

Residential Ethernet (RE) (a working paper)

The following paper represents an initial attempt to codify the content of multiple IEEE 802.3 Residential Ethernet (RE) Study Group slide presentations. The author has also taken the liberty to expand on various slide-based proposals, with the goal of triggering/facilitating future discussions.

For the convenience of the author, this paper has been drafted using the style of IEEE standards. The quality of the figures and the consistency of the notation should not be confused with completeness of technical content.

Rather, the formality of this paper represents an attempt by the author to facilitate review by interested parties. Major changes and entire clause rewrites are expected before consensus-approved text becomes available.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Residential Ethernet (RE) (a working paper)

Draft 0.136

Contributors:
See page 4.

Abstract: This working paper provides background and introduces possible higher level concepts for the development of Residential Ethernet (RE).

Keywords: residential, Ethernet, isochronous, real time

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Contributors

This working paper is based on contributions or review comments from the people listed below. Their listing doesn't necessarily imply they agree with the entire content or the author's interpretation of their input.

Jim Battaglia	Pioneer
Alexei Beliaev	Gibson
Dirceu Cavendish	NEC Labs America
George Claseman	Micrel
Feifei (Felix) Feng	Samsung Electronics
John Nels Fuller	Independent
Geoffrey M. Garner	Samsung Electronics
Kevin Gross	Cirrus Logic
Jim Haagen-Smit	HP
David V James	JGG
Dennis Lou	Pioneer
Michael D. Johas Teener	Broadcom
Fred Tuck	EchoStar

Version history

Version	Date	Author	Comments
0.082	2005Apr28	DVJ	<ul style="list-style-type: none"> Updates based on 2005Apr27 meeting discussions – Restructure document presentation order – Provide list of contributors, with appropriate disclaimer – Provide version history, for convenience of frequent reviewers – Fix page numbering for easy review (continuous count from start) – Fix clause numbering cross-reference bug (period after number) – Urban recording session (see 5.1.4) added for completeness – Conflicting traffic (see 5.1.5) added for completeness – Changed ‘ping’ to ‘refresh’, within the context of SRP – Changes the multicast addressing for classA frames – Refined state machines
0.085	2005May11	DVJ	<ul style="list-style-type: none"> – Updated front-page list of contributors – Updated book for continuous pages (Clause 1 discontinuity fixed) – Miscellaneous editing fixes – Initial pinging description added. – Previous Clause 9 (identifier assignments) moved to format clause. – The <i>subType</i> identifier assignments now specified. – The bunching annex (work in progress) now includes: <ul style="list-style-type: none"> A more typical age-based classA prioritization assumption. Other parameters of interest (idle and full-load durations). (Further thought on queue sizing, to avoid discards, is needed.)
0.088	2005Jun03	DVJ	<ul style="list-style-type: none"> – Application latency scenarios clarified. Generalized based on Norm Finn concerns. Clarified/corrected based on Kevin Gross comments. – Subscription revised, to converge with Felix presentation. – Bursting and bunching scenarios revised for applicability and clarity.
0.090	2005Jun06	DVJ	– Misc editorials in bursting and bunching annex.
0.092	2005Jun10	DVJ	– Extensive cleanup of Clause 5 subscription protocols, based on 2005Jun08 teleconference review comments.

Version	Date	Author	Comments
0.121	2005Jun24	DVJ	<ul style="list-style-type: none"> – Extensive cleanup of clock-synchronization protocols, base on 2005Jun22 teleconference review comments. Affected areas include: <ul style="list-style-type: none"> Subclause 5.1: Revised, based comments from Alexei Subclause 5.5: Time-synchronization overview updated Clause 7: Time-synchronization descriptions added Note that the state machines have now become obsolete. Annex J: Time-synchronization code added
0.125	2005Jun30	DVJ	<ul style="list-style-type: none"> – Grand-master description provided in 5.5.4. – Clock deviation moved from code to state machine. – Clock-synchronization code enhanced and split into distinct core (one per station) and port components. – Code cleanups, corresponding to the above.
0.127	2005Jul04	DVJ	<ul style="list-style-type: none"> – Pacing descriptions greatly enhanced. Miscellaneous error/clarity fixes, primarily clock related 5.7—A better overview provided Clause 9—Detailed state machines provided Synchronized time-of-day clock/Limitations of current approaches has been migrated to Annex D, where other alternatives are listed.
0.133	2005Jul12	DVJ	<ul style="list-style-type: none"> – Update of contributors list. – Pacing assumptions/objectives/strategies added to Clause 5. – Pacing state machine in Clause 8 simplified: migrated shapers into transmit code, eliminating shapers A,B,C merged Transmit100Mbs and Transmit1Gbs into TransmitTx
0.134	2005Jul17	DVJ	<ul style="list-style-type: none"> – Pacing disclaimers, based on preceding meeting discussions. classB could be found to be unnecessary 3-cycle to 2-cycle delay reduction may not be worth its complexity. – Miscellaneous fixups. – Additions to pacing, resisions of state machines. – Excess material deletions.
0.135	2005Jul30	DVJ	<ul style="list-style-type: none"> – Alternative rate-based pacing protocol provided. – The stream-addressing alternatives have been better organized.
—	TBD	—	—

Background

This working paper is highly preliminary and subject to changed. Comments should be sent to its editor:

David V. James
3180 South Ct
Palo Alto, CA 94306
Home: +1-650-494-0926
Cell: +1-650-954-6906
Fax: +1-360-242-5508
Email: dvj@alum.mit.edu

Formats

In many cases, readers may elect to provide contributions in the form of exact text replacements and/or additions. To simplify document maintenance, contributors are requested to use the standard formats and provide checklist reviews before submission. Relevant URLs are listed below:

General: <http://grouper.ieee.org/groups/msc/WordProcessors.html>
Templates: <http://grouper.ieee.org/groups/msc/TemplateTools/FrameMaker/>
Checklist: <http://grouper.ieee.org/groups/msc/TemplateTools/Checks2004Oct18.pdf>

Topics for discussion

Readers are encouraged to provide feedback in all areas, although only the following areas have been identified as specific areas of concern.

- a) Terminology. Is classA an OK way to describe the traffic within an RE stream?
Alternatives:
synchronous traffic? isochronous traffic? RE traffic? quasi-synchronous traffic?

TBDs

Further definitions are needed in the following areas:

- a) ClassA addressing models: review, select, and revise.
- b) Pacing models: review, select, and revise.

Contents		1
		2
List of figures.....	10	3
		4
List of tables.....	13	5
		6
1. Overview.....	15	7
		8
1.1 Scope and purpose.....	15	9
1.2 Introduction	15	10
		11
2. References.....	19	12
		13
3. Terms, definitions, and notation	20	14
		15
3.1 Conformance levels.....	20	16
3.2 Terms and definitions	20	17
3.3 Service definition method and notation.....	22	18
3.4 State machines	23	19
3.5 Arithmetic and logical operators	26	20
3.6 Numerical representation.....	26	21
3.7 Field notations	27	22
3.8 Bit numbering and ordering.....	28	23
3.9 Byte sequential formats	29	24
3.10 Ordering of multibyte fields	29	25
3.11 MAC address formats.....	30	26
3.12 Informative notes.....	31	27
3.13 Conventions for C code used in state machines	31	28
		29
4. Abbreviations and acronyms	32	30
		31
5. Architecture overview	33	32
		33
5.1 Latency constraints.....	33	34
5.2 Service classes	36	35
5.3 Architecture overview	37	36
5.4 Subscription.....	39	37
5.5 Synchronized time-of-day clocks	48	38
5.6 Formats	52	39
5.7 Pacing	56	40
		41
6. Frame formats.....	59	42
		43
6.1 ClassA frames.....	59	44
6.2 clockSync frame format	60	45
6.3 Subscription frame.....	63	46
6.4 Common <i>info</i> field format.....	64	47
6.5 Unique identifier values	65	48
		49
7. Clock synchronization	66	50
		51
7.1 Clock-synchronization overview.....	66	52
7.2 Terminology and variables.....	75	53
7.3 Clock synchronization state machines.....	76	54

8. Subscription state machines.....	80	1
		2
8.1 Terminology and variables.....	80	3
8.2 Subscription state machines	81	4
		5
9. Transmit state machines (proposal 1).....	92	6
		7
9.1 Pacing overview	92	8
9.2 Terminology and variables.....	99	9
9.3 Pacing state machines.....	100	10
		11
10. Transmit state machines (proposal 2).....	107	12
		13
10.1 Rate-based scheduling overview	107	14
10.2 Terminology and variables.....	112	15
10.3 Pacing state machines.....	113	16
		17
Annex A (informative) Bibliography.....	121	18
		19
Annex B (informative) Background material	122	20
		21
Annex C (informative) Encapsulated IEEE 1394 frames	127	22
		23
C.1 Hybrid network topologies.....	127	24
C.2 1394 isochronous frame formats	128	25
C.3 Frame mappings	130	26
C.4 CIP payload modifications	131	27
		28
Annex D (informative) Review of possible alternatives.....	134	29
		30
D.1 Clock-synchronization alternatives	134	31
D.2 Pacing alternatives.....	135	32
D.3 IEEE 1394 alternative.....	136	33
		34
Annex E (informative) Time-of-day format considerations	137	35
		36
E.1 Possible time-of-day formats.....	137	37
E.2 Time format comparisons.....	139	38
		39
Annex F (informative) Bursting and bunching considerations.....	140	40
		41
F.1 Topology scenarios.....	140	42
F.2 Bursting considerations	142	43
		44
Annex G (informative) Denigrated alternatives.....	160	45
		46
G.1 Stream frame formats	160	47
G.2 Subscription.....	162	48
		49
Annex H (informative) Frequently asked questions (FAQs)	169	50
		51
H.1 Unfiltered email sequences.....	169	52
H.2 Formulated responses	170	53
		54

Annex I (informative) Comment responses.....	171	1
		2
I.1 Recent review-comment resolutions	171	3
		4
Annex J (informative) C-code illustrations.....	175	5
		6
Index	185	7
		8
		9
		10
		11
		12
		13
		14
		15
		16
		17
		18
		19
		20
		21
		22
		23
		24
		25
		26
		27
		28
		29
		30
		31
		32
		33
		34
		35
		36
		37
		38
		39
		40
		41
		42
		43
		44
		45
		46
		47
		48
		49
		50
		51
		52
		53
		54

List of figures

	1
	2
Figure 1.1—Topology and connectivity	3
Figure 3.1—Service definitions	4
Figure 3.2—Bit numbering and ordering	5
Figure 3.3—Byte sequential field format illustrations	6
Figure 3.4—Multibyte field illustrations	7
Figure 3.5—Illustration of fairness-frame structure	8
Figure 3.6—MAC address format	9
Figure 3.7—48-bit MAC address format.....	10
Figure 5.1—Interactive audio delay considerations	11
Figure 5.2—Home recording session	12
Figure 5.3—Garage jam session.....	13
Figure 5.4—Urban recording session	14
Figure 5.5—Conflicting data transfers	15
Figure 5.6—Hierarchical control.....	16
Figure 5.7—Hierarchical flows	17
Figure 5.8—Controller activation.....	18
Figure 5.9—Agents on an established path	19
Figure 5.10—Periodic registration messages	20
Figure 5.11—Secondary registrations	21
Figure 5.12—Side-path deregistration.....	22
Figure 5.13—Final-path deregistration.....	23
Figure 5.14—Streaming data over registered paths.....	24
Figure 5.15—Insufficient bandwidth conditions	25
Figure 5.16—Periodic registration messages	26
Figure 5.17—Timer synchronization flows.....	27
Figure 5.18—Grand-master precedence flows	28
Figure 5.19—Time synchronization principles	29
Figure 5.20—Timer snapshot locations.....	30
Figure 5.21—Bridge PLL possibilities	31
Figure 5.22—Content framing methods	32
Figure 5.23—Plug addressing.....	33
Figure 5.24—ClassA frame format and associated data.....	34
Figure 5.25—ClassA frame format and associated data.....	35
Figure 5.26—ClassA frame formats	36
Figure 6.1—ClassA frame formats	37
Figure 6.2—clockSync frame format	38
Figure 6.3— <i>systemTag</i> subfields.....	39
Figure 6.4— <i>uniqueID</i> format	40
Figure 6.5—Complete seconds timer format.....	41
Figure 6.6—Subscription frame format.....	42
Figure 6.7—Common <i>info</i> field format.....	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

Figure 6.8— <i>protocolType</i> field value.....	65	1
Figure 7.1—Grand-master precedence	66	2
Figure 7.2—Hierarchical flows	67	3
Figure 7.3—Offset synchronization	69	4
Figure 7.4—Cascaded offsets (a possible scenario)	70	5
Figure 7.5—Rate synchronization	71	6
Figure 7.6—Cascaded rate differences (a possible scenario)	72	7
Figure 7.7—Rate-adjustment effects	73	8
Figure 7.8— <i>flexTimer</i> implementation example	74	9
Figure 7.9— <i>baseTimer</i> implementation example	75	10
Figure 9.1—Topology-dependent pacing delays.....	92	11
Figure 9.2—Paced 1 Gb/s classA flows	93	12
Figure 9.3—Cycle slippage	94	13
Figure 9.4—Paced 100 Mb/s classA flows.....	94	14
Figure 9.5—Cycle slippage	95	15
Figure 9.6—Transmit-port structure.....	96	16
Figure 9.7—Pacing at 1 Gb/s.....	97	17
Figure 9.8—Pacing at 100 Mb/s	98	18
Figure 9.9—Credit adjustments over time.....	98	19
Figure 10.1—Rate-based priorities	107	20
Figure 10.2—Reshaped bridge-traffic topology	109	21
Figure 10.3—Reshaped bridge-traffic timing.....	109	22
Figure 10.4—Transmit-queue structure.....	110	23
Figure 10.5—Credit-based shapers.....	111	24
Figure 10.6—Pacer credit adjustments over time.....	112	25
Figure B.1—SerialBus topologies	122	26
Figure B.2—Isochronous data transfer timing	123	27
Figure B.3—RPR rings	124	28
Figure B.4—RPR resilience	125	29
Figure B.5—RPR destination stripping	125	30
Figure B.6—RPR spatial reuse	126	31
Figure B.7—RPR service classes	126	32
Figure C.1—IEEE 1394 leaf domains	127	33
Figure C.2—IEEE 802.3 leaf domains	127	34
Figure C.3—IEEE 1394 isochronous packet format	128	35
Figure C.4—Encapsulated IEEE 1394 frame payload	128	36
Figure C.5—Conversions between IEEE 1394 packets and RE frames.....	130	37
Figure C.6—Multiframe groups	131	38
Figure C.7—Isochronous 1394 CIP packet format	131	39
Figure C.8—Time-of-day format conversions	132	40
Figure C.9—Grand-master precedence mapping	133	41
Figure 5.1—Complete seconds timer format.....	137	42
Figure E.2—IEEE 1394 timer format.....	137	43
		44
		45
		46
		47
		48
		49
		50
		51
		52
		53
		54

Figure E.3—IEEE 1588 timer format.....	138	1
Figure E.4—EPON timer format.....	138	2
Figure E.5—Compact seconds timer format	138	3
Figure E.6—Nanosecond timer format.....	138	4
Figure F.1—Bridge design models	140	5
Figure F.2—Three-source topology	141	6
Figure F.3—Six-source topology	141	7
Figure F.4—Three-source bunching timing; input-queue bridges	142	8
Figure F.5—Cumulative coincidental burst latencies.....	143	9
Figure F.6—Three-source bunching; input-queue bridges.....	144	10
Figure F.7—Six source bunching timing; input-queue bridges.....	145	11
Figure F.8—Cumulative bunching latencies; input-queue bridge.....	146	12
Figure F.9—Three-source bunching; output-queue bridges.....	147	13
Figure F.10—Six source bunching; output-queue bridges	148	14
Figure F.11—Cumulative bunching latencies; output-queue bridge.....	149	15
Figure F.12—Three-source bunching; variable-rate output-queue bridges.....	150	16
Figure F.13—Six source bunching; variable-rate output-queue bridges.....	151	17
Figure F.14—Cumulative bunching latencies; variable-rate output-queue bridge.....	152	18
Figure F.15—Three-source bunching; throttled-rate output-queue bridges	153	19
Figure F.16—Six source bunching; throttled-rate output-queue bridges	154	20
Figure F.17—Cumulative bunching latencies; throttled-rate output-queue bridge.....	155	21
Figure F.18—Three-source bunching; throttled-rate output-queue bridges.....	156	22
Figure F.19—Three-source bunching; throttled-rate output-queue bridges.....	157	23
Figure F.20—Six source bunching; classA throttled-rate output-queue bridges.....	158	24
Figure F.21—Cumulative bunching latencies; classA throttled-rate output-queue bridge	159	25
Figure G.1—classA frame formats	160	26
Figure G.2—ClassA frame formats	161	27
Figure G.3—Agents on an established path	163	28
Figure G.4—Controller activation	164	29
Figure G.5—Pinging the talker.....	164	30
Figure G.6—Path creation	165	31
Figure G.7—Side-path extensions	165	32
Figure G.8—Side-path demolition	166	33
Figure G.9—Released path.....	166	34
Figure G.10—Error responses	167	35
Figure G.11—Side-path demolition	168	36
		37
		38
		39
		40
		41
		42
		43
		44
		45
		46
		47
		48
		49
		50
		51
		52
		53
		54

List of tables

	1
	2
Table 3.1—State table notation example	3
Table 3.2—Called state table notation example	4
Table 3.3—Special symbols and operators.....	5
Table 3.4—Names of fields and sub-fields	6
Table 3.5— <i>wrap</i> field values	7
Table 5.1—Service classes and their quality-of-service relationships	8
Table 5.2—Tagged priority values	9
Table 6.1—Assigned <i>subType</i> identifiers.....	10
Table 7.1—External clock-synchronization pairs	11
Table 7.2—Clock-synchronization intervals	12
Table 7.3—ClockCore state table.....	13
Table 7.4—ClockPort state table.....	14
Table 8.1—AgentAction state table	15
Table 8.2—AgentTalker state table.....	16
Table 8.3—AgentTimer state table	17
Table 8.4—AgentListener state table	18
Table 9.1—ClockPort state table.....	19
Table 9.2—ReceiveRx state table	20
Table 9.3—TransmitTx state table	21
Table 10.1—Tagged priority values	22
Table 10.2—TransmitRx state table.....	23
Table 10.3—TransmitTx state table	24
Table C.1— <i>flag</i> field values	25
Table C.2— <i>counts</i> field values.....	26
Table E.1—Time format comparison	27
Table F.1—Cumulative bursting latencies	28
Table F.2—Cumulative bunching latencies; input-queue bridge	29
Table F.3—Cumulative bunching latencies; output-queue bridge	30
Table F.4—Cumulative bunching latencies; variable-rate output-queue bridge	31
Table F.5—Cumulative bunching latencies; throttled-rate output-queue bridge	32
Table F.6—Cumulative bunching latencies; classA throttled-rate output-queue bridge.....	33
	34
	35
	36
	37
	38
	39
	40
	41
	42
	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Residential Ethernet (RE) (a working paper)

This document and has no official status within IEEE or alternative SDOs.

Feedback to: dvj@alum.mit.edu

(See page 4 for the list of contributors.)

1. Overview

1.1 Scope and purpose

This working paper is intended to supplement Ethernet with real-time capabilities, with the scope and purpose listed below:

Scope: Residential Ethernet provides time-sensitive delivery between plug-and-play stations over reliable point-to-point full-duplex cable media. Time-sensitive data transmissions use admission control negotiations to guarantee bandwidth allocations with predictable latency and low-jitter delivery. Device-clock synchronization is also supported. Ensuring real-time services through routers, data security, wireless media, and developing new PMDs are beyond the scope of this project.

Purpose: To enable a common network for existing home Ethernet equipment and locally networked consumer devices with time-sensitive audio, visual and interactive applications and musical equipment. This integration will enable new applications, reduce overall installation cost/complexity and leverage the installed base of Ethernet networking products, while preserving Ethernet networking services. An appropriately enhanced Ethernet is the best candidate for a universal home network platform.

1.2 Introduction

1.2.1 Documentation status

This working paper is intended to identify possible architectures for Residential Ethernet (RE), the title currently assigned to an IEEE Study Group. Although this Study Group intends to become a formal IEEE 802 Working Group, the first step in this process (approval of a PAR) has not occurred.

This working paper attempts to represent opinions of its contributors (see page 4), although numerous others contributed to its content. The documented is formatted to minimize the difficulties associated with porting the text into a yet-to-be-defined standards document, although numerous changes and clause partitioning would be expected before that occurs.

1.2.2 Background

Ethernet has successfully propagated from the data center to the home, becoming the wired home computer interconnect of choice. However, insufficient support of real-time services has limited Ethernet's success as a consumer audio-video interconnects, where IEEE Std 1394 Serial Bus and Universal Serial Bus (USB) have dominated the marketplace.

This working paper for Residential Ethernet (RE) supports time-sensitive network traffic (called classA traffic), as well as legacy IEEE 1394 traffic, while associating the interconnect with Ethernet commodity pricing and relatively seamless frame-transport bridging.

1.2.3 Design objectives

Design objectives for Residential Ethernet (RE) protocols include the following:

- a) Scalable. Time-sensitive classA transfers can be supported over multiple speed links:
 - 1) 100 Mb/s. Normal (~1500 bytes, or 120 μ s) and classA frames coexist on 100 Mb/s links.
 - 2) 1 Gb/s. Jumbo (~8,200 bytes, or 66 μ s) and classA frames can coexist on 1 Gb/s links.
- b) Compatible. Existing devices and protocols are supported, as follows:
 - 1) Interoperable. Communications of existing 802.3 stations are not degraded by classA traffic.
 - 2) Heterogeneous. Existing 1394 A/V devices can be bridged over RE connections.
- c) Efficient. Time-sensitive transmissions are efficient as well as robust:
 - 1) Bandwidth is independently managed on non-overlapping paths.
 - 2) ClassA transmissions are limited to the links between talker and listener stations.
 - 3) Up to 75% of the link bandwidth can be allocated for classA transmissions.
- d) Applicable. Time-sensitive transmission characteristics are applicable to the marketplace.
 - 1) Precise. A common synchronous clock allows playback times to be precisely synchronized.
 - 2) Low latency. Talker and listener delays are less than human perceptible delays, for interactive home (see 5.1.2 and 5.1.3) and between-home (telephone or internet based) applications.
- e) Predictable. Subject to the (c3) constraint, classA traffic is unaffected by the network topology or the traffic loads offered by other stations.

1.2.4 Strategies

Strategies for achieving the aforementioned objectives include the following:

- a) Subscription. ClassA transmission bandwidths are limited to prenegotiated bandwidths.
- b) Pacing. ClassA transmissions are limited to subscription-negotiated per-cycle bandwidths. (The 125 μ s cycle is consistent with existing IEEE 1394 A/V and telecommunication systems.)
 - 1) Topology. Bandwidths can be guaranteed over arbitrary non-cyclical topologies.
 - 2) Presence. Subscription protocols can readily detect the presence/absence of talker streams.
- c) Simplicity. Simplicity is achieved by utilizing well behaved protocols:
 - 1) Only duplex point-to-point Ethernet links are supported.
 - 2) PLLs. Precise global clock synchronization eliminates the need for PLLs within bridges.
 - 3) Plugs. Self-administered stream identifiers are based on talker-managed plug identifiers. (This eliminates the need to define/provide/configure stream identifier servers.)
 - 4) RSVP. Subscription is based on a layer-2 simplification of the RSVP protocols, called SRP. (SRP allows listeners to autonomously/robustly adapt to spanning tree topology changes).

1.2.5 Interoperability

RE interoperates with existing Ethernet, but the scope of RE services is limited to the RE cloud, as illustrated in Figure 1.1; normal best-effort services are available everywhere else. The scope of the RE cloud is limited by a non-RE capable bridge or a half-duplex link, neither of which can support RE services.

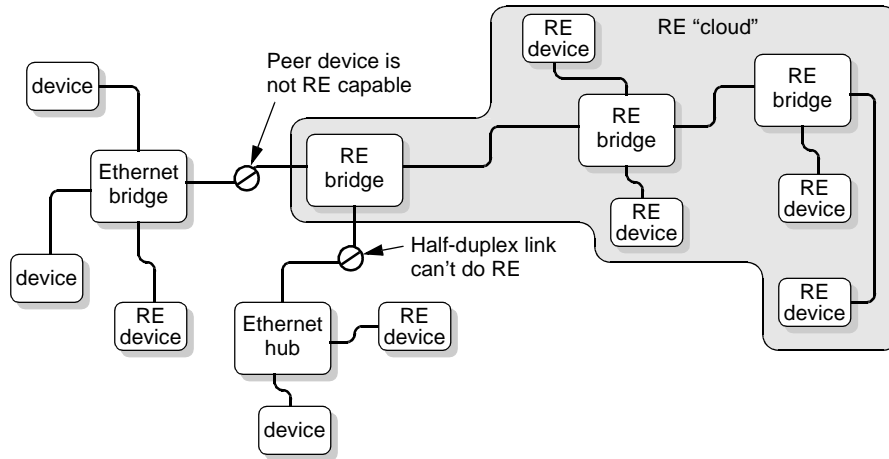


Figure 1.1—Topology and connectivity

Separation of RE devices is driven by the requirements of RE bridges to support subscription (bandwidth allocation), time-of-day clock-synchronization, and (preferably) of pacing of time-sensitive transmissions.

1.2.6 Document structure

The clauses and annexes of this working paper are listed below. The recommended reading order for first-time readers is Clause 5 (an overview), Clause F (critical considerations), Clause 7/8 (details of design). Other clauses provide useful background and reference material.

- Clause 1: Overview
- Clause 2: References
- Clause 3: Terms, definitions, and notation
- Clause 4: Abbreviations and acronyms
- Clause 5: Architecture overview
- Clause 6: Frame formats
- Clause 7: Clock synchronization
- Clause 8: Subscription state machines
- Annex A: Bibliography
- Annex B: Background material
- Annex C: Encapsulated IEEE 1394 frames
- Annex D: Review of possible alternatives
- Annex E: Time-of-day format considerations
- Annex G: Denigrated alternatives
- Annex F: Bursting and bunching considerations
- Annex H: Frequently asked questions (FAQs)
- Annex I: Comment responses
- Annex J: C-code illustrations

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

2. References

NOTE—This clause should be skipped on the first reading (continue with Clause 5).
This references list is highly preliminary, references will be added as this working paper evolves.

The following documents contain provisions that, through reference in this working paper, constitute provisions of this working paper. All the standards listed are normative references. Informative references are given in Annex A. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this working paper are encouraged to investigate the possibility of applying the most recent editions of the standards indicated below.

ANSI/ISO 9899-1990, Programming Language-C.^{1,2}

IEEE Std 802.1D-2004, IEEE Standard for Local and Metropolitan Area Networks: Media Access Control (MAC) Bridges.

¹Replaces ANSI X3.159-1989

²ISO documents are available from ISO Central Secretariat, 1 Rue de Varembe, Case Postale 56, CH-1211, Geneve 20, Switzerland/Suisse; and from the Sales Department, American National Standards Institute, 11 West 42 Street, 13th Floor, New York, NY 10036-8002, USA

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

3. Terms, definitions, and notation

NOTE—This clause should be skipped on the first reading (continue with Clause 5).

This text has been lifted from the P802.17 draft standard, which has a relative comprehensive list. Terms and definitions are expected to be added, revised, and/or deleted as this working paper evolves.

3.1 Conformance levels

Several key words are used to differentiate between different levels of requirements and options, as described in this subclause.

3.1.1 may: Indicates a course of action permissible within the limits of the standard with no implied preference (“may” means “is permitted to”).

3.1.2 shall: Indicates mandatory requirements to be strictly followed in order to conform to the standard and from which no deviation is permitted (“shall” means “is required to”).

3.1.3 should: An indication that among several possibilities, one is recommended as particularly suitable, without mentioning or excluding others; or that a certain course of action is preferred but not necessarily required; or that (in the negative form) a certain course of action is deprecated but not prohibited (“should” means “is recommended to”).

3.2 Terms and definitions

For the purposes of this working paper, the following terms and definitions apply. The Authoritative Dictionary of IEEE Standards Terms [B2] should be referenced for terms not defined in the clause.

3.2.1 audience: The set of listeners associated with a common streamID.

3.2.2 best-effort: Not associated with an explicit service guarantee.

3.2.3 bridge: A functional unit interconnecting two or more networks at the data link layer of the OSI reference model.

3.2.4 clock master: A bridge or end station that provides the link clock reference.

3.2.5 clock slave: A bridge or end station that tracks the link clock reference provided by the clock master.

3.2.6 cyclic redundancy check (CRC): A specific type of frame check sequence computed using a generator polynomial.

3.2.7 destination station: A station to which a frame is addressed.

3.2.8 frame: The MAC sublayer protocol data unit (PDU).

3.2.9 grand clock master: The clock master selected to provide the network time reference.

3.2.10 jitter: The variation in delay associated with the transfer of frames between two points.

3.2.11 latency: The time required to transfer information from one point to another.³

- 3.2.12 link:** A unidirectional channel connecting adjacent stations (half of a span). 1
- 3.2.13 listener:** A sink of a stream, such as a television or acoustic speaker. 2
- 3.2.14 local area network (LAN):** A communications network designed for a small geographic area, typically not exceeding a few kilometers in extent, and characterized by moderate to high data transmission rates, low delay, and low bit error rates. 3
- 3.2.15 MAC client:** The layer entity that invokes the MAC service interface. 4
- 3.2.16 management information base (MIB):** A repository of information to describe the operation of a specific network device. 5
- 3.2.17 maximum transfer unit (MTU):** The largest frame (comprising payload and all header and trailer information) that can be transferred across the network. 6
- 3.2.18 medium** (plural: **media**): The material on which information signals are carried; e.g., optical fiber, coaxial cable, and twisted-wire pairs. 7
- 3.2.19 medium access control (MAC) sublayer:** The portion of the data link layer that controls and mediates the access to the network medium. In this working paper, the MAC sublayer comprises the MAC datapath sublayer and the MAC control sublayer. 8
- 3.2.20 multicast:** Transmission of a frame to stations specified by a group address. 9
- 3.2.21 multicast address:** A group address that is not a broadcast address, i.e., is not all-ones, and identifies some subset of stations on the network. 10
- 3.2.22 network:** A set of communicating stations and the media and equipment providing connectivity among the stations. 11
- 3.2.23 pacer:** A credit-based entity that partitions residual bandwidths between two classes of frames. 12
- 3.2.24 packet:** A generic term for a PDU associated with a layer-entity above the MAC sublayer. 13
- 3.2.25 path:** A logical concatenation of links and bridges over which streams flow from the talker to the listener. 14
- 3.2.26 plug-and-play:** The requirement that a station perform classA transfers without operator intervention (except for any intervention needed for connection to the cable). 15
- 3.2.27 protocol implementation conformance statement (PICS):** A statement of which capabilities and options have been implemented for a given Open Systems Interconnection (OSI) protocol. 16
- 3.2.28 service discovery:** The process used by listeners or controlling stations to identify, control, and configure talkers. 17
- 3.2.29 shaper:** A credit-based entity that limits short-term transmission bandwidths to a specified rate. 18
- 3.2.30 simple reservation protocol (SRP):** The subscription protocol used to allocate and sustain paths for streaming classA traffic. 19

³Delay and latency are synonyms for the purpose of this working paper. Delay is the preferred term. 20

- 3.2.31 span:** A bidirectional channel connecting adjacent stations (two links). 1
- 3.2.32 source station:** The station that originates a frame. 2
- 3.2.33 station:** A device attached to a network for the purpose of transmitting and receiving information on that network. 3
- 3.2.34 stream:** A sequence of frames passed from the talker to listener(s), which have the same streamID. 4
- 3.2.35 subscription:** The process of establishing committed paths between the talker and one or more listeners. 5
- 3.2.36 talker:** The source of a stream, such as a cable box or microphone. 6
- 3.2.37 topology:** The arrangement of links and stations forming a network, together with information on station attributes. 7
- 3.2.38 transmit (transmission):** The action of a station placing a frame on the medium. 8
- 3.2.39 transparent bridging:** A bridging mechanism that is transparent to the end stations. 9
- 3.2.40 unicast:** The act of sending a frame addressed to a single station. 10

3.3 Service definition method and notation 11

The service of a layer or sublayer is the set of capabilities that it offers to a user in the next higher (sub)layer. Abstract services are specified in this working paper by describing the service primitives and parameters that characterize each service. This definition of service is independent of any particular implementation (see Figure 3.1). 12

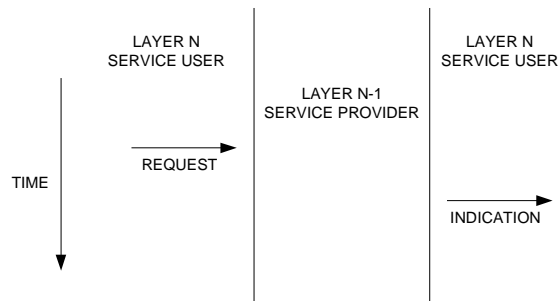


Figure 3.1—Service definitions 13

Specific implementations can also include provisions for interface interactions that have no direct end-to-end effects. Examples of such local interactions include interface flow control, status requests and indications, error notifications, and layer management. Specific implementation details are omitted from this service specification, because they differ from implementation to implementation and also because they do not impact the peer-to-peer protocols. 14

3.3.1 Classification of service primitives 15

Primitives are of two generic types. 16

- a) REQUEST. The request primitive is passed from layer N to layer N-1 to request that a service be initiated.
- b) INDICATION. The indication primitive is passed from layer N-1 to layer N to indicate an internal layer N-1 event that is significant to layer N. This event can be logically related to a remote service request, or can be caused by an event internal to layer N-1.

The service primitives are an abstraction of the functional specification and the user-layer interaction. The abstract definition does not contain local detail of the user/provider interaction. For instance, it does not indicate the local mechanism that allows a user to indicate that it is awaiting an incoming call. Each primitive has a set of zero or more parameters, representing data elements that are passed to qualify the functions invoked by the primitive. Parameters indicate information available in a user/provider interaction. In any particular interface, some parameters can be explicitly stated (even though not explicitly defined in the primitive) or implicitly associated with the service access point. Similarly, in any particular protocol specification, functions corresponding to a service primitive can be explicitly defined or implicitly available.

3.4 State machines

3.4.1 State machine behavior

The operation of a protocol can be described by subdividing the protocol into a number of interrelated functions. The operation of the functions can be described by state machines. Each state machine represents the domain of a function and consists of a group of connected, mutually exclusive states. Only one state of a function is active at any given time. A transition from one state to another is assumed to take place in zero time (i.e., no time period is associated with the execution of a state), based on some condition of the inputs to the state machine.

The state machines contain the authoritative statement of the functions they depict. When apparent conflicts between descriptive text and state machines arise, the order of precedence shall be formal state tables first, followed by the descriptive text, over any explanatory figures. This does not override, however, any explicit description in the text that has no parallel in the state tables.

The models presented by state machines are intended as the primary specifications of the functions to be provided. It is important to distinguish, however, between a model and a real implementation. The models are optimized for simplicity and clarity of presentation, while any realistic implementation might place heavier emphasis on efficiency and suitability to a particular implementation technology. It is the functional behavior of any unit that has to match the standard, not its internal structure. The internal details of the model are useful only to the extent that they specify the external behavior clearly and precisely.

3.4.2 State table notation

<p>NOTE—The following state machine notation was used within 802.17, due to the exactness of C-code conditions and the simplicity of updating table entries (as opposed to 2-dimensional graphics). Early state table descriptions can be converted (if necessary) into other formats before publication.</p>

Each row of the table is preferably provided with a brief description of the condition and/or action for that row. The descriptions are placed after the table itself, and linked back to the rows of the table using numeric tags.

3.4.2.1 Parallel-execution state tables

State machines may be represented in tabular form. The table is organized into two columns: a left hand side representing all of the possible states of the state machine and all of the possible conditions that cause transitions out of each state, and the right hand side giving all of the permissible next states of the state machine as well as all of the actions to be performed in the various states, as illustrated in Table 3.1. The syntax of the expressions follows standard C notation (see 3.13). No time period is associated with the transition from one state to the next.

Table 3.1—State table notation example

Current		Row	Next	
state	condition		action	state
START	sizeofMacControl > spaceInQueue	1	—	START
	passM == 0	2	—	START
	—	3	TransmitFromControlQueue();	FINAL
FINAL	SelectedTransferCompletes()	4	—	START
	—	5	—	FINAL

Row 3.1-1: Do nothing if the size of the queued MAC control frame is larger than the PTQ space.

Row 3.1-2: Do nothing in the absence of MAC control transmission credits.

Row 3.1-3: Otherwise, transmit a MAC control frame.

Row 3.1-4: When the transmission completes, start over from the initial state (i.e., START).

Row 3.1-5: Until the transmission completes, remain in this state.

Each combination of current state, next state, and transition condition linking the two is assigned to a different row of the table. Each row of the table, read left to right, provides: the name of the current state; a condition causing a transition out of the current state; an action to perform (if the condition is satisfied); and, finally, the next state to which the state machine transitions, but only if the condition is satisfied. The symbol “—” signifies the default condition (i.e., operative when no other condition is active) when placed in the condition column, and signifies that no action is to be performed when placed in the action column. Conditions are evaluated in order, top to bottom, and the first condition that evaluates to a result of TRUE is used to determine the transition to the next state. If no condition evaluates to a result of TRUE, then the state machine remains in the current state. The starting or initialization state of a state machine is always labeled “START” in the table (though it need not be the first state in the table). Every state table has such a labeled state.

Each row of the table is preferably provided with a brief description of the condition and/or action for that row. The descriptions are placed after the table itself, and linked back to the rows of the table using numeric tags.

3.4.2.2 Called state tables

A RETURN state is the terminal state of a state machine that is intended to be invoked by another state machine, as illustrated in Table 3.2. Once the RETURN state is reached, the state machine terminates execution, effectively ceasing to exist until the next invocation by the caller, at which point it begins execution again from the START state. State machines that contain a RETURN state are considered to be only instantiated when they are invoked. They do not have any persistent (static) variables.

Table 3.2—Called state table notation example

Current		Row	Next	
state	condition		action	state
START	sizeofMacControl > spaceInQueue	1	—	FINAL
	passM == 0	2		
	—	3	TransmitFromControlQueue();	RETURN
FINAL	MacTransmitError();	4	errorDefect = TRUE	RETURN
	—	5	—	

Row 3.2-1: The size of the queued MAC control frame is less than the PTQ space.

Row 3.2-2: In the absence of MAC control transmission credits, no action is taken.

Row 3.2-3: MAC control transmissions have precedence over client transmissions.

Row 3.2-4: If the transmission completes with an error, set an error defect indication.

Row 3.2-5: Otherwise, no error defect is indicated.

3.5 Arithmetic and logical operators

In addition to commonly accepted notation for mathematical operators, Table 3.3 summarizes the symbols used to represent arithmetic and logical (boolean) operations. Note that the syntax of operators follows standard C notation (see 3.13).

Table 3.3—Special symbols and operators

Printed character	Meaning
&&	Boolean AND
	Boolean OR
!	Boolean NOT (negation)
&	Bitwise AND
	Bitwise OR
^	Bitwise XOR
<=	Less than or equal to
>=	Greater than or equal to
=	Equal to
!=	Not equal to
=	Assignment operator
//	Comment delimiter

3.6 Numerical representation

NOTE—The following notation was taken from 802.17, where it was found to have benefits:

- The subscript notation is consistent with common mathematical/logic equations.
- The subscript notation can be used consistently for all possible radix values.

Decimal, hexadecimal, and binary numbers are used within this working paper. For clarity, decimal numbers are generally used to represent counts, hexadecimal numbers are used to represent addresses, and binary numbers are used to describe bit patterns within binary fields.

Decimal numbers are represented in their usual 0, 1, 2, ... format. Hexadecimal numbers are represented by a string of one or more hexadecimal (0-9,A-F) digits followed by the subscript 16, except in C-code contexts, where they are written as $0x123EF2$ etc. Binary numbers are represented by a string of one or more binary (0,1) digits, followed by the subscript 2. Thus the decimal number “26” may also be represented as “ $1A_{16}$ ” or “ 11010_2 ”.

MAC addresses and OUI/EUI values are represented as strings of 8-bit hexadecimal numbers separated by hyphens and without a subscript, as for example “01-80-C2-00-00-15” or “AA-55-11”.

3.7 Field notations

3.7.1 Use of italics

All field names or variable names (such as *level* or *myMacAddress*), and sub-fields within variables (such as *thisState.level*) are italicized within text, figures and tables, to avoid confusion between such names and similarly spelled words without special meanings. A variable or field name that is used in a subclause heading or a figure or table caption is also italicized. Variable or field names are not italicized within C code, however, since their special meaning is implied by their context. Names used as nouns (e.g., *subclassA0*) are also not italicized.

3.7.2 Field conventions

This working paper describes values that are packetized or MAC-resident, such as those illustrated in Table 3.2.

Table 3.4—Names of fields and sub-fields

Name	Description
<i>newCRC</i>	Field within a register or frame
<i>thisState.level</i>	Sub-field within field <i>thisState</i>
<i>thatState.rateC[n].c</i>	Sub-field within array element <i>rateC[n]</i>

Run-together names (e.g., *thisState*) are used for fields because of their compactness when compared to equivalent underscore-separated names (e.g., *this_state*). The use of multiword names with spaces (e.g., “This State”) is avoided, to avoid confusion between commonly used capitalized key words and the capitalized word used at the start of each sentence.

A sub-field of a field is referenced by suffixing the field name with the sub-field name, separated by a period. For example, *thisState.level* refers to the sub-field *level* of the field *thisState*. This notation can be continued in order to represent sub-fields of sub-fields (e.g., *thisState.level.next* is interpreted to mean the sub-field *next* of the sub-field *level* of the field *thisState*).

Two special field names are defined for use throughout this working paper. The name *frame* is used to denote the data structure comprising the complete MAC sublayer PDU. Any valid element of the MAC sublayer PDU, can be referenced using the notation *frame.xx* (where *xx* denotes the specific element); thus, for instance, *frame.serviceDataUnit* is used to indicate the *serviceDataUnit* element of a frame.

Unless specifically specified otherwise, reserved fields are reserved for the purpose of allowing extended features to be defined in future revisions of this working paper. For devices conforming to this version of this working paper, nonzero reserved fields are not generated; values within reserved fields (whether zero or nonzero) are to be ignored.

3.7.3 Field value conventions

This working paper describes values of fields. For clarity, names can be associated with each of these defined values, as illustrated in Table 3.5. A symbolic name, consisting of upper case letters with underscore separators, allows other portions of this working paper to reference the value by its symbolic name, rather than a numerical value.

Table 3.5—*wrap* field values

Value	Name	Description
0	STANDARD	Standard processing selected
1	SPECIAL	Special processing selected
2,3	—	Reserved

Unless otherwise specified, reserved values allow extended features to be defined in future revisions of this working paper. Devices conforming to this version of this working paper do not generate nonzero reserved values, and process reserved fields as though their values were zero.

A field value of TRUE shall always be interpreted as being equivalent to a numeric value of 1 (one), unless otherwise indicated. A field value of FALSE shall always be interpreted as being equivalent to a numeric value of 0 (zero), unless otherwise indicated.

3.8 Bit numbering and ordering

Data transfer sequences normally involve one or more cycles, where the number of bytes transmitted in each cycle depends on the number of byte lanes within the interconnecting link. Data byte sequences are shown in figures using the conventions illustrated by Figure 3.2, which represents a link with four byte lanes. For multi-byte objects, the first (left-most) data byte is the most significant, and the last (right-most) data byte is the least significant.

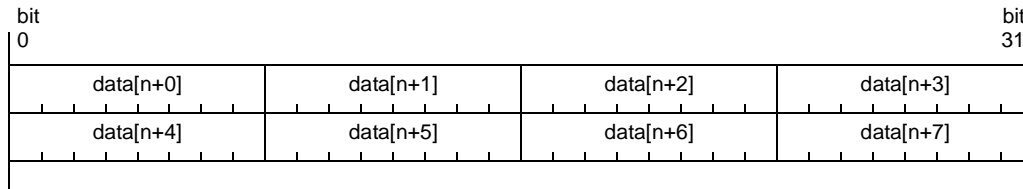


Figure 3.2—Bit numbering and ordering

Figures are drawn such that the counting order of data bytes is from left to right within each cycle, and from top to bottom between cycles. For consistency, bits and bytes are numbered in the same fashion.

NOTE—The transmission ordering of data bits and data bytes is not necessarily the same as their counting order; the translation between the counting order and the transmission order is specified by the appropriate reconciliation sublayer.

3.9 Byte sequential formats

Figure 3.3 provides an illustrative example of the conventions to be used for drawing frame formats and other byte sequential representations. These representations are drawn as fields (of arbitrary size) ordered along a vertical axis, with numbers along the left sides of the fields indicating the field sizes in bytes. Fields are drawn contiguously such that the transmission order across fields is from top to bottom. The example shows that *field1*, *field2*, and *field3* are 1-, 1- and 6-byte fields, respectively, transmitted in order starting with the *field1* field first. As illustrated on the right hand side of Figure 3.3, a multi-byte field represents a sequence of ordered bytes, where the first through last bytes correspond to the most significant through least significant portions of the multi-byte field, and the MSB of each byte is drawn to be on the left hand side.

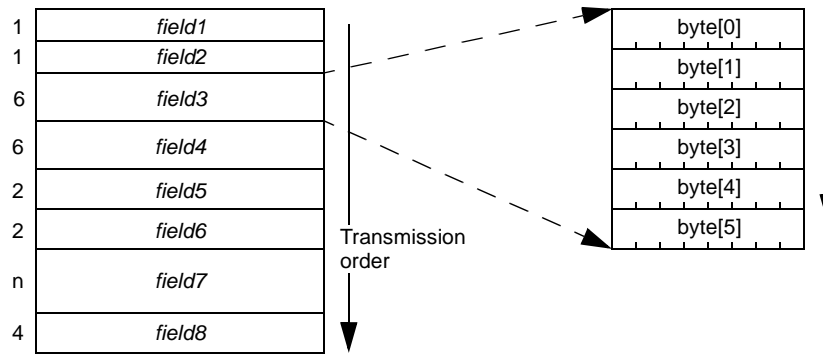


Figure 3.3—Byte sequential field format illustrations

NOTE—Only the left-hand diagram in Figure 3.3 is required for representation of byte-sequential formats. The right-hand diagram is provided in this description for explanatory purposes only, for illustrating how a multi-byte field within a byte sequential representation is expected to be ordered. The tag “Transmission order” and the associated arrows are not required to be replicated in the figures.

3.10 Ordering of multibyte fields

In many cases, bit fields within byte or multibyte objects are expanded in a horizontal fashion, as illustrated in the right side of Figure 3.4. The fields within these objects are illustrated as follows: left-to-right is the byte transmission order; the left-through-right bits are the most significant through least significant bits respectively.

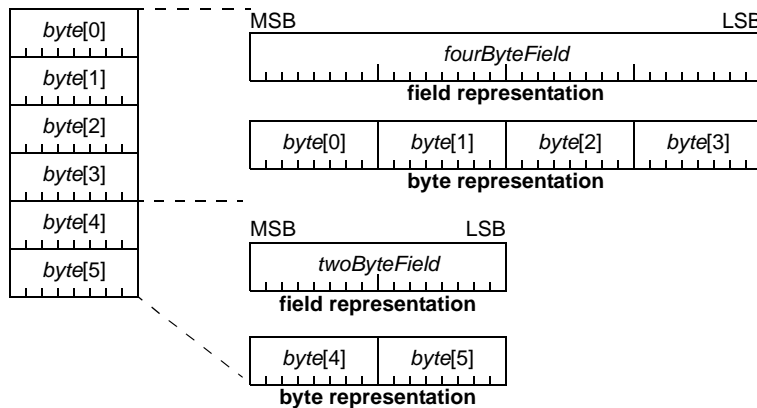


Figure 3.4—Multibyte field illustrations

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

The first *fourByteField* can be illustrated as a single entity or a 4-byte multibyte entity. Similarly, the second *twoByteField* can be illustrated as a single entity or a 2-byte multibyte entity.

