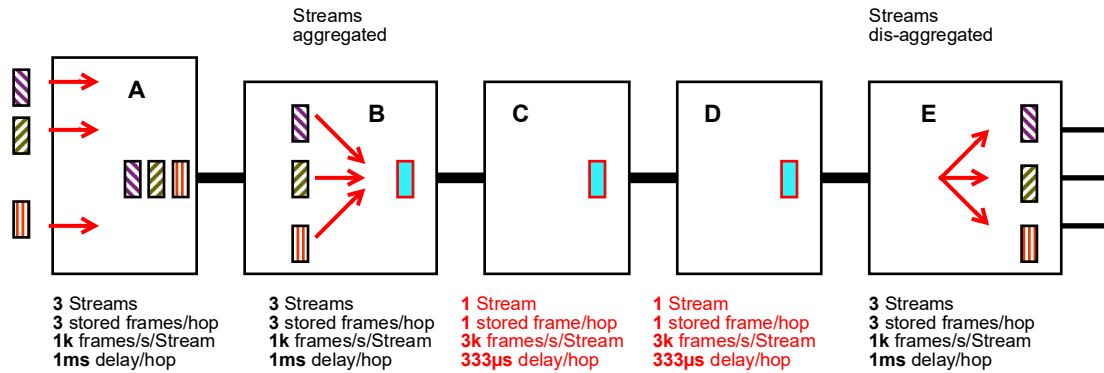


Figure T-5—CQF Stream Aggregation example

3 As indicated in the captions in Figure T-5, the aggregating Bridge B requires the same amount of buffer
4 space and state machines, and imposes the same forwarding delay on all three Streams, as does Bridge A
5 (see T.5.1). Once aggregated, however, Bridge C and Bridge D forward the Streams with less delay, and
6 using less buffer space and fewer state machines, than Bridges A and B. Finally, Bridge E dissolves the
7 aggregated Stream into its components (see T.5.2). Bridge E, again requires the same resources imposes the
8 same delay on the three Streams as the non-aggregated Bridge A. The case of Stream mixing, that is,
9 disaggregating Streams and then re-aggregating them in different combinations in the same Bridge, is
10 described in T.5.4.

11 T.5.1 CQF Stream aggregation

12 Stream aggregation usually includes changing the cycle time from a slower to a faster cycle. Changing the
13 cycle time is discussed in T.3.8. In general, count-based bin assignment is required, because the Streams
14 being aggregated can be arriving from different ports, and be distributed along the duration of a large cycle
15 time. On the port outputting the aggregated Stream, the bins are typically rotated at a higher rate, and often,
16 are on a medium with higher bandwidth than the media on which the Streams arrived. Typically, the input
17 Streams' frames can arrive in lumps, with intervening gaps, due to fan-in (multiple input ports to a single
18 output port). In the worst case, the same number of output bins are required on the output port as would have
19 been needed had there been no aggregation.

20 T.5.2 CQF Stream disaggregation

21 Stream disaggregation usually includes changing the cycle time from a faster to a slower cycle. Changing
22 the cycle time is discussed in T.3.8..

23 Whether count-based or time-based bin assignment can be used in the disaggregating bridge depends on the
24 degree to which the Bridges along the path of the aggregated Stream are able to maintain the timing of the
25 aggregating Bridge's inputs. In theory, if time-based CQF is used all along the aggregates Stream's path, and
26 if the disaggregating Bridge can align its input cycles with the aggregating Bridge's input cycles, in spite of
27 the intervening Bridges, then it would be possible to us time to select the output buffers in the disaggregating
28 Bridge. No mechanism is provided in this standard to achieve that synchrony.

29 Assuming that the disaggregating Bridge uses count-based bin assignment, the frames of the individual
30 disaggregated Streams are distributed into the slow output bins in the normal manner.

1 The number of bins required in the disaggregating Bridge may be slightly larger than the number of bins
2 required for normal forwarding of the equivalent disaggregated Streams, because the Bridges carrying the
3 aggregated Stream (e.g. Bridges C and D in Figure T-5) can introduce a small amount of delay variation.
4 Any such variation increases the buffer requirements in the disaggregating Bridge. See T.5.3.

5 **T.5.3 CQF delay variation and disaggregation buffers**

6 One of the major issues with aggregating Streams is the amount of buffer space required at the
7 disaggregation point. For example, suppose that Bridges B, C, and D in Figure T-5 use Asynchronous
8 Traffic Shaping, instead of CQF, to forward the aggregated Stream. Because the aggregated Stream is
9 allocated essentially the sum of the bandwidths of its component Streams, if some Streams in the
10 aggregation are flowing intermittently, then other Streams can be given the bandwidth not used by their
11 associates in the aggregation. This can easily result in bursts and gaps in the individual Streams when finally
12 delivered to the disaggregation Bridge (T.5.2) or mixing Bridge (T.5.4).