NOTE—The following text was taken from 802.17, where it was found to have benefits:
The details should, however, be revised to illustrate fields within an RE frame header *serviceDataUnit*.

To minimize potential for confusion, four equivalent methods for illustrating frame contents are illustrated in Figure 3.5. Binary, hex, and decimal values are always shown with a left-to-right significance order, regardless of their bit-transmission order.

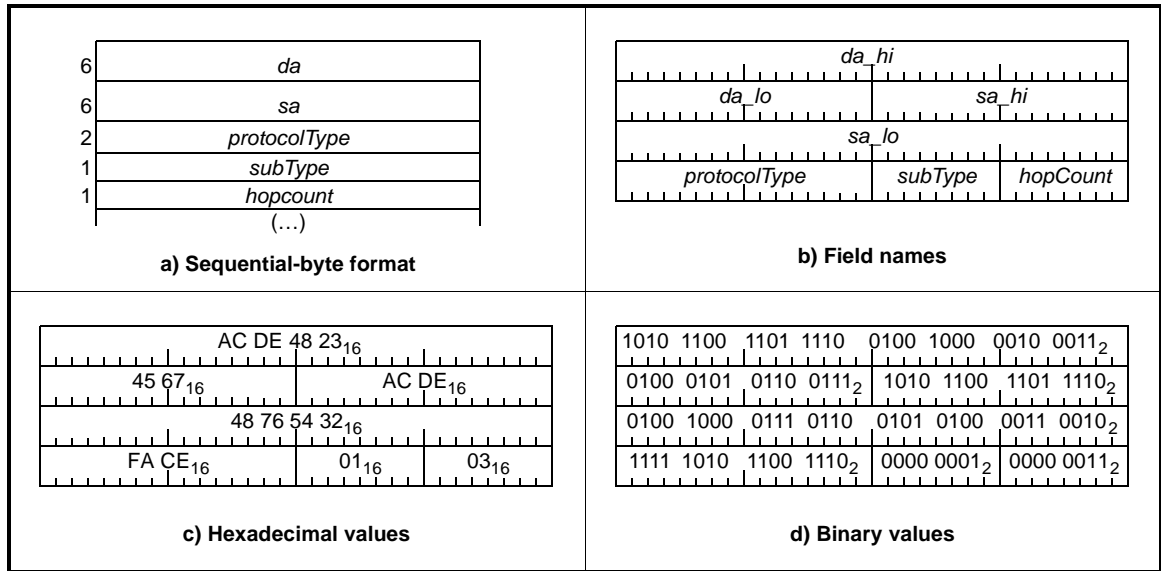


Figure 3.5—Illustration of fairness-frame structure

3.11 MAC address formats

The format of MAC address fields within frames is illustrated in Figure 3.6.

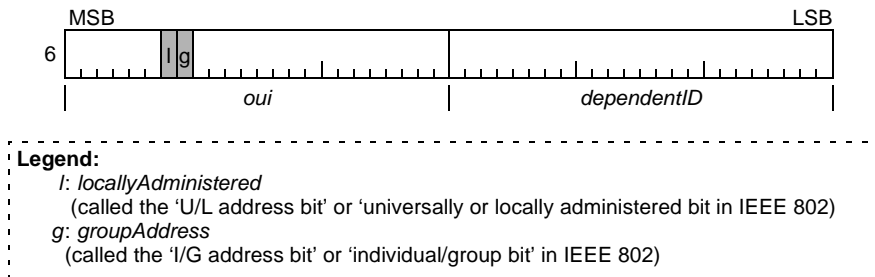


Figure 3.6—MAC address format

3.11.1 oui: A 24-bit organizationally unique identifier (OUI) field supplied by the IEEE/RAC for the purpose of identifying the organization supplying the (unique within the organization, for this specific context) 24-bit *dependentID*. (For clarity, the *locallyAdministered* and *groupAddress* bits are illustrated by the shaded bit locations.)

3.11.2 *dependentID*: An 24-bit field supplied by the *oui*-specified organization. The concatenation of the *oui* and *dependentID* provide a unique (within this context) identifier.

To reduce the likelihood of error, the mapping of OUI values to the *oui/dependentID* fields are illustrated in Figure 3.7. For the purposes of illustration, specific OUI and *dependentID* example values have been assumed. The two shaded bits correspond to the *locallyAdministered* and *groupAddress* bit positions illustrated in Figure 3.6.

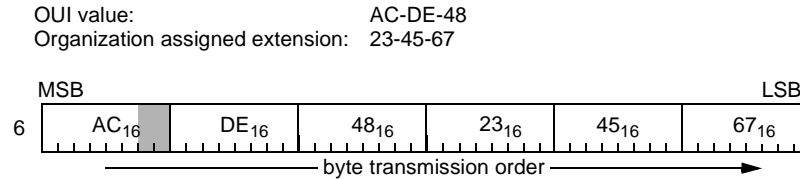


Figure 3.7—48-bit MAC address format

3.12 Informative notes

Informative notes are used in this working paper to provide guidance to implementers and also to supply useful background material. Such notes never contain normative information, and implementers are not required to adhere to any of their provisions. An example of such a note follows.

NOTE—This is an example of an informative note.

3.13 Conventions for C code used in state machines

Many of the state machines contained in this working paper utilize C code functions, operators, expressions and structures for the description of their functionality. Conventions for such C code can be found in Annex J.

4. Abbreviations and acronyms

NOTE—This clause should be skipped on the first reading (continue with Clause 5).

This text has been lifted from the P802.17 draft standard, which has a relative comprehensive list. Abbreviations/acronyms are expected to be added, revised, and/or deleted as this working paper evolves.

This working paper contains the following abbreviations and acronyms:

BER	bit error ratio
CRC	cyclic redundancy check
FCS	frame check sequence
FIFO	first in first out
GARP	Generic Attribute Registration Protocol
HEC	header error check
IEC	International Electrotechnical Commission
IEEE	Institute of Electrical and Electronics Engineers
IETF	Internet Engineering Task Force
ISO	International Organization for Standardization
ITU	International Telecommunication Union
LAN	local area network
LSB	least significant bit
MAC	medium access control
MAN	metropolitan area network
MIB	management information base
MSB	most significant bit
MTU	maximum transfer unit
OAM	operations, administration, and maintenance
OSI	open systems interconnect
PDU	protocol data unit
PHY	physical layer
RE	Residential Ethernet
RFC	request for comment
RPR	resilient packet ring
SRP	simple reservation protocol
TDM	time division multiplexing
VOIP	voice over internet protocol

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5. Architecture overview

5.1 Latency constraints

5.1.1 Interactive audio delay considerations

The latency constraints of the RE environment are based on the sensitivity of the human ear. To be comfortable when playing music, the delay between the instrument and the human ear should not exceed 10-to-15 ms, as illustrated in Figure 5.1. The individual hop delays must be considerably smaller, since instrument-sourced audio traffic may pass through multiple links and processing devices before reaching the ear, as illustrated in 5.1.2 and 5.1.3.

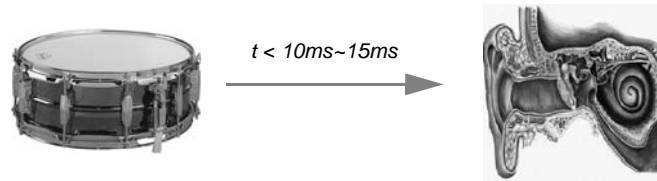


Figure 5.1—Interactive audio delay considerations

5.1.2 Home recording session

To illustrate hop-latency requirements, consider RE usage for a home recording session, as illustrated in Figure 5.2. The audio inputs (microphone and guitar) are converted, passed through a bridge, mixed within a laptop computer, converted at the speaker, and return to the performer's ear through the air.

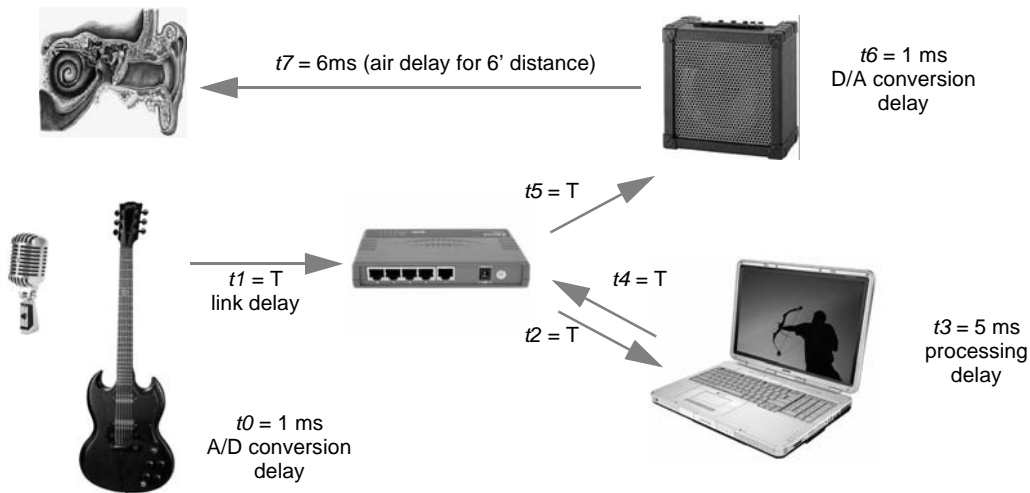


Figure 5.2—Home recording session

A fixed time T is assumed for each passage through a link, based on potential buffering and conflicting-traffic delays. Due to multiple link hops and the latency contributions, the constraints on the value of T are much less than the constraining 15ms instrument-to-ear latency, as illustrated in Equation 5.1.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

$$\begin{aligned} t_0 + t_1 + t_2 + t_3 + t_4 + t_5 + t_6 + t_7 &< 15 \text{ ms} \\ 1\text{ms} + T + T + 5\text{ms} + T + T + 1\text{ms} + 6\text{ms} &< 15\text{ms} \\ 4 \times T + 13\text{ms} &< 15\text{ms} \\ T &< 0.5 \text{ ms} \end{aligned} \tag{5.1}$$

To better understand the range of possibilities, consider an extremely aggressive implementation of end-point stations could reduce the link-latency requirements. For example, more aggressive end-point processing delays $\{t_0=0.25 \text{ ms}, t_3=2 \text{ ms}, t_6=0.25 \text{ ms}, t_7=6 \text{ ms}\}$ would yield a constraint of $T < 1.6 \text{ ms}$.

Note that these aggressive processor delays are unlikely to decrease as the MIPs rating of processors increase, due to the inherent delays associated with finite input response (FIR) filters and efficiencies achieved through block-processing. For example, 16-sample block processing of a 128-point FIR filter implies an inherent 80-cycle delay (16 for input block accumulation, 64 for filtering). With a 40 kHz sampling rate, this corresponds to a theoretical processing-latency limitation of 2 ms.

These numbers are only approximations; actual values (as determined by the marketplace) could vary substantially. For audiophiles, an overall processing latency of 5 ms may be desired; for discount shoppers, an overall latency of 15 ms may be tolerable. Larger ad-hoc networks of cascaded 4-port or 8-port bridges may be present. As with golden speaker cables, purchases may be based on perceptions of quality (the bridge latency specification), rather than reality (perceivable latencies).

5.1.3 Garage jam session

As another example, consider RE usage for a garage jam session, as illustrated in Figure 5.3. The audio inputs (microphone and guitar) are converted, passed through a guitar effects processor, two bridges, mixed within an audio console, return through two bridges, and return to the ear through headphones.

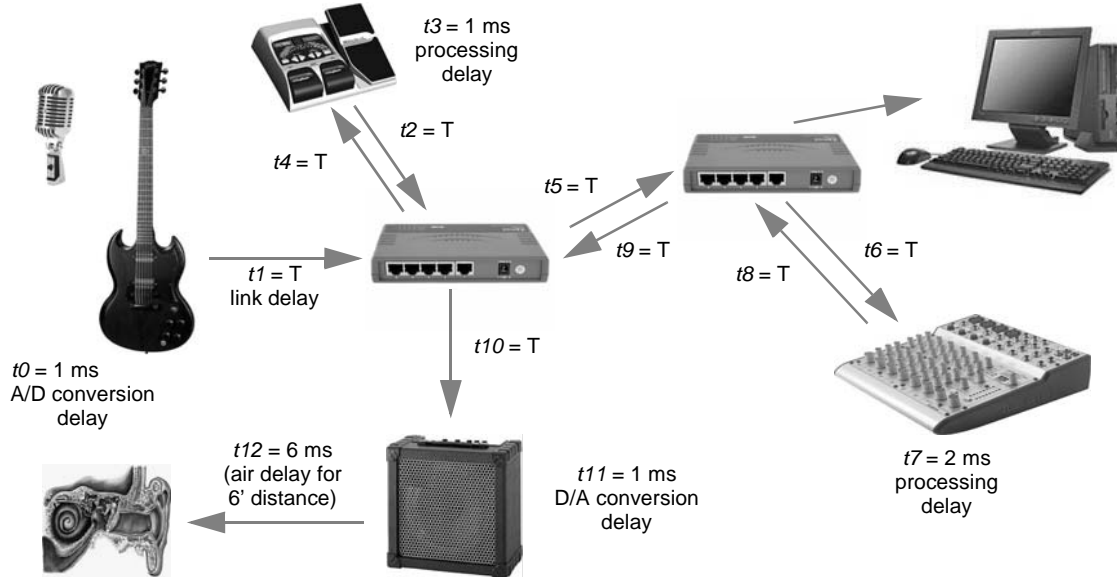


Figure 5.3—Garage jam session

Again, a fixed time T is assumed for each passage through a link, based on potential buffering and conflicting-traffic delays. Due to multiple hops and the latency contributions, the constraints yield a T value that is much less than the constraining 15ms instrument-to-ear latency (see Equation 5.2).

$$\begin{aligned}
 t_0 + t_1 + t_2 + t_3 + t_4 + t_5 + t_6 + t_7 + t_8 + t_9 + t_{10} + t_{11} + t_{12} &< 15 \text{ ms} & (5.2) \\
 1\text{ms} + T + T + 1\text{ms} + T + T + T + 2\text{ms} + T + T + T + 1\text{ms} + 6\text{ms} &< 15\text{ms} \\
 8 \times T + 11\text{ms} &< 15\text{ms} \\
 T &< 0.5 \text{ ms}
 \end{aligned}$$

To better understand the range of possible latencies, consider extremely aggressive implementations of end-point stations. For example, more aggressive end-point processing delays $\{t_0=0.25 \text{ ms}, t_3=0.25 \text{ ms}, t_7=2 \text{ ms}, t_{11}=0.25 \text{ ms}, t_{12}=6 \text{ ms}\}$ would yield a constraint of $T < 0.78 \text{ ms}$.

5.1.4 Urban home recording session

Within urban environments, headphones may be preferred to audio speakers, as illustrated in Figure 5.4 (a small modification of Figure 5.2). The audio inputs (microphone and guitar) are converted, passed through a bridge, mixed within a laptop computer, converted at the headphones, and near immediately presented to the performer's ear.

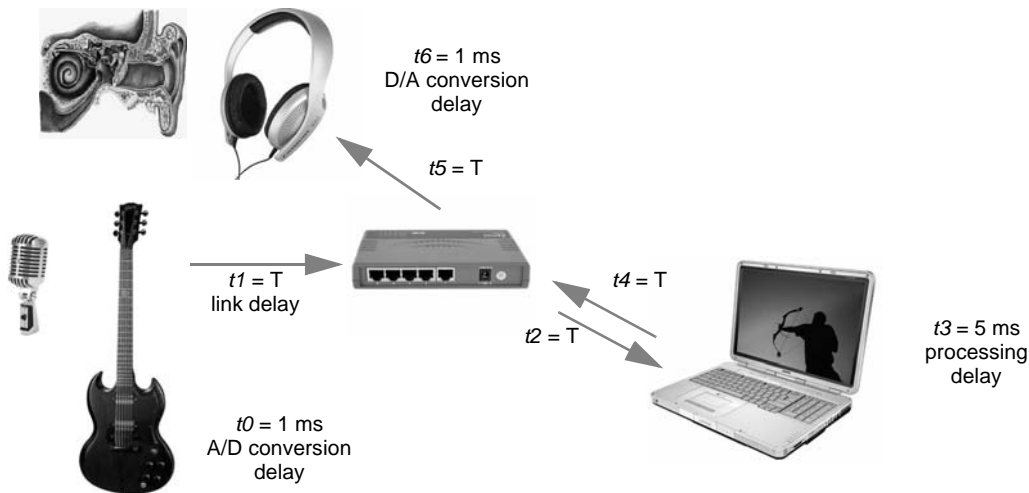


Figure 5.4—Urban recording session

While the earphones eliminate the air-to-ear hop-count delays, the sensitivity to delays is increased for the case of a vocal performer due to a comb filter formed by the interaction of headphone sound and sound conducted through the head. Remaining below the 0.5 to 5 ms range where comb filtering is prevalent is impractical, since the $\{t_0=1 \text{ ms}, t_3=5 \text{ ms}, t_6=1 \text{ ms}\}$ values already exceed the 0.5 ms limitation.

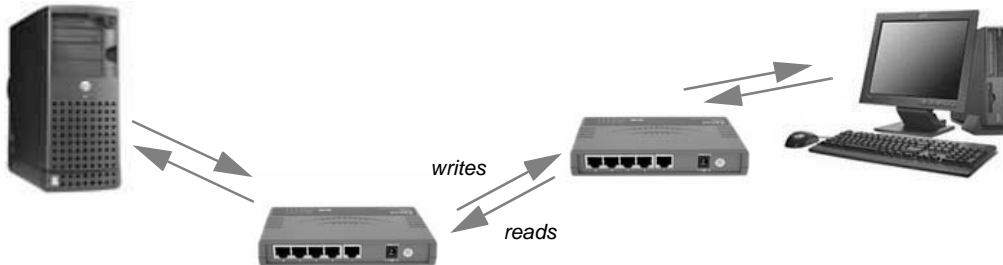
Professionals believe that increasing latency to 5 ms or more within such headphone-feedback environments is preferred over operation in the 0.5 to 5 ms range where comb filtering is prevalent. Again, due to multiple hops and the latency contributions, the constraints yield a T value that is much less than the constraining 15ms instrument-to-ear latency (see Equation 5.3).

$$\begin{aligned}
 t_0 + t_1 + t_2 + t_3 + t_4 + t_5 + t_6 &< 15 \text{ ms} & (5.3) \\
 1\text{ms} + T + T + 5\text{ms} + T + T + 1\text{ms} &< 15 \text{ ms} \\
 4 \times T + 7\text{ms} &< 15 \text{ ms} \\
 T &< 2\text{ms}
 \end{aligned}$$

To better understand the range of possible latencies, consider extremely aggressive implementations of end-point stations. For example, more aggressive end-point processing delays $\{t_0=0.25 \text{ ms}, t_3=2 \text{ ms}, t_6=0.25 \text{ ms}\}$ would yield a $T < 3.1 \text{ ms}$ constraint.

1 **5.1.5 Conflicting data transfers**
 2

3 Home networks may carry data traffic as well as time-sensitive traffic, as illustrated in Figure 5.3. During
 4 musical performances (or evening A/V screenings), high bandwidth computer-to-server transfers could
 5 occur over the same data-transfer links, as illustrated in Figure 5.5.



19 **Figure 5.5—Conflicting data transfers**

20 With the high data-transfer rates of disks and disk-array systems, the bandwidth capacity of residential
 21 Ethernet links could (if not otherwise limited) easily be reached. Thus, some form of prioritized bridging is
 22 necessary to ensure robust delivery of time-sensitive traffic.
 23

24 **5.2 Service classes**
 25
 26
 27

28 **Editors' Notes:** To be removed prior to final publication.
 29 The classA and classC service classes have consensus among the contributors to this working paper. The
 30 concept of classB services was included in IEEE Std 802.17-2004 and is being included for consideration
 31 by universal plug and play (UPnP), congestion management (CM), or legacy applications.
 32

33 This working paper defines three service classes (A, B, or C) with which the data transfer is associated, as
 34 summarized in Table 5.1. The classA service provides low-jitter transfer of traffic (and therefore lower
 35 worst-case delays) up to its allocated rate. Traffic above the allocated rate is rejected. The classB service
 36 provides bounded delay transfer of traffic. The classC service provides best-effort data-transfer services.
 37
 38

39 **Table 5.1—Service classes and their quality-of-service relationships**
 40

41
42
43
44
45
46
47
48
49
50

class of service		qualities of service		
class	examples of use	jitter	guaranteed bandwidth	type
A	real time	low	yes	allocated
B	near real time	bounded		
C	best effort	unbounded	no	opportunistic

51
 52 Link capacity required to support the classA and classB service is allocated via provisioning and these
 53 services can be characterized as allocated services. The provisioning activity is expected to ensure that the
 54

aggregate service commitment on each link does not exceed that link's capacity. The allocation rates distributed by provisioning regulates access to these guaranteed services.

Link capacity has to be ensured to support classA and classB service guarantees. This is done by allocating bandwidth through provisioning that prevents over-provisioning the links, using a subscription protocol (see 5.4).

5.3 Architecture overview

5.3.1 Abstract concepts

From the perspective of end-point stations, RE systems supports classA data-frame traffic, called streams. Each stream has one talker and one or more listeners, as illustrated in Figure 5.6-a.

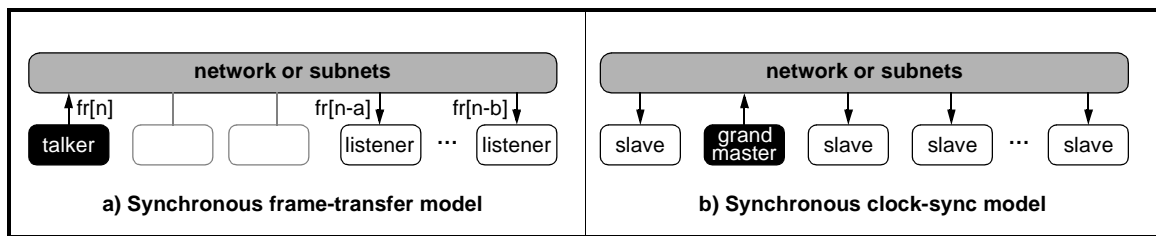


Figure 5.6—Hierarchical control

The delay between the talker and listener(s) is nominally a fixed number of 125 μ s cycles, although the number of cycles may be cable-length and/or bridge topology dependent. Additional delays can be inserted by the application(s), when synchronization between multiple listeners is required, since the talker's data can be time-stamped and all clocks are synchronized.

To reduce costs (and support GPS-inaccessible locations), synchronized clocks are provided by the interconnect. All classA talkers provide clock references, but only one of these stations is nominated to be the clock master; the others are called clock slaves (see Figure 5.6-b). The selected clock master is called the grand clock master, oftentimes abbreviated as "grand master".

Clock synchronization involves synchronizing the clock-slave clocks to the reference provided by the grand clock master. Tight accuracy is possible with matched-length duplex links, since bidirectional messages can cancel the cable-delay effects.

5.3.2 Detailed illustrations

In many cases, abstract illustrations (see Figure 5.6) are insufficient to illustrate expected behaviors. Thus, more detailed illustrations are oftentimes used to also show bridges and spans within the network cloud, as illustrated in Figure 5.7.

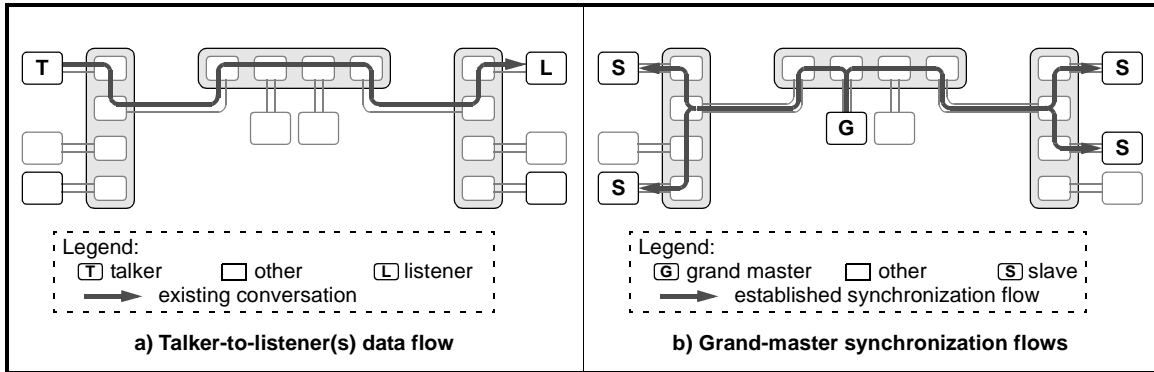


Figure 5.7—Hierarchical flows

5.3.3 Architecture components

The architecture of a home RE system involves the following protocols:

- a) Discovery (beyond the scope of this working paper).
A controller discovers the proper streamID/bandwidth parameters to allow the listener to subscribe to the desired talker-sourced stream.
- b) Subscription. The controller commands the listener to establish a path from the talker.
Subscription may pass or fail, based on availability of routing-table and link-bandwidth resources.
- c) Synchronization. The distributed clocks in talkers and listeners are accurately synchronized.
Synchronized clocks avoid cycle slips and playback-phase distortions.
- d) Pacing. The transmitted classA traffic is paced to avoid other classA traffic disruptions.

5.4 Subscription

5.4.1 Simple Reservation Protocol (SRP) overview

Subscription involves explicit negotiation for bandwidth resources, performed in a distributed fashion, flowing over the paths of intended communication. This subscription protocol are called the Simple Reservation Protocols (SRP). SRP represents an instance of the Generic Attribute Registration Protocol (GARP), with similar objectives to the layer-3 based Resource Reservation Protocol (RSVP). SRP shares many of the baseline RSVP and GARP features, including the following:

- SRP is simplex, i.e. reservations apply to unidirectional data flows.
- SRP is receiver-oriented, i.e., the receiver of a stream initiates and maintains the resource reservation used for that stream.
- SRP maintains “soft” state in bridges, providing graceful support for dynamic membership changes and automatic adaptations to changes in network topology.
- SRP is not a routing protocol, but depends on transparent bridging and STP routing protocols.

SRP simplicity is derived from its restricted layer-2 ambitions, as follows.

- SRP is symmetric, i.e. the listener-to-talker path is the inverse of the talker-to-listener path.
- SRP does not provide for transcoding; any stream is fully characterized by its streamID and bandwidth.

The viability of SRP is enhanced by basing its protocols on GARP, a protocol defined within IEEE Std 802.1D. Specifically, the RequestJoin and RequestLeave messages correspond to primitives defined within GARP.

SRP is defined to be a general 1-to-N resource-reservation scheme, although this discussion focuses on subscription of classA bandwidth resources. The SRP protocols could, however, be used to reserve other resource-limited resources, such as buffer allocations, latency targets, and frame-loss rates.

NOTE—SRP is thought to be applicable to N-to-N topologies, as well as 1-to-N topologies. However, the detailed review of N-to-N topologies (which would be necessary to verify the feasibility of such extensions) is beyond the scope of this working paper.

5.4.2 Soft reservation state

SRP takes a “soft state” approach to managing the reservation state in bridges. SRP soft state is created and periodically refreshed by listener generated RequestJoin messages; this state is deleted if no matching RequestJoin messages arrive before the expiration of a “cleanup timeout” interval. Listener’s may also force state deletions by generating an explicit RequestLeave message.

RequestJoin messages are idempotent. When a route changes, the next RequestJoin message will initialize the path state to the new route, and future RequestJoin messages will establish state there. The state on the now-unused segment of the route will be deleted after a timeout interval. Thus, whether a RequestJoin message is “new” or a “refresh” is determined separately by each station, depending upon the existence of state at that station.

SRP soft state is also deleted in the continued absence of associated talker-generated ConfirmJoin messages; the listener’s registration is discarded if no matching ConfirmJoin indication arrives before the expiration of a “cleanup timeout” interval. Thus, talker stations or agents may implicitly deregister by stopping its ConfirmJoin confirmations, or explicitly deregister by sending distinct ConfirmGone messages.

Editors' Notes: To be removed prior to final publication.
Additional discussions may be appropriate to discuss operation of the ConfirmGone messages.

SRP sends its messages as layer-2 datagrams with no reliability enhancement. Periodic transmissions by listener/talker stations and agents is expected to handle the occasional loss of an SRP message.

In the steady state, state is refreshed on a hop-by-hop basis to allow merging. Propagation of a change stops when and if it reaches a point where merging causes no resulting state change. This minimizes the SRP control traffic and is essential for scaling to large audiences.

5.4.3 Subscription bandwidth constraints

The SRP subscription protocols limit cumulative bandwidth allocations to a fixed percentage less than the capacity of the link, much like IEEE 1394 limits isochronous traffic to less than the capacity of its bus. This guarantees that high priority management information can be transmitted across the link. For RE systems, classA traffic is limited to 75% of the capacity of any RE link. Enforcement of such a limit is done in multiple ways:

- a) Subscription. Requests for establishing classA transmission paths are rejected if the cumulative bandwidths of all paths would consume more than 75% of the link bandwidth.
- b) Transmit queue hardware of RE stations (including bridges) discards classA content that (if transmitted) would cause classA traffic to exceed 75% of the transmit link capacity. Details are TBD.

Method (b) is desired to recovery from unexpected transient conditions (typically topology changes) that result in admission control violations, and is also useful for managing misbehaving devices

5.4.4 Controller entities

Subscription when a relative-intelligent controller discovers the need to establish a classA path between talker and listener entities. For example, user interactions with a television (called the controller) may cause streams flowing between the content source (called the talker) and speakers (the listeners), as illustrated in Figure 5.8.

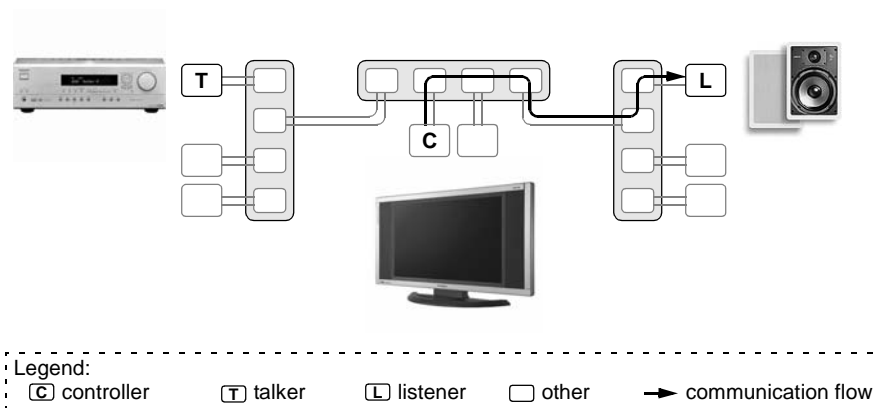


Figure 5.8—Controller activation

A controller can potentially simplify the listener by reducing the need to providing user interface and device-discovery capabilities. However, a controller could also reside within talker and/or listener components. However, actions between controllers and talker/listener stations are beyond the scope of this working paper.

5.4.5 Bridge-resident agents

Subscription facilities register classA communication paths from a talker to one or more listeners. Streams of time-sensitive data can then flow over these established paths, as illustrated by the dark arrow paths in Figure 5.9-a. Maintaining these established paths involves active participation of agents within the end-point talker, local listener, local talker, and end-point listener entities, as illustrated in Figure 5.9-b.

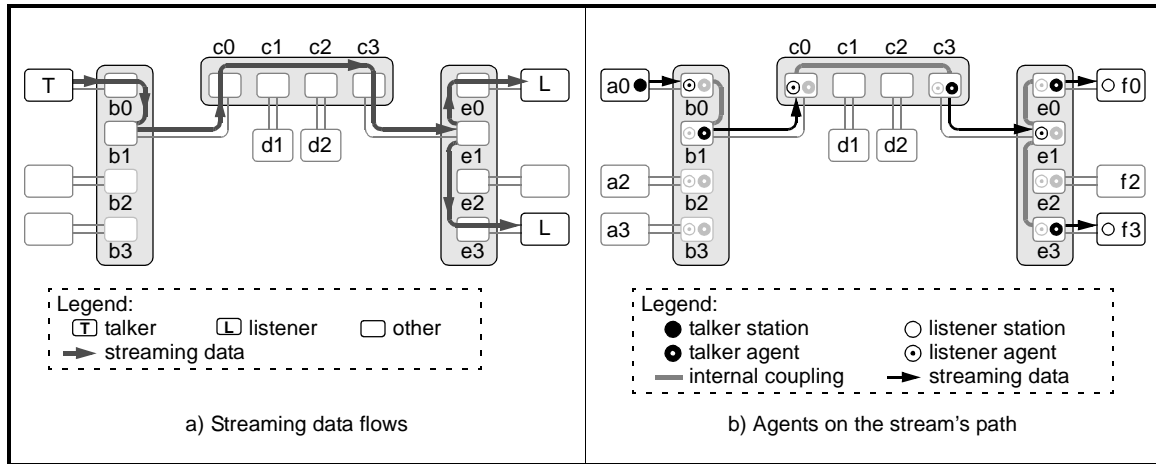


Figure 5.9—Agents on an established path

The talker stations/agents are responsible for maintaining an account consisting of {streamID, bandwidth} pairs, one for each of their distinct flows. Requests for additional link bandwidth are checked against these accounts and denied if the cumulative bandwidth would exceed 75% of the link capacity.

For each of the registered talker agents within a bridge, the listener agent remains active until all but the last talker agent registration is discarded. Thus, the talker agent in an upstream station receives its deregistration notice only after the last of the downstream listener stations has been deregistered.

The listener agent uses the same RequestJoin messages to establish and to maintain the path. This reduces design complexity and (most importantly) automatically re-routes stream flows after topology changes.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5.4.6 Registration

Registering a new listener and talker starts with a RequestJoin message sent from the listener $f0$ towards the talker $a0$, as illustrated by the dark arrow (1a) in Figure 5.10-a. These registration messages are not forwarded directly, but activate cooperative listener and talker agents with the bridge.

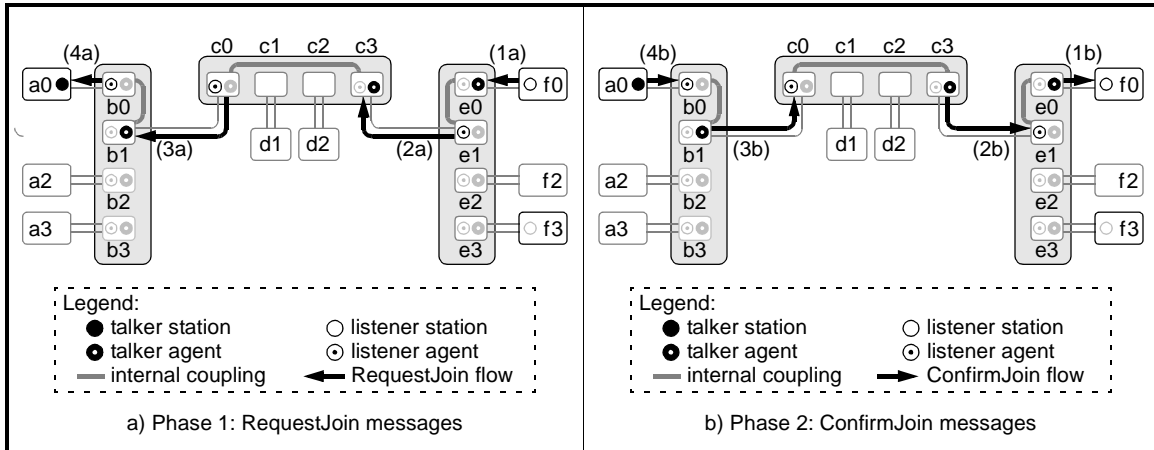


Figure 5.10—Periodic registration messages

In response to the received RequestJoin message (1a), bridgeE reserves talker-agent and listener-agent registration table entries in ports $e0$ and $e1$ respectively. A cascaded RequestJoin message (2a) is then sent towards talker station $a0$.

The cascaded forwarding continues through bridgeC. In response to the received RequestJoin message (2a), bridge C reserves talker-agent and listener-agent registration table entries in ports $c3$ and $c0$ respectively. A cascaded RequestJoin message (3a) is then sent towards talker station $a0$.

The cascaded forwarding continues through bridgeB. In response to the received RequestJoin message (3a), bridge B reserves talker-agent and listener-agent registration table entries in ports $b1$ and $b0$ respectively. A cascaded RequestJoin message (4a) is then sent towards talker station $a0$.

Referring now to Figure 5.10-b, the talker and talker agents are responsible for providing confirming ConfirmJoin messages, to confirm their continued presence. In this example, the RequestJoin messages {1a,2a,3a,4a} of Figure 5.10-a are continually confirmed by the ConfirmJoin messages {1b,2b,3b,4b} of Figure 5.10-b), respectively. In the continued absence of the expected ConfirmJoin messages, the talker (or talker-agent) assumes the listener (or listener-agent) is absent or has been deactivated.

Another timeout is associated with the absence of periodic RequestJoin messages. In the continued absence of these expected messages, the talker assumes the listener is absent or has been deactivated. Based on this assumption, the associated talker (station or agent) registration resources are released.

5.4.7 Secondary listener registrations

A second listener registers by sending a RequestJoin message towards the talker, as illustrated by the dark-arrow path in Figure 5.11-a. When an established registration is discovered, the bridge (not the talker) processes the message. Thus, the registration is expanded to include a new-listener side path, as illustrated in Figure 5.11-b.

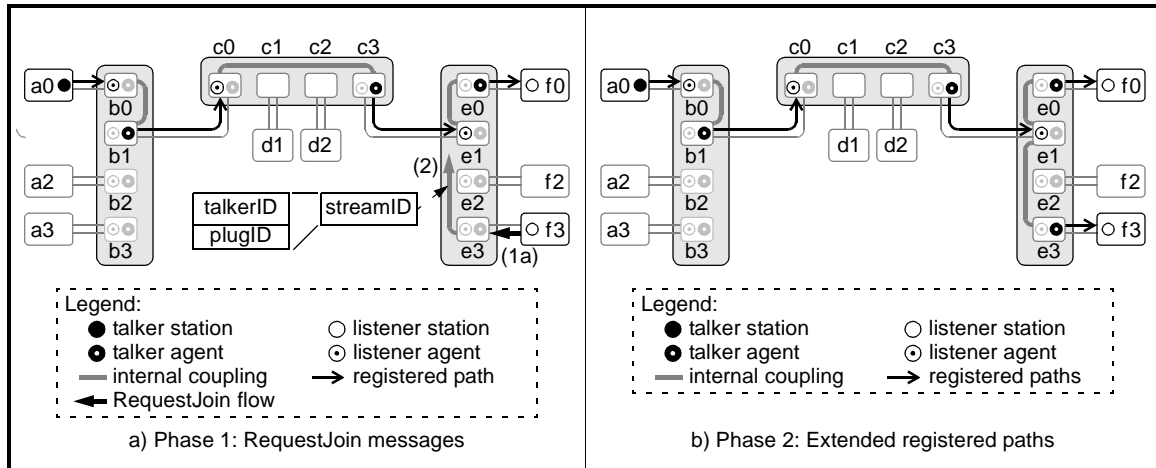


Figure 5.11—Secondary registrations

Each talker and listener agent maintains separate registration state, so that only active paths are registered. Maintaining distinct registrations also allows the bridge to detect when the last listener disconnects, so that its previously shared upstream span can be deregistered appropriately.

Each path is uniquely identified by its associated streamID. The streamID consists of a {talkerId, plugID} information that uniquely identifies the associated talker resource), as illustrated by the rectangle inserts within Figure 5.11-a. The talkerID represents the MAC address of the talker and the plugID distinguishes between possible streaming sources within the talker.

The multicast address used to route the classA multicast frames, as well as the allocated classA bandwidth, are returned to the listeners within ResponseForm messages.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5.4.8 Secondary listener deregistration

A retiring secondary listener normally leaves an established registration by sending a RequestLeave message towards the talker. That RequestLeave message (1a) propagates to the nearest merging bridge connection, as illustrated in Figure 5.12-a. When an established/merged registration is discovered, the bridge (not the talker) deregisters the listener, as illustrated by the disappearance of external path e0-to-f0 and internal path e1-to-e0 within Figure 5.12-b.

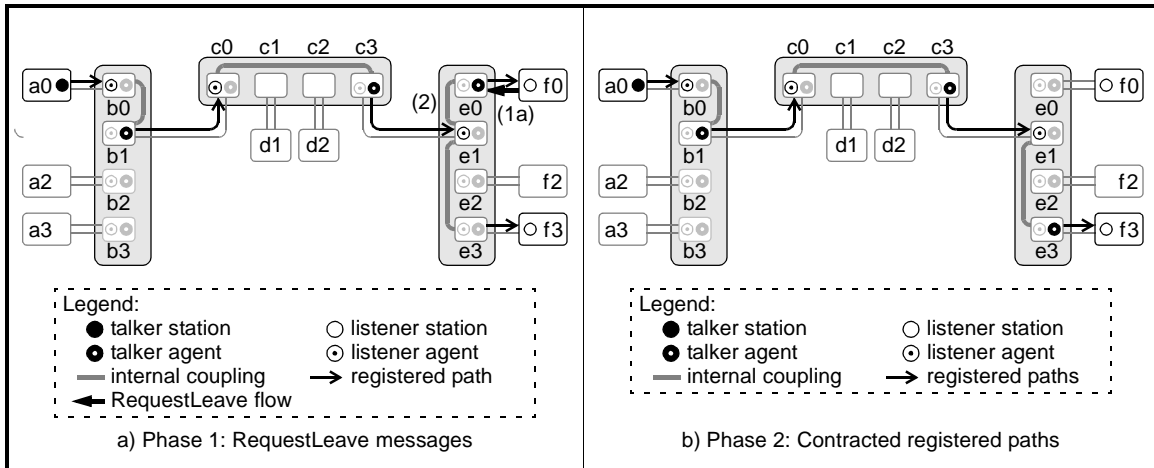


Figure 5.12—Side-path deregistration

5.4.9 Final deregistration

The final retiring listener also sends a RequestLeave message (1a) towards the talker. In this case, variants of that message {2a,3a,4a} eventually propagate to the talker, as illustrated in Figure 5.13-a. No listeners remain registered after this cascaded propagation of RequestLeave messages, as illustrated in Figure 5.13-b.

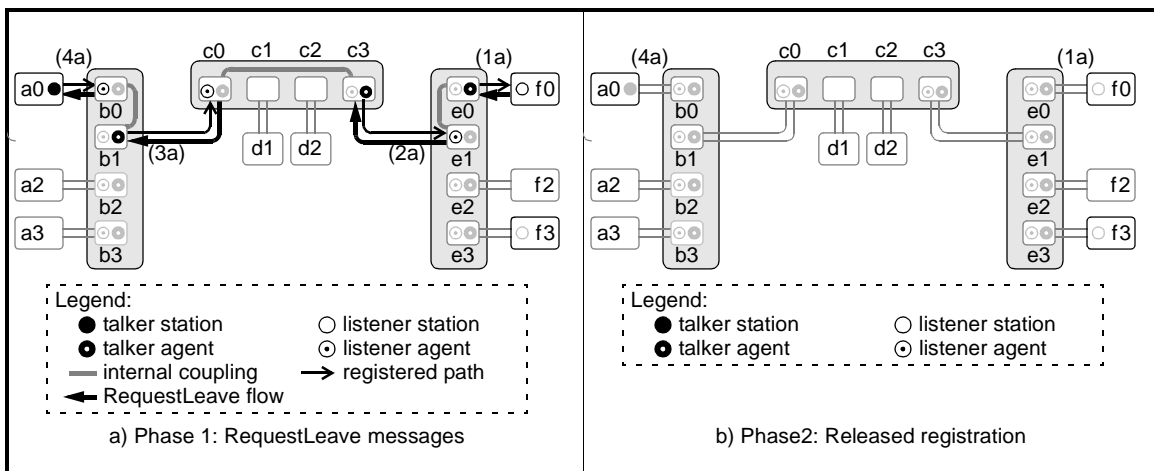


Figure 5.13—Final-path deregistration

5.4.10 Stream transmissions

Once listeners are registered (see Figure 5.14-a), a talker communicates critical parameters within the ConfirmPath message (instead of the initial ConfirmJoin messages) and starts its stream transmissions over the registered paths, as illustrated by the arrows in Figure 5.14-b.

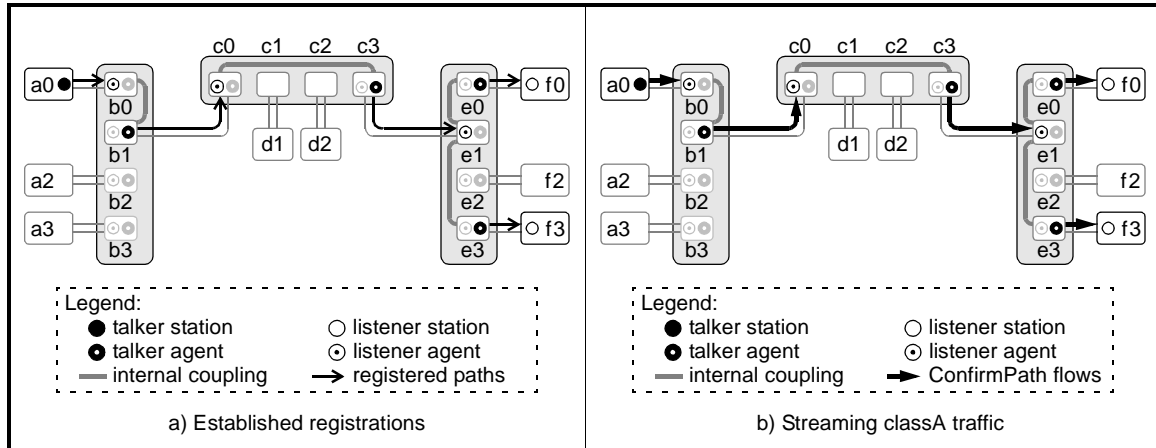


Figure 5.14—Streaming data over registered paths

The ConfirmPath message could be a variant of the ConfirmJoin message with a distinct command-code value. Like the baseline ConfirmJoin message, the ConfirmPath message is also sufficient to sustain the talker’s registration. This simplifies the talkers (and talker agents) by eliminating the need to concurrently transmit two distinct periodic registration-sustaining messages.

5.4.11 Insufficient bandwidth conditions

The available link bandwidths can sometimes be insufficient when the talker starts its stream transmissions. For example, bandwidths may be sufficient to sustain listener *f0* but not listener *f3*, as illustrated by the *e0-to-f0* and *e3-to-f3* paths in Figure 5.15-a, respectively.

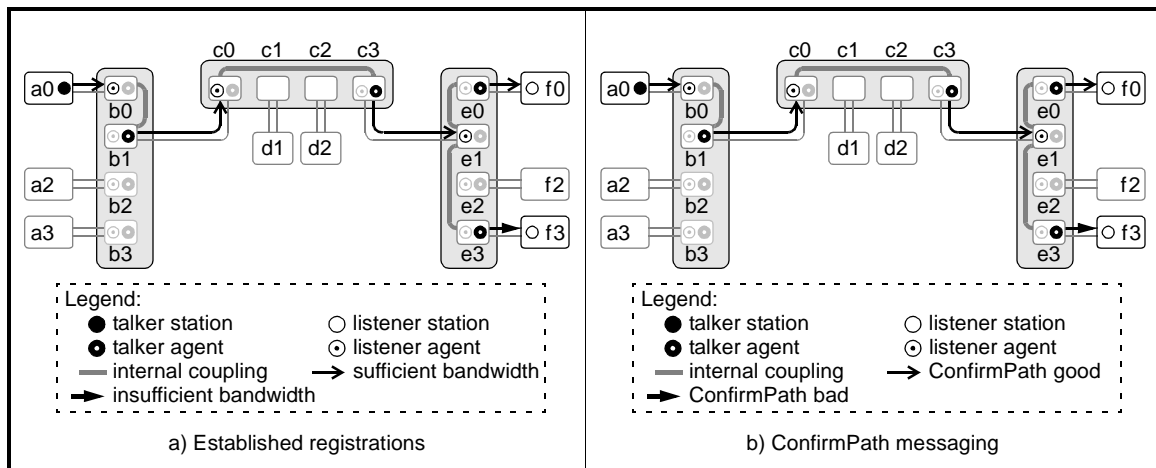


Figure 5.15—Insufficient bandwidth conditions

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1 In this case, listener *f3* does not receive the talker's streaming classA traffic. However, listener *f3* continues
2 to receive its ConfirmJoin messages, each of which contains an error indication code. Listener *f3* is thus
3 informed of the insufficient-bandwidth error condition, allowing corrective/reporting actions to be initiated
4 by higher level protocols.
5

6 **5.4.12 Errors conditions**

7
8 Errors may be associated with a variety of failure conditions, including (but not limited to) those listed
9 below.

- 10 a) Resources. Insufficient resources are available within the bridge.
11 (These insufficient-resource errors are handled by GARP specified mechanisms, see TBD.)
12 1) Insufficient registration-table storage is available in the bridge's downstream talker agent.
13 2) Insufficient registration-table storage is available in the bridge's upstream listener agent.
14
15 b) Bandwidth. Insufficient bandwidths are available within the bridge.
16 (These insufficient-bandwidth errors are handled by ConfirmJoin error codes, see 5.4.11.)
17 1) Insufficient bandwidth is available on the link from the talker agent to its adjacent listener.
18 2) Insufficient link or memory bandwidth is available with the bridge.
19

20 **5.4.13 Heartbeat timeouts**

21
22 Talker agents/stations are responsible for periodically polling locally registered listener agents/stations, to
23 demonstrate their continued presence. In the absence of these polling updates, the listeners assume the talker
24 is absent and deregister the inactive path (or inactive branch from the path). These talker-absent timeouts are
25 performed independently on each span.
26

27 Listener agents/stations are responsible for periodically reregistering with locally registered talker
28 agents/stations, to confirm their continued presence. In the absence of these reregistration updates, the
29 talkers assume the listener is absent and deregister the inactive path (or inactive branch from the path).
30 These listener-absent timeouts are performed independently on each span.
31

32 These periodic heartbeat-based timeouts handle a variety of error conditions, including the following:

- 33 a) A RequestJoin, RequestLeave, ConfirmJoin, or ConfirmPath is (corrupted and) not delivered.
34 b) The physical topology is changed, causing changes in the paths of streaming classA traffic.
35 c) A talker or listener is decommissioned and thus is no longer functionally present.
36 d) A flooded RequestJoin message reaches a non-talker end station or subnet.
37 e) After the talker's port is learned, a bridge discontinues flooding extraneous RequestJoin messages.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5.4.14 Untended flooding

Registering a new listener normally involves cascaded RequestJoin message sent from the listener $f0$ towards the talker $a0$, as illustrated in Figure 5.10-a. In some cases, the talker's address may be unlearned and flooding may be necessary. Thus, BridgeB could sometimes be forced to flood the RequestJoin to stations $\{a0, a2, a3\}$, when an unlearned address can't be directed to station $a0$, as illustrated in Figure 5.10-b.

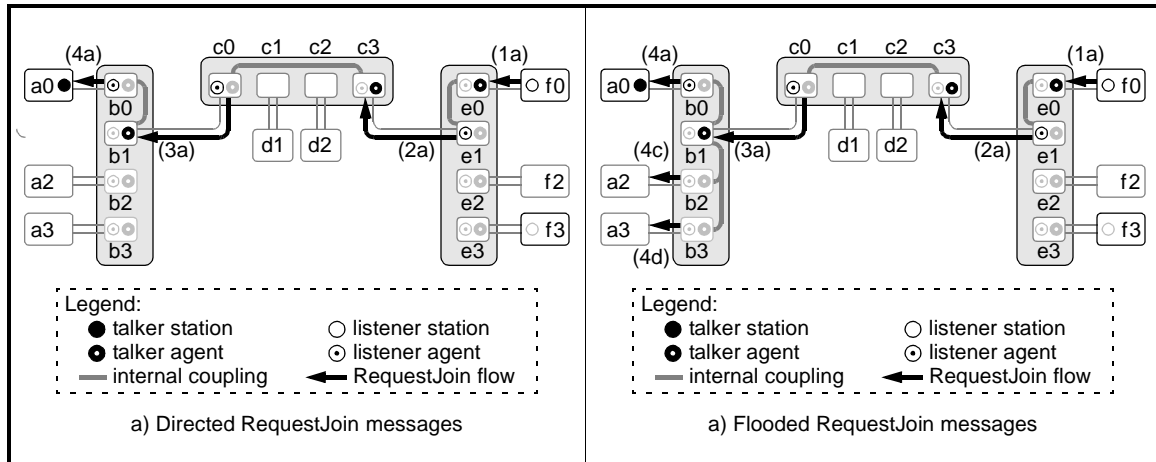


Figure 5.16—Periodic registration messages

In this example, talker $a0$ is present and its ConfirmJoin messages will soon propagate back to bridgeB, where the address of talker station $a0$ is learned. When this occurs, the flooding to stations $\{a2, a3\}$ stops.

Editors' Notes: To be removed prior to final publication.

Additional discussions may be appropriate to discuss what happened when the talker address is absent, as simply summarized below.

As noted previously (see 5.4.13), the talker agent is responsible for providing confirming ResponseJoin messages, so that the absence of a talker station can be readily detected. Allocated registration-table entries within bridges can be released after the talker-station absence is detected. Thus, flooding causes no harm.

5.4.15 GARP primitives

This subclause was intended to clarify the higher level SRP functionality. Thus, names of primitives were chosen for clarity, rather than consistency with the expected GARP messages. For the benefit of experienced GARP users, a sketch of the intended mappings of primitives is provided within this subclause.

The RequestJoin and RequestLeave messages correspond to like-names primitives within GARP. The ConfirmJoin and ConfirmPath messages correspond to variants of the leave-all messages within GARP.

5.5 Synchronized time-of-day clocks

5.5.1 Assumptions

This working paper specifies a protocol to synchronize independent timers running on separate stations of a distributed networked system, based on concepts specified within IEEE Std 1588-2002. Although a high degree of accuracy and precision is specified, the technology is applicable to low-cost consumer devices. The protocols are based on the following design assumptions:

- a) Each end station and intermediate bridges provide independent clocks.
- b) All clocks are accurate, typically to within ± 100 PPM.
- c) Point-to-point transmit/receive duplex connections are provided.
- d) Transmit/receive propagation delays within duplex cables are well matched.

5.5.2 Objectives

With these assumptions in mind, the time synchronization objectives include the following:

- a) Precise. Multiple timers can be synchronized to within 10's of nanoseconds.
- b) Inexpensive. For consumer A/V devices, the costs of synchronized timers are minimal. (GPS, atomic clocks, or 1PPM clock accuracies would be inconsistent with this criteria.)
- c) Scalable. The protocol is independent of the networking technology. In particular:
 - 1) Cyclical physical topologies are supported.
 - 2) Long distance links (up to 2 km) are allowed.
- d) Plug-and-play. The system topology is self-configuring; no system administrator is required.

5.5.3 Strategies

Strategies used to meet these objectives include the following:

- a) Precision is achieved by calibrating and adjusting *timeOfDay* clocks.
 - 1) Offsets. Offset value adjustments eliminate immediate clock-value errors.
 - 2) Rates. Rate value adjustments reduce long-term clock-drift errors.
- b) Simplicity is achieved by the following:
 - 1) Concurrence. Most configuration and adjustment operations are performed concurrently.
 - 2) Feed-forward. PLLs are unnecessary within bridges, but possible within applications.
 - 3) Symmetric. Clock-master/clock-slave computations are similar (only slave results are saved).
 - 4) Periodic. Messages are sent periodically, rather than in timely response to other requests.
 - 5) Frequent. Frequent (typically 1 kHz) interchanges reduces needs for precise clocks.
- c) Balanced functionality.
 - 1) Low-rate. Complex computations are infrequent and can be readily implemented in firmware.
 - 2) High-rate. Frequent computations are simple and can be readily implemented in hardware.

5.5.4 Grand-master selection

5.5.4.1 Grand-master selection

Clock synchronization involves streaming of clock-synchronization information from a grand-master timer to one or more slave timers. Although primarily intended for non-cyclical physical topologies (see Figure 5.17a), the synchronization protocols also function correctly on cyclical physical topologies (see Figure 5.17b), by activating only a non-cyclical subset of the physical topology.

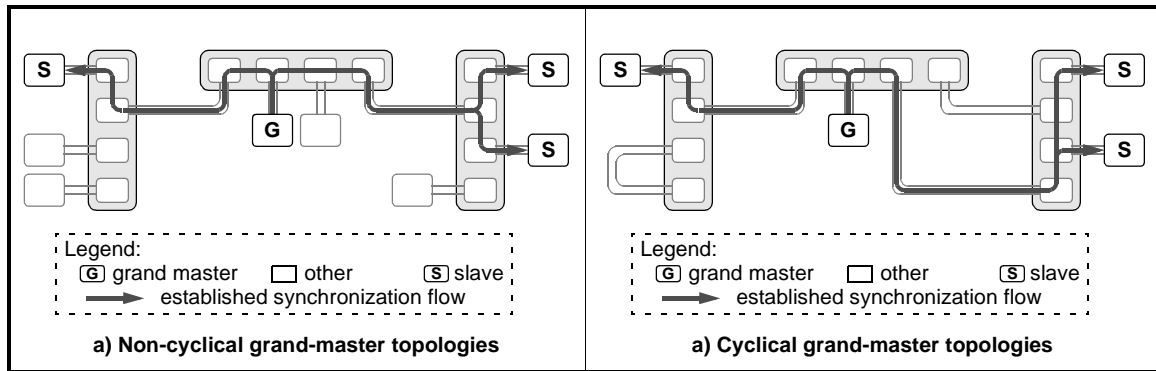


Figure 5.17—Timer synchronization flows

In concept, the clock-synchronization protocol starts with the selection of the reference-timer station, called a grand-master station (oftentimes abbreviated as grand-master). Each station is associated with a distinct preference value; the grand-master is the station with the “best” preference values.

Stations forward the best of their observed preference values to neighbor stations, allowing the overall best-preference value to be ultimately selected and known by all. The station whose preference value matches the overall best-preference value ultimately becomes the grand-master.

5.5.4.2 Communicated preference values

The grand-master station observes that its precedence is better than values received from its neighbors, as illustrated in Figure 5.18a. A slave stations observes its precedence to be worse than one of its neighbors and forwards the best-neighbor precedence value to adjacent stations, as illustrated in Figure 5.18b. To avoid cyclical behaviors, a *hopsCount* value is associated with preference values and is incremented before the best-precedence value is communicated to others.

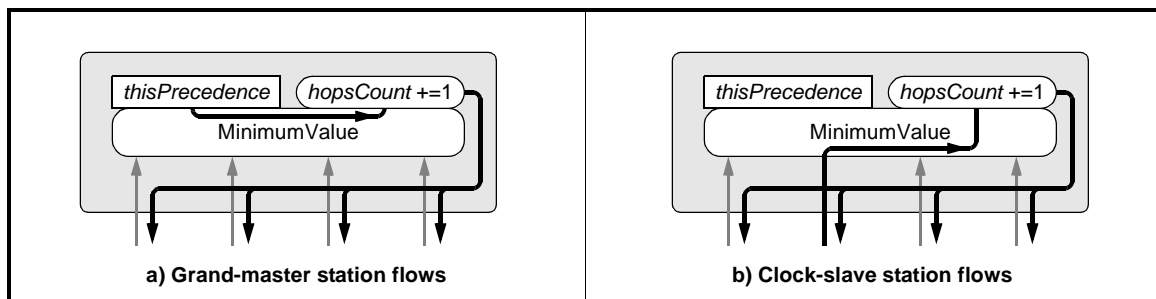


Figure 5.18—Grand-master precedence flows

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

The grand-master selection precedence includes multiple components, listed and described below (see 7.1.2). The *portTag* value is only needed within a bridge and is therefore not transmitted between stations.

- a) *systemTag*. A changeable value that is associated with each grand-master capable station. This value is can specify grand-master preferences (e.g., a home gateway may be preferred).
- b) *uniqueID*. A unique value associated with each station, typically based on its MAC address. This value is used as a tie breaker, when two contenders have identical *systemTag* values.
- c) *hopsCount*. A value that is incremented when passing through stations. This is the tie breaker, when two ports receive identical *systemTag:uniqueID* values.
- d) *portTag*. A changeable value that is associated with each port on a grand-master capable station. This is the tie breaker, when two ports receive identical *systemTag:uniqueID:hopsCount* values.

5.5.5 Synchronization principles

Timer synchronization is based on the concept of free-running local times (*localD*, *localE*, and *localF*) with compensating offset values (*offsetD*, *offsetE*, and *offsetF*), as illustrated in Figure 5.19. Updates involve changes to the offset values, not the free-running local timer values. In this example, we assume that: StationE is synchronized to its adjacent StationD; StationF is synchronized to its adjacent StationE. As a result, StationF is indirectly synchronized to StationD (through StationE).

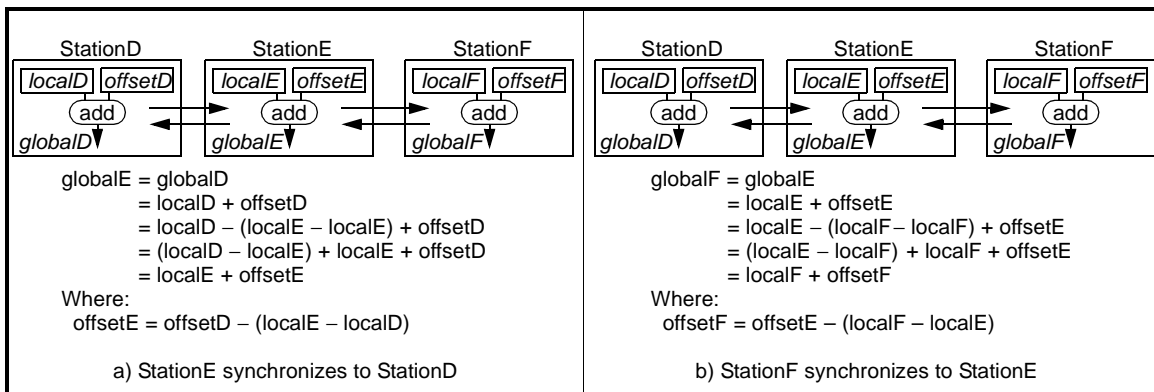


Figure 5.19—Time synchronization principles

The formulation of the *offsetE* value begins the assumption that the *globalE* and *globalD* times are identical. Addition of (*localE*–*localE*) and regrouping of terms leads to the formulation of the desired *offsetE* value, based on *offsetD* and (*localE*–*localD*) time difference values, as illustrated in Figure 5.19-a. Synchronization is thus possible using periodic transfers of *offsetD* values and computations of (*localE*–*localD*) timer

The formulation of the *offsetF* value begins the assumption that the *globalF* and *globalE* times are the identical. Addition of (*localF*–*localF*) and regrouping of terms leads to the formulation of the desired *offsetF* value, based on *offsetE* and (*localF*–*localE*) time difference values, as illustrated in Figure 5.19-b. Synchronization is thus possible using periodic transfers of *offsetE* values and computations of (*localF*–*localE*) timer differences.

In concept, the *offsetE* value is adjusted first; its adjusted value is then used to compute the desired *offsetF* value. In actuality, the periodic computations of *offsetE* and *offsetF* values are performed concurrently.

5.5.6 Timer snapshot locations

Mandatory jitter-error accuracies are sufficiently loose to allow transmit/receive snapshot circuits to be located with the MAC, as illustrated in Figure 5.20a. Vendors may elect to further reduce timing jitter by latching the receive/transmit times within the PHY, where the uncertain FIFO latencies can be best avoided.

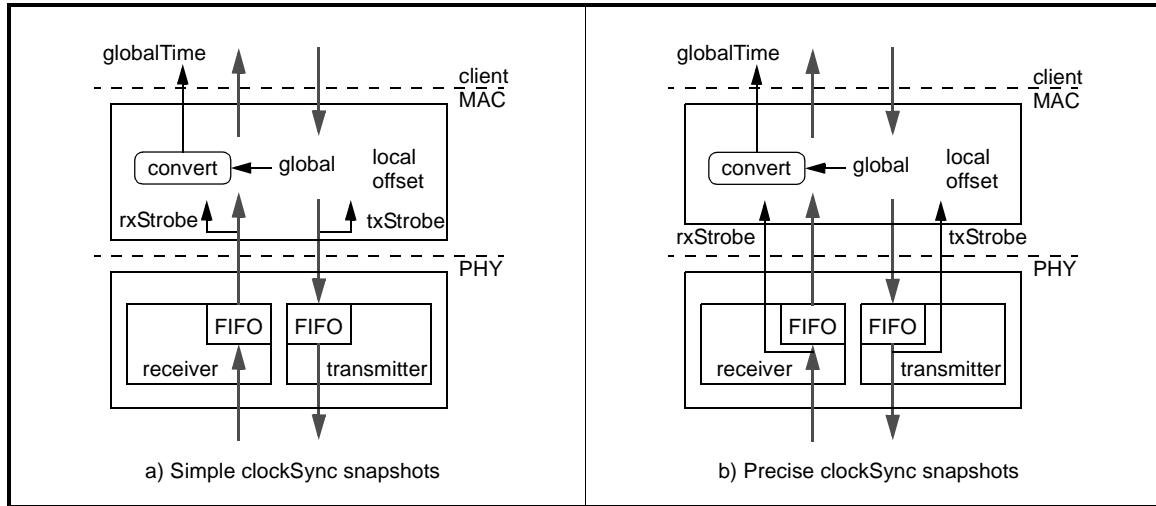


Figure 5.20—Timer snapshot locations

5.5.7 Bridge PLL possibilities

In addition to other valuable properties, the precise low-latency time-of-day synchronization protocols reduce jitter sufficiently to eliminate the needs for PLLs within bridges, as illustrated in Figure 5.21a. Elimination of such PLLs (illustrated in Figure 5.21b) simplifies the bridge design, while allowing each end-point application to independently optimize the effective capture-time and jitter-magnitude requirements of its PLL.

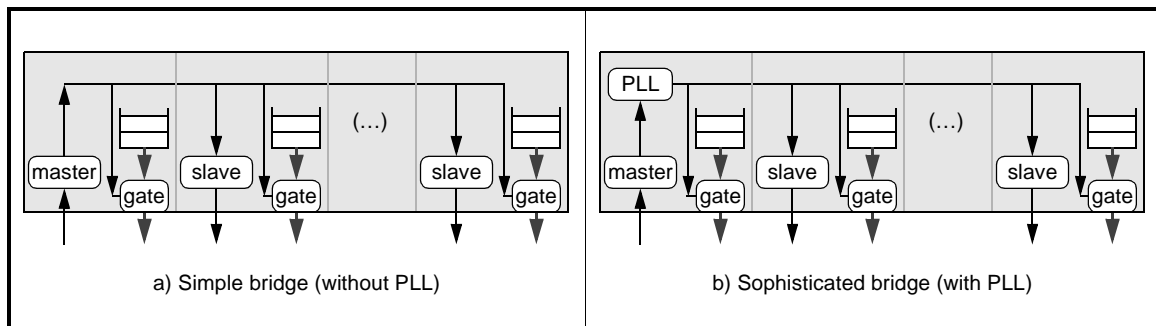


Figure 5.21—Bridge PLL possibilities

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5.6 Formats

5.6.1 Content framing

ClassA content is the client supplied per-cycle classA information, transferred from a talker to one or more listeners. The content within each cycle can be small or large; stereo audio stream transfers involve only approximately 20 bytes per cycle. Uncompressed 32-bits/pixel frame buffers (2 megapixels, 30Hz) would transmit 30 kilobytes per cycle. Framing of this content must be efficient for small sizes and sufficient for large sizes, as illustrated in Figure 5.22.

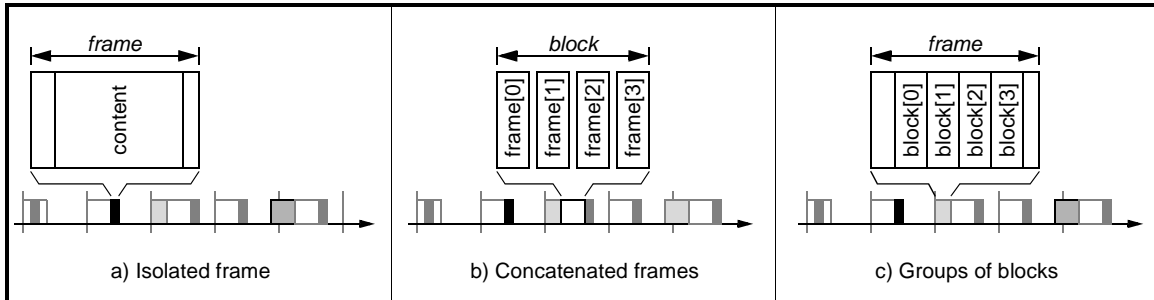


Figure 5.22—Content framing methods

For low bandwidth transmissions, each frame transports distinct classA content, as illustrated in Figure 5.22-a. For high bandwidth transmissions, the content can span multiple frames, as illustrated in Figure 5.22-b (see also C.3.2).

As an alternative improved-efficiency alternative, low bandwidth content could be encapsulated into blocks, where multiple blocks are included within each frame transmission, as illustrated in Figure 5.22-c. This allows the per-frame overhead (the inter-packet gap, header, and trailer fields) to be amortized over multiple blocks. For example, the eight inputs from a guitar may be packed together into the same frame. However, the packing of multichannel content is beyond the scope of this working paper.

Another approach would be to reduce the need for concatenated frames by using the (defacto standard) jumbo-frame sizes, which are approximately 9,000 bytes in size. However, support of the jumbo frame size is not ensured, and (when supported) is considerably less than 2^{16} -byte maximum size of an IEEE 1394 isochronous frame, or the 118 kilobyte size implied by 75% utilization of a 10Gb/s link.

5.6.2 Station plug addressing

Stream addressing is based on the concept of plugs, as illustrated in Figure 5.23. Streams are identified by their 48-bit talker-station identifier concatenated with that talker's 16-bit *plugId*. Each talker station may have up to 2^{16} streams, via logical plugs, in addition to the station's hardwired connections. Stations are expected to provide higher level commands for connecting/mixing/amplifying/converting/etc. data between combinations of hardwired and logical plugs. However, the details of such commands are beyond the scope of this working paper.

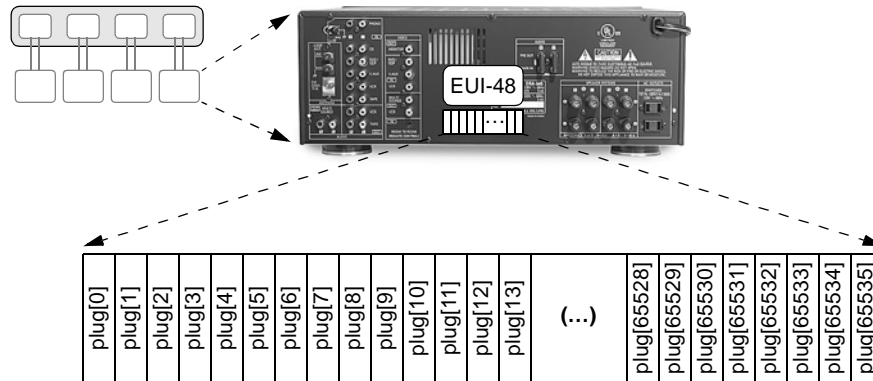


Figure 5.23—Plug addressing

5.6.3 Stream frame formats (alternative 1)

Streaming classA frames are no different than other multicast Ethernet frames. The distinction is that each of these multicast addresses is assumed to have associated *streamID* and bandwidth information saved within each forwarding bridge, as illustrated in Figure 5.24.

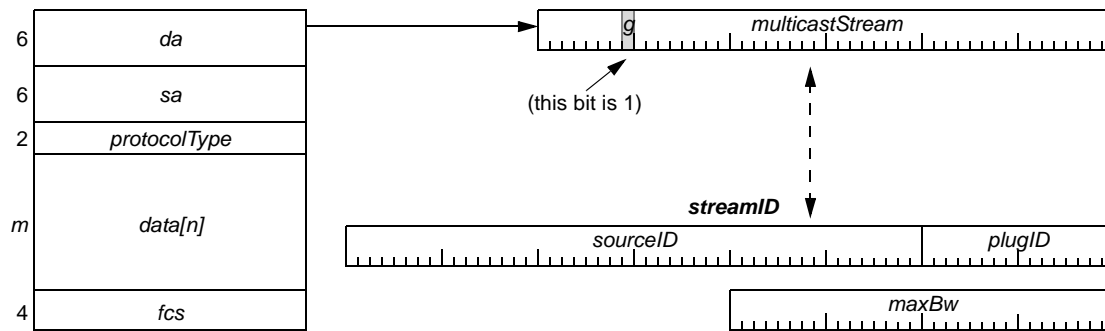


Figure 5.24—ClassA frame format and associated data

The *streamID* consists of two components: *sourceID* and *plugID*. The 48-bit *sourceID* identifies the source and usually equals the *sa* value; the *plugID* identifies the resource within that source. A distinct *maxBw* (maximum bandwidth) field identifies the negotiated maximum for classA bandwidth.

5.6.4 Stream frame formats (alternative 2)

Streaming classA frames are no different than other Ethernet frames. The distinction is that each of these frames supplies a nonzero user_priority field, as illustrated in Figure 5.25.

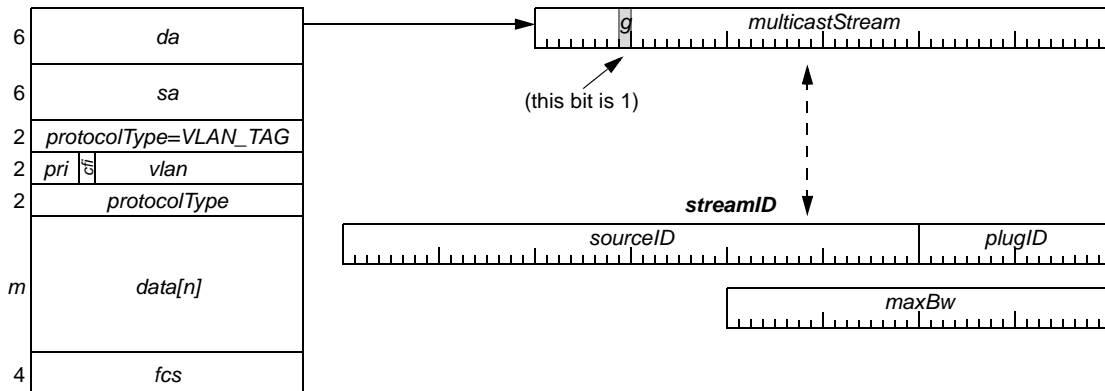


Figure 5.25—ClassA frame format and associated data

The *streamID* consists of two components: *sourceID* and *plugID*. The 48-bit *sourceID* identifies the source and usually equals the *sa* value; the *plugID* identifies the resource within that source. A distinct *maxBw* (maximum bandwidth) field identifies the negotiated maximum for classA bandwidth.

5.6.5 Stream frame formats (alternative 3)

Frames within a stream are no different than other Ethernet frames, with the exception of their distinct *da* (destination address) field, as illustrated in Figure 5.26. The most significant 32-bit portion of the *da* classifies the frame as an classA frame. The less significant 16-bit portion specifies the *plugID* portion of the *streamID* associated with the frame.

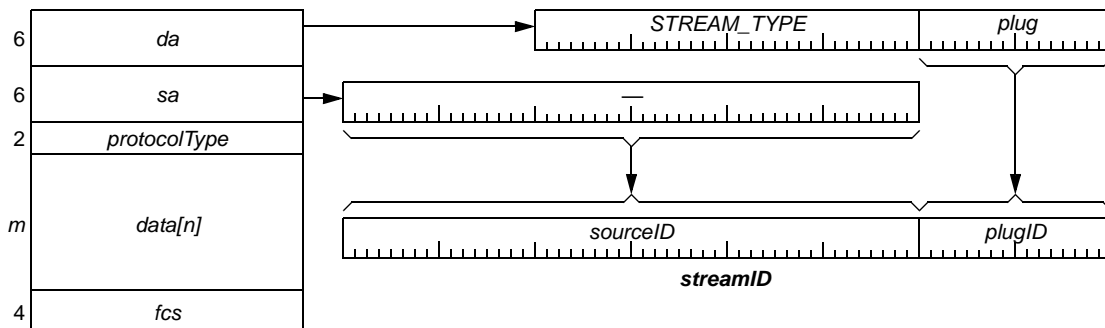


Figure 5.26—ClassA frame formats

5.6.6 Stream frame format alternatives comparison

Quality of service is thus specified by the user_priority field parameter within VLAN-tagged frames, as listed in Table 5.2.

Table 5.2—Tagged priority values

Alternative	Compact	Similar	Multicast server
1: DA-multicast	good	good	poor
2: VLAN-priority	poor	best	poor
3: SA-multicast	good	poor	good

The DA-multicast header is the compact, its forwarding mechanism are similar to those now supported, but a multicast server is required to provide unique multicast-stream addresses.

The VLAN-priority header is the 4 bytes larger, its forwarding mechanism is nearly identical to those now supported, but a multicast server is required to provide unique multicast-stream addresses.

The SA-multicast header is the compact, its forwarding mechanism is quite different than those now supported by bridges, but has the advantage that no multicast server is not required.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

5.7 Pacing

5.7.1 Assumptions

This working paper specifies a protocols for pacing classA traffic streams as they pass through multiple bridges. Although a high degree of scalability is implied, the technology is applicable to inexpensive consumer devices. The protocols are based on the following design assumptions:

- a) Sizes. The maximum frame size is assumed to be 2 kB, for consistency with established 802.3 frame-extension working group directions.
- b) Speeds. Only full-duplex 100 Mb/s, 1 Gb/s, and 100 Gb/s 100-meter links must be supported.
- c) Limits. The classA traffic transmissions are limited to 75% of the available link bandwidth.

5.7.2 Objectives

With these assumptions in mind, the time synchronization objectives include the following:

- a) Reliable. The worst-case delay between talker and listener stations is small, deterministic, and not effected by operating conditions, including the following:
 - 1) Loading. Arbitrary talker-station and listener-station traffic patterns can be supported.
 - 2) Scaling. Any 802.1 supported spanning tree topologies can be supported.
- b) Plug-and-play. Manual provisioning of the system is not required.
- c) Compatible. The pacing of high-class frames cannot disrupt legacy or lower-class transmissions.
- d) Friendly. Some higher-class traffic that cannot be reliably paced, due to legacy sources or bridges; retains precedence over lower-class traffic.
- e) Robust. Higher-class traffic never starves the forwarding of lower-class control traffic.
- f) Efficient. Unused higher-class bandwidth can be readily reclaimed lower-class traffic.

5.7.3 Strategies

Strategies used to meet these objectives include the following:

- a) Buckets. Higher-class traffic is grouped into buckets; buckets are forwarded every 125 μ s cycle.
- b) Limits. The levels of higher-class traffic are limited to 75% of the link bandwidths.
 - 1) Excess classA traffic above this 75% limit is discarded.
 - 2) Excess classB traffic above this 75% limit is temporarily processed as classC traffic.
- c) Reuse. Unused higher-precedence bandwidths are reused if not consumed as intended.
 - 1) Unused classA traffic within the 75% limit is available for classB traffic.
 - 2) Unused classA/classB traffic within this 75% limit is available for classC traffic.
- d) Downgrade. When passing through unsupportive bridges, classA traffic is downgraded to classB. The classB traffic is no longer paced, but retains its precedence over classC traffic.

5.7.4 Service classes

Pacing is intended to ensure timely delivery of pre-subscribed classA traffic, in the presence of arbitrary classB and classC loading conditions. Interactions between these three service classes is summarized below:

- a) ClassA. A pre-subscribed paced time-sensitive service with guaranteed latency. The classA traffic is paced and (at its scheduled transmit time) has priority over classB traffic.
- b) ClassB. A pre-subscribed time-sensitive service with guaranteed bandwidth. The classB traffic is shaped and has priority over classC traffic. Two type of classB traffic are expected, as follows:
 - 1) Legacy. Time-sensitive multicast traffic sourced by non-supportive talker stations.
 - 2) Hybrid. ClassA traffic that has passed through a nonsupportive bridge. Such previously-classA traffic can no longer be paced and therefore is downgraded to classB.
- c) ClassC. A best-effort service that utilizes bandwidths not consumed by classA or classB traffic. The classA and classB subscription/shaping restrictions ensure a minimum 25% of link bandwidths are available for classC transmissions.

5.7.5 Fine-grained pacing

Pacing involves the throttling of classA streams so that their average bandwidth can be guaranteed over small averaging intervals. Such fine-grained pacing has the following advantages:

- a) Latency. Talker-to-listener delays are small, deterministic, and link-utilization independent.
- b) Jitter. Delay variations between a talker and listeners are bounded and topology independent.
- c) Intervals. Short intervals simplify the detection/enforcement of maximum classA bandwidths. (A goal is to limit classA bandwidths to no more than 75% of the link capacity, see 1.2.3.)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

6. Frame formats

6.1 ClassA frames

6.1.1 ClassA frame fields

A classA frame differs from other frames in the format of its multicast *da* (destination address), as illustrated in Figure 6.1.

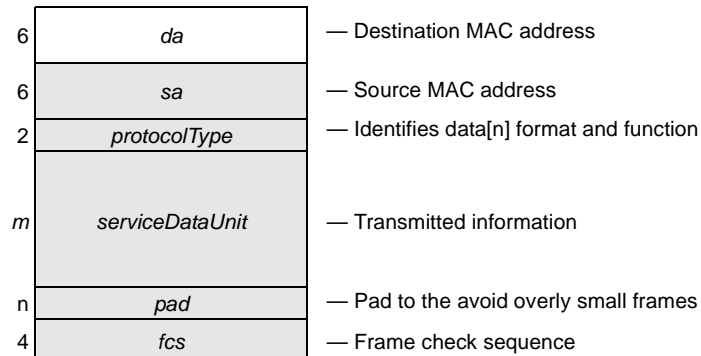


Figure 6.1—ClassA frame formats

6.1.1.1 *da*: A 6-byte (destination address) field that specifies a multicast address associated with the stream.

6.1.1.2 *sa*: A 48-bit (source address) field that specifies the local station sending the frame. The *sa* field contains an individual 48-bit MAC address (see 3.11) as specified in 9.2 of IEEE Std 802-2001.

6.1.1.3 *protocolType*: A 16-bit field contained within the payload. When the value of *protocolType* is greater than or equal to 1536 (600_{16}) the *protocolType* field indicates the nature of the MAC client protocol (type interpretation), selecting from values designated by the IEEE Type Field Register. When less than 1536 (0_{16} – $5FF_{16}$), the *protocolType* is interpreted as the length of the frame (length interpretation). The length and type interpretations of this field are mutually exclusive.

6.1.1.4 *serviceDataUnit*: An *m*-byte field the contains the service data unit provided by the client.

6.1.1.5 *pad*: If the sum of the other field lengths is less than 64 bytes, then the number of zero-valued *pad* bytes are sufficient to make a 64-byte frame. Otherwise, the *pad* field is not present.

6.1.1.6 *fcs*: A 4-byte (frame check sequence) field whose 32-bit CRC covers the frame's content.

6.2 clockSync frame format

6.2.1 clockSync fields

Clock synchronization (clockSync) frames facilitate the synchronization of neighboring clock span-master and clock span-slave stations. The frame, which is normally sent once each isochronous cycle, includes time-snapshot information and the identity of the network's clock master, as illustrated in 6.2. The gray boxes represent physical layer encapsulation fields that are common across all Ethernet frames.

6	<i>da</i>	— Destination MAC address
6	<i>sa</i>	— Source MAC address
2	<i>protocolType</i>	— Distinguishes RE frames from others (see 6.5.1)
1	<i>subType</i>	— Distinguishes clockSync from other RE frames (see 6.5.2)
1	<i>hopsCount</i>	— Hop count from the grand master
1	<i>syncCount</i>	— Sequence number for clockSync frames
1	<i>cycleCount</i>	— Cycle count for pacing
2	<i>systemTag</i>	— More-significant grand-master election precedence
8	<i>uniqueID</i>	— Less-significant grand-master election precedence
8	<i>lastFlexTime</i>	— Incoming link's frame transmission time (1 cycle delayed)
8	<i>deltaTime</i>	— Outgoing link's frame propagation time
8	<i>offsetTime</i>	— Offset time within the neighbor
4	<i>diffRate</i>	— Cumulative rates from the grand-master
4	<i>lastBaseTime</i>	— Incoming link's frame transmission time (1 cycle delayed)
4	<i>fcs</i>	— Frame check sequence

Figure 6.2—clockSync frame format

6.2.1.1 *da*: A 48-bit (destination address) field that specifies the station(s) for which the frame is intended. The *da* field contains either an individual or a group 48-bit MAC address (see 3.11), as specified in 9.2 of IEEE Std 802-2001.

6.2.1.2 *sa*: A 48-bit (source address) field that specifies the local station sending the frame. The *sa* field contains an individual 48-bit MAC address (see 3.11), as specified in 9.2 of IEEE Std 802-2001.

6.2.1.3 *protocolType*: A 16-bit field contained within the payload that identifies the format and function of the following fields (see 6.5.1).

6.2.1.4 *subType*: A 16-bit field that identifies the format and function of the following fields (see 6.5.2).

6.2.1.5 *hopsCount*: An 8-bit field that identifies the maximum number of hops between the talker and associated listeners.

6.2.1.6 *syncCount*: An 8-bit field that is incremented on each clockSync frame transmission.

6.2.1.7 *cycleCount*: A 7-bit field that equals $(cycle \% 125)$, where *cycle* represents units of 125 μ s within the transmitting station's *timeOfDay* value.

6.2.1.8 *systemTag*: A 16-bit field that has highest precedence in the grand-master selection protocols.

6.2.1.9 *uniqueID*: A 64-bit field that specifies the precedence of the grand clock master, specified in 6.2.3.

6.2.1.10 *lastFlexTime*: A 64-bit field that specifies the time within the source station when the previous clockSync frame was transmitted. The format of this field is specified in 6.2.4.

6.2.1.11 *deltaTime*: A 64-bit field that specifies the differences between clockSync receive and transmit times, as measured on the opposing link. The format of this field is specified in 6.2.4.

6.2.1.12 *offsetTime*: A 64-bit field that specifies the offset time within the source station. The format of this field is specified in 6.2.4.

6.2.1.13 *diffRate*: A 32-bit field that specifies the *diffRate* value within the source station.

6.2.1.14 *lastBaseTime*: A 32-bit field that specifies the *timer1* value within the source station when the previous clockSync frame was transmitted.

6.2.1.15 *fcs*: A 32-bit (frame check sequence) field that is a cyclic redundancy check (CRC) of the frame.

6.2.2 *systemTag* subfields

The format of the 16-bit *systemTag* field is based on the format of the spanning tree protocol precedence value, as illustrated in Figure 6.3.

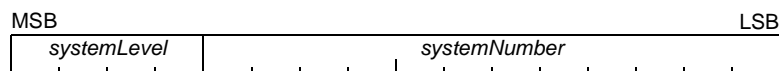


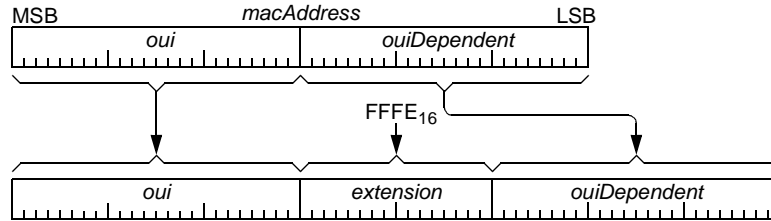
Figure 6.3—*systemTag* subfields

6.2.2.1 *systemLevel*: A 4-bit field that comprise a settable priority component that permits the relative priority of bridges to be managed.

6.2.2.2 *systemNumber*: A 12-bit field that comprise a locally assigned system identifier extension. (The term *systemID* is equivalent to 'system ID', as specified within IEEE Std 802.1D-2004.)

1 **6.2.3 uniqueID fields**

2
 3 The format of the 64-bit *uniqueID* field is a unique identifier. For stations that have a uniquely assigned
 4 48-bit *macAddress*, the 64-bit *uniqueID* field is derived from the 48-bit MAC address, as illustrated in
 5 Figure 6.4.



16 **Figure 6.4—uniqueID format**

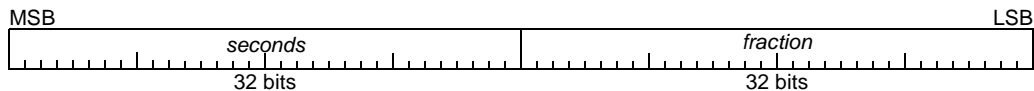
17
 18 **6.2.3.1 oui:** A 24-bit field assigned by the IEEE/RAC (see xx).

19
 20 **6.2.3.2 extension:** A 16-bit field assigned to encapsulated EUI-48 values.

21
 22 **6.2.3.3 ouiDependent:** A 24-bit field assigned by the owner of the *oui* field (see xx).

23
 24 **6.2.4 Time field formats**

25
 26 Time-of-day values within a frame are specified by 64-bit values, consistent with IETF specified NTP[B8]
 27 and SNTP[B9] protocols. These 64-bit values consist of two components: a 32-bit *seconds* and 32-bit
 28 *fraction* fields, as illustrated in Figure 6.5.



34 **Figure 6.5—Complete seconds timer format**

35
 36 **6.2.4.1 seconds:** A 32-bit field that specifies time in seconds.

37
 38 **6.2.4.2 fraction:** A 32-bit field that specified time offset within the second, in units of 2^{-32} second.

39
 40 The concatenation of 32-bit *seconds* and 32-bit *fraction* field specifies a 64-bit *time* value, as specified by
 41 Equation 6.1.

42
 43
$$time = seconds + (fraction / 2^{32}) \tag{6.1}$$

44 Where:

45 *seconds* is the most significant component of the time value (see Figure 6.5).

46 *fraction* is the less significant component of the time value (see Figure 6.5).

6.3 Subscription frame

6.3.1 Subscription frame fields

Subscription frames contain channel-acquisition information, as illustrated in Figure 6.6.

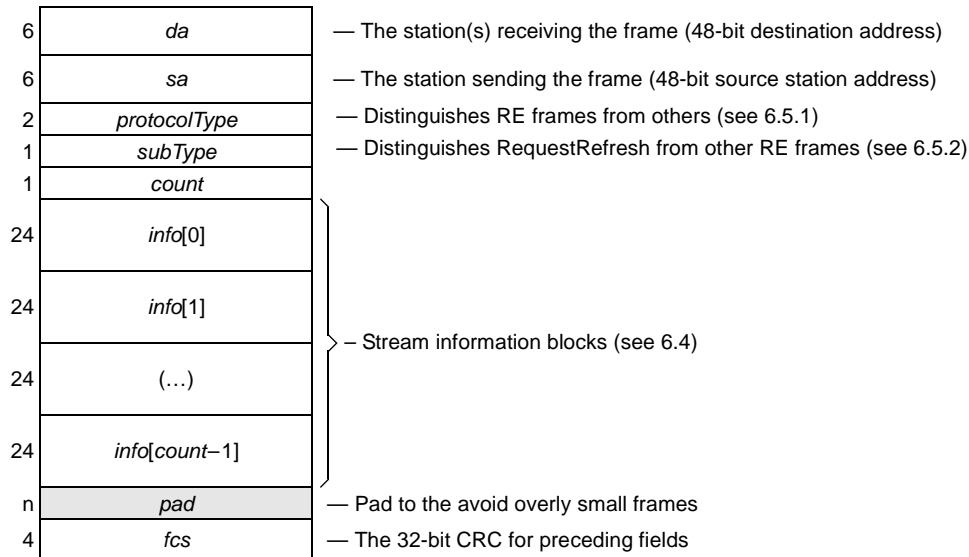


Figure 6.6—Subscription frame format

6.3.1.1 *da*: A 6-byte (destination address) field that normally specifies the destination address for the frame transmission, with unicast and multicast forms.

6.3.1.2 *sa*: A 6-byte (source address) field that normally specifies the source address for the frame transmission. If a bridge is present between the frame and its associated listener, the *sa* value identifies the bridge.

6.3.1.3 *protocolType*: A 2-byte field that normally specifies the frame length, or the format and function of the following fields (excluding the 4-byte *fcs* field). This RE assigned value distinguishes its frame formats from others (see 6.5.1).

6.3.1.4 *subType*: A 1-byte field that distinguishes the ResponseError frame from other frames defined within this working paper.

6.3.1.5 *count*: A 1-byte field that specifies the number of elements within the following *info*-block array.

6.3.1.6 *info*: A 24-byte array element that provides listener subscription information (see 6.4).

6.3.1.7 *pad*: If the sum of the other field lengths is less than 64 bytes, then the number of zero-valued *pad* bytes are sufficient to make a 64-byte frame. Otherwise, the *pad* field is not present.

6.3.1.8 *fcs*: The 4-byte (frame check sequence) field whose 32-bit CRC covers the frame's content. For RE content frames, the standard definition applies.

6.4 Common *info* field format

Many frame transports an array of one or more *info*[] fields, whose content is illustrated in Figure 6.7.

2	<i>command</i>	— Database action command
6	<i>talkerID</i>	— Multicast talker identifier
2	<i>plugID</i>	— Resource within the talker
6	<i>mcastID</i>	— Multicast destination label
2	<i>maxCycles</i>	— Delay from the talker
4	<i>maxBw</i>	— Maximum required bandwidth
2	<i>reserved</i>	— Reserved

Figure 6.7—Common *info* field format

6.4.1 *command*: A 2-byte field that differentiates between database-update actions.

6.4.2 *talkerID*: A 6-byte field that identifies the stream’s talker.

6.4.3 *plugID*: A 16-bit field that specifies the plug identifier within the talker.

The concatenation of the 48-bit *talkerID* and 16-bit *plugID* fields forms a 64-bit *streamID* that uniquely identifies the classA multicast stream.

6.4.4 *mcastID*: A 6-byte (multicast identifier) field that routes frames between the talker and audience.

6.4.5 *maxCycles*: A 2-byte field that is updated by bridges, as the RequestRefresh flows from the talker to the listener, allowing the maximum number of delay cycles between the talker and listener stations to be known to the talker.