13 Count-based CQF makes a critical reduction in the delay variation of the aggregated Stream, even when
14 components of the stream are intermittent or missing, as long as the number of bins in each hop is
15 minimized. Time-based CQF adds no delay variation; the aggregating Bridge's bin assignments are
16 maintained all the way to the disaggregating Bridge. The disaggregating Bridge performs count-based bin
17 assignment, but requires excess buffer space only in amount equal to the size of the input bins on the
18 aggregating Bridge (Bridge A in Figure T-5), as its output cycles are not necessarily in phase with the
19 aggregating Bridge's input cycles.

20 **T.5.4 CQF Stream mixing**

21 Stream mixing occurs when Streams are disaggregated and then reaggregated, perhaps in different
22 combinations, in the same Bridge. If we neglect the actual operations of wrapping and unwrapping the
23 frames (if needed), this operation is very similar to either aggregation or disaggregation; mixing does *not*
24 require a set of disaggregation buffers followed by a set of aggregation buffers. Rather,

25 In general, any Stream that does not pass intact through a Bridge has to be split into its component Streams.
26 (Split at the topmost level—aggregations of aggregations need not be split all the way down the stack). Each
27 component Stream passes through a count-based bin assignment state machine, and is assigned an output bin
28 appropriate for its output port, whether or not it is also being reaggregated. As for the case for
29 disaggregation (T.5.2) the output buffer space required depends on the cycle times of the original,
30 component Streams at the time they entered the aggregation.

31 **T.6 Additional considerations**

32 **T.6.1 Computing the CQF reservation**

33 At the lowest level, e.g. in a count-base bin assignment state machine, CQF reservations are in terms of bit
34 times per cycle, with an implied cycle time. However, Streams are not usually characterized in this manner.

35 << Editor's note: This section is clearly incomplete. It will include a description of how to convert TSN and
36 CBS specs to CQF, and also how to combine multiple Streams' TSN and/or CBS specs into a single CQF
37 spec. >>

1 T.6.2 Frame size problem

2 Stream data does not always consist of frames that are all the same size. The advantage of uniform frame
3 size is that, in the ideal case, one can allocate a Stream one frame per cycle, and choose the cycle time and/or
4 the Stream's bandwidth reservation so that there is no wasted bandwidth. Similarly, if we imagine that a
5 Stream alternates frames of 4 000 bit times and 800 bit times, we can allocate 4800 bit times per T_C and still
6 get perfect results.

7 But, in a service provider situation where we are allocating a certain bandwidth per customer, but the frame
8 sizes are essentially random, things are not so simple. Let us suppose that the maximum frame for a Stream
9 is 13 000 bit times, which is approximately equal to a maximum-length Ethernet frame, and that the cycle
10 time $T_C = 100\mu\text{s}$. $13\,000/100\mu\text{s} = 130\text{ Mb/s}$. But, allocating a bandwidth of 13 000 bits/ T_C will not give
11 the Stream 130 Mb/s. In the worst case, one 13 000 bit frame followed by one minimum-length frame = 672
12 bits, the Stream gets $(13\,000+672)/(200\mu\text{s}) = 68.36\text{ Mb/s}$.

13 We could overprovision the Stream by a factor of almost 2, keep the same T_C , and get minimal latency.
14 However, we could also assign the Stream to a longer T_C . In the worst case, there are $(13\,000-8)$ wasted bits
15 in each cycle. Therefore, we can guarantee 130 Mb/s using a cycle time of $500\mu\text{s}$ by provisioning $(5*13\,000$
16 $+ 13\,000 - 8)/(5*100\mu\text{s})$, or 156 Mb/s, which is a 20% overprovisioning, rather than a 90%
17 overprovisioning, at the cost of five times the per-hop latency.

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

22 Another approach to increased resource utilization efficiency is to run each CQF Stream, on ingress to the
23 TSN network, through a "sausage maker". That is, frames can be encapsulated using a scheme that
24 combines and/or splits frames into uniform-sized chunks (sausages), either small or large, that can be carried
25 end-to-end through the TSN network, then split out into their original form. This means that
26 overprovisioning due to the mix of frame sizes is reduced to that required by the encapsulation, itself. (In
27 fact, that overhead can be negative, if small frames are aggregated² into large transmission units.)