6.4.6 *maxBw*: A 4-byte field that specifies the level of negotiated classA bandwidth, measured in bytes of per-cycle content.

6.4.7 *reserved*: A 2-byte zero-valued field that is ignored.

6.5 Unique identifier values

6.5.1 *protocolType* identifier

NOTE—The following *protocolType*-assignment text will ultimately be updated with assigned values.

The clockSync (see 6.2) and subscription (see 6.3) frames are distinguished from other frames by their 16-bit distinct *protocolType* value, as illustrated in Figure 6.8. The following 1-byte *subType* field further distinguishes between these uses (see 6.5.2).

Assigned *protocolType* value:
QR-ST

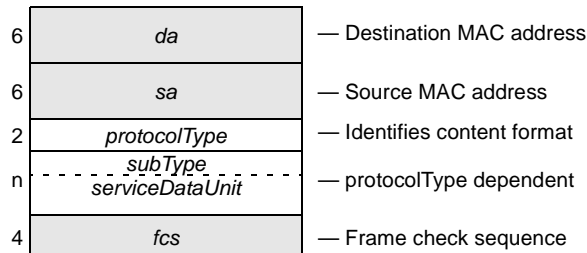


Figure 6.8—*protocolType* field value

6.5.2 *subType* identifier

Distinct *subType* identifiers distinguish between RE frame types, as specified by Table 6.1.

Table 6.1—Assigned *subType* identifiers

Value	Name	Row	See	Description
TBD	CLOCK_SYNC	1	6.2	Demarcates boundaries between isochronous cycles.
192-255	E1394	2	C.2.2	Encapsulated IEEE 1394 packet (or portion of 1394 packet)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

7. Clock synchronization

7.1 Clock-synchronization overview

7.1.1 Clock synchronization services

Clock synchronization involves the transmission and reception of clockSync frames interchanged between adjacent-span stations, using the state machines defined within this clause. When considered as a whole, these provide the following services:

- a) Election. The grand clock master is elected from among the grand-clock-master capable stations.
- b) Isolation. Timeouts identify the boundaries, beyond which RE services are not supported.
- c) Clock-sync. Clock-slave stations are synchronized to the grand master station's time reference.

7.1.2 Grand-master precedence

Grand-master precedence is based on the concatenation of multiple fields, as illustrated in Figure 7.1. The *portTag* value is used within bridges, but is not transmitted between stations.

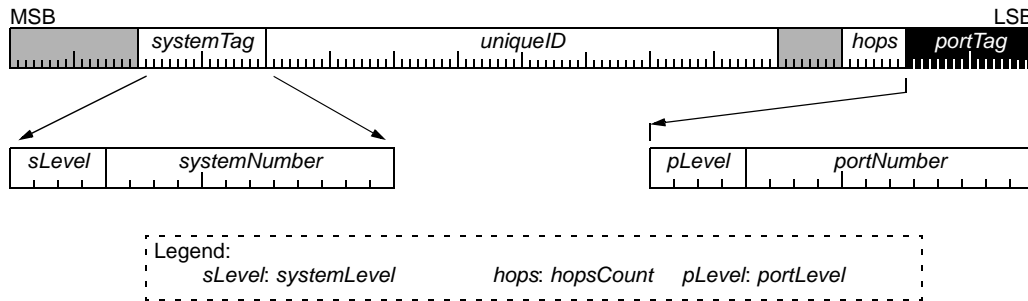


Figure 7.1—Grand-master precedence

7.1.3 Clock-synchronization agents

Clock-synchronization information conceptually flows from a grand-master station to clock-slave stations, as illustrated in Figure 7.2a. A more detailed illustration shows pairs of synchronized clock-master and clock-slave components, as illustrated in Figure 7.2b.

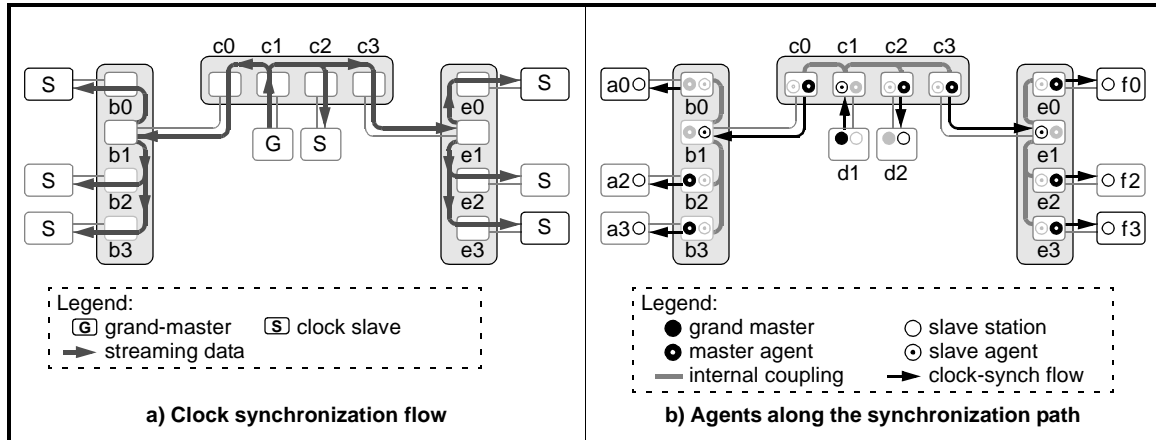


Figure 7.2—Hierarchical flows

7.1.4 Clock-synchronized pairs

Each bridge port provides clock-master and clock-slave agents, although both are never simultaneously active. External communications (see 7.2b) synchronize clock-slaves to clock-masters, as listed in Table 7.1.

Table 7.1—External clock-synchronization pairs

Grand master	Clock master agent	Clock slave agent	Clock slave	Type of synchronization
d1	—	c1	—	Station-to-bridge
—	c0	b1	—	Bridge-to-bridge
—	c3	e1	—	
—	b0	—	a0	Bridge-to-station
—	b2	—	a2	
—	b3	—	a3	
—	c2	—	d2	
—	e0	—	f0	
—	e2	—	f2	
—	e3	—	f3	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1 Internal communications distribute synchronized time from clock-slave agents b1, c1, and e1 to the other
2 clock-master agents on bridgeB, bridgeC, and bridgeE respectively. However, bridge-internal port-to-port
3 synchronization protocols are implementation-dependent and beyond the scope of this working paper.
4

5 Within a clock-slave, precise time synchronization involves adjustments of timer offset and rate values. The
6 adjustments of the timer's offset is called offset synchronization (see 7.1.6); the adjustments of the timer's
7 rate is called rate synchronization (see 7.1.8). Both involve calibration of local clock-master/clock-slave dif-
8 ferences and the propagation of cumulative differences in the clock-slave direction, as described by the C
9 code of Annex J.

10
11 Time synchronization yields distributed but closely-matched *timeOfDay* values within stations and bridges.
12 No attempt is made to eliminate intermediate jitter with bridge-resident jitter-reducing phase-lock loops
13 (PLLs.) but application-level phase locked loops (not illustrated) are expected to filter high-frequency jitter
14 from the supplied *timeOfDay* values
15

16 **7.1.5 Clock-synchronization intervals**

17
18 Clock synchronization involves the processing of periodic events. Three distinct time periods are involved,
19 as listed in Table 7.2. The clock-period events trigger the update of free-running timer values; the period
20 affects the timer-synchronization accuracy and is therefore constrained to be small.
21

22
23 **Table 7.2—Clock-synchronization intervals**

24
25

Name	Time	Description
clock-period	< 20 ns	Time between timer-register value updates
send-period	10 ms	Time between sending of periodic clockSync frames between adjacent stations
slow-period	100 ms	Time between computation of clock-master/clock-slave rate differences

26
27
28
29
30
31

32
33 The send-period events trigger the interchange of clockSync frames between adjacent stations. While a
34 smaller period (1 ms or 100 μs) could improve accuracies, the larger value is intended to reduce costs by
35 allowing computations to be executed by inexpensive (but possibly slow) bridge-resident firmware.
36

37 The slow-period events trigger the computation of timer-rate differences. The timer-rate differences are
38 computed over two slow-period intervals, but recomputed every slow-period interval. The larger 100 ms (as
39 opposed to 10 ms) computation interval is intended to reduce errors associated with sampling of
40 clock-period-quantized slow-period-sized time intervals.
41
42
43
44
45
46
47
48
49
50
51
52
53
54

7.1.6 Offset synchronization

Offset synchronization involves a subset of the time-synchronization components, as illustrated by white-colored boxes in Figure 7.5. Each clock consists of a progressing *timeOfDay* value, whose offset and rate are periodically adjusted. The free-running *flexTimer* timer is never reset; synchronization of stationE (with respect to stationD) is accomplished by adjustments to the *flexOffset* and *flexRate* values within stationE.

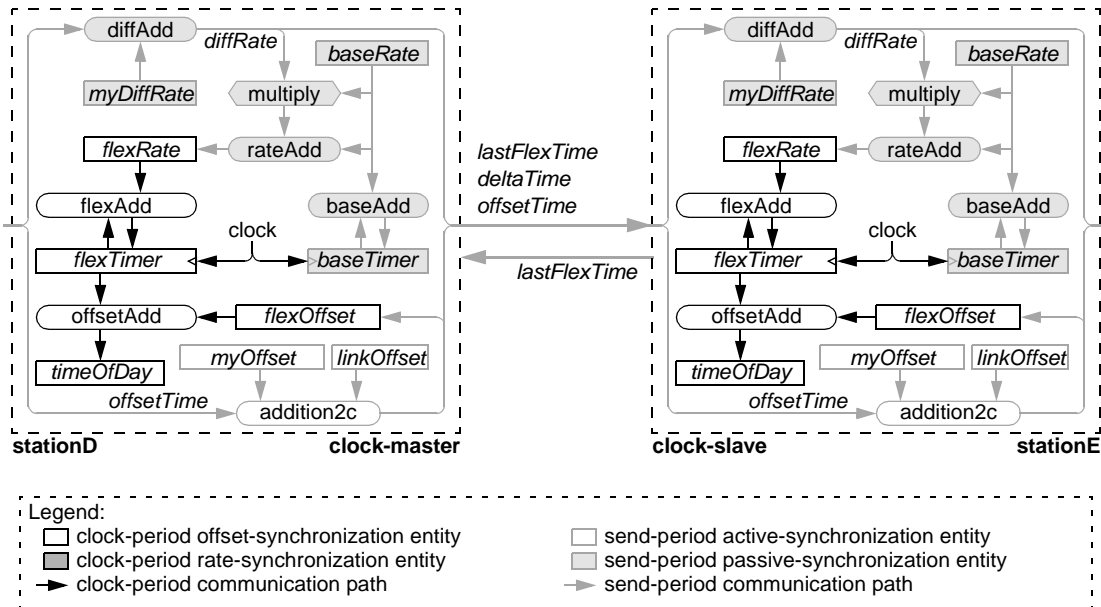


Figure 7.3—Offset synchronization

The offset-synchronization protocols interchange parameters periodically, possibly every 10 ms. The *lastFlexTime*, *deltaTime*, and *offsetTime* values are sent periodically from the clock-master to the clock-slave. The *lastFlexTime* is sent periodically from the clock-slave to the clock-master, providing information necessary for the clock-master to produce a *deltaTime* value for the clock-slave.

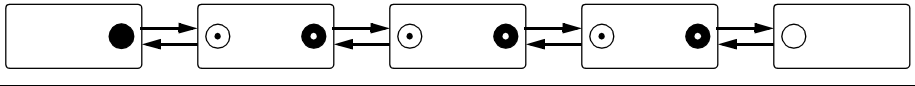
The offset-compensation protocols for stationE adjust its *myOffset* value so that the instantaneous values of *stationE.timeOfDay* and *stationD.timeOfDay* are the same. Computations are performed on clockStrobe reception and clockStrobe transmission.

As an option, an additional *linkOffset* value is available to manually compensate for mismatched transmit-link/receive-link duplex-cable delays on the clock-master side. The *linkOffset* value is expected to be manually set when the cable mismatch is known through other mechanisms, such as specialized cable-characterization equipment.

The station's *offsetTime* value is constructed by adding the received *clockStrobe.offsetTime*, local *myOffset*, and local *linkOffset* values. This revised *clockStrobe.offsetTime* value is used within each station and is passed to the downstream neighbor (when such a neighbor is present).

7.1.7 Cascaded offsets

The concept of cascaded offset values can be better understood by considering a simple 3-bridge example, as illustrated in Figure 7.4. The slave-agent in bridgeB is synchronized to its neighbor grand-master via clockSync frames sent on the connecting bidirectional span. Within bridgeB, the clock-slave agent passes the time directly to the clock-master agent. The slave-agent in bridgeC is synchronized to its neighbor clock-master via clockSync frames sent on the connecting bidirectional span. Other ports are similarly synchronized, thus synchronizing the right-most clock-slave station to the left-most grand-master station.



Parameter	grand-master	bridgeB	bridgeC	bridgeD	clock-slave
name	grand-master	bridgeB	bridgeC	bridgeD	clock-slave
number	1	2	3	4	5
<i>flexTimer</i>	100	500	-300	200	400
<i>myOffset</i>	10	-400	800	-500	-200
<i>flexOffset</i>	10	-390	410	-90	-290
<i>timeOfDay</i>	110				

Representing:
 $myOffset[k+1] = flexTimer[k] - flexTimer[k+1];$
 $flexOffset[k+1] = flexOffset[k] + myOffset[k+1];$
 $timeOfDay[k] = flexTimer[k] + flexOffset[k];$

Figure 7.4—Cascaded offsets (a possible scenario)

To simplify this illustration, consider only the seconds portion of the *flexTimer* value within each station or bridge. These values may differ dramatically, based (perhaps) on the power-cycling or topology formation sequence. Thus, the grand-master could have a *flexTimer* value of 100 while its bridgeB neighbor has a *flexTimer* value of 500.

The *myOffset* value within bridgeB will converges to the value of -400, representing the differences between grand-master and bridgeB *flexTimer* values. The *flexOffset* value received from the grand-master is added to this *myOffset* value, so that bridgeB's *flexOffset* becomes -390. The *flexTimer* and *flexOffset* values are added, to yield a resultant bridgeB *timeOfDay* value of 110, properly synchronized to the identical grand-master's value.

Similarly, bridgeC is synchronized to bridgeB, bridgeD to bridgeC, and the clock-slave to bridgeD.

7.1.8 Rate synchronization

Rate synchronization involves a subset of the time-synchronization components, as illustrated by white-colored boxes in Figure 7.5. The free-running *baseTimer* timer facilitate the determination of rate differences between the clock-master and clock-slave stations.

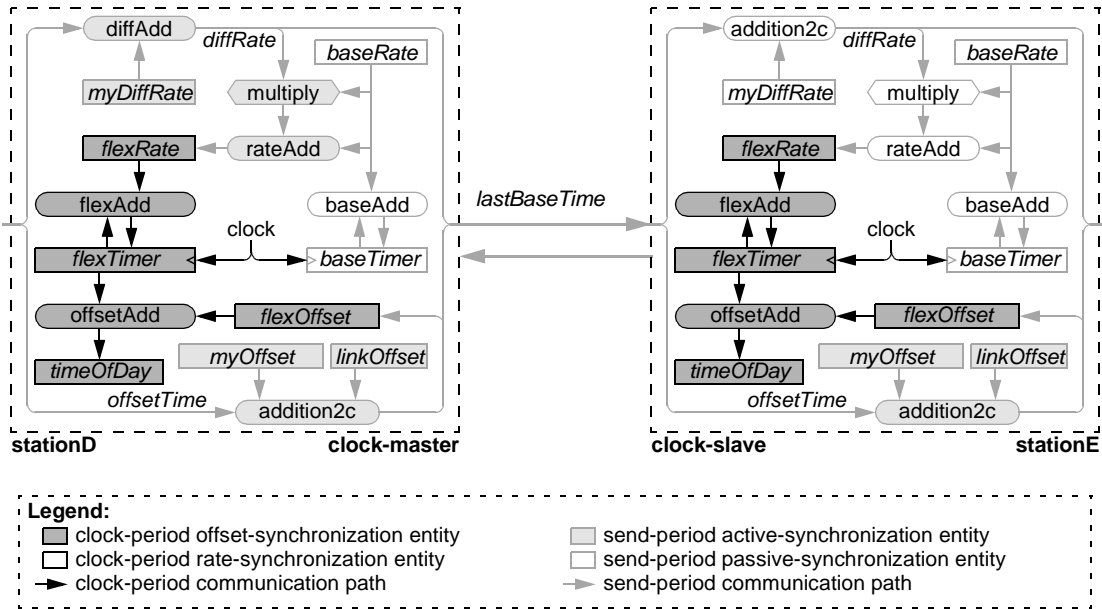


Figure 7.5—Rate synchronization

The rate-synchronization protocols interchange parameters periodically, but less frequently than the offset-synchronization protocols, possibly every 100 ms. The *lastBaseTime* value is sent periodically from the clock-master to the clock-slave. Nothing is returned from the clock-slave station.

The rate-compensation protocols for stationE adjust its *myDiffRate* value to accommodate for differences between the *stationD.baseTimer* and *stationE.baseTimer* rates. Computations are performed on clockStrobe reception and clockStrobe transmission.

The station's *diffRate* value is constructed by adding the received *clockStrobe.diffRate* and local *myDiffRate* values. This revised *clockStrobe.diffRate* value is used within each station and is passed to the clock-slave side neighboring station (if present).

7.1.9 Cascaded rate differences

The concept of cascaded rate values can be better understood by considering a simple 3-bridge example, as illustrated in Figure 7.6. Within this figure, the *myDiffRateN* and *diffRateN* represent parts-per-million (PPM) normalized values of *myDiffRate* and *diffRate* respectively.

Parameter					
name	grand-master	bridgeB	bridgeC	bridgeD	clock-slave
number	1	2	3	4	5
crystal deviation	+10 PPM	+100 PPM	-100 PPM	-75 PPM	+75 PPM
<i>myDiffRateN</i>	0 PPM	-90 PPM	200 PPM	-25 PPM	-150 PPM
<i>diffRateN</i>	0 PPM	-90 PPM	110 PPM	+85 PPM	-65 PPM
<i>flexTimer</i> deviation	10 PPM				

Representing:
 $myDiffRateN[k+1] = flexRate[k] - flexRate[k+1];$
 $diffRate[k+1] = diffRate[k] + myDiffRate[k+1];$
 $flexTimerDeviation[k] = crystalDeviation[k] + diffRate[k];$

Figure 7.6—Cascaded rate differences (a possible scenario)

The slave-agent in bridgeB is synchronized to its neighbor grand-master via clockSync frames sent on the connecting bidirectional span. Within bridgeB, the clock-slave agent passes the time directly to the clock-master agent. The slave-agent in bridgeC is synchronized to its neighbor clock-master via clockSync frames sent on the connecting bidirectional span. Other ports are similarly synchronized, thus synchronizing the right-most clock-slave station to the left-most grand-master station.

To simplify this illustration, consider only the parts-per-million (PPM) normalized rate values within each station or bridge. These values may differ significant, based (perhaps) on the nominal value or ambient temperature. Thus, the grand-master could have a crystal deviation of +10 while its bridgeB neighbor has a crystal deviation of +100.

The *myDiffRate* value within bridgeB will converges to the value of -90 PPM, representing the differences between grand-master and bridgeB crystal accuracies. The *diffRate* value received from the grand-master is added to the *myDiffRate* value, so that bridgeB's *diffRate* becomes -90 PPM. The *diffRate* and crystal deviation values are additive, yielding a resultant bridgeB *flexTimer* deviation of 10 PPM, properly synchronized to the identical grand-master's value.

Similarly, the rate of bridgeC is synchronized to bridgeB, bridgeD to bridgeC, and the clock-slave to bridgeD.

7.1.10 Rate-difference effects

If the absence of rate adjustments, significant *timeOfDay* errors can accumulate between send-period updates, as illustrated on the left side of Figure 7.7. The 2 ms deviation is due to the cumulative effect of clock drift, over the 10 ms send-period interval, assuming clock-master and clock-slave crystal deviations of -100 PPM and $+100$ PPM respectively.

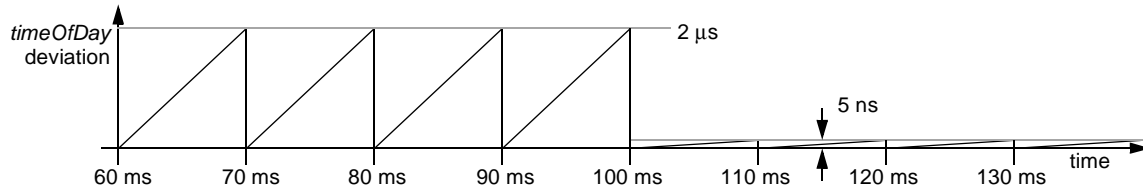


Figure 7.7—Rate-adjustment effects

While this regular sawtooth is illustrated as a highly regular (and thus perhaps easily filtered) function, irregularities could be introduced by changes in the relative ordering of clock-master and clock-slave transmissions, or transmission delays invoked by asynchronous frame transmissions. Tracking peaks/valleys or filtering such irregular functions are thought unlikely to yield similar *timeOfDay* deviation reductions.

The differences in rates could easily be reduced to less than 1 PPM, assuming a 200 ms measurement interval (based on a 100 ms slow-period interval) and a 100 ns arrival/departure sampling error. A clock-rate adjustment at time 100 ms could thus reduce the clock-drift related errors to less than 5 ns. At this point, the timer-offset measurement errors (not clock-drift induced errors) dominate the clock-synchronization error contributions.

7.1.11 flexTimer implementation example

The selection of the best time-of-day format is oftentimes complicated by the desire to equate the clock format granularity with the granularity of the implementation's 'natural' clock frequency. Unfortunately, the 'natural' frequency within a multimodal {1394, 802-100Mb/s, 802.3 1Gb/s} implementation is uncertain, and may vary based on vendors and/or implementation technologies.

The difficulties of selecting a 'natural' clock-frequency can be avoided by realizing that any clock with sufficiently fine resolution is acceptable. Flexibility involves using the most-convenient clock-tick value, but adjusting the timer advance rate associated with each clock-tick occurrence.

The same mechanism easily supports both near-arbitrary clocking rates and fine-grained rate-adjustments, needed to support timer-synchronization protocols, as illustrated in Figure 7.8. Within this figure, the shaded bytes represent values that can safely be hardwired to zero with insignificant loss of accuracy.

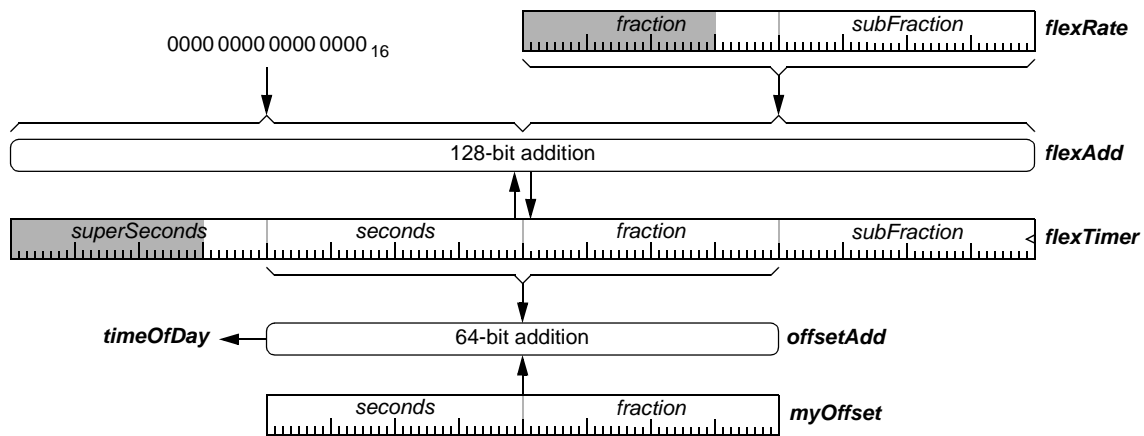


Figure 7.8—flexTimer implementation example

This illustration is not intended to constrain implementations, but to illustrate how the system's clock and timer formats can be optimized independently. This allows the *timeOfDay* timer format to be based on arithmetic convenience, timing precision, and years-before-overflow characteristics (see Annex E).

7.1.12 An alternative *baseTimer* implementation

An alternative implementation could implement the *baseTimer*-related circuitry in hardware. For such implementations, the associated firmware can be simplified, since the multiplies are eliminated from the most frequently executed loop (see Annex J).

A possible *baseTimer* hardware implementation is much simpler than the fully adjustable timer implementation, due to the absence of offset-compensation, rate-compensation, and seconds-accumulation hardware, as illustrated in Figure 7.9. Within this figure, the shaded bytes represent values that can safely be hardwired to zero with insignificant loss of accuracy.

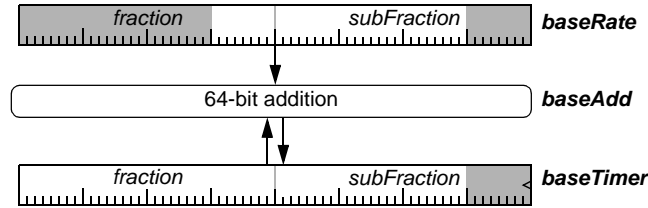


Figure 7.9—*baseTimer* implementation example

7.2 Terminology and variables

7.2.1 Common state machine definitions

The following state machine inputs are used multiple times within this clause.

CYCLES

The number of isochronous cycles within each second; defined to be 8,000.

NULL

Indicates the absence of a value and (by design) cannot be confused with a valid value.

queue values

Enumerated values used to specify shared queue structures.

Q_CRX_SYNC—The identifier associated with the received clockSync frames.

Q_CTX_SYNC—The identifier associated with the transmitted clockSync frames.

Q_ARX_REQ*—The identifier associated with the received subscription request frames.

Q_ATX_REQ*—The identifier associated with the transmitted subscription request frames.

Q_ATX_RES*—The identifier associated with the transmitted ResponseError frames.

Q_ARX_STR*—The identifier associated with the talker agent's streaming input.

Q_ATX_STR*—The identifier associated with the talker agent's streaming output.

NOTE—Those queue identifiers with an "*" are used in other clauses, but are described above. This allows all queue identification values in one location, rather than interleaving their definitions throughout this working paper.

7.2.2 Common state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

currentTime

A value representing the current time.

7.2.3 Common state machine routines

ClockSyncArrived(stationInfoPtr, portInfoPtr)

Snapshots the clockSync frame arrival time, on specified station and port (see Annex J).

ClockSyncDeparted(stationInfoPtr, portInfoPtr)

Snapshots the clockSync frame departure time, on specified station and port (see Annex J).

ClockSyncTransmit(stationInfoPtr, portInfoPtr, clockSyncPtr)

Forms a clockSync frame for transmission (see Annex J).

ClockSyncReceive(stationInfoPtr, portInfoPtr, clockSyncPtr, rateAdjust)

Processes a clockSync frame after reception (see Annex J).

Dequeue(queue)

Returns the next available frame from the specified queue.

frame—The next available frame.

NULL—No frame available.

Enqueue(queue, frame)

Places the frame at the tail of the specified queue.

Min(value1, value2)

Returns the numerically smaller of two values.

QueueEmpty(queue)

Indicates when the queue has emptied.

TRUE—The queue has emptied.

FALSE—(Otherwise.)

TimerTick(stationInfoPtr)

Updates *flexTimer* (and *baseTimer*) entities on each clock tick (see Annex J).

7.2.4 Variables and literals defined in other clauses

This clause references the following parameters, literals, and variables defined in Clause TBD:

TBDs

7.3 Clock synchronization state machines

7.3.1 ClockCore state machine

7.3.1.1 ClockCore state machine definitions

The following state machine inputs are used multiple times within this clause:

None.

7.3.1.2 ClockCore state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

clockPeriod

The duration of a synchronized timer update interval.

$$clockPeriod < 20 \text{ ns}$$

currentTime

See 7.2.2.

clockDeviation

The deviation from nominal frequency of the station-local crystal-stabilized clock.

msCount

A count that is incremented at the end of each 1 millisecond interval.

msTime

The start time of the current 1 millisecond timing interval.

nominalFrequency

The nominal frequency of the station-local crystal-stabilized clock.

tickTime

A time snapshot taken at the start of each *clockPeriod* interval.

7.3.1.3 ClockCore state machine routines*TimerTick(stationInfoPtr)*

See 7.2.3.

7.3.1.4 ClockCore state table

The ClockAgent state machine calls other C-code routines, as specified in Table 7.3. A purpose of the ClockAgent state machine is to ensure correctness of the other routines, by ensuring their indivisible executions. The notation used in the state table is described in 3.4.

Table 7.3—ClockCore state table

Current		Row	Next	
state	condition		action	state
START	—	1	$clockPeriod = 1.0 / (\text{nominalFrequency} * (1.0 + (\text{deviation} / 1000000.)))$	START
FIRST	$(currentTime - tickTime) \geq clockPeriod$	2	TimerTick(siPtr); tickTime = currentTime;	FIRST
	$(currentTime - msTime) \geq .001$	3	msTime = currentTime; msCount += 1;	
	—	4	—	

Row 7.3-1: Compute the *clockPeriod*, based on the nominal frequency and deviation.

Row 7.3-2: Update the *flexTimer* and *baseTimer* once every *clockPeriod* interval.

Row 7.3-3: Update the millisecond counter once every millisecond.

Row 7.3-4: Otherwise, no operations are performed.

7.3.2 ClockPort state machine

7.3.2.1 ClockPort state machine definitions

The following state machine inputs are used multiple times within this clause.

None.

7.3.2.2 ClockPort state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

frame

The contents of a clockSync frame.

lastInterval

A saved value of *rateInterval*, when the last rate-interval update was scheduled to occur.

rateInterval

A counter that increments on transitions of 100 ms rate-update intervals.

rateCount

A milliseconds-snapshot taken during the clockSync receive processing.

The *rateCount* value paces the relatively infrequent rate-update computations.

rxClockLast

The previously observed value of *rxClockSync*, used to detect changes in this toggling value.

rxClockSync

An indication whose value is toggled on the PHY-sensed arrival of each clockSync frame.

This value is toggled before a frame can be dequeued from the Q_CRX_SYNC queue.

rxCount

A milliseconds-snapshot taken during clockSync receive processing.

The *rxCount* value paces the detection of clockSync-silence timeouts.

sendCount

A milliseconds-snapshot taken during the clockSync transmission processing.

The *sendCount* value paces the normal clockSync frame transmissions.

selectCount

A value that tracks *siPtr*→*selectCount*, to facilitate detection of station-precedence changes.

sinkCount

A milliseconds-snapshot taken during the clockSync reception and timeout processing.

The *sinkCount* value paces the infrequent clockSync-reception timeout processing.

txClockSync

An indication whose value is toggled on the PHY-sensed departure of each clockSync frame.

This value is toggled shortly after a frame has departed from the Q_CTX_SYNC queue.

txClockLast

The previously observed value of *txClockSync*, used to detect changes in this toggling value.

7.3.2.3 ClockPort state machine routines

ClockSyncArrived(stationInfoPtr, portInfoPtr)

ClockSyncDeparted(stationInfoPtr, portInfoPtr)

ClockSyncTransmit(stationInfoPtr, portInfoPtr, clockSyncPtr)

ClockSyncReceive(stationInfoPtr, portInfoPtr, clockSyncPtr, rateAdjust)

ClockSyncTransmit(stationInfoPtr, portInfoPtr, clockSyncPtr)

See 7.2.3.

Dequeue(queue)

Enqueue(queue, frame)

See 7.2.3.

7.3.2.4 ClockPort state table

The ClockPort state machine calls other C-code routines, as specified in Table 7.4. A purpose of the ClockPort state machine is to ensure correctness of the other routines, by ensuring their timely and indivisible executions. The notation used in the state table is described in 3.4.

Table 7.4—ClockPort state table

Current		Row	Next	
state	condition		action	state
START	(frame = Dequeue(Q_CRX_SYNC)) != NULL	1	rxCount = sendCount;	NEAR
	rxClockSync != rxClockLast	2	ClockSyncArrived(siPtr, piPtr); rxClockLast = rxClockSync	START
	txClockSync != txClockLast	3	ClockSyncDeparted(siPtr, piPtr); txClockLast = txClockSync	
	selectCount != siPtr->selectCount && (msCount - sendCount) >= 1	4	selectCount = siPtr->selectCount; sendCount = siPtr->msCount; ClockSyncTransmit(siPtr, piPtr, &frame);	
	(siPtr->msCount - sendCount) >= 10	5	Enqueue(Q_CTX_SYNC, frame);	
	(siPtr->msCount - rateCount) >= 100	6	rateInterval += 1;	
	(siPtr->msCount - sinkCount) >= 50	7	ClockSyncReceive(siPtr, piPtr, NULL, 0); sinkCount = siPtr->msCount;	
	—	8	—	
NEAR	lastInterval != rateInterval	9	ClockSyncReceive(siPtr, piPtr, &frame, 1); lastInterval = rateInterval;	FINAL
	—	10	ClockSyncReceive(siPtr, piPtr, &frame, 0);	
FINAL	—	11	rxCount = sinkCount = siPtr->msCount;	START

Row 7.3-1: When a clock-sync frame arrives, mark its arrival time and process.

Row 7.3-2: Process the PHY-generated signal to determine when the clockSync frame arrived.

Row 7.3-3: Process the PHY-generated signal to determine when the clockSync frame departed.

Row 7.3-4: Transmit quickly when the grand-master selection is changing.

Row 7.3-5: Transmit routinely when the grand-master selection has stabilized.

Row 7.3-6: Trigger the rate adjustments on approximate 100 ms intervals.

Row 7.3-7: A port timeout occurs in the continued absence of clockSync frame arrivals.

Row 7.3-8: Otherwise, wait for the next event to occur.

Row 7.3-9: Restart the rate interval condition after the last rate-measurement completion.

Row 7.3-10: Otherwise, process the received clockSync frame without rate-interval measurements.

Row 7.3-11: Restart the receive-timeout counter after processing each clockSync frame.

8. Subscription state machines

NOTE—This clause should be skipped on the first reading (continue with Annex B).
The following state machines were previously highly preliminary and subject to change.
They have not yet been updated to track on recent changes to the SRP, so they are also obsolete.
Thus, the structure and formatting is useful but the details should be ignored.

Subscription state machines are responsible for performing talker-agent and listener-agent duties.

8.1 Terminology and variables

8.1.1 Common state machine definitions

The following state machine definitions are used multiple times within this clause.

NULL

Indicates the absence of a value and (by design) cannot be confused with a valid value.

subtype specifiers

ST_ERROR—A control response that provides an SRP refresh-operation error indication.

ST_FRESH—A control request that provides blocks of SRP refresh parameters.

ST_LEAVE—A control request that provides a block of SRP leave parameters.

8.1.2 Common state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

localTimer

A 64-bit timer representing the current 64-bit internal free-running time-of-day value.

myMacAddress

MAC address of the bridge.

refreshFlag

A variable that is toggled periodically; each change activates refresh interval activities.

srpState

The information associated with an element of talker-agent state. This includes:

maxBw—The maximum bandwidth of the associated stream.

maxCycles—The maximum cycles to the attached listener.

refreshTime—The time of the last observed RequestRefresh frame.

srcPortID—The port identifier of the assumed source.

srcMac—The address of the downstream bridge.

state—The connectivity state, one of the following:

IS_JOINING—Stream communications are now using this path.

IS_LEAVING—Stream communication are no longer using this path.

IS_FAILED—Stream communications have failed; message must be sent.

IS_ACTIVE—Stream communications remain active.

IS_PASSIVE—The SRP state is queued for deletion, behaving as though nonexistent.

streamTime—The time of the last observed stream flow.

streamID—The streamID of the associated stream.

subCode—The error subcode associated with the IS_FAILED state.

8.1.3 Common state machine routines*StateSearch(streamID)*

Returns the talker-state information associated with the specified stream value.

srpState—matching talker-agent state

NULL—no matching state found

8.1.4 Variables and literals defined in other clauses

This clause references the following parameters, literals, and variables defined in Clause 7

*Dequeue(queue)**Enqueue(queue, frame)**localTimer*

Q_ARX_REQ

Q_ATX_REQ

Q_ARX_STR

Q_ATX_STR

Q_ATX_RES

8.2 Subscription state machines**8.2.1 AgentAction state machine**

The AgentAction state machine controls the sequencing of AgentTalker, AgentTimer, and AgentListener state machines. There are multiple instances of these state machine, one per bridge port, each of which is invoked. A refresh flag is also complemented at a regular interval.

The following subclasses describe parameters used within the context of this state machine.

8.2.1.1 AgentAction state machine definitions

–none–

8.2.1.2 AgentAction state machine variables*localTimer**refreshFlag*

See 8.1.2.

refreshTime

The time when the last refresh was performed.

refreshTimeout

The time interval between successive refresh operations.

8.2.1.3 AgentAction state machine routines*AgentListeners()*

A routine that calls all of the AgentListener state machines (one for each bridge port).

AgentTalkers()

A routine that calls all of the AgentTalker state machines (one for each bridge port).

AgentTimers()

A routine that calls all of the AgentTimer state machines (one for each bridge port).

1 **8.2.1.4 AgentAction state table**

2
 3 The AgentAction state machine is specified in Table 8.1.

4
 5
 6 **Table 8.1—AgentAction state table**

7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18

Current		Row	Next	
state	condition		action	state
START	—	1	AgentTalkers(); AgentTimers(); AgentListeners();	LOOP
TIMER	(localTimer – refreshTime) >= refreshTimeout	2	refreshTime = localTimer; refreshFlag ^= 1;	FINAL
	—	3	—	

19
 20 **Row 8.1-1:** Execute each of the AgentTalker, AgentTimer, and AgentListener state machines.

21
 22 **Row 8.1-2:** Complement the refresh flag at the end of each refresh interval.

23
 24 **Row 8.1-3:** Otherwise, wait until the arrival of the next refresh interval.

25
 26 **8.2.2 AgentTalker state machine**

27
 28 The AgentTalker state machine monitors received RequestRefresh and RequestLeave frames. There are
 29 multiple AgentTalker state machines per bridge, one for each of the bridge ports.

30
 31 The following subclauses describe parameters used within the context of this state machine.

32
 33 **8.2.2.1 AgentTalker state machine definitions**

34 IS_FAILED

35 IS_JOINING

36 IS_LEAVING

37 See 8.1.2.

38 NULL

39 Indicates the absence of a value and (by design) cannot be confused with a valid value.

40 Q_ARX_REQ

41 Q_ARX_STR

42 Q_ATX_STR

43 See 8.1.4.

44 ST_REFRESH

45 ST_LEAVE

46 See 8.1.1.

47 *subCode* field values

48 SC_DA_LOST—No route to the specified destination is present.

49 SC_DA_MINE—The route to the specified destination loops back.

50 SC_BAD_HERE—This port’s SRP state has different parameters than the refresh request.

51 SC_BW_LIMIT—The requested stream bandwidth would exceed 75% of the link capacity.

52 SC_BAD_THERE—Another port’s SRP state has different parameters than the refresh request.

53 SC_UP_FULL—The associated listener port has insufficient space to support the refresh request.

8.2.2.2 AgentTalker state machine variables	1
<i>block</i>	2
A data structure representing the contents of a RequestRefresh info block.	3
<i>frame</i>	4
The received RequestRefresh/RequestLeave control frame (see 6.3).	5
<i>linkCapacity</i>	6
A variable representing the operational bandwidth of the link. (This can be affected by autonegotiation protocols and capabilities of the span partners.)	7
<i>localTimer</i>	8
See 8.1.4.	9
<i>matching</i>	10
A variable representing the presence of matching SRP state within another talker-agent port.	11
<i>myMacAddress</i>	12
See 8.1.2.	13
<i>oldState</i>	14
The information associated with a closely matching element of another talker-agent state.	15
<i>refreshTime</i>	16
A variable representing the arrival time of the preceding RequestRefresh message.	17
<i>srpState</i>	18
See 8.1.2.	19
<i>tstState</i>	20
The information associated with a closely matching element of this talker-agent state.	21
<i>stream</i>	22
A variable representing a stream identifier.	23
	24
	25
	26
8.2.2.3 AgentTalker state machine routines	27
	28
<i>Dequeue(queue)</i>	29
See 8.1.4.	30
<i>FullSearch(srpState, info)</i>	31
Searches through other talker agents searching for an entry with matching <i>info</i> parameters.	32
The search starts at the <i>srpState</i> -specified entry and returns each matching entry at most once.	33
The search ignores the <i>srpState</i> entries with a phase of IS_FAILED or IS_PASSIVE.	34
<i>tstState</i> —Another talker agent has the same <i>streamID</i> and matching state.	35
NONE—Another talker agent has the same <i>streamID</i> , but different state.	36
NULL—No more other-talker agents have the same <i>streamID</i> .	37
<i>InfoSelect(frame, i)</i>	38
Returns the <i>streamID</i> -specified information block within the RequestRefresh frame.	39
<i>info</i> —selected frame parameters	40
NULL—no matching parameters found	41
<i>LinkBandwidth()</i>	42
Returns the cumulative link bandwidth associated with the talker agent.	43
(This excludes bandwidths associated with entries in the IS_FAILED phase.)	44
<i>ListenerListing(srpState)</i>	45
Publishes the <i>srpState</i> information in the associated listener agent registry.	46
<i>srpState</i> —Completes successfully.	47
NULL—(Otherwise).	48
<i>SrcRoute(da)</i>	49
Returns the port identifier passed through when routed to the <i>da</i> -specified MAC.	50
positive—matching <i>portID</i> value	51
negative—no matching port found	52
<i>StateSearch(streamID)</i>	53
See 8.1.3.	54

StateForm(streamID, bandwidth)

Allocates and initializes the talker-state information associated with the argument values.

srpState—matching talker-agent state

NULL—no state-space available

8.2.2.4 AgentTalker state table

The AgentTalker state machine is responsible for establishing and demolishing paths, as specified in Table 8.2. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 8.2—AgentTalker state table

Current		Row	Next	
state	condition		action	state
START	(frame = Dequeue(Q_ARX_REQ)) != NULL	1	—	PARSE
	—	2	—	RETURN
PARSE	frame.subtype == ST_FRESH	3	info = NULL;	LOOP
	frame.subtype == ST_LEAVE	4	tstState = StateSearch((info.talkerID<<16) info.portID);	LEAVE
	—	5	—	RETURN
LOOP	(info = InfoSelect(frame, info)) != NULL	6	tstState = StateSearch((info.talkerID<<16) info.portID);	TEST
	—	7	—	RETURN
TEST	tstState == NULL	8	—	FORM
	tstState.phase == IS_FAILED	9	—	LOOP
	tstState.mcastID != block.mcastID	10	—	FORM
	tstState.maxCycles != block.maxCycles	11	—	
	tstState.maxBw != block.maxBw	12	—	
	tstState.phase == IS_LEAVING	13	tstState.phase = IS_ACTIVE	POKE
—	14	—		
POKE	—	15	tstState.refreshTime = localTimer;	LOOP
FORM	(srpState = StateForm()) != NULL	16	srpState.mcastID = info.mcastID; srpState.talkerID = info.talkerID; srpState.plugID = info.plugID; srpState.maxCycle = info.maxCycles; srpState.maxBw = info.maxBw; oldState = FullSearch(NULL, info);	CHECK
	—	17	—	LOOP

Table 8.2—AgentTalker state table

Current		Row	Next	
state	condition		action	state
CHECK	tstState != NULL	18	srpState.subCode = SC_BAD_HERE;	NACK
	port < 0	19	srpState.subCode = SC_DA_NONE;	
	port == myPortID	20	srpState.subCode = SC_DA_MINE;	
	LinkBandwidth() > 0.75 * linkCapacity	21	srpState.subCode = SC_BW_LIMIT;	
	oldState == DIFF	22	srpState.subCode = SC_BAD_THERE	
	—	23	srpState.refreshTime = localTimer; srpState.streamTime = localTimer;	PEEK
NACK	—	24	srpState.phase = IS_FAILED	LOOP
PEEK	oldState != NULL	25	srpState.phase = IS_ACTIVE;	TOSS
	ListenerListing(srpState) == NULL	26	srpState.subCode = SC_UP_FULL;	NACK
	—	27	srpState.phase = IS_JOINING;	LOOP
TOSS	oldState.phase == IS_LEAVING	28	oldState.phase = IS_PASSIVE;	LAST
	—	29	—	
LAST	(oldState = FullSearch(oldState, info)) != NULL	30	—	TOSS
	—	31	—	LOOP
LEAVE	tstState == NULL	32	—	RETURN
	tstState.phase == IS_FAILED	33	—	
	FullSearch(NULL, info) == NULL	34	tstState.phase = IS_LEAVING;	
	—	35	Release(tstState);	

Row 8.2-1: Dequeue a received subscription-request message, if available.

Row 8.2-2: Otherwise, wait for the next subscription-request message.

Row 8.2-3: Process received RequestRefresh messages.

Row 8.2-4: Process received RequestLeave messages.

Row 8.2-5: Discard unrecognized refresh messages.

Row 8.2-6: Find state associated with the selected blocks within the RequestRefresh messages.

Row 8.2-7: Stop processing after the last RequestRefresh block has been processed.

Row 8.2-8: If a matching entry cannot be found, a new one must be formed.

Row 8.2-9: The refresh is ignored while the matching entry is dedicated to error reporting.

Row 8.2-10: If the matching entry has a distinct multicast identifier, the refresh is erroneous.

Row 8.2-11: If the matching entry has a distinct *maxCycles* count, the refresh is erroneous.

Row 8.2-12: If the matching entry has a distinct maximum bandwidth, the refresh is erroneous

Row 8.2-13: If the state was leaving, it changes to active.

Row 8.2-14: Otherwise, the state (joining or active) remains unchanged.

1 **Row 8.2-15:** Update the refresh timeout when a matching entry is observed.

2
3 **Row 8.2-16:** If storage is available, update the new state based on the supplied *info* field parameters.

4 **Row 8.2-17:** If no storage is available, nothing can be done and the *info* state is discarded.

5 (A timeout is necessary to detect this discard, since no storage state is available for error reporting purposes.)

6
7 **Row 8.2-18:** With a matching/inconsistent same-port state, the appropriate error-status code is returned.

8 **Row 8.2-19:** If no upstream port can be found, the appropriate error-status code is returned.

9 **Row 8.2-20:** If the upstream port is one's self, the appropriate error-status code is returned.

10 **Row 8.2-21:** If the cumulative bandwidth limit is exceeded, the appropriate error-status code is returned.

11 **Row 8.2-22:** With a matching/inconsistent other-port state, the appropriate error-status code is returned.

12 **Row 8.2-23:** Otherwise, the timeouts are reset before the refresh is accepted.

13
14 **Row 8.2-24:** The SRP state is marked to communicate the failure condition.

15
16 **Row 8.2-25:** If matching state is found on another talker agent, this port's state is set to active.

17 **Row 8.2-26:** Otherwise, this port's state is set to joining.

18 (This triggers the near-immediate transmission of a limited refresh message, to first establish the stream.)

19
20 **Row 8.2-28:** If an existing entry is marked as leaving, its state is changed to passive to ensure removal.

21 (This talker agent is joining, so the connection remains and there is no need to announce another's leaving.)

22 **Row 8.2-29:** Otherwise, the existing entry is ignored.

23
24 **Row 8.2-30:** Check to confirm the presence an another existing entry.

25 **Row 8.2-31:** Or, terminate the search in the absence of another existing entry.

26
27 **Row 8.2-32:** If no matching to the leaving request is found, the leave request is ignored.

28 **Row 8.2-33:** If a matching error response is found, the leave request is ignored.

29 **Row 8.2-34:** If no other port has an active request, the leave request is accepted.

30 **Row 8.2-35:** If another port has an active request, this leave request can be safely ignored.

31 32 **8.2.3 AgentTimer state machine**

33
34 The AgentTimer state machine monitors received RequestRefresh and RequestLeave frames. There are
35 multiple AgentTimer state machines per bridge, one for each of the bridge ports.

36
37 The following subclauses describe parameters used within the context of this state machine.

38 39 **8.2.3.1 AgentTimer state machine definitions**

40
41 IS_ACTIVE

42 IS_FAILED

43 See 8.1.2.

44 NULL

45 Indicates the absence of a value and (by design) cannot be confused with a valid value.

46 Q_ATX_RES

47 Q_ARX_STR

48 Q_ATX_STR

49 See 8.1.4.

50 ST_ERROR

51 See 8.1.1.

52 A *subtype* specifier that distinguishes the ResponseError frame from other RE frames.

53
54

8.2.3.2 AgentTimer state machine variables	1
<i>frame</i>	2
The received streaming classA frame or generated SRP ResponseError frame (see 6.1).	3
<i>info</i>	4
A data structure representing the contents of a RequestRefresh/RequestLeave info block.	5
<i>localTimer</i>	6
See 8.1.4.	7
<i>myMacAddress</i>	8
See 8.1.2.	9
<i>refreshTime</i>	10
A variable representing the arrival time of the preceding RequestRefresh message.	11
<i>refreshTimeout</i>	12
A variable representing a timeout interval for RequestRefresh messages.	13
<i>srpState</i>	14
See 8.1.2.	15
<i>stream</i>	16
A variable representing a stream identifier.	17
8.2.3.3 AgentTimer state machine routines	18
<i>CastSearch(mcastID)</i>	19
Returns the talker-state information associated with the specified multicast identifier.	20
<i>srpState</i> —matching talker-agent state	21
NULL—no matching state found	22
<i>Dequeue(queue)</i>	23
<i>Enqueue(queue, frame)</i>	24
See 8.1.4.	25
<i>QueueHasSpace(index)</i>	26
Indicates whether space is available for frame transmissions.	27
TRUE—Space is available.	28
FALSE—(Otherwise.)	29
<i>StateSearch(streamID)</i>	30
See 8.1.3.	31
<i>StateSelect(index)</i>	32
Returns the talker-agent state associated with the specified <i>index</i> .	33
<i>info</i> —matching talker-agent state	34
NULL—no state-space available	35
<i>StateToss(index)</i>	36
Discards talker-state information associated with the argument value.	37
	38
	39
	40
	41
	42
	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

8.2.3.4 AgentTimer state table

The AgentTimer state machine is responsible for reporting timeout and upstream-communicated errors, as specified in Table 8.3. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 8.3—AgentTimer state table

Current		Row	Next	
state	condition		action	state
START	(frame = Dequeue(Q_ARX_STR)) != NULL	1	srpState = CastSearch(frame.da);	FLOW
	(frame = Dequeue(Q_ARX_RES)) != NULL	2	info = frame.info; tstState = StateSearch((info.talkerID<<16) info.portID);	SERVE
	—	3	srpState = NULL	LOOP
FLOW	srpState == NULL	4	—	START
	—	5	Enqueue(Q_ATX_STR, frame); srpState.streamTime = localTimer;	
SERVE	tstState != NULL	6	tstState.phase = IS_FAILED; tstState.subCode = frame.subCode;	START
	—	7	—	
LOOP	(srpState = StateSelect(srpState)) != NULL	8	—	TIMES
	—	9	—	RETURN
TIMES	srpState.phase == IS_FAILED	10	—	NEAR
	srpState.phase == IS_JOINING	11	—	LOOP
	srpState.phase == IS_LEAVING	12	StateToss(srpState);	
	srpState.phase == IS_PASSIVE	13		
	(localTimer – srpState.refreshTime) >= refreshTimeout	14		
	(localTimer – srpState.streamTime) >= dataTimeout	15		
—	16	—		
NEAR	QueueHasSpace(Q_ATX_RES)	17	frame.da = srpState.srcMac; frame.sa = myMacAddress; frame.subType = ST_ERROR; frame.subCode = srpState.subCode; frame.streamId = srpState.streamID; frame.maxBw = srpState.maxBw; frame.cycles = srpState.maxCycles; Enqueue(Q_ATX_RES, frame); StateToss(srpState);	LOOP
	—	18	—	

Row 8.3-1: Monitor the received stream flow, as frames pass through.	1
Row 8.3-2: Process received error messages, when they become available.	2
Row 8.3-3: Otherwise, aging timeouts are invoked.	3
	4
Row 8.3-4: Stream flows are not forwarded in the absence of matching state.	5
Row 8.3-5: Otherwise, stream flows are monitored and flow downstream.	6
	7
Row 8.3-6: In the presence of matching talker-agent state, the stream passes through.	8
Row 8.3-7: In the absence of matching talker-agent state, the stream passes through.	9
	10
Row 8.3-8: Select each talker-state element associated with the port.	11
Row 8.3-9: Stop when all talker-state elements have been processed.	12
	13
Row 8.3-10: A failed entry is processed distinctively.	14
Row 8.3-11: The joining phase indications has no timeout.	15
Row 8.3-12: The leaving phase indications has no timeout.	16
Row 8.3-13: The passive phase indication has been effectively discarded, so discard it immediately.	17
Row 8.3-14: In the absence of sustained refresh messages, the active SRP state is discarded.	18
Row 8.3-15: In the absence of sustained stream flows, the active SRP state is discarded.	19
Row 8.3-16: Otherwise, no timeout actions are required.	20
	21
Row 8.3-17: In the presence of a failed phase indication, a ResponseError is sent downstream.	22
Row 8.3-18: Otherwise, no action is taken.	23
	24
	25
	26
	27
	28
	29
	30
	31
	32
	33
	34
	35
	36
	37
	38
	39
	40
	41
	42
	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

8.2.4 AgentListener state machine

The AgentListener state machine generates RequestRefresh and RequestLeave control frames. There are multiple AgentListener state machines on each bridge, one is associated with each of the bridge ports.

The following subclauses describe parameters used within the context of this state machine.

8.2.4.1 AgentListener state machine definitions

Q_ATX_REQ

See 8.1.4.

IS_PASSIVE

See 8.1.2.

NULL

Indicates the absence of a value and (by design) cannot be confused with a valid value.

8.2.4.2 AgentListener state machine variables

frame

An SRP control frame.

localTimer

See 8.1.4.

myMacAddress

See 8.1.2.

refreshTime

A variable representing the transmission time of the preceding RequestRefresh message.

refreshTimeout

A variable representing a timeout interval for RequestRefresh messages.

refreshList

A list of *srpState* entries prepared for upstream transmission.

srpState

See 8.1.2.

8.2.4.3 AgentListener state machine routines

Enqueue(queue, frame)

See 8.1.4.

EnqueueList(queue, list)

Transfers content from the *rpState* lists into one or more frames.

Each of these frames is then placed into the specified queue.

JoiningList()

Forms a list of the joining-phase entries from the listener agent's state array.

JoiningToActive(list)

Within all listed entries, each phase value of *IS_JOINING* is changed to *IS_ACTIVE*.

QueueHasSpace(index)

Indicates whether space is available for frame transmissions.

TRUE—Space is available.

FALSE—(Otherwise.)

RefreshList()

Forms a list of the joining-phase and active-phase entries from the listener agent's state array.

ReviseListenerList()

Revises the listener list entries to ensure consistency with distributed AgentTalker state content.

8.2.4.4 AgentListener state table

The AgentListener state machine is responsible for generating upstream RequestRefresh and RequestLeave frames, as specified in Table 8.4. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 8.4—AgentListener state table

Current		Row	Next	
state	condition		action	state
START	—	1	ReviseListenerList();	FIRST
FIRST	QueueHasSpace(Q_ARX_REQ)	2	—	TIMER
	—	3	—	RETURN
CHECK	localTimer >= (refreshTime + refreshTimeout) && ((refreshList= RefreshList()) != NULL)	4	refreshTime = localTimer;	FRESH
	srpState = QueueHasLeave()	5	frame.da = upstreamAddress; frame.sa = myMacAddress; frame.info = srpState.info; EnqueueFrame(Q_ATX_REQ, frame); srpState.phase = IS_PASSIVE;	START
	(refreshList = JoiningList()) != NULL	6	—	FRESH
	—	7	—	RETURN
FRESH	—	8	EnqueueList(Q_ATX_REQ, refreshList); JoinToActive(refreshList);	START

Row 8.4-1: Refresh the listener list, ensuring consistency with distributed AgentTalker state content.

Row 8.4-2: In the presence of transmission-queue storage, transmissions are enabled.

Row 8.4-3: Otherwise, transmissions are inhibited.

Row 8.4-4: When periodically enabled, the list of joining and active states is sent.

Row 8.4-5: Leave requests are checked; distinct ones cause a RequestListen frame to be sent.

Row 8.4-6: When entries are found, the list of joining states is sent.

Row 8.4-7: Otherwise, no talker-agent refresh/leave messages are transmitted.

Row 8.4-8: Enqueue the refresh-list entries for eventual transmission.

Afterwards, change the phase from joining to active, to inhibit unnecessary future transmissions.

9. Transmit state machines (proposal 1)

NOTE—Multiple bunch-avoiding pacing protocols are presented for consideration:
 a) Clause 9 (this clause) presents a pseudo-synchronous transmission model.
 b) Clause 10 presents a cross-flow shaper transmission model.

9.1 Pacing overview

9.1.1 Delays

The preferred topologies consists entirely of paced bridges, as illustrated in Figure 9.1a. Within such topologies, a frame transmitted by station a0 in cycle[n] incurs fixed nominal delays while passing through bridges. Thus, this frame nominally departs bridgeB in cycle[n+2], bridge C in cycle[n+4], and bridgeE in cycle[n+6].

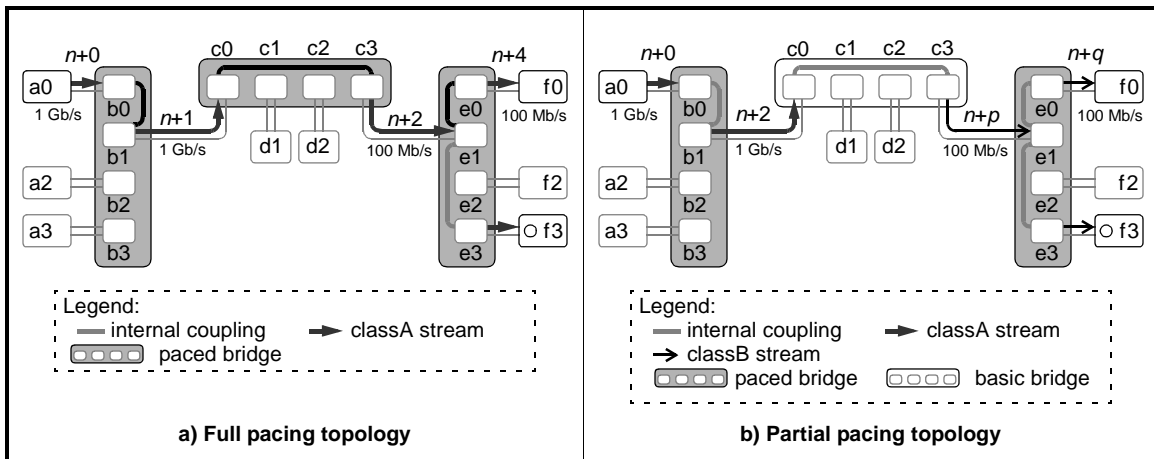


Figure 9.1—Topology-dependent pacing delays

Within Figure 9.1a, the actual transmission times can vary from their nominal targets, due to contention with other traffic. Each bridge compensates for early and late arrivals, so that the extent of deviations from nominal on link b1-to-c0 are the same as those on link e0-to-f0.

Within Figure 9.1b, an intermediate basic bridge is assumed. Output from bridgeC is therefore downgraded from classA to classB, to avoid degradation of well-paced traffic. Thus, the fully-paced properties of bridgeE still apply to possible f3-to-f0 traffic (not illustrated).

The uncertainty of cycle p and q cycle delays in Figure 9.1b are due to passing through the non-paced bridgeC. Although much of this traffic would arrive earlier, some of the traffic could be delayed up to the nominal delays of Figure 9.1a. In more complex topologies, such delays could exceed the nominal delays through paced bridges, due to bunching effects (see Annex F).

To support such topology, this working paper mandates that compliant end stations provide larger elasticity buffers (see TBD) than required within fully paced topologies. However, defining topology restrictions to ensure elasticity-buffer sufficiency is beyond the scope of this working paper.

9.1.2 Paced 1 Gb/s classA flows

Pacing involves sending accumulated classA traffic once every isochronous cycle, rather than allowing larger (typically an MTU) frames to be accumulated. After each cycle's classA traffic has been sent, the remaining time is available for sending classB/classC traffic. This provides low-jitter bandwidth guarantees, as does time division multiplexing (TDM), while allowing unused classA bandwidths to be utilized by classB/classC traffic.

A pacing bridge maintains this pacing behavior, thus avoiding problems normally associated with bunching (see Annex F). For a bridge between 1 Gb/s link1 and 1 Gb/s link2 (see Figure 9.2a), paced frames can be forwarded with a nominal 1-cycle delay (see Figure 9.2b). The 1-cycle delay is necessary to account for offset migration and store-and-forward processing delays.

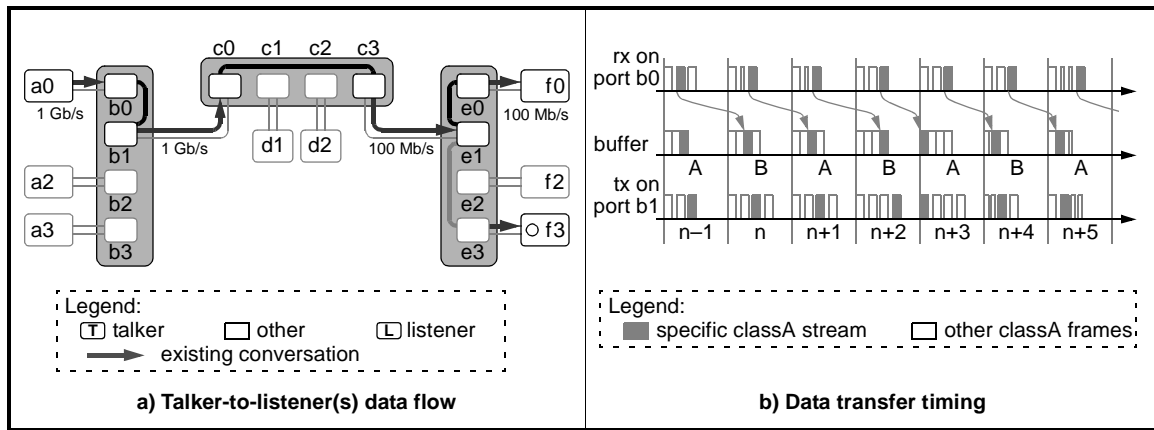


Figure 9.2—Paced 1 Gb/s classA flows

Offset migration refers to changes in a classA frame's within-cycle placement on (for example) link1 and link2. Depending on the timing of unrelated events, the offset of the classA-data frame within the cycle can migrate over time, as other conversations are started, ended, advanced, delayed, joined, or routed elsewhere.

A possible implementation could utilize double output buffers, processed as follows:

cycle $[n+2 \times k+0]$: classA traffic is saved in buffer[A] and transmitted from buffer[B].

cycle $[n+2 \times k+1]$: classA traffic is saved in buffer[B] and transmitted from buffer[A].

The boundaries between cycles are marked by a distinct set of cycleSync markers (not illustrated), rather than relying on precise time-synchronization and deadbands to imply their temporal placement.

The classA transmissions within each cycle are shaped, to allow for unrelated asynchronous frame transmissions. The shaper allows a higher-than 75% transmission rate, to ensure transmission completion well before before the next cycle begins, even in the presence of conflicting non-classA transmissions.

To better understand the minimal buffer requirements, consider frame transfers that are momentarily disrupted by an MTU-sized classC transmission, started near the end of link1's classA transmissions. For the receive-side slippage scenario of 9.3a, data[n] arrives in cycle[n] and fills buffer[A]. Since buffer A is not destined for transmission until cycle[n+1], conflicts are avoided.

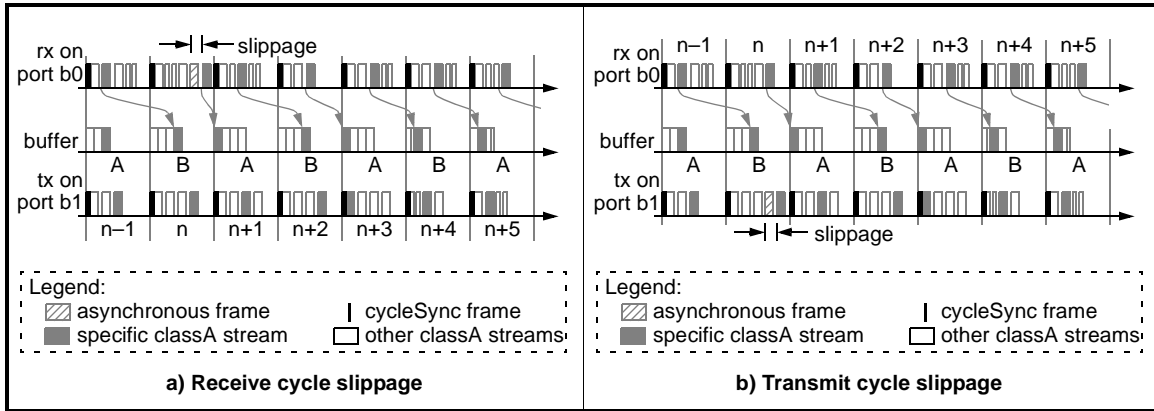


Figure 9.3—Cycle slippage

For the transmit-side slippage scenario of Figure 9.3b, buffer[B] is fully emptied in cycle[n]. Since buffer[B] is not destined for filling until cycle n+1, conflicts are avoided.

9.1.3 Paced 100 Mb/s flows

Editors' Notes: To be removed prior to final publication.
 A two-cycle delay is illustrated, although the protocols can be simplified by assuming a three cycle delay. The tradeoff between protocol simplicity and a passthrough latency has not been carefully reviewed.

A 100 Mb/s pacing bridge also maintains this pacing behavior, thus avoiding problems normally associated with bunching (see Annex F). For a bridge between 100 Mb/s link3 and 100 Mb/s link4 (see Figure 9.4a), paced frames can be forwarded with a nominal 2-cycle delay (see Figure 9.4b).

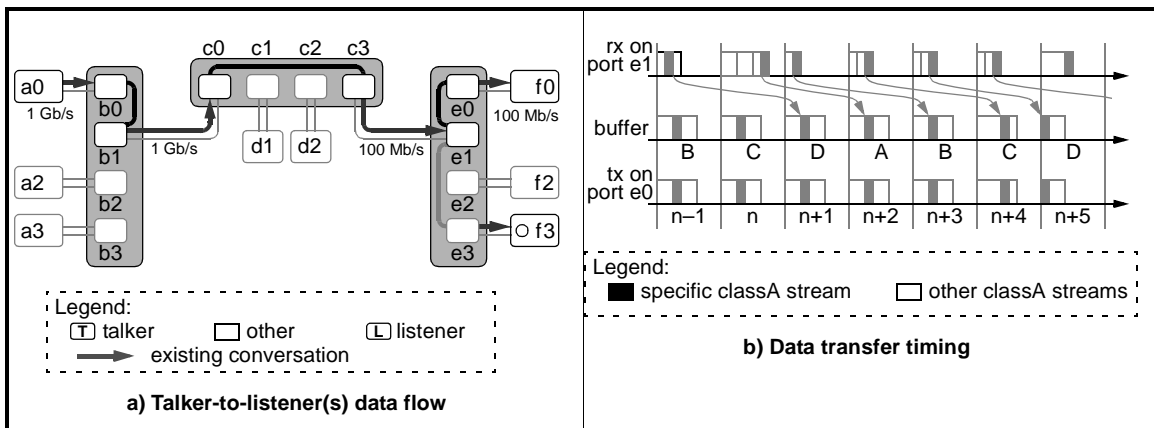


Figure 9.4—Paced 100 Mb/s classA flows

A possible implementation would involved six output buffers, processed as follows:

- cycle $[n+4 \times k+0]$: classA traffic is saved in buffer[A] and transmitted from buffer[C].
- cycle $[n+4 \times k+1]$: classA traffic is saved in buffer[B] and transmitted from buffer[D].
- cycle $[n+4 \times k+2]$: classA traffic is saved in buffer[C] and transmitted from buffer[A].
- cycle $[n+4 \times k+3]$: classA traffic is saved in buffer[D] and transmitted from buffer[B].

To better understand the minimal buffer requirements, consider frame transfers that are momentarily disrupted by an MTU-sized classC transmission, started near the end of link3 classA transmissions. For the receive-side slippage scenario of Figure 9.5a, data $[n]$ arrives in cycle $[n+1]$ and fills buffer[A]. Since buffer A is not destined for transmission until cycle $[n+2]$, conflicts are avoided.

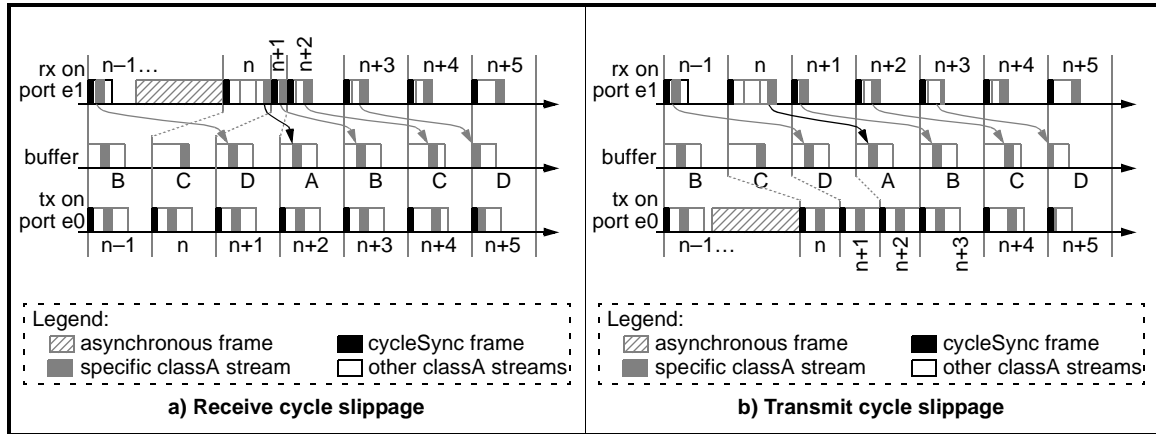


Figure 9.5—Cycle slippage

For the transmit-side slippage scenario of Figure 9.5b, buffer[D] is fully emptied in cycle $[n+1]$ and in cycle $[n+2]$. Since buffer[D] is not destined for filling until cycle $[n+3]$, conflicts are avoided.

To achieve a robust 2-cycle latency objective, restrictions are placed on non-classA transmissions. These restrictions are as follows:

- a) An MTU (or sequence of frames not exceeding an MTU) may be appended to the last classA frame within any cycle whose cycleSync frame transmission was not delayed.
- b) Within any cycle, any non-classA frame may be transmitted after the last classA frame, but only if this frame transmission would not delay the transmission of the next cycleSync frame.

Condition (a) is sufficient to ensure that all transmissions occur within the intended or following cycle, assuming a 100 Mb/s span, 2000 byte MTU, 125 μ s cycle, and 75% classA loading. With these assumptions, the worst-case delay from the start of the intended cycle, as specified by Equation 9.1, is well within the 2-cycle 250 μ s constraint.

$$delay \geq (MTU - 0.25 \times cycle) + 0.75 \times cycle \quad (9.1)$$

$$delay \geq 2000 \times ((8 \text{ bits/byte}) \times (1 \text{ second}) / (100 \text{ Mb/s})) + 0.50 \times (125 \mu\text{s})$$

$$delay \geq (160 \mu\text{s}) + (62.5 \mu\text{s})$$

$$delay \geq 222.5 \mu\text{s}$$

9.1.4 Transmit port structure

An end station and bridge have functionally distinct transmit queues for classA, classB, and classC traffic, allowing each to be managed separately, as illustrated in Figure 9.6. The transmit port is responsible for pacing classA/classB traffic and shaping classB/classC traffic, so as to limit the high-class traffic to 75% of the link bandwidth. The transmit-port structure is slightly different for 100 Mb/s and 1 Gb/s transmit ports, due to the distinct times associated with an MTU transmission.

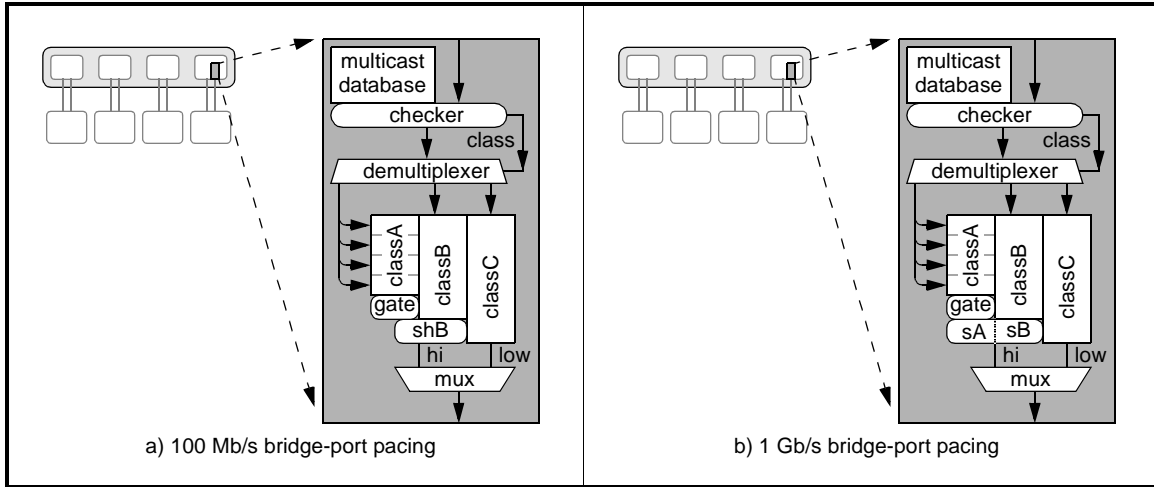


Figure 9.6—Transmit-port structure

Although classA frames have the highest priority, the classA frames are gated to prevent their early departure. Gating involves blocking classA frames that arrived with $sourceCycle=n$, until the start of cycle $n+p$. After the start of cycle $n+p$, the transmitter waits for the completion of preceding non-classA frames (or residual cycle $n+p-1$ classA frames), then transmits these arrived-in-cycle- n frames with $sourceCycle=n+p$. As noted previously, p is a design-dependent integer constant, preferably no more than 4 cycles (see 5.1.2 and 5.1.3).

A bridge has to cope with frame-reception uncertainties (due to preceding frame-transmission uncertainties), in addition to its own frame-transmission uncertainties. As such, the values of p are expected to be slightly larger in bridges than in talker-station or listener-station designs.

Within bridges, the distinction between service classes is based on the multicast addresses within frames. These multicast addresses are checked against the multicast database, which supplies $class$ information in addition to the normal multicast routing (forward or not-forward) information. This $class$ information controls the demultiplexer, which routes to the appropriate classA, classB, or classC output queues.

The cycle slippage on a 100 Mb/s link mandates the use of four 3/4-cycle output buffers, which incur a 2-cycle pass-through delay. The classA traffic is gated to avoid wrong-cycle transmissions and excessive consumption, but is not otherwise not shaped. The overlapping shB shaper of Figure 9.6a is intended to illustrate the use of classA transmission counts and the classB shaper, not the shaping of classA traffic.

On such 1 Gb/s transmitter ports, the classA traffic is shaped to reduce lower-class blockage, as well as gated to avoid wrong-cycle transmissions and excessive consumption. The adjacency of shA/shB shapers in Figure 9.6b is intended to illustrate distinct classA/classB shaping functions, but sharing of classA transmission counts between shapers.

Achievable delays through a bridge depend only on the speed of the input-link speed, as summarized in Table 9.1. These numbers are slightly misleading, since transmissions on a 100 Mb/s link have implied additional delays incurred when passing through its adjacent 100 Mb/s receiver.

Table 9.1—ClockPort state table

Link type		Delay	
Input	Output	Cycles	Time
100 Mb/s	—	2	250 μs
1 Gb/s	—	1	125 μs

9.1.5 Pacing at 1 Gb/s

Pacing at 1 Gb/s, as illustrated in Figure 9.7. For ontime cycles, a residual amount of classB/classC traffic is allowed throughout the cycle, as illustrated in Figure 9.7a. For slipped cycles, a residual amount of classB/classC traffic becomes available after the delay effects have been overcome, as illustrated in Figure 9.7b.

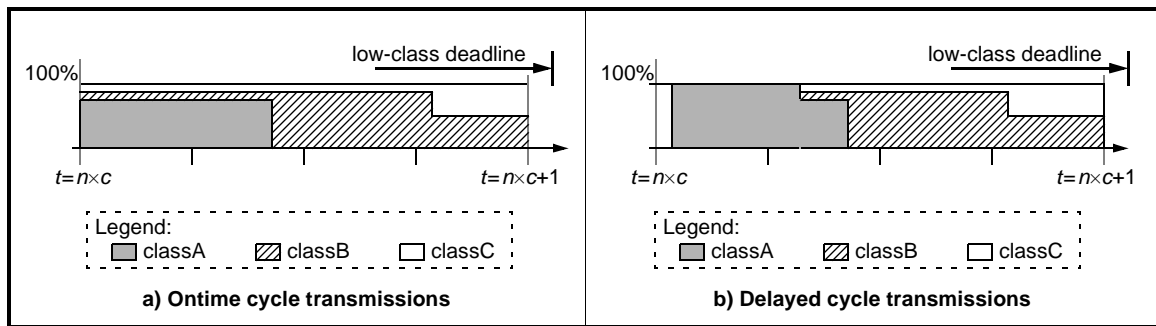


Figure 9.7—Pacing at 1 Gb/s

9.1.6 Pacing at 100 Mb/s

Pacing at 100 Mb/s, as illustrated in Figure 9.8. For ontime cycles, a residual amount of classB/classC traffic is allowed throughout the cycle, as illustrated in Figure 9.8a. For delayed cycles, a residual amount of classB/classC traffic becomes available after the delay effects have been overcome, as illustrated in Figure 9.8b.

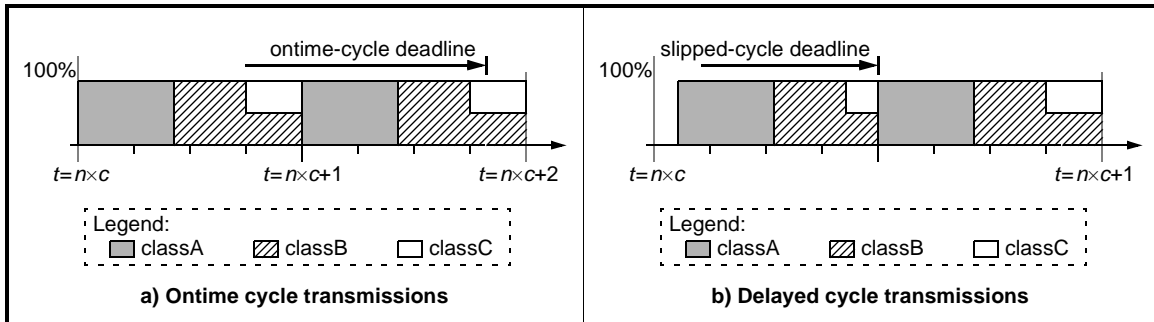


Figure 9.8—Pacing at 100 Mb/s

9.1.7 Shaper behavior

Although multiple shaper are specified within this working paper, the behavior of most shapers can be characterized by a common algorithm and instance-specific parameters (as done within RPR[B5]). The shapers' credits are adjusted down or up, as illustrated in Figure 9.9. The decrement and increment values typically represent sizes of a transmitted frame and of credit increments in each update interval, respectively.

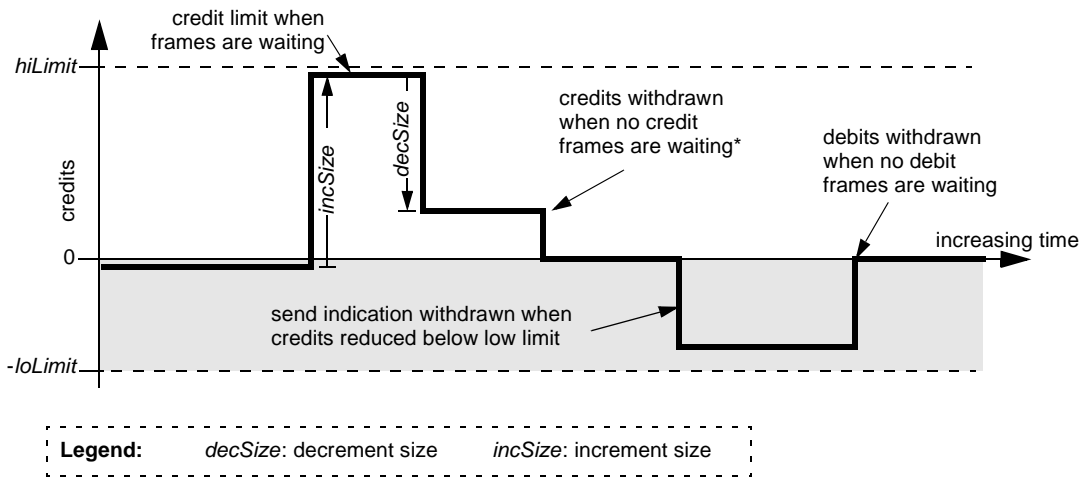


Figure 9.9—Credit adjustments over time

Crossing below the zero threshold generates a rate-limiting indication (the removal of a send indication), so that offered traffic can stop. By design, the credit value never goes below the $-loLimit$ extreme. To bound the burst traffic after inactivity intervals, when no frames are ready for transmission, credits are reduced to zero (if currently higher than zero) and can accumulate to no more than the zero-value limit.

The $hiLimit$ threshold limits the positive credits, to avoid overflow. When frames are ready for transmission (and are being blocked by transit traffic), credits can accumulate to no more than this $hiLimit$ value.

In concept, the shapers consist of a token bucket. The credits in the token bucket are incremented by the size of each debit-frame when it is being transmitted. The number of credits in a token bucket is decremented by the size of each credit-frame when it is being transmitted. When a credit-frame is waiting, it is transmitted only if the number of credits in the token bucket is positive; When a debit-frame is waiting, it is transmitted only if the number of credits in the token bucket is negative.

9.2 Terminology and variables

9.2.1 Common state machine definitions

The following state machine inputs are used multiple times within this clause.

queue values

Enumerated values used to specify shared queue structures.

QP_TX_PUSH—The input port's receive-from-ports queue.

QP_TX_CA—The first of the output port's classA buffers.

QP_TX_CB—The output port's classB queue.

QP_TX_CC—The output port's classC queue.

QP_TX_LINK—The output port's transmit-PHY queue.

QP_TX_SYNC—The port's queue that provides clockSync frames.

9.2.2 Common state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

currentTime

A value representing the current time.

mtuSize

The size of the maximum transfer unit (MTU).

Value: 2000 bytes

NOTE—The specified *mtuSize* is larger than currently supported by IEEE Std 802.3, but consistent with expected near-term frame-extension revisions of this standard.

speedIs100Mbps

A value that communicates the operating speed of the link.

TRUE—The port is operating at a speed of 100 Mb/s.

FALSE—The port is operating at speeds of 1 Gb/s or above.

thisCycle

A cycle counter derived from *thisTime*, as defined by Equation 9.2.

$$\text{Floor}(\text{thisTime} * 8000); \quad (9.2)$$

thisTime

A normalized time-of-day counter derived from *timeOfDay*, as defined by Equation 9.3.

$$(\text{timeOfDay} / (4.0 * (1 << 30))) \quad (9.3)$$

9.2.3 Common state machine routines

–none–

9.2.4 Routines defined in other clauses

This clause references the following routines defined in Clause 7:

Dequeue(queue)

Enqueue(queue, frame)

Min(value1, value2)

See 7.2.3.

9.3 Pacing state machines

9.3.1 ReceiveRx state machine

The ReceiveRx state machine is responsible for receiving pacing classA traffic, shaped classB traffic, and best-effort classC traffic. An intent is to transfer each to the appropriate output queue.

The following subclasses describe parameters used within the context of this state machine.

9.3.1.1 ReceiveRx state machine definitions

CYCLE_SYNC

An assigned *subType* value that distinguishes a clockSync from other Residential Ethernet frames.

GROUP_BIT

A constant value derived from IEEE Std 802-2001 and specified by Equation 9.4.

$((\text{macAddress} \ \& \ \text{GROUP_BIT}) \neq 0)$ (9.4)

queue values

Enumerated values used to specify shared queue structures.

QP_TX_CA, QP_TX_CB, QP_TX_CC

QP_TX_PUSH

See 9.2.2.

RES_ETHER

The *protocolType* code value assigned to Residential Ethernet.

9.3.1.2 ReceiveRx state machine variables

class

A value that represents the results of a forwarding database search.

delta

A value that represents the difference between frame-signaled and computed cycle values.

frame

The contents of a received frame.

myCycle

The two least-significant bits of the *thisCycle* value.

queueA

The selected classA queue identifier, based on *delta*-selected locations.

speedIs100Mbs

thisCycle

thisTime

See 9.2.2.

9.3.1.3 ReceiveRx state machine routines*DataBaseClass(macAddress, port)*Provides a forwarding database indication of how the *macAddress* is routed to the specified *port*.

CLASS_A—The associated multicast frame is forwarded as classA traffic.

CLASS_B—The associated multicast frame is forwarded as classB traffic.

CLASS_C—The associated multicast frame is forwarded as classC traffic.

BLOCKED—The associated multicast frame is not forwarded.

Dequeue(queue)

See 9.2.4.

*EnqueuePort(port, queue, frame)*Places the *frame* at the tail of the specified *queue* within the specified *port*.*ForwardUnicast(frame)*

Forwards a unicast frame to the selected output port, if any.

This routine mimics existing standards, which remain unaffected by this working paper.

Multicast(macAddress)

Indicates whether the supplied address is a multicast (or broadcast) address, as specified by Equation 9.5.

TRUE—The address is a multicast (or broadcast) address.

FALSE—(Otherwise.)

$$((\text{macAddress} \ \& \ \text{GROUP_BIT}) \ != \ 0) \quad (9.5)$$

9.3.1.4 ReceiveRx state table

The ReceiveRx state machine is specified in Table 9.2. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 9.2—ReceiveRx state table

Current		Row	Next	
state	condition		action	state
START	(frame = Dequeue(QP_TX_PUSH)) != NULL	1	—	FIRST
	—	2	myCycle = (thisCycle % 4); delta = (4 + myCycle - rxCycle) % 4;	PLACE
PLACE	delta == 3	3	delta = 0;	PLUS
	—	4	—	
PLUS	speedIs100Mbs	5	queueA = QP_TX_CA + (4 + 2 - delta) % 4;	START
	—	6	queueA = QP_TX_CA + (4 + 1 - delta) % 4;	
FIRST	frame.protocolType==RES_ETHER && frame.subType == CYCLE_SYNC	7	rxCycle = (frame.cycleCount % 4);	START
	Multicast(frame.da)	8	class = DataBaseClass(frame.da, port);	CAST
	—	9	ForwardUnicast(frame)	START
PUSH	class == CLASS_A	10	EnqueuePort(port, queueA, frame);	START
	class == CLASS_B	11	EnqueuePort(port, QP_TX_CB, frame);	
	class == CLASS_C	12	EnqueuePort(port, QP_TX_CC, frame);	
	—	13	—	

Row 9.2-1: If a frame has arrived, process that frame.

Row 9.2-2: Otherwise, compute the cycle offset for later classA queue placement.

Row 9.2-3: Frames that arrive early are processed as though they arrived within this cycle.

Row 9.2-4: Otherwise, the difference between labeled and actual cycles determines the frame's placement.

Row 9.2-5: Frames arriving from a 100 Mb/s link are placed 2-cycles ahead, to allow for cycle slips.

Row 9.2-6: Frames arriving from a 1 Gb/s link are placed 1-cycle ahead, since cycle slips are avoided.

Row 9.2-7: The cycleSync frames identify the cycle number, despite cycle-slip possibilities.

Row 9.2-8: Multicast frames are sent to all enabled ports.

Row 9.2-9: Unicast frames are processed normally.

Row 9.2-10: Multicast classA frames are forwarded to the appropriate cycle-sensitive classA queue.

Row 9.2-11: Multicast classB frames are forwarded to the classB queue.

Row 9.2-12: Multicast classC frames are forwarded to the classC queue.

Row 9.2-13: If no class is specified, multicast frames are not routed through this port.

9.3.2 TransmitTx state machine

The TransmitTx state machine is responsible for pacing/shaping classA traffic and shaping classB traffic destined for 1 Gb/s links. An intent is to support projected MTU-sized transfers and interleaved lower-class traffic, without exceeding the 1-cycle delay inherent with cycle-synchronous bridge-forwarding protocols.

The following subclauses describe parameters used within the context of this state machine.

9.3.2.1 TransmitTx state machine definitions

BPS

Represents a bound on the number of transmitted bytes per second, as defined by Equation 9.6.

$$(\text{speedIs100Mbps} ? 12500000 : 125000000) \quad (9.6)$$

CAP

Represents a bound on the number of transmitted bytes, as defined by Equation 9.7.

$$\begin{aligned} &((\text{speedIs100Mbps} \ \&\& \ \text{phase} \neq \text{MORE}) ? \\ &((\text{cycle} + 1) * 8000. - \text{thisTime}) * \text{BPS} : \text{MTU}) \end{aligned} \quad (9.7)$$

queue values

Enumerated values used to specify shared queue structures.

QP_TX_CA

QP_TX_CB

QP_TX_CC

QP_TX_LINK

QP_TX_SYNC

See 9.2.2.

9.3.2.2 TransmitTx state machine variables

creditA

A shaper credit whose positive value enables classA/classB primary transmissions.

creditB

A shaper credit whose positive/negative values enable secondary classB/classC transmissions.

cycle

The cycle whose classA data is being transmitted.

cycleSize

The number of bytes included within a 125 μ s cycle.

Value:

1562.5—for 100 Mb/s links

15625—for 1 Gb/s links

frame

The contents of a to-be-transmitted frame.

hiLimitB

A value that limits the cumulative *creditB* credits.

Value: MTU.

limit

A value that limits the amount of transmitted primary classA/classB bandwidth.

loLimitB

A value that limits the cumulative *creditB* debits.

Value: MTU.

phase

An indication of what remains to be transferred within the cycle.

HEAD—The cycleSync frame are to be sent.

MORE—Other classA/classB frames are to be sent.

DONE—All classA frames have been sent.

<i>queue</i>	1
A variable that identifies the appropriate classA queue for this cycle's transmissions.	2
<i>speedIs100Mbs</i>	3
See 9.2.2.	4
<i>thisCycle</i>	5
<i>thisTime</i>	6
See 9.2.2.	7

9.3.2.3 TransmitTx state machine routines

<i>Cap(speedIs100Mbs, phase, creditA, cycle, thisTime)</i>	10
Provides a cap on the lengths of classB and classC transmissions.	11
<pre> if (speedIs100Mbs) { if (phase == MORE) return(0); near = (cycle + 1.0) * 8000; safe = (cycle + 0.8) * 8000; return(thisTime <= safe ? MTU : near * BPS); } else { if (phase == MORE) return(-creditsA/16); near = (cycle + 1.05) * 8000; return((near - thisTime) * BPS); } </pre>	(9.8) 12
<i>Dequeue(queue)</i>	13
See 9.2.4.	14
<i>DequeueSize(queue, size)</i>	15
Returns the next available frame from the specified queue, from frames no larger than <i>size</i> .	16
<i>Enqueue(queue, frame)</i>	17
See 9.2.4.	18
<i>QueueEmpty(queue)</i>	19
Returns the an indication of whether the queue is empty.	20
0—The specified queue is not empty.	21
1—The specified queue is empty.	22
<i>Size(frame)</i>	23
Returns the size of the specified frame.	24

9.3.2.4 TransmitTx state table

The TransmitTx state machine is specified in Table 9.3. The link-speed independent rows are white; the link-speed dependent rows are shaded. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 9.3—TransmitTx state table

Current		Row	Next	
state	condition		action	state
START	cycle > (thisCycle + 1)	1	cycle = thisCycle; phase = HEAD;	PREP
	cycle < (thisCycle - 1)	2		
	cycle < thisCycle && phase == DONE	3	cycle += 1; phase = HEAD;	
	!QueueEmpty(QP_TX_LINK)	4	—	START
	phase == HEAD	5	queue = QP_TX_CA + (cycle % 4); frame = Dequeue(QP_TX_SYNC); limit = 0.75 * cycleSize; creditA = 16 * BPS * (thisTime - cycle*8000.); phase = MORE;	POST
	—	6	cap = Cap(speedIs100Mbs, phase, creditA, cycle, thisTime);	PLUS
PLUS	creditB >= 0 && (frame = DequeueSize(QP_TX_CB, cap)) != NULL	7	creditB = Max(loLimitB, creditB - Size(frame)); creditA += 16 * Size(frame);	FINAL
	creditB <= 0 && (frame = DequeueSize(QP_TX_CC, cap)) != NULL	8	creditB = Min(hiLimitB, creditB + Size(frame)); creditA += 16 * Size(frame);	
	(frame = DequeueSize(QP_TX_CB, cap)) != NULL	9	creditA += 16 * Size(frame);	
	(frame = DequeueSize(QP_TX_CC, cap)) != NULL	10		
	phase != MORE	11	—	START
	(frame = DequeueSize(queue, limit)) != NULL	12	—	POST
	(frame = Dequeue(queue)) != NULL	13	—	START
	(frame = DequeueSize(QP_TX_CB, limit)) != NULL	14	limit -= Size(frame);	FINAL
—	15	phase = DONE; creditB = Min(hiLimitB, limit+creditB);	START	
POST	—	16	limit -= Size(frame); creditA -= Size(frame);	FINAL
FINAL	—	17	Enqueue(QP_TX_LINK, frame);	START

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Row 9.3-1: If <i>cycle</i> has advanced two-beyond <i>thisCycle</i> , something is in error.	1
(The <i>cycle</i> value can advance one-beyond <i>thisCycle</i> , due to small <i>timeOfDay</i> update discontinuities.)	2
Row 9.3-2: If <i>cycle</i> has dropped two-behind <i>thisCycle</i> , something is in error.	3
(Large <i>timeOfDay</i> update discontinuities can cause cycle to advance or retreat beyond normal bounds.)	4
	5
Row 9.3-3: The phase is initialized to HEAD at the start of each cycle.	6
Row 9.3-4: Wait for the queue to be emptied, so something can be transmitted.	7
Row 9.3-5: When the next cycle starts, a clockSync frame is transmitted.	8
The <i>limit</i> value is set to limit classA transmissions to no more than 75% of the link bandwidth.	9
The <i>creditA</i> value initialized to account for cycleSync frame slippage (for 1 Gb/s ports only).	10
Row 9.3-6: Set caps on the maximum transmission size of classB/classC transmissions.	11
	12
Row 9.3-7: If enabled and available, a classB frame is transmitted.	13
The <i>creditB</i> values is decremented by the transmitted frame size, to effect a classB shaper.	14
Row 9.3-8: If enabled and available, a classC frame is transmitted.	15
The <i>creditB</i> values is incremented by the transmitted frame size, to effect a classB shaper.	16
Row 9.3-9: If available, a classB frame is transmitted.	17
Row 9.3-10: If available, a classC frame is transmitted.	18
Row 9.3-11: Otherwise, no frame is transmitted.	19
	20
Row 9.3-12: An enabled, available, and properly sized classA frame is readied for transmission.	21
Row 9.3-13: An enabled, available, and improperly sized classA frame is discarded.	22
Row 9.3-14: An enabled, available, and properly sized classB frame is readied for transmission.	23
Row 9.3-15: If enabled but unavailable, this cycle's primary frame transmissions have completed.	24
	25
Row 9.3-16: The shaper's <i>creditA</i> value is decremented to lightly throttle primary transmissions.	26
The <i>limit</i> value is also decremented, to enforce the 75% cycle classA/classB transmission limitation.	27
	28
Row 9.3-17: Transmission is affected by placing the frame in the port's transmit queue.	29
	30
	31
	32
	33
	34
	35
	36
	37
	38
	39
	40
	41
	42
	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

10. Transmit state machines (proposal 2)

NOTE—Multiple bunch-avoiding pacing protocols are presented for consideration:
 a) Clause 9 presents a pseudo-synchronous transmission model.
 b) Clause 10 (this clause) presents cross-flow shaper transmission model.