28 T.6.3 Tailored bandwidth offerings

29 We can note that, in a service provider environment, overprovisioning can be almost eliminated by a
30 combination of 1) Stream Aggregation (T.5) and 2) offering the customer only a specific set of choices for a
31 bandwidth contract, corresponding to the values of T_C implemented in the provider's network.

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

36 T.6.4 Overprovisioning to improve latency

37 A minimum of network resources is consumed when a Stream with constant frame size is allocated just
38 enough bit times per cycle for a single frame, and the cycle time is $1/(\text{the Stream's frame rate})$. One frame is
39 delivered per cycle.

²Not to be confused with Link Aggregation

1 If that same number of bit times is reserved for that Stream in a class of service with a T_C that is, for
2 example, five times shorter than the minimum-resource T_C , then that Stream will experience shorter end-to-
3 end delay through the network. This lower delay comes at a cost; four out of five cycles have unused
4 allocable transmission time. That bandwidth is available to best-effort traffic, but not to other Streams with
5 reserved resources.

6 In a network that is not saturated with Stream traffic, this can be a viable trade-off.

7 **T.6.5 CQF and credit-based shaper**

8 Looking at Figure T-2, we see that, once the major cycle at priority level 4 begins transmitting, the best-
9 effort traffic is interrupted until all of the CQF level 4 data is transmitted. At some point, as the amount of
10 traffic in a very slow CQF cycle increases, the burstiness of the best-effort transmission opportunities could,
11 in theory, become a problem. This can be mitigated by applying a credit-based shaper function to the slowest
12 CQF cycle(s). However, the parameters of this shaper must be adjusted as the load on the slow CQF cycle(s)
13 changes, because a Bridge must always finish transmitting all of the data in a bin. Thus, adding a credit-
14 based shaper would detract from a significant advantage of time-based CQF—its freedom from requiring
15 reconfiguring a Bridge each time a Stream is added.

16 **T.6.6 Interactions among CQF, ATS, and control traffic**

17 For time-based CQF to function correctly—in particular, for it to guarantee no congestion loss—all of the
18 frames in a bin have to be transmitted before the beginning of the cycle. The formulas for buffer allocation
19 and end-to-end delay depend on this. Therefore, to configure queues with ATS shapers at priority levels more
20 important than a CQF queue, one must be sure that, in the worst case, the CQF queue can still empty its bin
21 within the allotted period (T_A in T.3.2.1).

22 These control protocols can include both link-local protocols, such as LLDP or BPDUs (<insert reference>),
23 and protocols likely to be forwarded as data by a Bridge, such as Layer 3 routing protocols.

24 Similarly, steps have to be taken by the implementer and/or network manager to understand high-priority
25 control traffic such as bridging or routing protocols. Typically, this means creating one or more Stream
26 reservations to control the impact of control protocols on Stream data. Of course, the impact of the other
27 Streams on the control protocols also has to be analyzed. Experience with TSN tends to show that such
28 protocols, not having been regulated carefully before the advent of deterministic networking, are quite
29 tolerant of the levels of delay imposed by transforming them from conflicting, highest-priority frames, to
30 them reasonably high-priority, but bandwidth limited, Streams.

31 ‘

¹ **Annex ZY**

² (informative)

³ **Bibliography**

⁴ This chapter should include the .1Q bibliography in total, with the proper conditional text tags.

⁵ *Insert the following references in the appropriate collating sequence and renumber accordingly:*

⁶ [B1] Seaman, Mick, “Paternoster policing and scheduling” [http://www.ieee802.org/1/files/public/](http://www.ieee802.org/1/files/public/docs2019/cr-seaman-paternoster-policing-scheduling-0519-v04.pdf)
⁷ [docs2019/cr-seaman-paternoster-policing-scheduling-0519-v04.pdf](http://www.ieee802.org/1/files/public/docs2019/cr-seaman-paternoster-policing-scheduling-0519-v04.pdf),

⁸ [B2] Finn, Norman, [https://mentor.ieee.org/802.1/dcn/21/1-21-0056-00-ICne-input-synchronization-for-](https://mentor.ieee.org/802.1/dcn/21/1-21-0056-00-ICne-input-synchronization-for-cyclic-queueing-and-forwarding.pdf)
⁹ [cyclic-queueing-and-forwarding.pdf](https://mentor.ieee.org/802.1/dcn/21/1-21-0056-00-ICne-input-synchronization-for-cyclic-queueing-and-forwarding.pdf)

¹⁰

¹ **Annex ZZ**

² (informative)

³ **Commentary**

⁴ << Editor's Note: This is a temporary Annex intended to record issues and their resolutions as the project
⁵ proceeds. It will be removed prior to Standards Association ballot. >>

⁶