10.1 Rate-based scheduling overview

The clause describes a rate-based scheduling technique. The rate-based scheduling concepts are similar to those within rate monotonic scheduling protocols, commonly used within real-time systems. Objectives associated with time-sensitive forwarding alternatives include the following:

- a) Multiple time-sensitive transmission rates are supported, including:
 - 1) High rate 8 kHz traffic, such as the traffic generated by simple bridges between RE and existing IEEE 1394[B6] A/V devices.
 - 2) Higher -rate 64 kHz traffic, such as the traffic generated by highly interactive latency-sensitive video game video and/or sensors.
 - 3) Lower rate traffic, such as voice over internet protocol (VOIP) traffic, without forcing this traffic to be reblocked into smaller (and therefore less efficient) frame sizes.
- b) Frame forwarding should not be dependent on successful time-of-day synchronization between the bridge and adjacent stations. Frame forwarding should succeed before the grand clock-master station has been selected, or when the selected grand-master clock station changes.
- c) Frame-forwarding protocols should leverage existing bridge queue and service models, although specification of abstract rate shaper details is expected.

Rate-based scheduling involves associating a priority with frame transmissions, where the priority is a monotonic function of the frame transmission frequency, as illustrated in Figure 10.1. Assuming the cumulative traffic is limited to less than the link capacity, the latency of each traffic class is guaranteed (the latency guarantee is approximately an MTU more than an inter-arrival period).

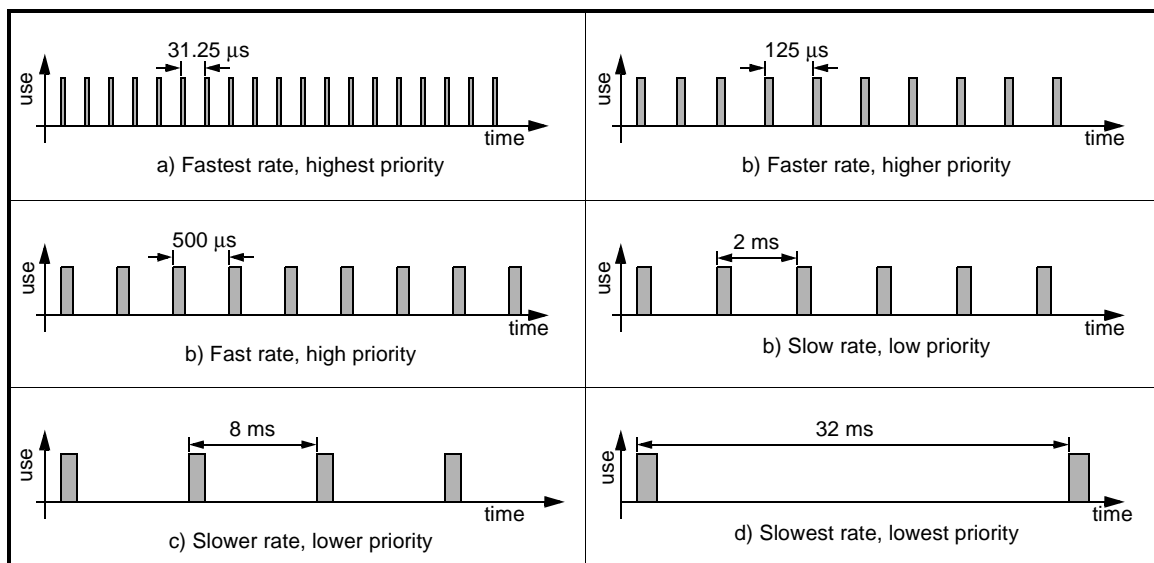


Figure 10.1—Rate-based priorities

10.1.1 Rate-based priorities

Quality of service is based on the availability of user_priority field parameter associated within transmitted time-sensitive frames, as listed in Table 10.1.

Table 10.1—Tagged priority values

Code	Interval (ms)	Name	Description
0	n/a	CLASS_C	Best effort, with minimal guaranteed BW
1	n/a	CLASS_B	Preferred, with minimal guaranteed BW
2	32	CLASS_A5	Guaranteed BW over longest interval
3	8	CLASS_A4	Guaranteed BW over longer interval
4	2	CLASS_A3	Guaranteed BW over long interval
5	0.5	CLASS_A2	Guaranteed BW over short interval
6	0.125	CLASS_A1	Guaranteed BW over shorter interval
7	0.03125	CLASS_A0	Guaranteed BW over shortest interval

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

10.1.2 Port-to-port reshaping

The concept of rate-based scheduling assumes shaped talkers and reshaped talker agents within bridges, as illustrated in Figure 10.2 (only the components associated with specific flows are illustrated). In this illustration, classA0 traffic flows between points (a, b, c, d, e), exhibits bunched and reshaped behaviors, as illustrated in Figure 10.3.

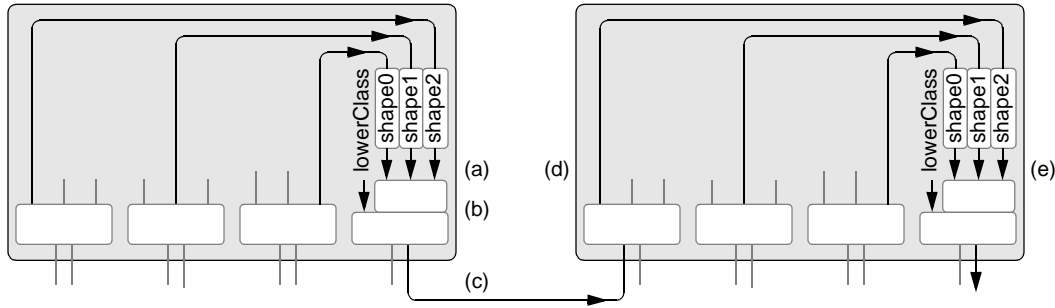


Figure 10.2—Reshaped bridge-traffic topology

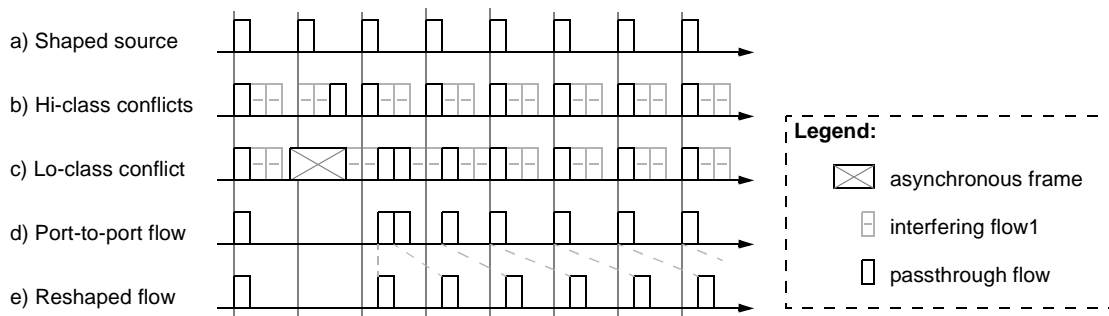


Figure 10.3—Reshaped bridge-traffic timing

The (a) through (e) time lines represent the flow of frames from within one talker-or-bridge into another bridge-or-listener, described as follows:

- a) A properly shaped source stream is originally generated within a talker, or a port-to-port flow (consisting of multiple streams) within a bridge.
- b) Forwarding of multiple sources to a shared transmission link can produce jitter, due to slight differences in frame-to-frame spacings.
- c) Forwarding of multiple sources to a shared transmission link can produce additional jitter, when higher-class traffic waits for the completion of previously initiated lower-class transmissions.
- d) Bunching becomes apparent in the port-to-port flow, representing the portion of the received (c) traffic that is forwarded to a specific transmitter port.
- e) A shaper delays the forwarding of bunched frames, so that the port-to-port flow is properly shaped. Delays can be invoked by time stamping frames with an in-the-future transmission time.

The reshaped flow (e) retains the properly shaped properties of the preceding flow (a), while incurring a maximum delay d through the bridge. These properties ensure a linear maximum delay of $n \times d$, for streams that flow through N bridges.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

10.1.3 Transmit ports

10.1.3.1 Transmit port structure

The transmit port is responsible for shaping classA traffic (to avoid bunching) and pacing classB/classC traffic (to avoid classC traffic starvation). Pacing and shaping algorithms assume functionally distinct queues within each transmit port, as illustrated in Figure 10.4.

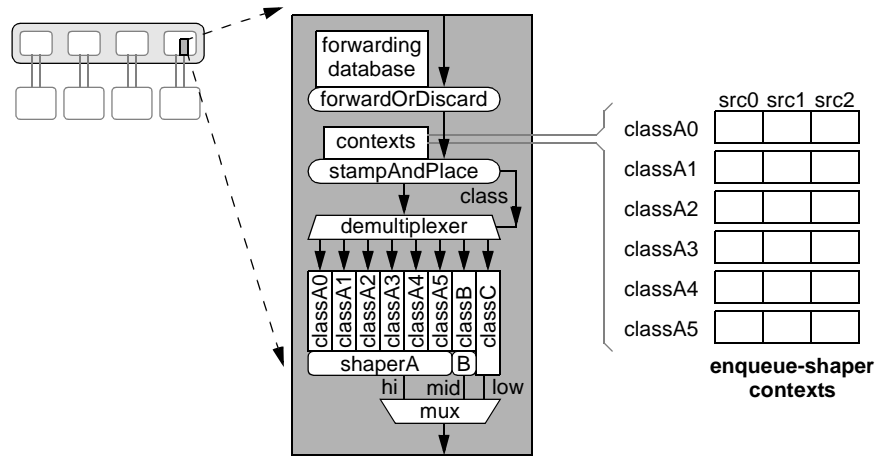


Figure 10.4—Transmit-queue structure

The intent of per-class shapers is to avoid priority inversions, wherein higher-class frames are delayed by the presence of concurrent lower-class traffic. Independent per-class shapers and queues allow enqueued higher-class and lower-class frames to be forwarded independently, thus avoiding priority inversions within queues.

The intent of per-source shapers is to avoid increasingly large cumulative bunching delays. The per-source reshaping eliminates bunches before merging, so that the pass-through bunching severity for 1-bridge and *n*-bridge flows are the same.

10.1.3.2 Enqueue reshaping contexts

The desired per-class latencies could not be guaranteed in the presence of classA traffic bunching. To avoid bunching, frames are shaped before being placed into classified transmit queues.

A shaper is responsible for attaching a time-stamp label to frames. One time-stamp shaper is logically associated with each source port and each classA traffic subclass (classA0, classA1, classA2, classA3, classA4, classA5). E.g, a four-port switch (which has three possible source ports) would have 18 time-stamp shapers.

The purpose of a time-stamp shaper is to associate a time-stamp label with each queued frame. The time-stamp label represents a time in the future; the frame’s transmission is deferred until the current time reaches the frame’s time-stamp value. This facilitates the delayed forwarding of successive frames within each bunch, thus suppressing the bunching effects found on receive-link transmission.

The context for each time-stamp shaper is based on the frame’s receive port and traffic class. While the context is considerably larger than that associated with strict per-port shapers, only one shaper (within each port) is ever active. Thus, context-switching per-port shaper instances represent a viable implementation technology

10.1.3.3 Dequeue shaping and pacing

Transmit ports utilize a shaper and pacer, as illustrated as shaperA and B components within Figure 10.4. The purpose of these is to ensure forward progress of best-effort control traffic. In concept, this involves a two-step bandwidth partitioning mechanism:

- The shaperA limits the cumulative classA and primary classB traffic to 75% of the link bandwidth. The intent is to ensure that 25% residual bandwidth remains available for lower-class traffic.
- Pacer B partitions the residual 25% traffic equally between classB and classC traffic. This ensures that classB traffic is never starved, in the presence of 75% classA traffic. This ensures that classC traffic is not starved, in the presence of excess classB traffic.

10.1.4 Credit-based shapers and pacers

10.1.4.1 Credit-based shapers

Although multiple shapers are specified within this working paper, the behavior of most shapers can be characterized by a common algorithm and instance-specific parameters (as done within RPR[B5]). The shaper's credits are adjusted down or up, as illustrated in Figure 10.5. The decrement and increment values typically represent sizes of a transmitted frame and of credit increments in each update interval, respectively.

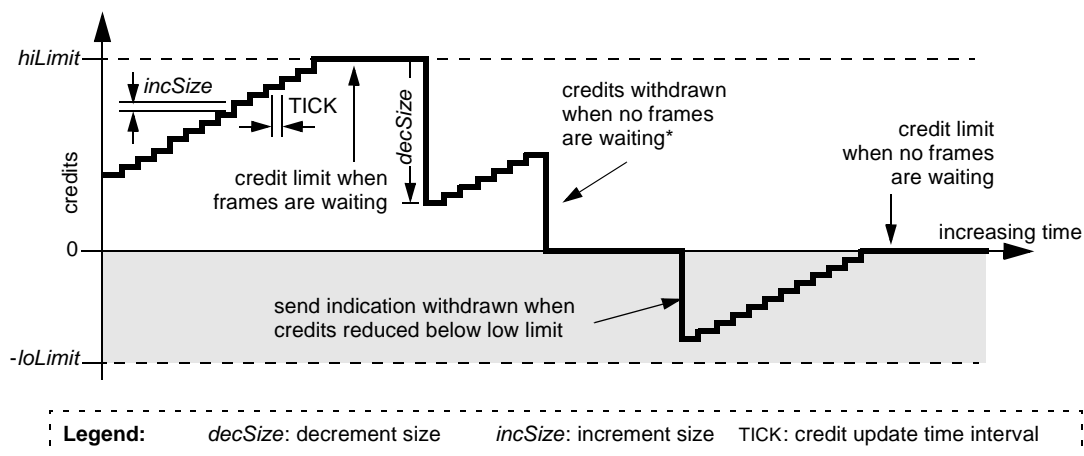


Figure 10.5—Credit-based shapers

Crossing below the zero threshold generates a rate-limiting indication, so that offered traffic can stop. By design, the credit value never goes below the $-loLimit$ extreme. To bound the burst traffic after inactivity intervals, when no frames are ready for transmission, credits are reduced to zero (if currently higher than zero) and can accumulate to no more than the zero-value limit.

The $hiLimit$ threshold limits the positive credits, to avoid overflow. When frames are ready for transmission (and are being blocked by transit traffic), credits can accumulate to no more than this $hiLimit$ value.

In concept, the shaper consist of a token bucket. The number of credits in a token bucket is decremented by the size of each transmitted frame. The credits in the token bucket are incremented at the end of every credit update interval. A frame is only transmitted when the credits are positive.

10.1.4.2 Credit-based pacers

Although multiple pacers are specified within this working paper, the behavior of most pacers can be characterized by a common algorithm and instance-specific parameters (as done within RPR[B5]). The pacer's credits are adjusted down or up, as illustrated in Figure 10.6. The decrement and increment values typically represent sizes of debit and credit frames in each update interval, respectively.

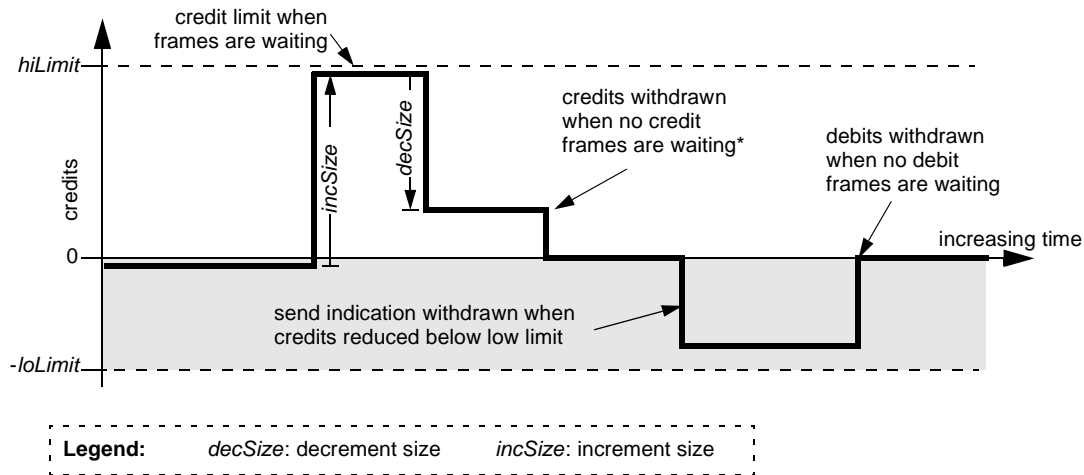


Figure 10.6—Pacer credit adjustments over time

Crossing below the zero threshold generates a rate-limiting indication, so that offered traffic can stop. By design, the credit value never goes below the $-loLimit$ extreme. To bound the burst traffic after inactivity intervals, when no frames are ready for transmission, credits are reduced to zero (if currently higher than zero) and can accumulate to no more than the zero-value limit.

The $hiLimit$ threshold limits the positive credits, to avoid overflow. When frames are ready for transmission (and are being blocked by transit traffic), credits can accumulate to no more than this $hiLimit$ value.

In concept, the pacer consists of a token bucket. The credits in the token bucket are incremented by the size of each transmitted debit-frame. The number of credits in a token bucket is decremented by the size of each transmitted credit-frame. A credit-frame is only transmitted when the credits are positive; a debit-frame is only transmitted when the credits are negative.

10.2 Terminology and variables

10.2.1 Common state machine definitions

The following state machine inputs are used multiple times within this clause.

queue values	1
Enumerated values used to specify shared queue structures.	2
QP_TX_PUSH—The transmit port’s internal queue, where received frames are placed.	3
QP_TX_A0—The first of the output port’s classA buffers.	4
QP_TX_A1—The second of the output port’s classA buffers.	5
QP_TX_A2—The third of the output port’s classA buffers.	6
QP_TX_A3—The fourth of the output port’s classA buffers.	7
QP_TX_A4—The second of the output port’s classA buffers.	8
QP_TX_A5—The third of the output port’s classA buffers.	9
QP_TX_BP—The output port’s classB queue.	10
QP_TX_CP—The output port’s classC queue.	11
QP_TX_LINK—The output port’s transmit-PHY queue.	12

10.2.2 Common state machine variables

One instance of each variable specified in this clause exists in each port, unless otherwise noted.

currentTime

A value representing the current time.

framed

The contents of a received frame, with supplemental information, as follows:

frame—The contents of a frame.

sourcePort—The source port that received the frame.

txTime—A time-stamp value representing the intended (bunching delayed) transmission time.

10.2.3 Common state machine routines

Max(value1, value2)

Returns the numerically larger of two values.

10.2.4 Variables and routines defined in other clauses

This clause references the following variables and routines defined in Clause 7:

currentTime

See 7.2.2.

Dequeue(queue)

Enqueue(queue, frame)

Min(value1, value2)

See 7.2.3.

10.3 Pacing state machines

10.3.1 TransmitRx state machine

The TransmitRx state machine is responsible for enqueueing traffic (received on other ports and broadcast to all possible transmitter ports) for possible forwarding. An intent is to transfer each to the appropriate output queue.

The following subclasses describe parameters used within the context of this state machine.

10.3.1.1 TransmitRx state machine definitions	1
queue values	2
Enumerated values used to specify shared queue structures.	3
QP_TX_A0, QP_TX_A1, QP_TX_A2, QP_TX_A3, QP_TX_A4, QP_TX_A5	4
QP_TX_BP, QP_TX_CP	5
QP_TX_PUSH	6
See 10.2.1.	7
	8
	9
10.3.1.2 TransmitRx state machine variables	10
	11
<i>class</i>	12
A value that represents the frame's priority class.	13
<i>count</i>	14
A value that represents the current credits, while minimum and maximum limits are being applied.	15
<i>currentTime</i>	16
See 10.2.4.	17
<i>delay</i>	18
A value that represents the time delay assigned by the frame's shaper.	19
<i>framed</i>	20
See 10.2.2.	21
<i>sPtr</i>	22
Represents a pointer to shaper values.	23
	24
10.3.1.3 TransmitRx state machine routines	25
	26
<i>ContextCheck(sourcePort, class)</i>	27
Returns a pointer to the associated pacer context, with the following fields:	28
<i>credit</i> —The cumulative credit from past pacer activities.	29
<i>lastTime</i> —The last time the pacer was invoked.	30
<i>loLimit</i> —The low limit for shaper credits.	31
<i>rate</i> —The highest allowed rate of the paced traffic, in bytes-per-second.	32
<i>Dequeue(queue)</i>	33
See 10.2.4.	34
<i>Enqueue(queue, frame)</i>	35
Places the <i>frame</i> at the tail of the specified <i>queue</i> within the assumed port.	36
<i>ForwardClass(framed)</i>	37
The forwarding database is checked. If forwarding is enabled, the priority class is returned.	38
Otherwise, a NULL class value is returned. The following enumerated values are returned:	39
CLASS_A0—The associated multicast frame is forwarded as classA traffic.	40
CLASS_A1—The associated multicast frame is forwarded as classA traffic.	41
CLASS_A2—The associated multicast frame is forwarded as classA traffic.	42
CLASS_A3—The associated multicast frame is forwarded as classA traffic.	43
CLASS_A4—The associated multicast frame is forwarded as classA traffic.	44
CLASS_A5—The associated multicast frame is forwarded as classA traffic.	45
CLASS_B—The associated multicast frame is forwarded as classB traffic.	46
CLASS_C—The associated multicast frame is forwarded as classC traffic.	47
<i>Max(value1, value2)</i>	48
See 10.2.3.	49
<i>Min(value1, value2)</i>	50
See 10.2.4.	51
	52
	53
	54

10.3.1.4 TransmitRx state table

The TransmitRx state machine is specified in Table 9.2. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.

Table 10.2—TransmitRx state table

Current		Row	Next	
state	condition		action	state
START	(framed = Dequeue(QP_TX_PUSH))!=NULL	1	—	FIRST
	—	2	—	START
FIRST	(class = ForwardClass(framed)) != NULL	3	sPtr = ContextCheck(framed.sourcePort, class);	NEXT
	—	4	—	START
NEXT	class == CLASS_A0	5	queue = QP_TX_A0;	PACE
	class == CLASS_A1	6	queue = QP_TX_A1;	
	class == CLASS_A2	7	queue = QP_TX_A2;	
	class == CLASS_A3	8	queue = QP_TX_A3;	
	class == CLASS_A4	9	queue = QP_TX_A4;	
	class == CLASS_A5	10	queue = QP_TX_A5;	
	class == CLASS_B	11	queue = QP_TX_BP;	
	—	12	queue = QP_TX_CP;	
PACE	—	13	count = sPtr->credit + sPtr->rate * (currentTime - sPtr->lastTime) - Size(frame); count = Max(0, Min(sPtr->loLimit, count)); sPtr->credit = count; sPtr->lastTime = currentTime; delay = Max(0, -count / sPtr->rate); framed.txTime = currentTime + delay;	FINAL
FINAL	—	14	Enqueue(queue, frame);	START

Row 10.2-1: If a frame has arrived, process that frame.

Row 10.2-2: Otherwise, wait for the next frame to arrive.

Row 10.2-3: When forwarded frames, the shaper context is based on the source port and class.

Row 10.2-4: The non-forwarded frames are discarded.

Row 10.2-5: The classA0 frames are forwarded to the appropriate time-sensitive classA0 queue.

Row 10.2-6: The classA1 frames are forwarded to the appropriate time-sensitive classA1 queue.

Row 10.2-7: The classA2 frames are forwarded to the appropriate time-sensitive classA2 queue.

Row 10.2-8: The classA3 frames are forwarded to the appropriate time-sensitive classA3 queue.

Row 10.2-9: The classA4 frames are forwarded to the appropriate time-sensitive classA4 queue.

Row 10.2-10: The classA5 frames are forwarded to the appropriate time-sensitive classA5 queue.

Row 10.2-11: The classB frames are forwarded to the appropriate time-sensitive classB queue.	1
Row 10.2-12: The classC frames are forwarded to the appropriate time-sensitive classC queue.	2
	3
Row 10.2-13: ClassA frames are time stamped by shapers.	4
Shaper pacer parameters for each distinct { <i>class</i> , <i>source</i> } pair constrains bunching within each class.	5
Shaper parameters are updated as in 10.1.4.1: decremented on transmissions and incremented over time.	6
High and low limits are applied to the updated credits and the last-updated time is updated.	7
Negative credits correspond to transmission-delay values, which are attached to output-port queued frames.	8
	9
Row 10.2-12: The received frames are placed into the appropriate queue.	10
	11
10.3.2 TransmitTx state machine	12
	13
The TransmitTx state machine is responsible for pacing/shaping classA traffic and shaping classB traffic destined for 1 Gb/s links. An intent is to support projected MTU-sized transfers and interleaved lower-class traffic, without exceeding the 1-cycle delay inherent with cycle-synchronous bridge-forwarding protocols.	14
	15
	16
The following subclasses describe parameters used within the context of this state machine.	17
	18
	19
10.3.2.1 TransmitTx state machine definitions	20
	21
BPS	22
The nominal link transmission rate, in bytes per second.	23
MTU	24
The maximum frame size, in bytes.	25
queue values	26
Enumerated values used to specify shared queue structures.	27
QP_TX_A0, QP_TX_A1, QP_TX_A2, QP_TX_A3	28
QP_TX_BP, QP_TX_CP	29
QP_TX_LINK	30
See 10.2.1.	31
TICK	32
The amount of time between shaper updates.	33
Range: [1 bytes transmit time, 16-bit transmit time]	34
Default: 1 byte transmit time	35
	36
10.3.2.2 TransmitTx state machine variables	37
	38
<i>best</i>	39
A value that represents the weight and identify of the next-best classA queue.	40
<i>goodness</i> —The smallest <i>weight</i> × <i>wait</i> value associated with alternate classA transmissions.	41
<i>queue</i> —The queue associated with the best futuristic encapsulated frame.	42
<i>countA</i>	43
A speculative value of creditA, used only when the frame is qualified for transmission.	44
<i>countB</i>	45
A speculative value of creditB, used only when the frame is qualified for transmission.	46
<i>creditA</i>	47
A shaper credit whose positive value enables classA/classB primary transmissions.	48
<i>creditB</i>	49
A shaper credit value whose positive and negative values enable secondary classB and classC transmissions respectively.	50
	51
<i>currentTime</i>	52
See 10.2.4.	53
	54

<i>frame</i>	1
The contents of a to-be-transmitted frame.	2
<i>framed</i>	3
See 10.2.2.	4
<i>hiLimitA</i>	5
A value that limits the cumulative <i>creditA</i> credits.	6
Value: MTU.	7
<i>hiLimitB</i>	8
A value that limits the cumulative <i>creditB</i> credits.	9
Value: MTU.	10
<i>limit</i>	11
A value that limits the amount of transmitted primary classA/classB bandwidth.	12
<i>loLimitA</i>	13
A value that limits the cumulative <i>creditA</i> debits.	14
Value: MTU.	15
<i>loLimitB</i>	16
A value that limits the cumulative <i>creditB</i> debits.	17
Value: MTU.	18
<i>tickTime</i>	19
A value that defines when the time-tick interval ends.	20
	21
10.3.2.3 TransmitTx state machine routines	22
	23
<i>Dequeue(queue)</i>	24
See 10.2.4.	25
<i>Unqueue(queue, weight, &best, currentTime)</i>	26
Dequeues and returns the most overdue frame from the specified <i>queue</i> , excluding those frames whose scheduled transmission time is after the specified <i>currentTime</i> value.	27
<i>framed</i> —The oldest of the overdue frame.	28
NULL—No frame available.	29
In the presence of only futuristic frames, a <i>test</i> = $weight \times (txTime - currentTime)$ value is computed.	30
If <i>best.queue</i> is NULL or <i>test</i> < <i>best.goodness</i> , the <i>best.queue</i> and <i>best.goodness</i> components are updated to reflect the best alternate classA transmission queue.	31
	32
	33
<i>Enqueue(queue, frame)</i>	34
See 10.2.4.	35
<i>Size(frame)</i>	36
Returns the size of the specified frame.	37
<i>StaleFrame(frame, queue)</i>	38
Indicates whether the specified frame is stale and discardable, as specified by Equation 10.1.	39
0—The specified frame is not stale.	40
1—(Otherwise.)	41
	42
<pre>// The value of "internal" depends on the class, as specified in Table 10.1. (10.1) (index = (QP_TX_A0 - queue), (currentTime - framed.txTime) > (2 * (MTU + interval[index] * BPS)))</pre>	43
	44
	45
10.3.2.4 TransmitTx state table	46
	47
The TransmitTx state machine is specified in Table 9.3. In the case of any ambiguity between the text and the state machine, the state machine shall take precedence. The notation used in the state table is described in 3.4.	48
	49
	50
	51
Row 10.3-1: Update the classA credits after each tick interval.	52
Row 10.3-2: Wait for the queue to be emptied, so that something can be transmitted.	53
	54

Table 10.3—TransmitTx state table

Current		Row	Next	
state	condition		action	state
START	(currentTime - tickTime) >= TICK;	1	creditA = Min(hiLimitA, creditA + 0.75 * TICK * BPS); tickTime = currentTime ;	START
	!QueueEmpty(QP_TX_LINK)	2	—	
	creditA < 0	3	—	FAIR
	—	4	best.queue = NULL;	BEST
BEST	(framed = Unqueue(queue= QP_TX_A0, 32, &best, currentTime)) != NULL	5	countA = Min(loLimitA, creditA - Size(framed));	NEAR
	(framed = Unqueue(queue= QP_TX_A1, 16, &best, currentTime)) != NULL	6		
	(framed = Unqueue(queue= QP_TX_A2, 8, &best, currentTime)) != NULL	7		
	(framed = Unqueue(queue= QP_TX_A3, 4, &best, currentTime)) != NULL	8		
	(framed = Unqueue(queue= QP_TX_A4, 2, &best, currentTime)) != NULL	9		
	(framed = Unqueue(queue= QP_TX_A5, 1, &best, currentTime)) != NULL	10		
	best.queue != NULL && (framed = Dequeue(queue= best.queue)) != NULL	11		
	(framed = Dequeue(QP_TX_BP)) != NULL	12		
	—	13	creditA = 0;	START
FAIR	creditB >= 0 && (framed = Dequeue(QP_TX_BP)) != NULL	14	creditB = creditB - Size(framed);	FINAL
	creditB <= 0 && (framed = Dequeue(QP_TX_CP)) != NULL	15	creditB = creditB + Size(framed);	
	(framed = Dequeue(QP_TX_BP)) != NULL	16	creditB = 0;	
	(framed = Dequeue(QP_TX_CP)) != NULL	17		
	—	18	creditB = 0;	START
NEAR	StaleFrame(framed, queue)	19	—	START
	—	20	creditA = countA;	FINAL
FINAL	—	21	Enqueue(QP_TX_LINK, framed.frame);	START

Row 10.3-3: In the absence of classA credits, fairly transmit enqueued classB and classC frames.

Row 10.3-4: Fairly service classB/classC when the classA/classB transmissions are disallowed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Row 10.3-5: If enabled and available, a classA0 frame is transmitted.	1
Row 10.3-6: If enabled and available, a classA1 frame is transmitted.	2
Row 10.3-7: If enabled and available, a classA2 frame is transmitted.	3
Row 10.3-8: If enabled and available, a classA3 frame is transmitted.	4
Row 10.3-9: If enabled and available, a classA4 frame is transmitted.	5
Row 10.3-10: If enabled and available, a classA5 frame is transmitted.	6
Row 10.3-11: If available, a scheduled-for-the-future classA frame is transmitted.	7
Row 10.3-12: If enabled and available, a classB frame is transmitted.	8
Row 10.3-13: Since nothing is ready to be sent, the classA credits are cleared.	9
	10
Row 10.3-14: If enabled and available, a classB frame is transmitted.	11
The <i>creditB</i> values is decremented by the transmitted frame size, to avoid classC starvation.	12
Row 10.3-15: If enabled and available, a classC frame is transmitted.	13
The <i>creditB</i> values is incremented by the transmitted frame size, to avoid classB starvation.	14
Row 10.3-16: If available, a classB frame is transmitted.	15
Row 10.3-17: If available, a classC frame is transmitted.	16
Row 10.3-18: Otherwise, no frame is transmitted.	17
	18
Row 10.3-19: Stale frames, whose delivery times cannot be guaranteed, are discarded.	19
Row 10.3-20: Non-stale frames are not discarded.	20
	21
Row 10.3-21: The next frame is transmitted and credits are updated accordingly.	22
	23
	24
	25
	26
	27
	28
	29
	30
	31
	32
	33
	34
	35
	36
	37
	38
	39
	40
	41
	42
	43
	44
	45
	46
	47
	48
	49
	50
	51
	52
	53
	54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Annexes

Annex A

(informative)

Bibliography

NOTE—This clause should be skipped on the first reading (continue with Annex B).
Although not finalized, this bibliography provides useful material for understanding this working paper.

- [B1] IEEE 100, The Authoritative Dictionary of IEEE Standards Terms, Seventh Edition.¹
- [B2] IEEE Std 802-2002, IEEE Standards for Local and Metropolitan Area Networks: Overview and Architecture.
- [B3] IEEE Std 801-2001, IEEE Standard for Local and Metropolitan Area Networks: Overview and Architecture.
- [B4] IEEE Std 802.1D-2004, IEEE Standard for Local and Metropolitan Area Networks: Media Access Control (MAC) Bridges.
- [B5] IEEE Std 802.17-2004, IEEE Standard for Local and Metropolitan Area Networks: Resilient packet ring (RPR) access method and physical layer specifications.
- [B6] IEEE Std 1394-1995, High performance serial bus.
- [B7] IEEE Std 1588-2002, IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems.
- [B8] IETF RFC 1305: Network Time Protocol (Version 3) Specification, Implementation and Analysis, David L. Mills, March 1992²
- [B9] IETF RFC 2030: Simple Network Time Protocol (SNTP) Version 4 for IPv4, IPv6 and OSI, D. Mills, October 1996.
- [B10] IETF RFC 2205: Resource Reservation Protocol (RSVP), R. Braden, L. Zhang, S. Berson, and S. Herzog, S. Jamin, October 1996.

¹IEEE publications are available from the Institute of Electrical and Electronics Engineers, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331, USA (<http://standards.ieee.org/>).

²IETF publications are available via the World Wide Web at <http://www.ietf.org>.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Annex B

(informative)

Background material

B.1 Related standards

B.1.1 IEEE 1394 Serial Bus

As background, real-time features of an existing (and widely adopted on PCs) serial interface standard are summarized in this subclause: IEEE 1394-1995 High Performance Serial Bus. To avoid confusion with other serial buses (serial ATA, etc.), the term “SerialBus” is used within this annex to refer to this specific IEEE standard.

B.1.1.1 SerialBus topologies

Since its conception, SerialBus evolved from being a shared bus (like Ethernet) to a collection of point-to-point duplex links, as illustrated in Figure B.1. Arbitrary hierarchical topologies can be supported, but dotted-line redundant looping connections are only allowed in recent upgrades of the standard.

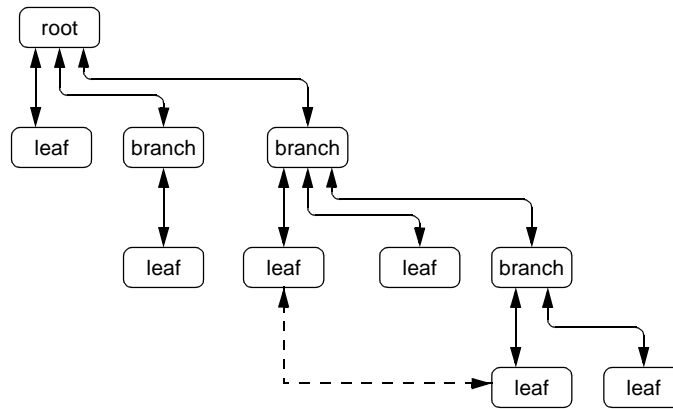


Figure B.1—SerialBus topologies

This physical duplex-link topology could, in concept, support concurrent non-overlapping data transfers. SerialBus only partially utilizes these capabilities (arbitration and data transfers can be overlapped), because its arbitration protocols were inherited from its initial conception as an arbitrated shared broadcast bus.

B.1.1.2 Isochronous data transfers

SerialBus isochronous traffic is transmitted at a 8 kHz rate, as illustrated by the 125 μ s cycles within Figure B.2.

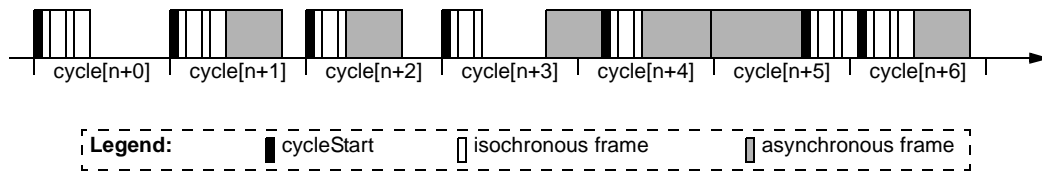


Figure B.2—Isochronous data transfer timing

In the absence of conflicting traffic, an 8kHz cycle starts with the transmission of a cycleStart frame, as illustrated in cycle[n+0]. The cycleStart frame triggers the sending of the isochronous frames that have been queued for cycle[n+0] transmission; these continue until all isochronous traffic has been sent.

After a cycle's isochronous traffic has been sent, one or more asynchronous transmissions are allowed, as illustrated in cycle[n+1].

Devices can be paused, compression rates can be variable, and connections can fail. For such reasons, the amounts of isochronous traffic within each cycle can vary below its scheduled limits, as illustrated in cycle[n+2].

The asynchronous traffic is not constrained to start at the end of a cycle, but can start at anytime that the frame is available and isochronous transfers are idle, as illustrated near the end of cycle[n+3]. If started near the end of a cycle, the isochronous transfer can be forced to start within the following cycle[n+4].

A large late-starting asynchronous frame can extend the start of isochronous transfers, so that spill-over into the next cycle is possible, as illustrated in cycle[n+5]. Since isochronous transfers have priority, the delay in the next isochronous cycle is reduced, and the isochronous traffic completes within the boundaries of cycle[n+6].

B.1.1.3 Isochronous reservations

Even the best of isochronous transfers fails when the offered load exceeds the link capacity. To eliminate this possibility, isochronous bandwidth is reserved before being consumed. On a single bus (of up to 64 stations), reservations are controlled through access to compare&swap register, which all isochronous stations provide, although only one is selected to be used (based on the largest populated device address).

On a multiple bus topology (buses interconnected through bridges), reservations management is more complex. In this case, frames are passed from the source to its desired-to-be-connected destination(s), reserving reservations along the data-transmission path. As is true on a single bus, reservation requests are rejected when insufficient bandwidth capacity remains. This is not described in the baseline 1394 specification, but is described in a follow-on P1394.1 draft (currently progressing through Sponsor ballot).

B.1.1.4 SerialBus experiences

Experiences, as follows:

- a) Cycle slip. Cycle-slip reduces design complexity, permits transmissions of large asynchronous frames, and improves asynchronous traffic throughput. Transmission precision is unnecessary: error in the cycleStart transmission time is encoded within that frame, allowing clock-slave devices to accurately adjust their phase-lock-loops, regardless of observed cycleStart transmission times.
- b) Cycle time. An 8 kHz cycle rate represents a good trade-off between efficiency (the overhead is less, when cycle times are longer) and latency (the latency is less, when cycle times are longer).
- c) Pseudo frames. The SerialBus isochronous frames have a distinct (6-bit channel number) addressing scheme. In hindsight, using a standard frame header (destination address and source address) would have many benefits, including the simplification of bridges between segments.
- d) Service classes. SerialBus has evolved to support three classes of traffic: isochronous, prioritized asynchronous, and baseline asynchronous. These are roughly equivalent to the classA, classB, and classC service classes defined for RPR (see B.1.2).

B.1.2 Resilient packet ring (RPR)

As background, the time-sensitive capabilities associated with IEEE P802.17 Resilient packet ring (RPR) are summarized in this subannex. RPR is a metropolitan area network (MAN) that can be transparently bridged to Ethernet.

B.1.2.1 RPR rings

RPR employs a ring structure using unidirectional, counter-rotating ringlets. Each ringlet is made up of links with data flow in the same direction. The ringlets are identified as ringlet0 and ringlet1, as shown in Figure B.3.

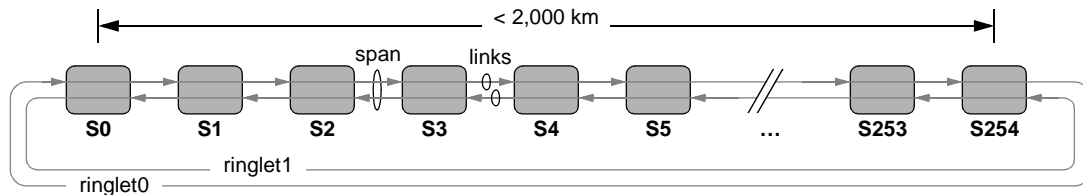


Figure B.3—RPR rings

Stations on the ring are identified by an IEEE 802 48-bit MAC address. All links on the ring operate at the same data rate, but may exhibit different delay properties. Ring circumference of less than 2,000 kilometers are assumed.

The portion of a ring bounded by adjacent stations is called a span. A span is composed of unidirectional links transmitting in opposite directions.

B.1.2.2 RPR resilience

RPR stations are resilient, in that communications can continue in that operations continue in the presence of single-point failures, as illustrated in Figure B.4. Resilient features can recover from failed links by bypassing the frame-manipulation portions of a partially failed station (see Figure B.4-b), thus avoiding a failed station (see Figure B.4-c and Figure B.4-d) or a failed span (see Figure B.4-e and Figure B.4-f).

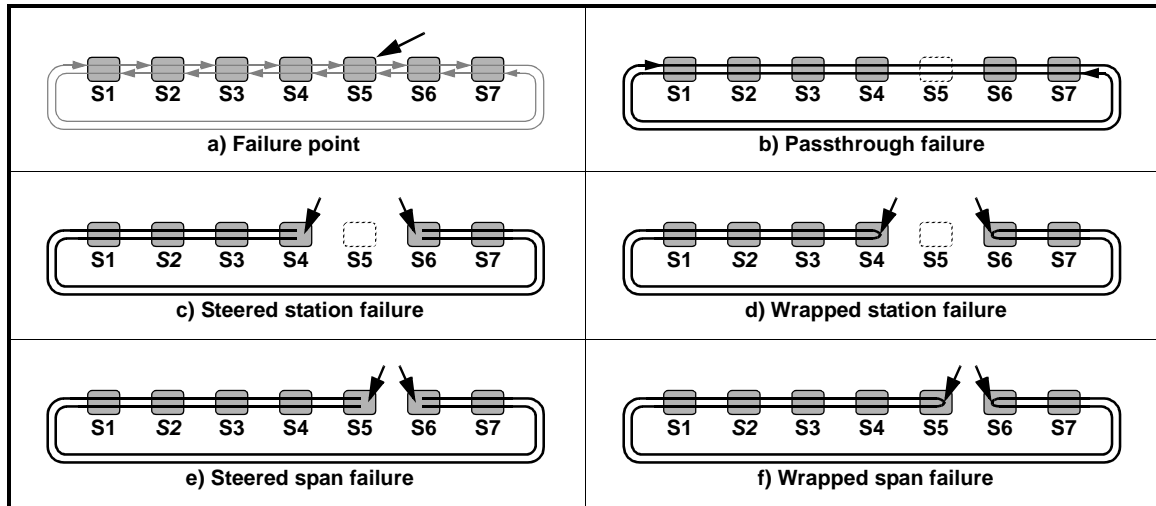


Figure B.4—RPR resilience

B.1.2.3 RPR spatial reuse

RPR efficiently strips local unicast frames at their destination, so that bandwidth on unaffected links is available for other frame transfers, as illustrated in Figure B.5-a. A unicast frame is added by the source station, and is stripped at the destination station. The frame is normally copied at the destination station for delivery to the local MAC client or MAC control entity. If ringlet selection is based on shortest hop-count, a response frame is likely to take an opposing ringlet path, as illustrated in Figure B.5-b.

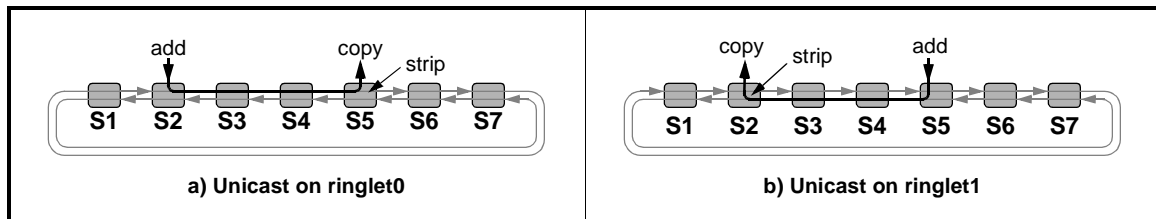


Figure B.5—RPR destination stripping

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1 The RPR frame transmissions on one link are largely independent of frame transmissions on other link. This
2 allows per-link bandwidths to be utilized beyond that possible with IEEE Std 802.5-1998 Token Ring or
3 ANSI FDDI ring based LAN technologies. Spatial reuse is illustrated in Figure B.6.

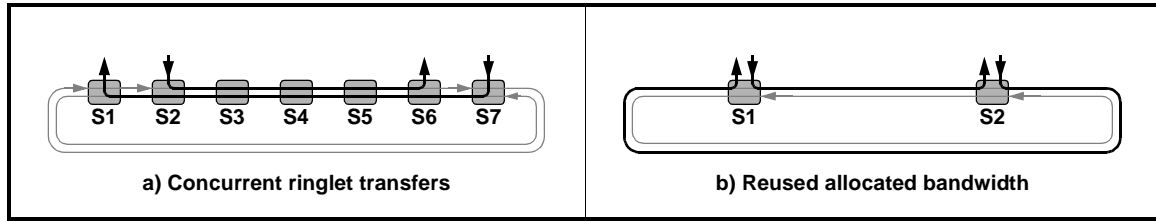


Figure B.6—RPR spatial reuse

16 Concurrent per-ringlet transmissions (see Figure B.6-a) allow stations bandwidths to exceed individual link
17 capacities. The effective bandwidths of non-overlapping transfers (see Figure B.6-b) are similarly improved.

19 B.1.2.4 RPR service classes

21 RPR provides transit queues, which allow received traffic to be queued during a station's frame
22 transmission, as illustrated in Figure B.7. The highest priority frames are classA and have their own bypass
23 buffer; the lower priority frames are classB and classC, and share the use of a distinct bypass buffer. To
24 minimize the classA latencies, servicing of the classA buffer has precedence over servicing of the
25 classB/classC buffer.

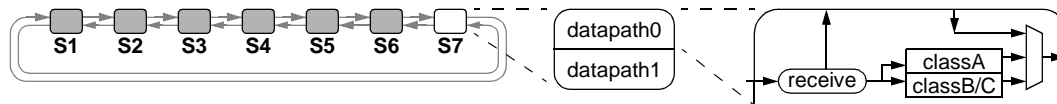


Figure B.7—RPR service classes

35 During the initial phases of investigation, techniques for allowing newly-arrived classA traffic to preempt an
36 active classB/classC frame transmission were considered. While such techniques are practical, the metro-
37 politan area networks (MANs) environments limits the effectiveness of such techniques; at these longer
38 distances, the link delays can often exceed the retransmission-blocked delays within individual stations.

Annex C

(informative)

Encapsulated IEEE 1394 frames

To illustrate the sufficiency and viability of the RE isochronous services, the transformation of IEEE 1394 packets is illustrated. A connection between an IEEE 1394 talker, IEEE 1394 adapter, intermediate Ethernet links, IEEE 1394 adapter, and an IEEE 1394 listener is assumed.

C.1 Hybrid network topologies

C.1.1 Supported IEEE 1394 network topologies

This annex focuses on the use of RE to bridge between IEEE 1394 domains, as illustrated in Figure C.1. The boundary between domains is illustrated by a dotted line, which passes through a SerialBus adapter station.

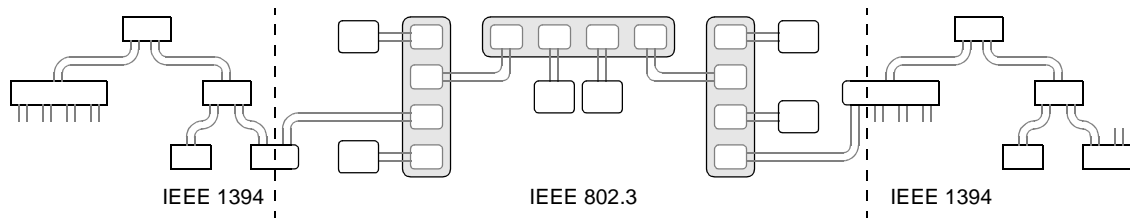


Figure C.1—IEEE 1394 leaf domains

C.1.2 Unsupported IEEE 1394 network topologies

Another approach would be to use IEEE 1394 to bridge between IEEE 802.3 domains, as illustrated in Figure C.2. While not explicitly prohibited, architectural features of the topology-supportive adapters and encapsulated-frame formats are beyond the scope of this working paper.

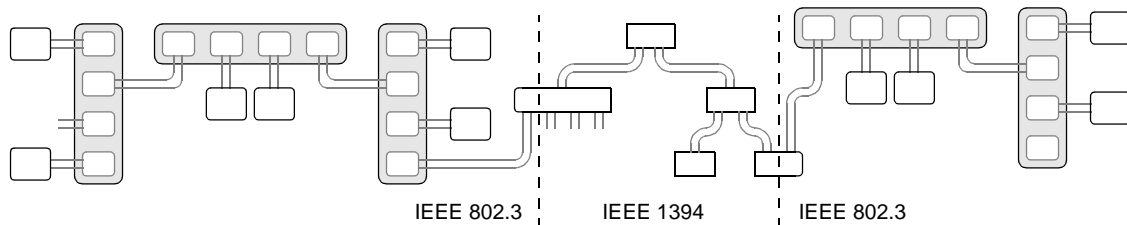


Figure C.2—IEEE 802.3 leaf domains

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

C.2 1394 isochronous frame formats

C.2.1 1394 isochronous frame formats

An IEEE 1394 isochronous frame contains header and payload components, as illustrated by Figure C.3. While all components could be encapsulated into an Ethernet frame, some of these fields would be redundant (with fields in the encapsulating frame) or unnecessary.

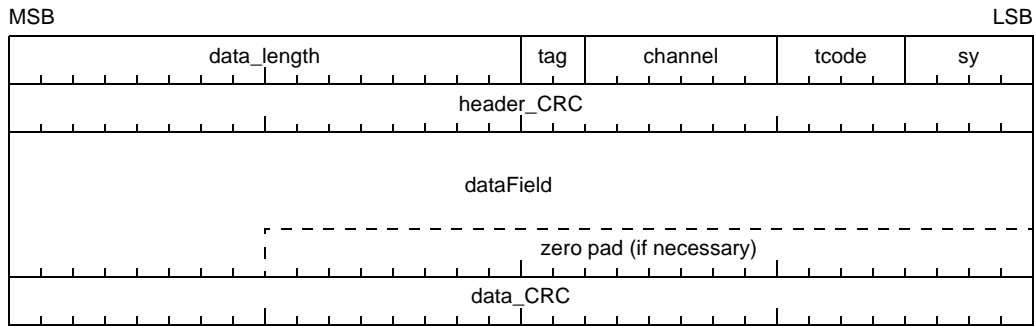


Figure C.3—IEEE 1394 isochronous packet format

C.2.2 Encapsulated IEEE 1394 frame payload

For uniframe groups, the IEEE 1394 isochronous frames are modified slightly and placed within an Ethernet *serviceDataUnit*. The format of this *serviceDataUnit* is illustrated by Figure C.4.

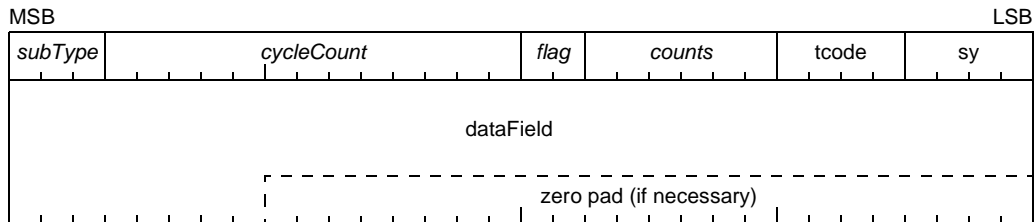


Figure C.4—Encapsulated IEEE 1394 frame payload

C.2.2.1 subType: A 3-bit field that distinguishes encapsulated 1394 frames from other formats with the same *protocolType* specifier.

C.2.2.2 cycleCount: A 13-bit field that identifies the isochronous cycle during which this frame was transmitted. For the first frame within any group, this information is needed to perform CIP header updates (see C.4). These fields also provide error-detecting consistency checks.

C.2.2.3 flag: A 2-bit field that distinctively identifies the frame type, as specified in Table C.1.

Table C.1—*flag* field values

Value	Name	Description
0	ONLY	Only frame within a uniframe group
1	LAST	Final frame within a multiframe group
2	CORE	Intermediate frame within an multiframe group
3	LEAD	First frame within a multiframe group

C.2.2.4 counts: A 6-bit field that identifies additional frame-group parameters, as specified in Table C.2. When interpreted as a *partCount* value, this effectively identifies the number of zero-pad bytes. When interpreted as a *frameCount* value, the values of $\{n-1, n-2, \dots, 1\}$ label the first through next-to-last frames of an n -frame multiframe group.

Table C.2—*counts* field values

flag	Name	Description
ONLY	<i>partCount</i>	The LSBs of the residual data_length field.
LAST		
CORE	<i>frameCount</i>	A sequence identifier for frames within the group
LEAD		

C.2.2.5 dataField: For a uniframe group, the contents of the SerialBus ‘data field’ bytes.

C.3 Frame mappings

C.3.1 Synchronous frame mappings

Adapters are required to manage differences between IEEE 1394 isochronous packets and RE frames, as illustrated in Figure C.5.

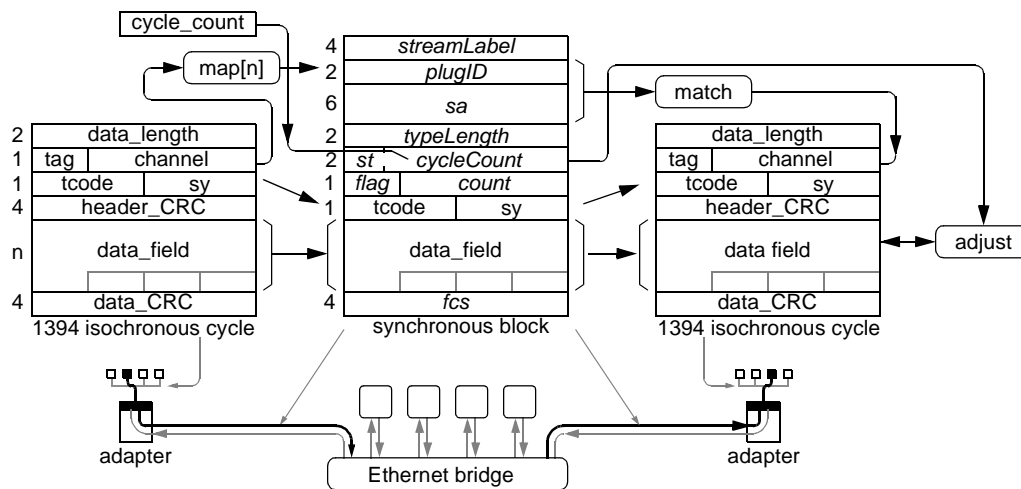


Figure C.5—Conversions between IEEE 1394 packets and RE frames

The IEEE 1394 to Ethernet frame translation involves the following:

- a) The IEEE 1394 `data_length` field is discarded (The `data_length` information can be reconstructed from the length of the received frame.)
- b) The IEEE 1394 `tag` field is ignored (this connection context is known to higher layer software).
- c) The IEEE 1394 `channel` field becomes an index into an array of communication contexts. The selected context provides the `plugID` value, the least-significant portion of the Ethernet `da`.
- d) The IEEE 1394 isochronous transmission cycle number is copied to the Ethernet `cycleCount` field. (The cycle number is the `cycle_time_data.cycle_count` field from the preceding cycle-start packet.)
- e) The IEEE 1394 `tcode` and `sy` fields are copied to the corresponding Ethernet fields.
- f) The `data_length`, `header_CRC`, and `data_CRC` fields are checked; if any are found to be inconsistent, no RE frame is created (the presumed to be corrupted frame is dropped).

NOTE — Unlike IEEE 1394, no synchronous frame transformations are required when passing through bridges. This is consistent with 802.3 specifications, which leave frames unmodified when passing through bridges.

The Ethernet to IEEE 1394 frame translation involves the following:

- a) Invalid Ethernet frames (multicast `sa` address, too-short or too-long, or bad `fcs`) are discarded.
- b) The IEEE 1394 `data_length` field is derived from the Ethernet frame length.
- c) The context with the matching `streamId` (`sa` concatenated with `plug`) values is selected. This context provides the provides the channel field value.
- d) The IEEE 1394 `tag` and `tcode` fields are set to identify isochronous IEEE 1394 packets.
- e) The IEEE 1394 `tcode` and `sy` fields are copied from the Ethernet frame.
- f) The IEEE 1394 `data_field` is directly mapped to the RE content field. (IEC61883-type content may have its synchronization fields updated as needed, see C.4.)
- g) The IEEE 1394 `header_CRC` and `data_CRC` fields are computed.

C.3.2 Multiframe groups

To avoid exceeding the maximum Ethernet frame size, large frames are decomposed into multiframe groups. The initial frames within the multiframe group are distinctively identified by their *counts* values, as illustrated in Figure C.6.

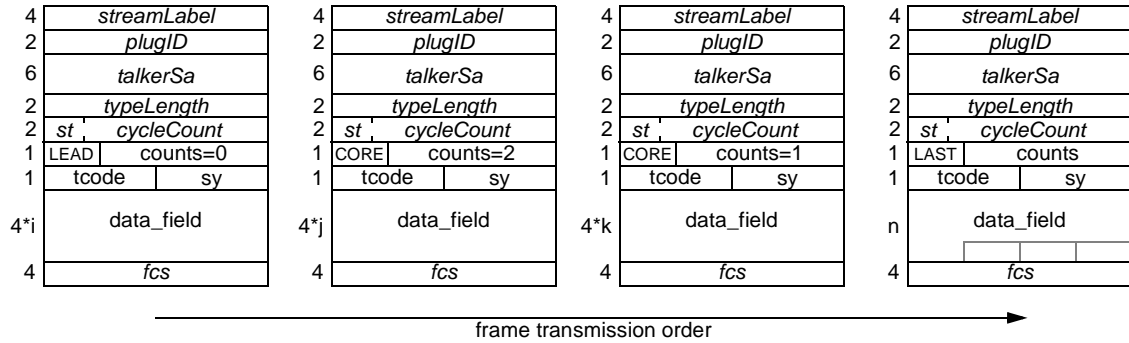


Figure C.6—Multiframe groups

The final frame within the group is identified by its distinctive *flag*=LAST identifier. For this frame, the *counts* field specifies the number of data bytes within the frame, modulo 64.

C.4 CIP payload modifications

Isochronous 1394 data packets may conform to a common isochronous packet (CIP) format, as defined by IEC 61883/FIS. The presence of a CIP format is indicated by a tag=1 bit in the Serial Bus isochronous packet header, as illustrated in Figure C.7. The white shading identifies those fields (when present and valid) are modified when passing through a RE-to-1394 adapter.

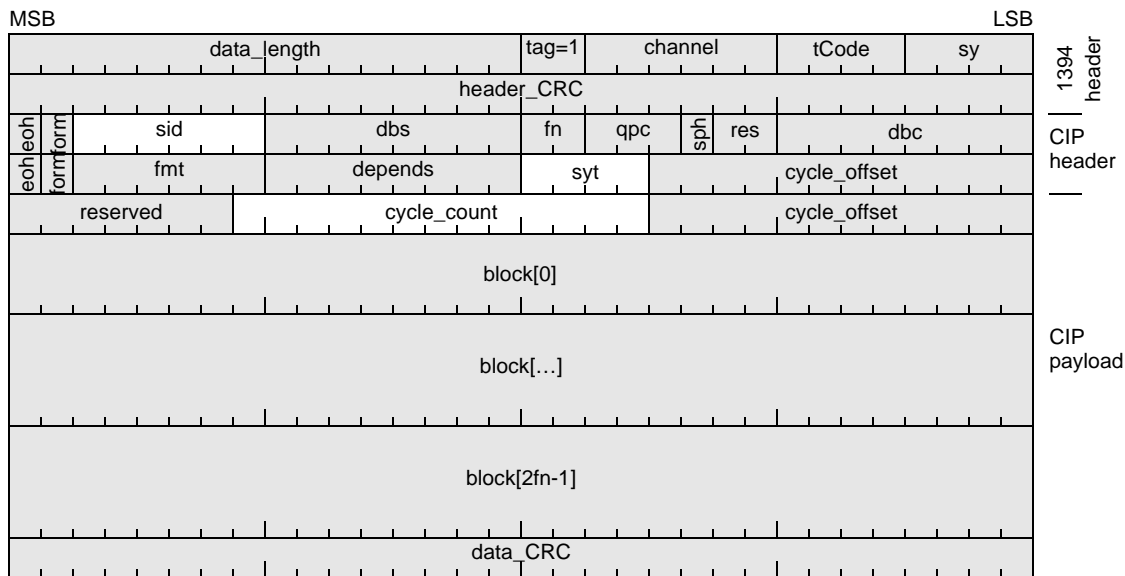


Figure C.7—Isochronous 1394 CIP packet format

The *sid* field must be set to the physical ID of the talking portal. This allows the listener to identify the bridge's talker portal.

Two-quadlet CIP headers may also contain absolute time stamp information or indicate its presence elsewhere in the packet's data payload. Absolute time stamps may be found in one or more places in isochronous:

- the *syt* field of the second quadlet of the CIP header if the *fmt* field in that quadlet has a value between zero and $1F_{16}$, inclusive; and
- the *cycle_count* and *cycle_offset* fields of all of the source packet headers (SPH) within the isochronous subaction.

Both of these time stamps are specified as absolute values that specify a future cycle time. Since isochronous subactions experience delays when routed over RE, these time stamps must be adjusted by the difference in cycle times between the talker and the RE-to-1394 bridge. The delay, in units of cycles, is the difference between the talker and 1394 adapter's transmission times, as specified in Equation 3.2.

$$\text{latency} = (\text{adapter.sendCycle} - \text{syncBock.talkerCycle}); \tag{3.1}$$

When the *syt* or *cycle_count* fields are present, their adjustments are specified by Equation 3.2. Because IEEE 1394 constrains *cycle_count* to the range zero to 7999, inclusive, the time stamp adjustments must be performed modulus 8000

$$\text{transmitted.syt} = (\text{received.syt} + \text{latency}) \% 8000; \tag{3.2}$$

$$\text{transmitted.cycle_count} = (\text{received.cycle_count} + \text{latency}) \% 8000; \tag{3.3}$$

C.4.1 Time-of-day format conversions

The difference between RE and IEEE 1394 time-of-day formats is expected to require conversions within the RE-to-1394 adapter. Although multiplies are involved in such conversions, multiplications by constants are simpler than multiplications by variables. For example, a conversion between RE and IEEE 1394 involves no more than two 32-bit additions and one 16-bit addition, as illustrated in Figure C.8.

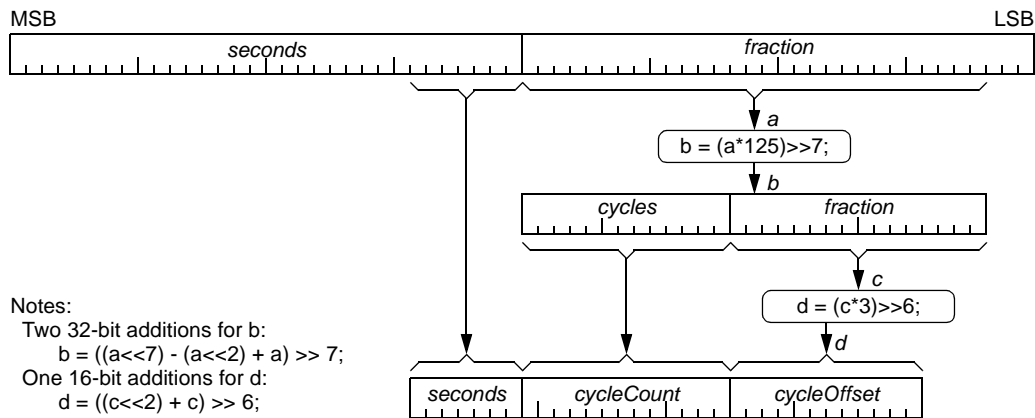


Figure C.8—Time-of-day format conversions

C.4.2 Grand-master precedence mappings

Compatible formats allow either an IEEE 1394 or IEEE 802.3 stations to become the network's grand-master station. While difference in format are present, each format can be readily mapped to the other, as illustrated in Figure C.9:

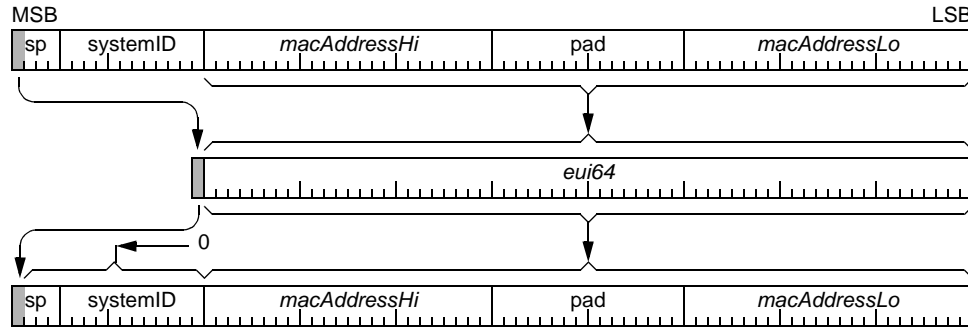


Figure C.9—Grand-master precedence mapping

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Annex D

(informative)

Review of possible alternatives

D.1 Clock-synchronization alternatives

D.1.1 Statistical averaging

Wide-area network based protocols distribute time by enclosing time-stamp values in specialized calibration frames. Higher level frame-processing protocols are responsible for determining the average transmission delays through the interconnect, so that calibration-frames can be used for accurate time-synchronization purposes.

The frame transmission latency is highly variable, based on delays incurred when waiting behind other previously-queue frames. Long-term averaging is typically used to cope with nonrandom delays, whether they be periodic, biased, or time-of-day dependent.

The use of long-time averages has limited applicability within the home, where small numbers of streams can exhibit very non-random statistical behaviors. Furthermore, long-term averaging intervals restricts transient-event response times, such as the insertion or removal of associated clock-synchronized devices.

D.1.2 Phase-locked synchronization

Local-area network based protocols, such as IEEE Std 1588, specify communication protocols for communicating timer-difference errors from a local clock-master station to its neighboring clock-slave station. However, this standard does not define how the clock-slave station compensates its values to track the time reference of the neighboring clock-master station.

The most common method of synchronizing clock-master and clock-slave devices involves phase-lock-loop (PLL) circuits. Such circuits integrate sensed differences between the clock-master and clock-slave devices, using these integrated values to adjust the clock-slave operating frequency.

The clock-slave resident PLLs are useful for reducing the transmission-induced timing-error jitters. However, the response time of a cascaded set of PLLs degrades as the number of cascaded devices increases. Also, the dynamics of more-responsive (gain peaking) cascaded PLL can be undesirable, causing the deviations of later stages to exponentially increase with their distance from the source, a characteristic commonly called the whip-lash effect.

D.1.3 Offset-locked synchronization

Another possible IEEE 1588 synchronization technique involves adding an offset value to the clock-slave device, where the value of that offset is based on the time differences sensed between the clock-master and clock-slave stations.

Constantly updated offsets ensures tracking of the clock-slave to the clock-master, without the response-time and whiplash effects normally associated with PLLs. However, since the clock rates remain unchanged, clock drifts can cause significant forward or backward jumps of the synchronized clock-slave timer. These discontinuities and transmit-time uncertainties can limit the accuracies of the slave-resident timer values.

D.2 Pacing alternatives

D.2.1 Higher level flow control

Higher layer protocols (such as the flow-control mechanisms of TCP) throttle the source to the bandwidth capabilities of the destination or intermediate interconnect. With the appropriate excess-traffic discards and rate-limiting recovery, such higher layer protocols can be effective in fairly distributing available bandwidth.

For real-time applications, however, the goal is to limit the number of talkers (so they can each have sufficient bandwidth), not to distribute the insufficient bandwidth fairly.

D.2.2 Over-provisioning

Over-provisioning involves using only a small portion of the available bandwidth, so that the cumulative bandwidth of multiple applications rarely exceeds that of the interconnect. This technique works well when frame losses are expected (voice over IP delays and gaps are similar to satellite-connected long distance phone calls) or when large levels of cumulative bandwidth ensure a tight statistical bound for maximum bandwidth utilization.

For most streaming applications within the home, however, frame losses are viewed as equipment defects (stutters in video or audio streams), which correspond to eventual loss of brand name values. Also, the existing kinds of transfers in a home (disk-to-disk, memory-to-display, tuner-to-display, multi-station games, etc.) do not (nor should not) have bandwidth limits.

D.2.3 Strict priorities

Existing networks can assign priority levels to different classes of traffic, effectively ensuring delivery of one before delivery of the other. One could provide the highest priority to the video traffic (with large bandwidth requirements), a high priority to the audio traffic (lower bandwidth, but critical), and the lowest priority level to file transfers. A typical number of priorities is eight.

Strict priority protocols are deficient in that the priorities are statically assigned, and the assignments (based on traffic class) often do not correspond to the desires of the consumer (my PBS show, rather than my teenager's games, perhaps). For example, priorities could result in transmission of two video streams, but not the audio associated with either.

Strict priority protocols usually assign fixed application-dependent priorities, assigning one priority to video and another to audio, for example. Mixed traffic (such as video streams with encapsulated audio) are not easily classified in this manner.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

D.3 IEEE 1394 alternative

Isochronous data transfers are well supported by the IEEE 1394 Serial Bus family of standards. This IEEE standards family (also called FireWire and iLink) is herein referred to simply as IEEE 1394.

Existing consumer equipment (digital camcorders, current generation high-definition televisions (HDTVs), digital video cassette recorders (DVCRs), digital video disk (DVD) recorders, set top boxes (STBs), and computer equipment intended for media authoring) support the IEEE 1394 interconnect. While some versions limit cable lengths to 4.5 meters, other physical layers support considerably longer lengths. A hub-like connection of IEEE 1394 devices supports seamless real-time services.

Although IEEE 1394 supports longer-reach physical layers, not all devices are compatible with these physical layers, or the distinct connectors associated with distinct physical layers. The RE protocols are based on Ethernet connections, a vast majority of which are based on 100 meter cables and the RJ-45 connector.

The IEEE 1394 isochronous packet addressing was designed with single-bus topologies in mind, which complicates the design of such bus bridges. The RE synchronous frames are designed with multiple stations and bridges in mind.

IEEE 1394 packets are differentiated by bus-local channel identifier, which must be allocated from a central per-bus resources and updated when isochronous packets pass through bridges. Mechanism must therefore be defined to agree upon the central per-bus resource, from among multiple available resources, and to renegotiate that agreement when any of the current central per-bus resources are removed.

Furthermore, absolute time stamps within some IEEE 1394 isochronous packets must be adjusted when passing through bridges. Such data-format dependent adjustments complicate bridge designs; their data-format dependent nature would most likely inhibit their successful adoption within an Ethernet bridge standard.

Annex E

(informative)

Time-of-day format considerations

To better understand the rationale behind the ‘extended binary’ timer format, other formats are evaluated and compared within this annex.

E.1 Possible time-of-day formats

E.1.1 Extended binary timer formats

The extended-binary timer format is used within this working paper and summarized herein. The 64-bit timer value consist of two components: a 32-bit *seconds* and 32-bit *fraction* fields, as illustrated in Figure 5.1.



Figure 5.1—Complete seconds timer format

The concatenation of 32-bit *seconds* and 32-bit *fraction* field specifies a 64-bit *time* value, as specified by Equation E.1.

$$time = seconds + (fraction / 2^{32}) \quad (E.1)$$

Where:

seconds is the most significant component of the time value (see Figure 5.1).

fraction is the less significant component of the time value (see Figure 5.1).

E.1.2 IEEE 1394 timer format

An alternate “1394 timer” format consists of *secondCount*, *cycleCount*, and *cycleOffset* fields, as illustrated in Figure E.2. For such fields, the 12-bit *cycleOffset* field is updated at a 24.576MHz rate. The *cycleOffset* field goes to zero after 3171 is reached, thus cycling at an 8kHz rate. The 13-bit *cycleCount* field is incremented whenever *cycleOffset* goes to zero. The *cycleCount* field goes to zero after 7999 is reached, thus restarting at a 1Hz rate. The remaining 7-bit *secondCount* field is incremented whenever *cycleCount* goes to zero.

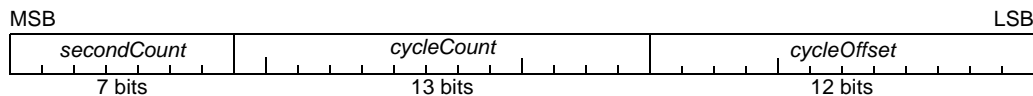


Figure E.2—IEEE 1394 timer format

E.1.3 IEEE 1588 timer format

IEEE 1588 timer format consists of seconds and nanoseconds fields components, as illustrated in Figure E.3. The nanoseconds field must be less than 10^9 ; a distinct *sign* bit indicates whether the time represents before or after the epoch duration.

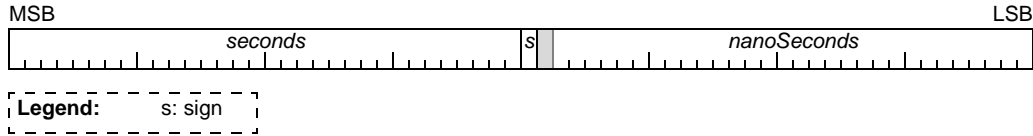


Figure E.3—IEEE 1588 timer format

E.1.4 EPON timer format

The IEEE 802.3 EPON timer format consists of a 32-bit scaled nanosecond value, as illustrated in Figure E.4. This clock is logically incremented once each 16 ns interval.

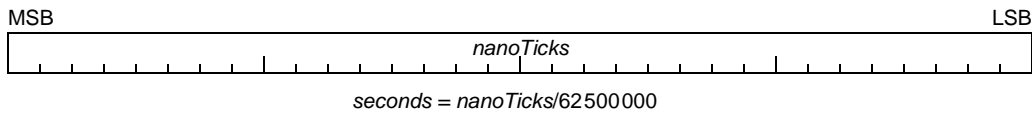


Figure E.4—EPON timer format

E.1.5 Compact seconds timer format

An alternate “compact seconds” format could consist of 8-bit *seconds* and 24-bit *fraction* fields, as illustrated in Figure E.5. This would provided similar resolutions to the IEEE 1394 timer format, without the complexities associated with its binary coded decimal (BCD) like encoding.

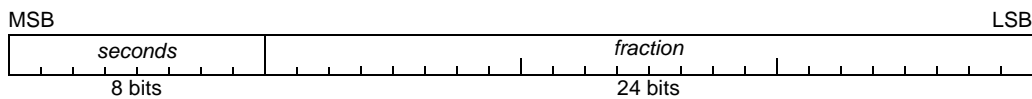


Figure E.5—Compact seconds timer format

E.1.6 Nanosecond timer format

An alternate “nanosecond” format could consists of 2-bit *seconds* and 30-bit *nanoSeconds* fields, as illustrated in Figure E.6.

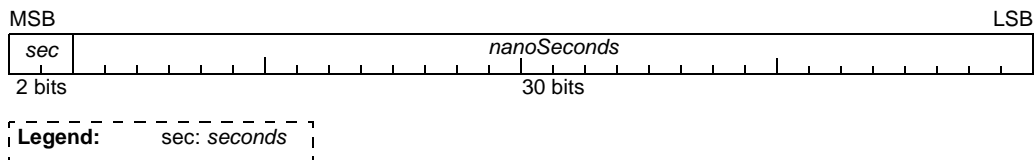


Figure E.6—Nanosecond timer format

E.2 Time format comparisons

To better understand the relative benefits of different time formats, the relevant properties are summarized in Table E.1. Counter complexity is not included in the comparison, since the digital logic complexity (see 7.1.11) is comparable for all formats.

Table E.1—Time format comparison

Name	Subclause	Range	Precision	Arithmetic	Seconds	Defined standards
Column	—	1	2	3	4	5
extended binary	TBD	136 years	232 ps	Good	Good	RFC 1305 NTP, RFC 2030 SNTIPv4
IEEE 1394	E.1.2	128 s	30 ns	Poor	Good	IEEE 1394
IEEE 1588	E.1.3	272 years	1 ns	Fair	Good	IEEE 1588
IEEE 802 (EPON)	E.1.4	69 s	16 ns	Good	Poor	IEEE 802.3
compact seconds	E.1.5	256 s	60 ns	Best	Good	—
nanoseconds	E.1.6	4 s	1 ns	Best	Poor	—

Column 1: A desirable property is the support of a wide range of second values, to eliminate the need for defining/coordinating/implementing auxiliary seconds-synchronization protocols. The 136-year range of the extended binary format is sufficient for this purpose.

Column 2: A desirable property is a fine-grained resolution, sufficient to measure each bit-transmission times. The ‘extended binary’ provides the most precision; exceeds the resolution of expected cost-effective time-capture circuits.

Column 3: Computation of time differences involves the subtraction of two timer-snapshot values. Subtraction of ‘extended binary’ numbers involving standard 64-bit binary arithmetic; no special field-overflow compensations are required. Only the less precise ‘compact seconds’ and nanoseconds formats are simpler, due to the reduced 32-bit size of the timer values.

Column 4: Time values must oftentimes be compared to externally provided values (e.g., timers extracted from GPS or stratum-clock sources). For these purposes, the availability of a seconds component is desired. The ‘extended binary’ format provides a seconds component that can be easily extracted or such purposes.

Annex F

(informative)

Bursting and bunching considerations

F.1 Topology scenarios

F.1.1 Bridge design models

The sensitivity of bridges to bursting and bunching is highly dependent on the queue management protocols within the bridge. To better understand these effects, a few bridge design models are evaluated, as illustrated in Figure F.1.

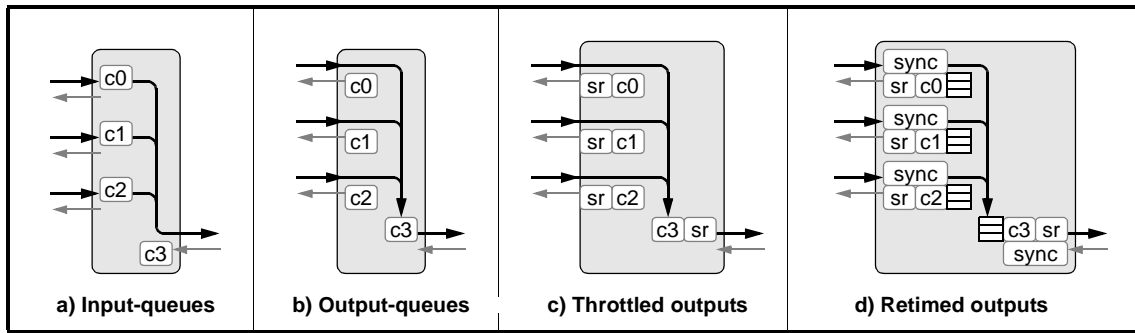


Figure F.1—Bridge design models

The input-queue design (see Figure F.1-a) assumes that frames are queued in receive buffers. The transmitter accepts frames from the receivers, based on service-class precedence. In the case of a tie (two receivers can provide same-class frames), the lowest numbered receive port has precedence. This model best illustrates nonlinear bunching problems.

The output-queue design (see Figure F.1-b) assumes that received frames are queued in transmit buffers. Within each service class, frames are forwarded in FIFO order. This model best illustrates linear bunching problems (for steady flows), but also exhibits nonlinear bunching (for nonsteady flows).

The throttled-output design (see Figure F.1-c) is an enhanced output-queue model, with an output shaper to limit transmission rates. The purpose of the output shaper is to ensure sufficient nonreserved bandwidth for less time-sensitive control and monitoring purposes. The model illustrates how shapers can worsen the output-queue bridge's bunching behaviors.

The retimed-outputs design (see Figure F.1-d) reduces (and can eliminate) bunching problems by detecting late-arrival frames at the receivers. Several synchronous-cycle buffers are provided at the transmitters, to compensate for transmission delays in the received data.

F.1.2 Three-source hierarchical topology

A hierarchical topology best illustrate potential problems with bunching, as illustrated in Figure F.2. Traffic from talkers {a0,a1,a2} flows into bridge B. Bridge B concentrates traffic received from three talkers, with the cumulative b3 traffic sent to c3. Identical traffic flows are assumed at bridge ports {c0,c1,c3}, although only one of these sources is illustrated. Bridges {C,D,E,F,G,H} behave similarly.

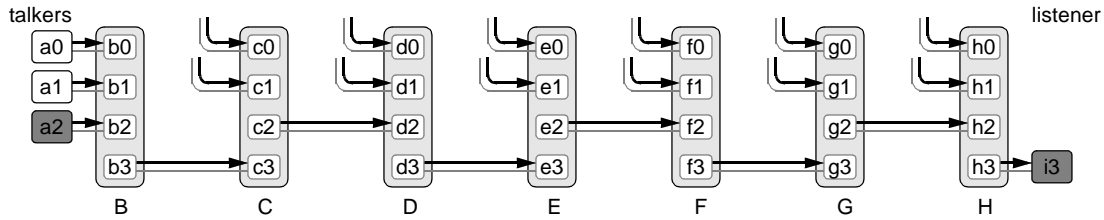


Figure F.2—Three-source topology

F.1.3 Six-source hierarchical topology

Spreading the traffic over multiple sources, as illustrated in Figure F.3, exasperates bursting and bunching problems. Traffic from talkers {a0,a1,a2,a3,a4,a5} flows into ports on bridge B. Bridge B concentrates traffic received from six talkers, with the cumulative b6 traffic sent to c6. Identical traffic flows are assumed at bridge ports {c0,c1,c3,c3,c4,c6}, although only one of these sources is illustrated. Bridges {C,D,E,F,G,H} behave similarly.

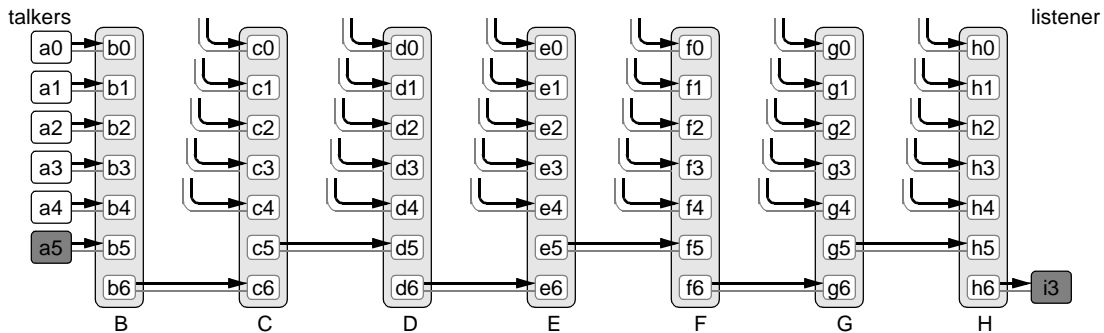


Figure F.3—Six-source topology

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

F.2 Bursting considerations

F.2.1 Three-source bursting scenario

A troublesome bursting scenario on a 100 Mb/s link can occur when small bandwidth streams coincidentally provide their infrequent 1500 byte frames concurrently, as illustrated in Figure F.4. Even though the cumulative bandwidths are considerably less than the capacity of the 100 Mb/s links, significant delays are incurred when passing through multiple bridges.

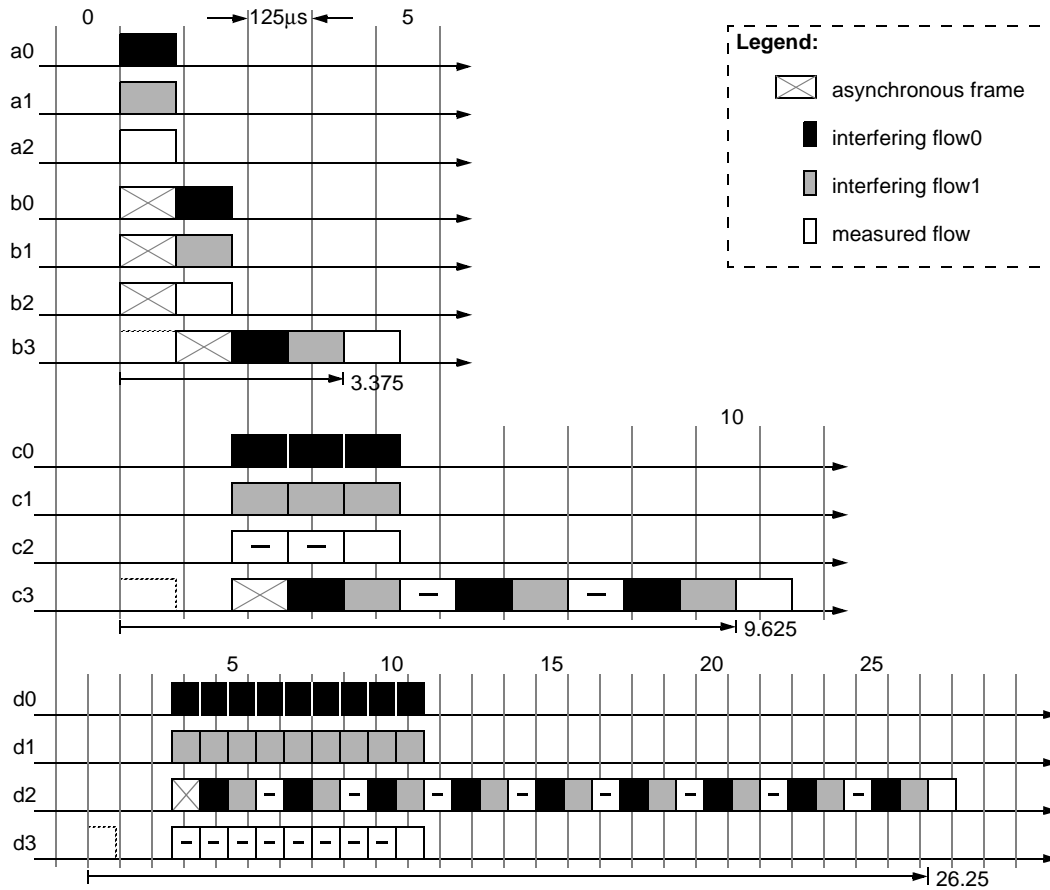


Figure F.4—Three-source bunching timing; input-queue bridges

F.2.1.1 Cumulative bunching latencies

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.1 and plotted in Figure F.5.

Table F.1—Cumulative bursting latencies

Topology	Units	Measurement point							
		A	B	C	D	E	F	G	H
3-source (see F.2.2.1)	mtu	1	4	11	30	85	248	735	2194
	ms	.120	.480	1.32	3.6	10.2	29.6	88.2	263
6-source (see F.2.2.2)	mtu	1	7	38	219	1300	7781	46662	229943
	ms	.120	.840	4.56	26.3	156	934	5600	27600

The values within this table are computed based on Equation F.1.

$$delay[n] = mtu \times (n + p^n) \tag{F.1}$$

Where:

- mtu* (maximum transfer unit) is the maximum frame size
- n* is the number of hops from the source
- p* is the number of receive ports in each bridge.

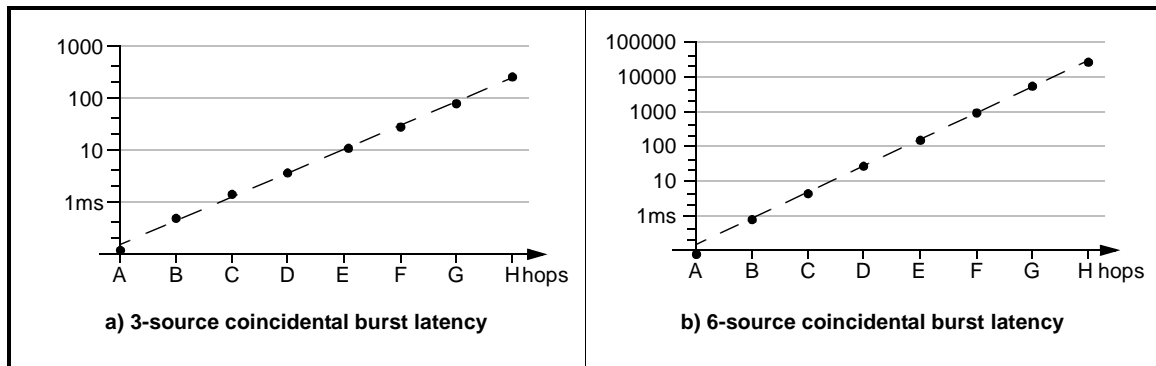


Figure F.5—Cumulative coincidental burst latencies

Conclusion: The classA traffic bandwidths should be enforced over a time interval that is on the order of an MTU size (120μs), so as to avoid excessive delays caused by coincidental back-to-back large-block transmissions.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

F.2.2 Bunching scenarios; input-queue bridges

F.2.2.1 Three-source bunching; input-queue bridges

To illustrate the effects of worst case bunching on input-queue bridges, specific flows are illustrated in Figure F.6. Bridge ports {b0,b1,b2} concentrates traffic from three talkers; one third of the cumulative traffic is forwarded through b3. Each stream consumes 25% of the link bandwidth; 25% is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c3},...,{e0,e1,e3}, only illustrate the passing-through listener traffic; the remainder of the traffic is assumed to be routed elsewhere.

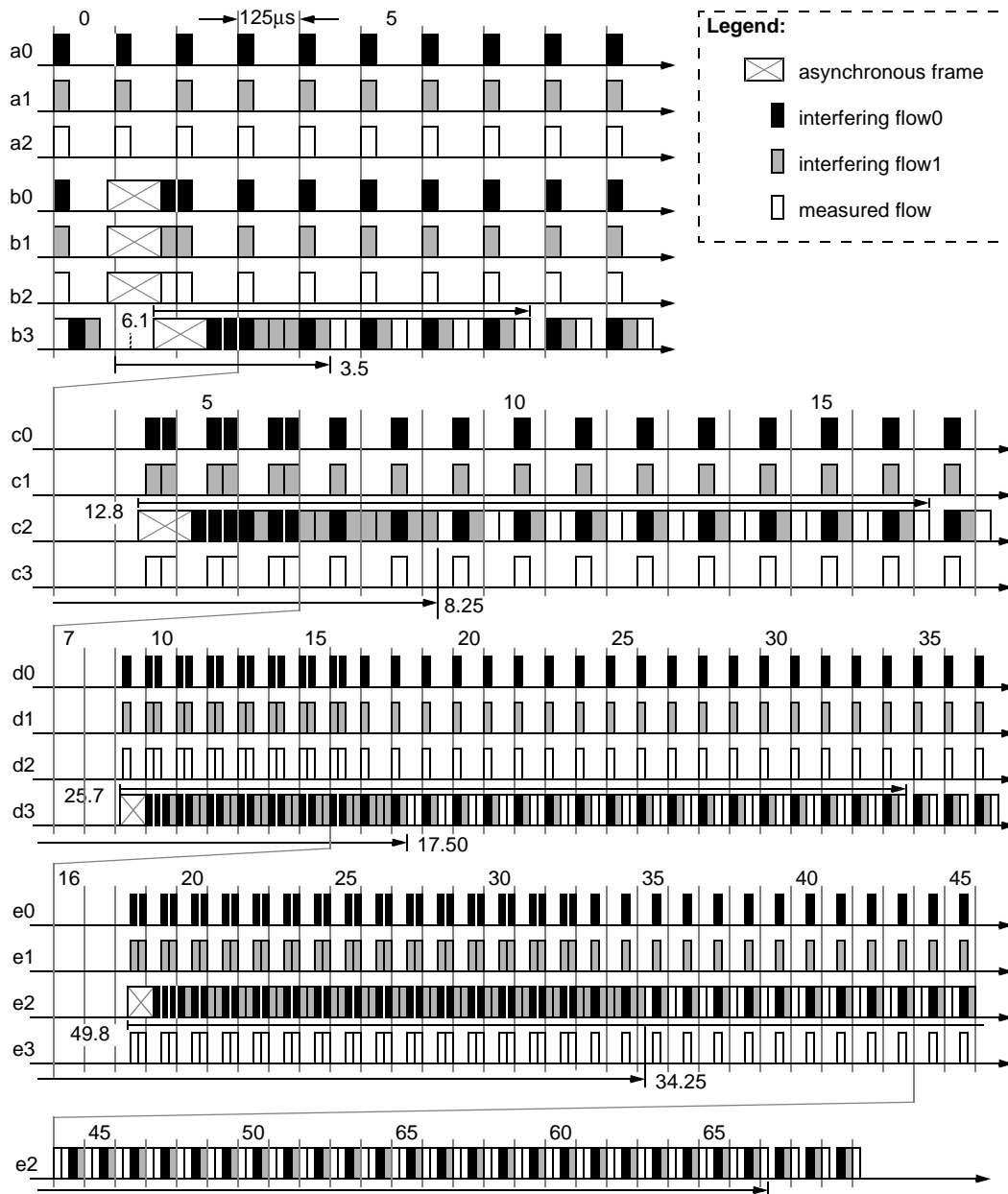


Figure F.6—Three-source bunching; input-queue bridges

F.2.2.2 Six-source bunching; input-queue bridges

To better illustrate the effects of worst case bunching on input-queue bridges, specific flows are illustrated in Figure F.7. Bridge ports {b0,b1,b2,b3,b4,b5} concentrates traffic from three talkers; one sixth of the cumulative traffic is forwarded through b6. Each of six streams consumes 12.5% of the link bandwidth, so that 25% is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c2,c3,c4,c6} only illustrate passing-through traffic; the remainder of the traffic is routed elsewhere.

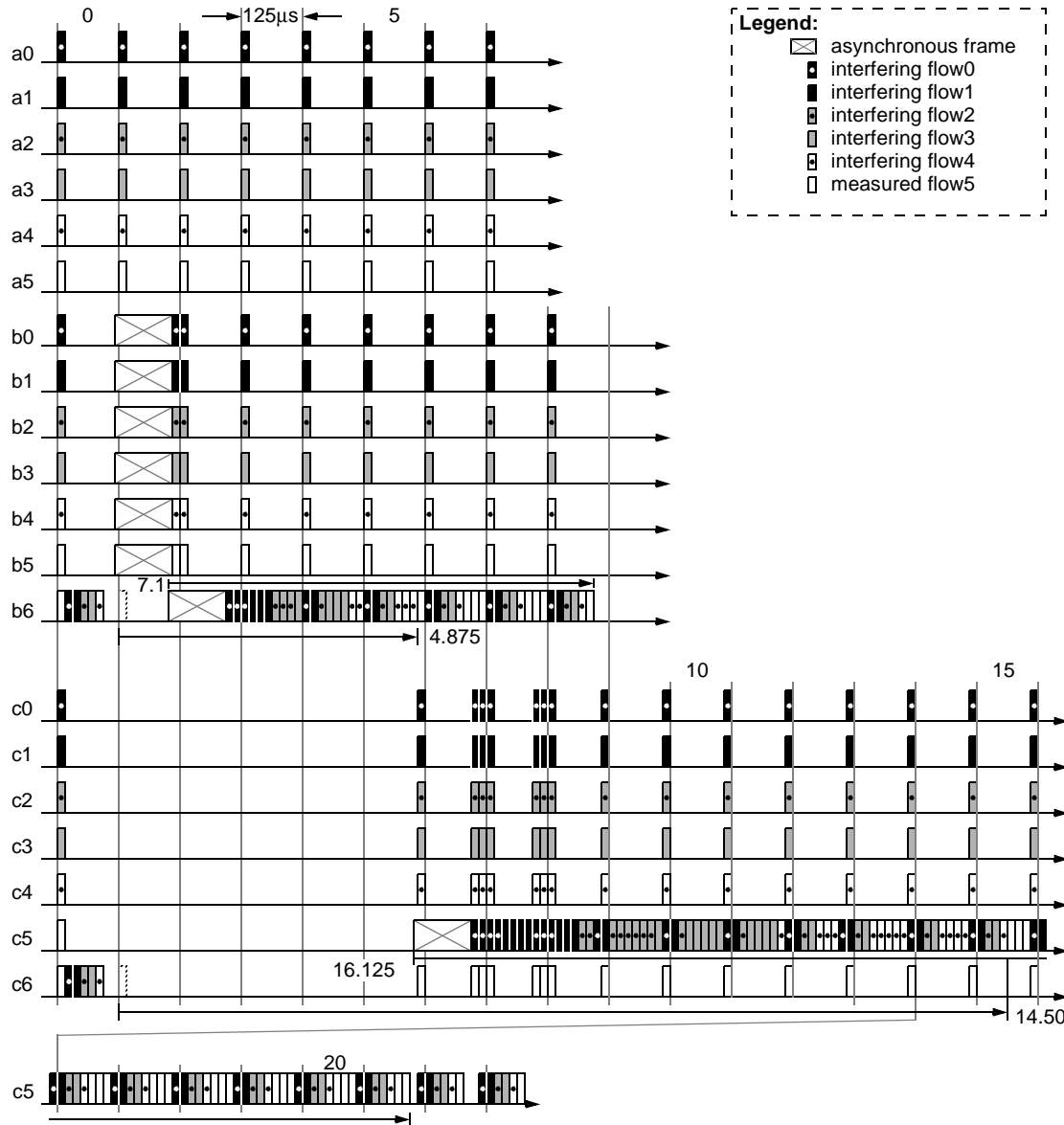


Figure F.7—Six source bunching timing; input-queue bridges

F.2.2.3 Cumulative bunching latencies, input-queue bridge

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.2 and plotted in Figure F.8.

Table F.2—Cumulative bunching latencies; input-queue bridge

Topology	Units	Measurement point							
		A	B	C	D	E	F	G	H
3-source (see F.2.2.1)	cycles	0.125	3.5	8.25	17.5	34.25	(70.75)	(143.2)	(288.2)
	ms	0.01	0.44	1.03	2.19	4.28	8.84	17.9	36.0
6-source (see F.2.2.2)	cycles	0.125	4.875	14.50	(39.33)	(107.2)	(288.2)	(771)	2058
	ms	0.01	0.61	1.81	4.92	13.4	36.0	96.4	257

The first few numbers are generated using graphical techniques, as illustrated in Figure F.2.2.2. The following numbers are estimated, based on Equation F.2.

$$delay[n+1] = (mtu + delay[n]) \times (1 / (1 - 0.75 \times (p-1) / p)) \tag{F.2}$$

Where:

- mtu* (maximum transfer unit) is the maximum frame size
- rate* is the fraction of the bandwidth reserved for class A traffic, assumed to be 0.75
- n* is the number of hops from the source
- p* is the number of receive ports in each bridge.

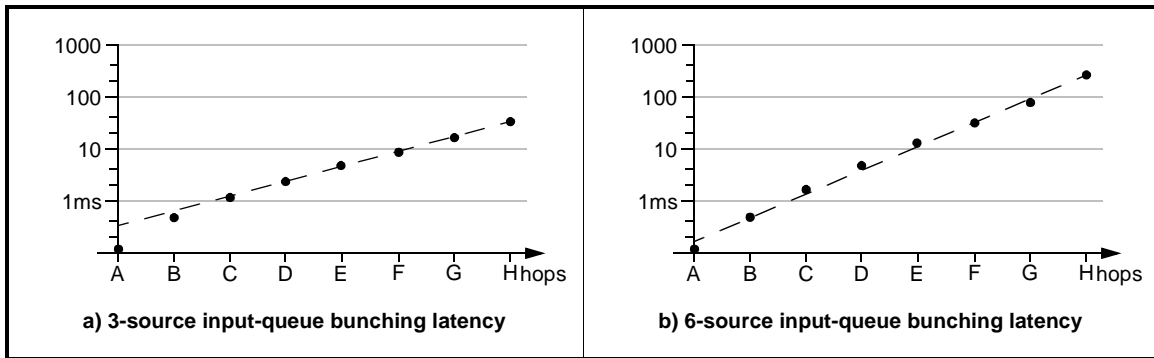


Figure F.8—Cumulative bunching latencies; input-queue bridge

Conclusion: A FIFO based output-queue bridge should be used. Alternatively (if input queuing is used), received frames should be time-stamped to ensure FIFO like forwarding.

F.2.3 Bunching topology scenarios; output-queue bridges

F.2.3.1 Three-source bunching timing; output-queue bridges

To illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.9. Bridge ports {b0,b1,b2} concentrates traffic from three talkers; one third of the cumulative traffic is forwarded through b3. Each stream consumes 25% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {b0,b1,b2},...,{e0,e1,e3} only illustrate the passing-through listener traffic; the remainder of the traffic is assumed to be routed elsewhere.

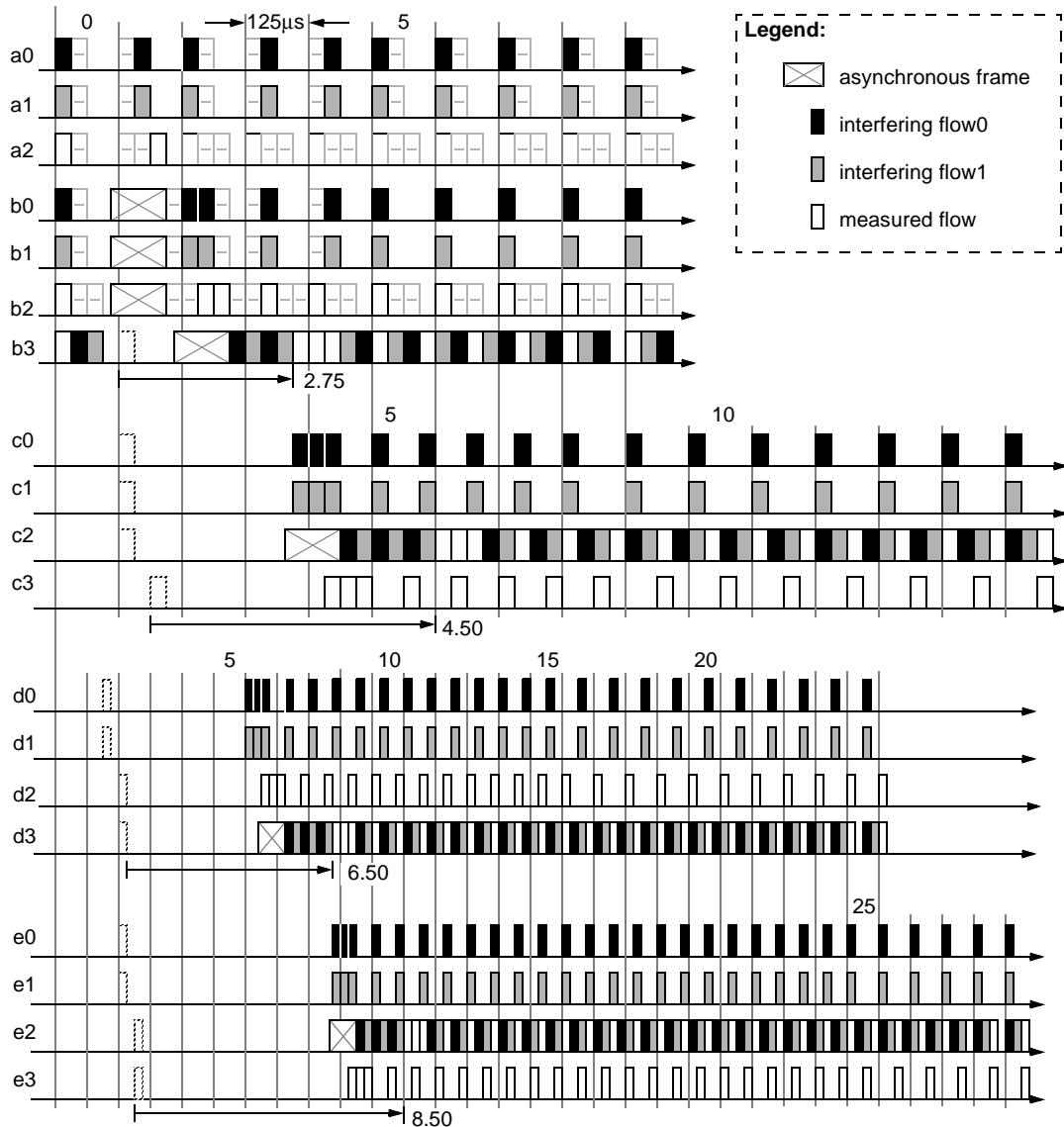


Figure F.9—Three-source bunching; output-queue bridges

F.2.3.2 Six-source bunching; output-queue bridges

To better illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.10. Bridge ports {b0,b1,b2,b3,b4,b5} concentrates traffic from six talkers; one sixth of the cumulative traffic is forwarded through port b6. Each of six streams consumes 12.5% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c2,c3,c4,c6} and {d0,d1,d2,d3,d4,d5} only illustrate passing-through traffic; the remainder of the traffic is routed elsewhere.

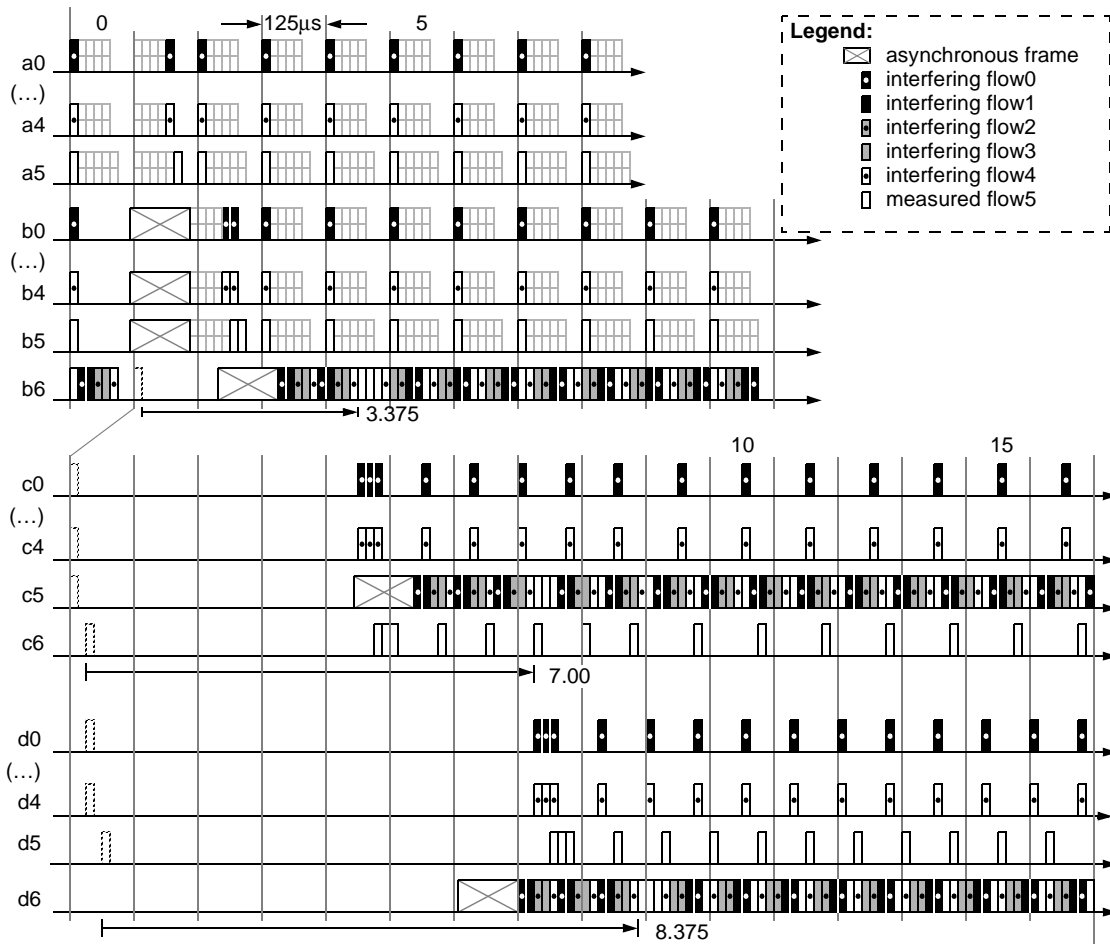


Figure F.10—Six source bunching; output-queue bridges

F.2.3.3 Cumulative bunching latencies; output-queue bridge

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.3 and plotted in Figure F.11.

Table F.3—Cumulative bunching latencies; output-queue bridge

Topology	Units	Measurement point							
		B	C	D	E	F	G	H	I
3-source (see F.2.2.1)	cycles	.875	2.75	4.5	6.5	8.5	–	–	–
	ms	0.10	0.34	0.56	0.81	1.6	–	–	–
6-source (see F.2.2.2)	cycles	.875	3.375	7.00	8.375	–	–	–	–
	ms	0.10	0.42	.875	1.05	–	–	–	–

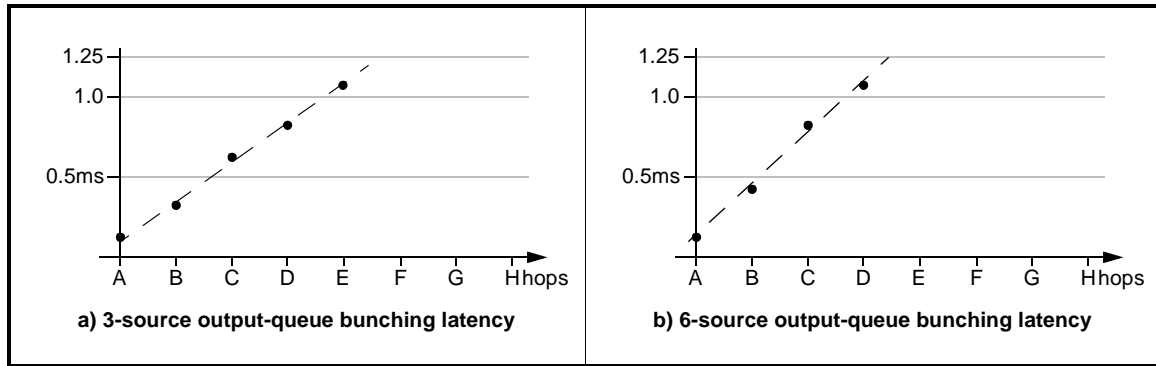


Figure F.11—Cumulative bunching latencies; output-queue bridge

Conclusion: For steady-state classA traffic, acceptably small linear latencies are introduced by output-queue bridges on 75% loaded links. Unfortunately, the nonsteady-state nature of variable-rate traffic makes this conclusion suspect (see F.2.4).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

F.2.4 Bunching topology scenarios; variable-rate output-queue bridges

F.2.4.1 Three-source bunching; variable-rate output-queue bridges

To illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.12. Bridge ports {b0,b1,b2} concentrates traffic from three talkers; one third of the cumulative traffic is forwarded through port b3. Each stream consumes 25% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c3},...,{e0,e1,e3} only illustrate the passing-through listener traffic; the remainder of the traffic is assumed to be routed elsewhere.



Figure F.12—Three-source bunching; variable-rate output-queue bridges

F.2.4.2 Six-source bunching; variable-rate output-queue bridges

To better illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.13. Bridge ports {b0,b1,b2,b3,b4,b5} concentrates traffic from six talkers; one sixth of the cumulative traffic is forwarded through port b6. Each of six streams consumes 12.5% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c2,c3,c4,c6}, {d0,d1,d2,d3,d4,d5}, and {e0,e1,e2,e3,e4,e6} only illustrate passing-through traffic; the remainder of the traffic is routed elsewhere.

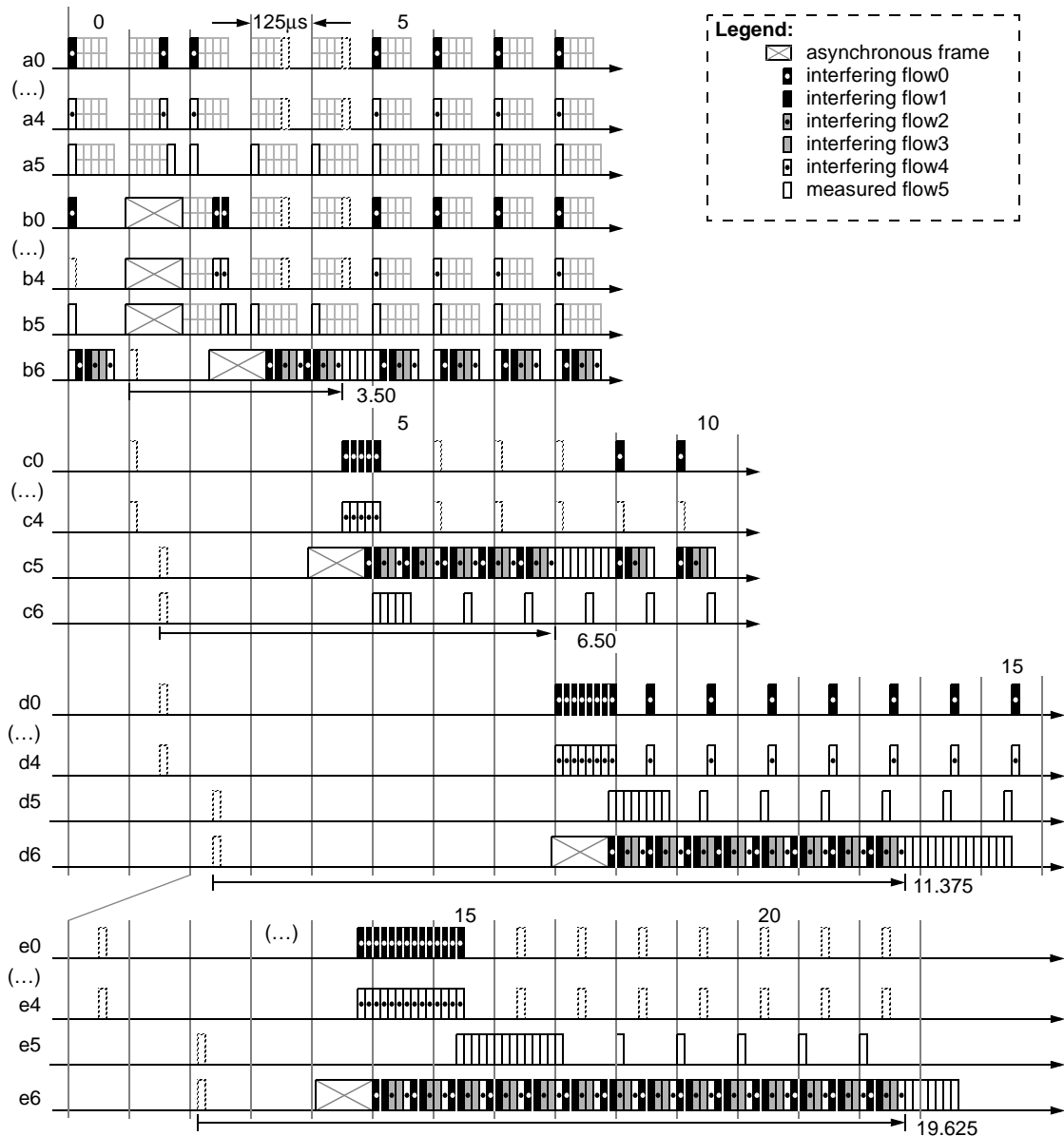


Figure F.13—Six source bunching; variable-rate output-queue bridges

F.2.4.3 Cumulative bunching latencies; variable-rate output-queue bridge

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.4 and plotted in Figure F.14.

Table F.4—Cumulative bunching latencies; variable-rate output-queue bridge

Topology	Units	Measurement point							
		A	B	C	D	E	F	G	H
3-source (see F.2.2.1)	cycles	0.75	2.75	4.75	7.25	10.75	–	–	–
	ms	0.10	0.34	0.59	0.90	1.34	–	–	–
6-source (see F.2.2.2)	cycles	0.75	3.50	6.50	11.38	19.63	–	–	–
	ms	0.10	0.44	0.81	1.42	2.45	–	–	–

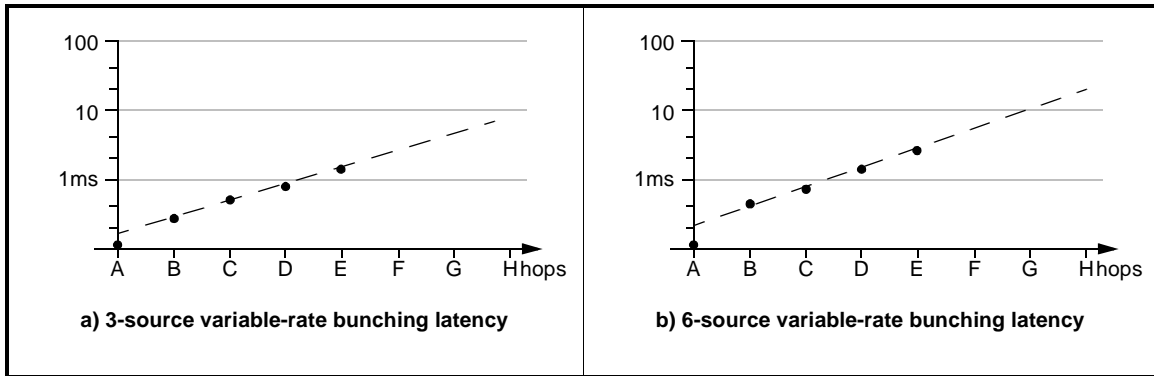


Figure F.14—Cumulative bunching latencies; variable-rate output-queue bridge

Conclusion: For nonsteady-state classA traffic, significant expedient latencies are introduced by output-queue bridges on 75% loaded links. Unfortunately, throttled outputs further exasperates this latency (see F.2.4).

F.2.5 Bunching topology scenarios; throttled-rate output-queue bridges

F.2.5.1 Three-source bunching; throttled-rate output-queue bridges

To illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.15. Bridge ports {b0,b1,b2} concentrates traffic from three talkers; one third of the cumulative traffic is forwarded through port b3. Each stream consumes 25% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c3}, {d0,d1,d2}, and {e0,e1,e3} only illustrate the passing-through listener traffic; the remainder of the traffic is assumed to be routed elsewhere.

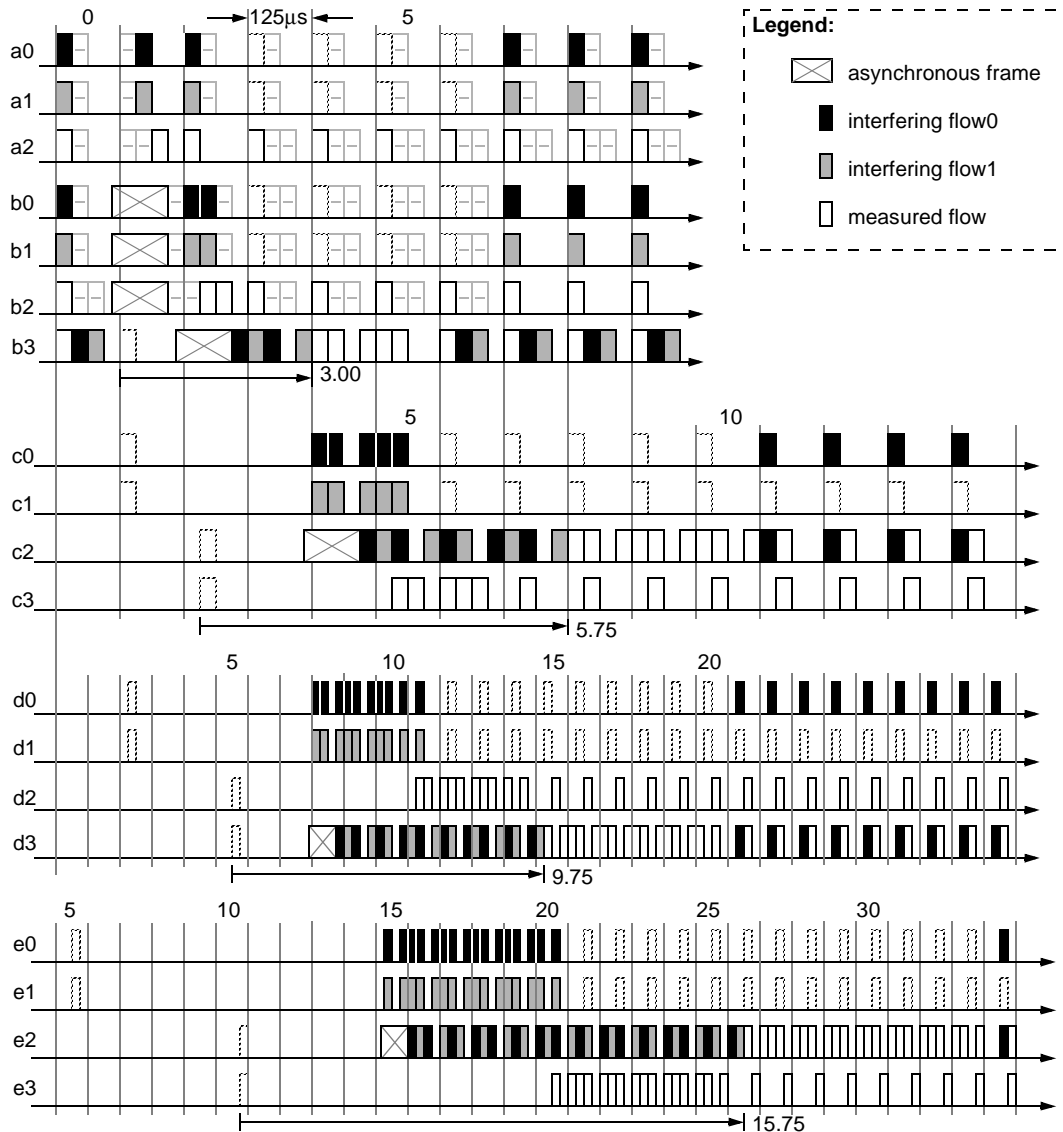


Figure F.15—Three-source bunching; throttled-rate output-queue bridges

F.2.5.2 Six-source bunching; throttled-rate output-queue bridges

To better illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.16. Bridge ports {b0,b1,b2,b3,b4,b5} concentrates traffic from six talkers; one sixth of the cumulative traffic is forwarded through port b6. Each of six streams consumes 12.5% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c2,c3,c4,c5},...,{e0,e1,e2,e3,e4,e6} only illustrate passing-through traffic; the remainder of the traffic is routed elsewhere.

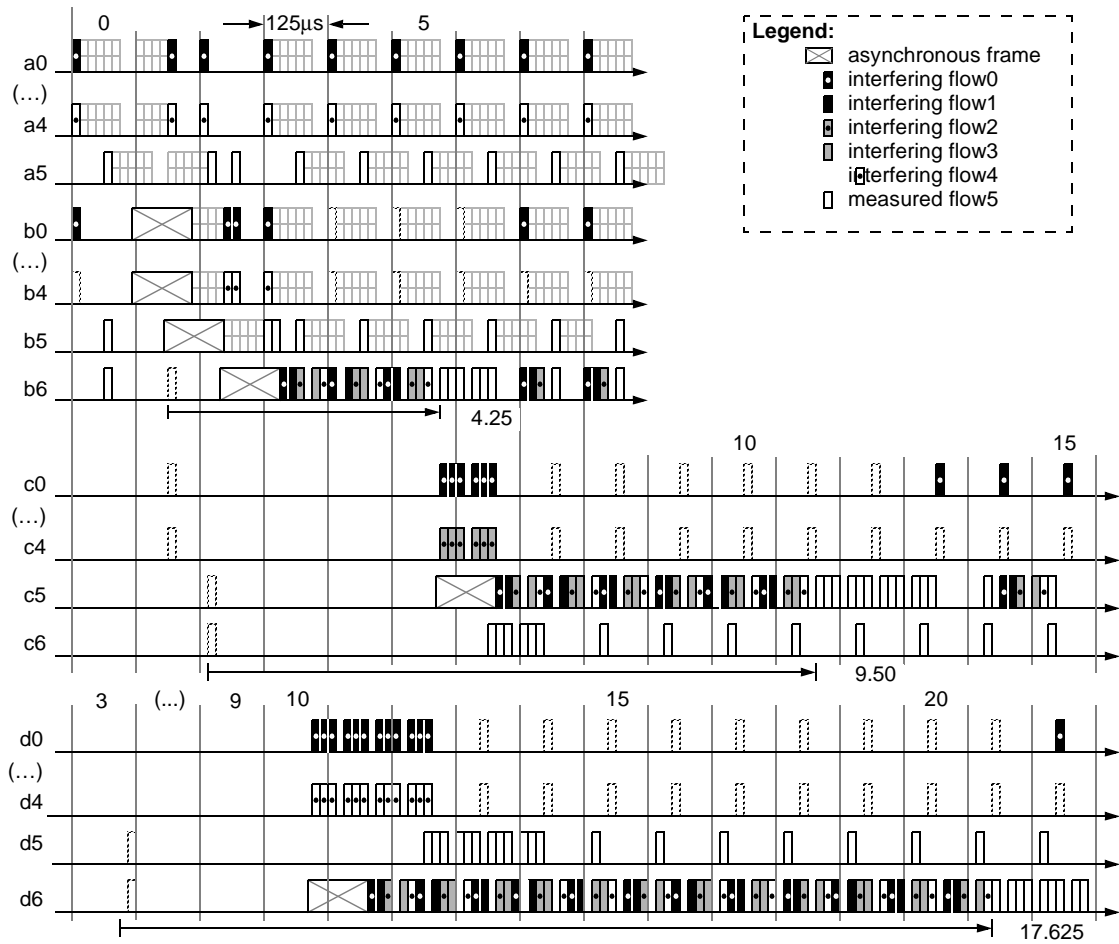


Figure F.16—Six source bunching; throttled-rate output-queue bridges

F.2.5.3 Cumulative bunching latencies; throttled-rate output-queue bridge

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.5 and plotted in Figure F.17.

Table F.5—Cumulative bunching latencies; throttled-rate output-queue bridge

Topology	Units	Measurement point							
		A	B	C	D	E	F	G	H
3-source (see F.2.2.1)	cycles	0.75	3.00	5.75	9.75	15.75	–	–	–
	ms	0.09	0.38	0.73	1.21	1.97	–	–	–
6-source (see F.2.2.2)	cycles	0.75	4.25	9.5	17.63	–	–	–	–
	ms	0.09	0.53	1.19	2.20	–	–	–	–

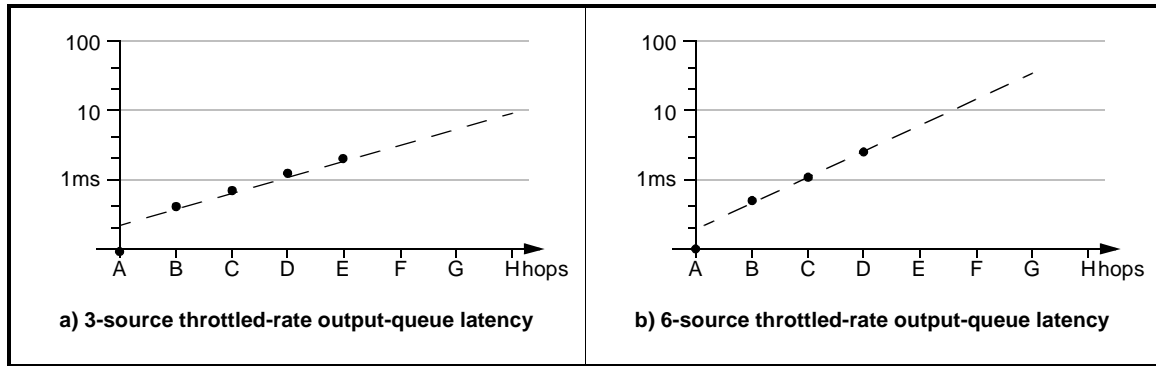


Figure F.17—Cumulative bunching latencies; throttled-rate output-queue bridge

Conclusion: On large topologies, the classA traffic latencies can accumulate beyond acceptable limits. Some form of receiver retiming may therefore be desired.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

F.2.6 Bunching topology scenarios; classA throttled-rate output-queue bridges

The extent of bunching extent is worst when large classC frames are present. However, bunching can also occur in the absence of large classC frames, as described in the remainder of this subannex.

F.2.6.1 Three-source bunching; classA throttled-rate output-queue bridges

To illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.18 and Figure F.19. Bridge ports {b0,b1,b2} concentrates traffic from three talkers; one third of the cumulative traffic is forwarded through port b3. Each stream consumes 25% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c3}, {c0,d1,d2}, and {e0,e1,e3} only illustrate the passing-through listener traffic; the remainder of the traffic is assumed to be routed elsewhere.

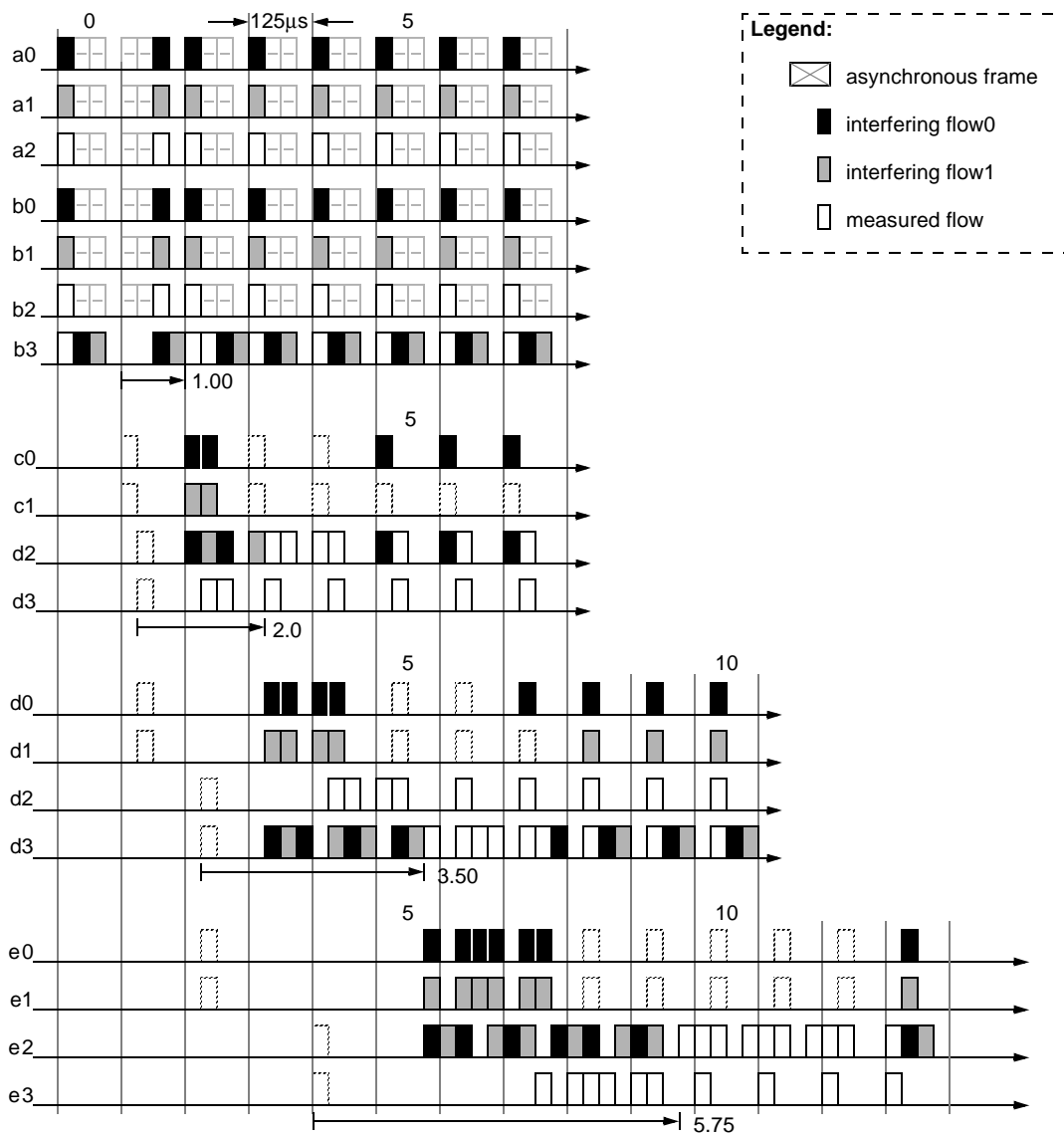


Figure F.18—Three-source bunching; throttled-rate output-queue bridges

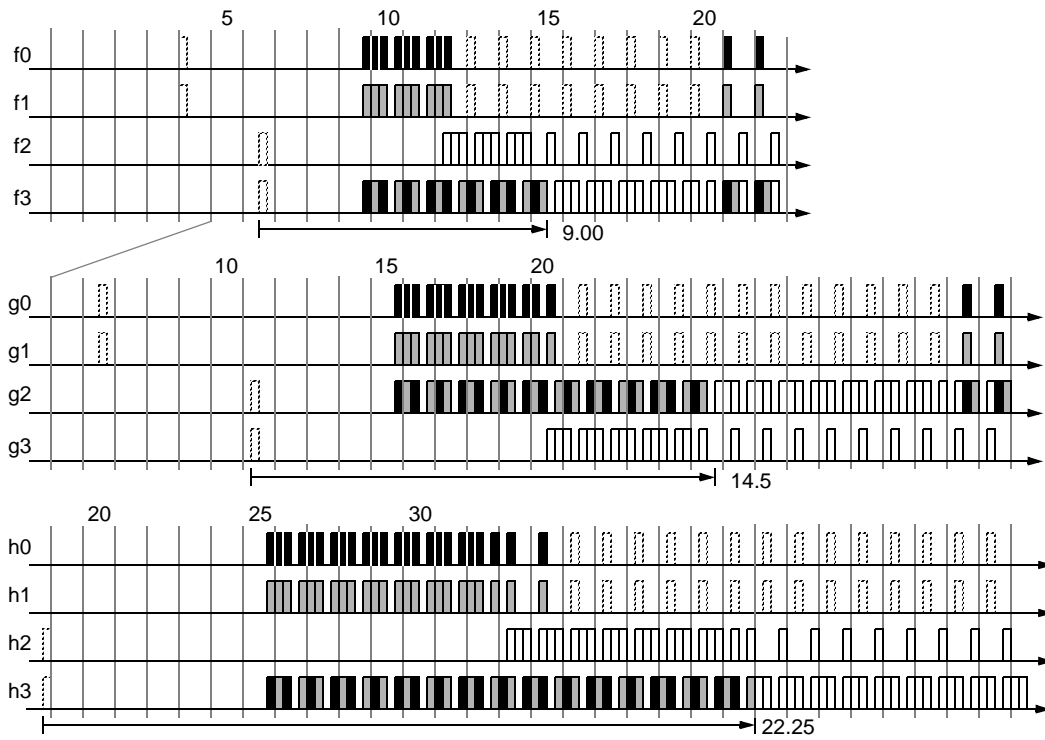


Figure F.19—Three-source bunching; throttled-rate output-queue bridges

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

F.2.6.2 Six-source bunching; classA throttled-rate output-queue bridges

To better illustrate the effects of worst case bunching, specific flows are illustrated in Figure F.20. Bridge ports {b0,b1,b2,b3,b4,b5} concentrates traffic from six talkers; one sixth of the cumulative traffic is forwarded through port b6. Each of six streams consumes 12.5% of the link bandwidth; 25% of the link bandwidth is available for asynchronous traffic.

For clarity, the traces for input traffic on ports {c0,c1,c2,c3,c4,c5},...,{d0,d1,d2,d3,d4,d6} only illustrate passing-through traffic; the remainder of the traffic is routed elsewhere.

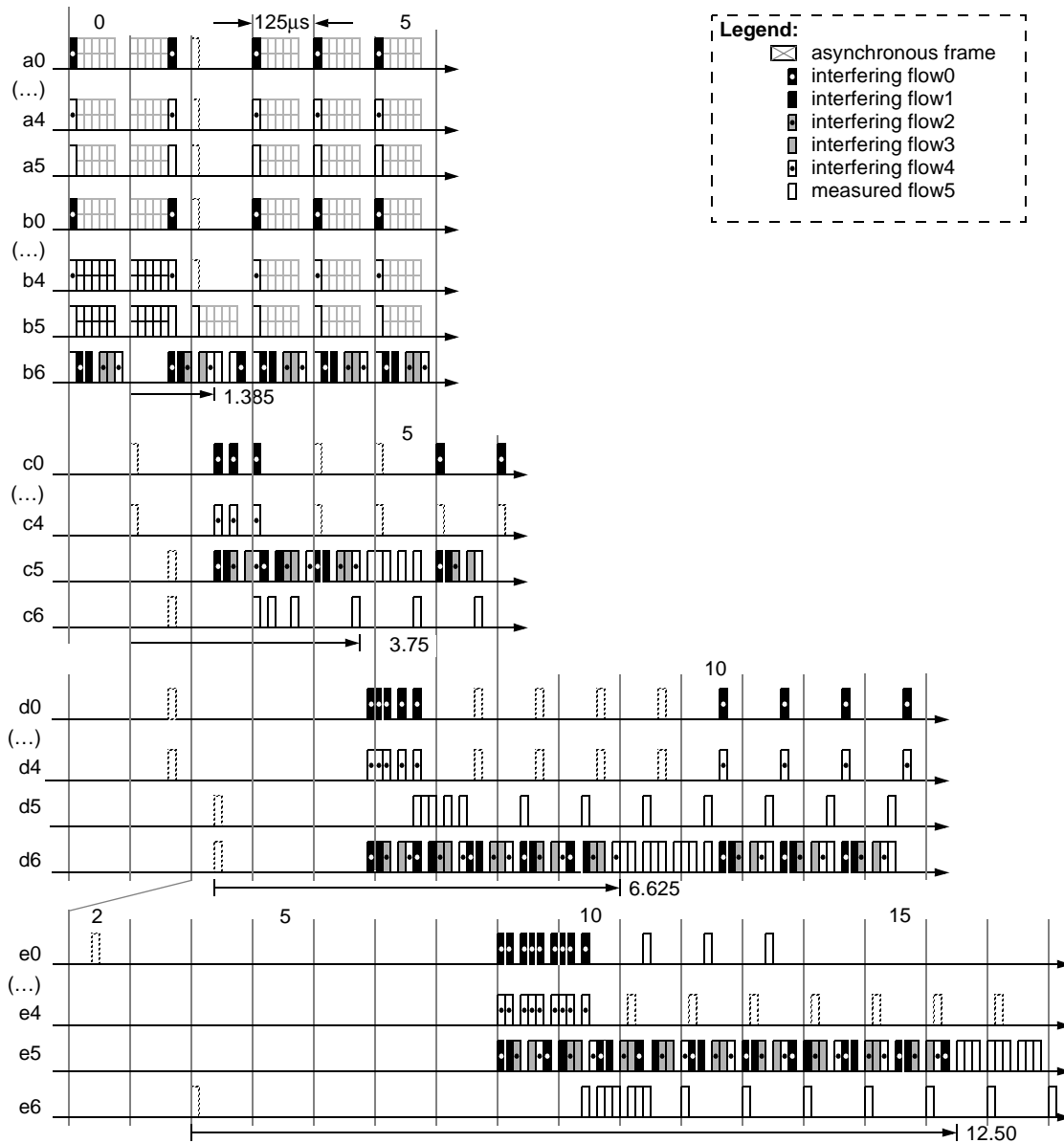


Figure F.20—Six source bunching; classA throttled-rate output-queue bridges

F.2.6.3 Cumulative bunching latencies; classA throttled-rate output-queue bridge

The cumulative worst-case latencies implied by coincidental bursting are listed in Table F.6 and plotted in Figure F.21.

Table F.6—Cumulative bunching latencies; classA throttled-rate output-queue bridge

Topology	Units	Measurement point							
		A	B	C	D	E	F	G	H
3-source (see F.2.2.1)	cycles	–	1.00	2.00	3.5	5.75	9.00	14.5	22.5
	ms	–	0.125	0.25	0.44	0.72	1.13	1.81	2.81
6-source (see F.2.2.2)	cycles	–	1.385	3.75	6.625	12.50	–	–	–
	ms	–	0.17	0.47	0.83	1.56	–	–	–

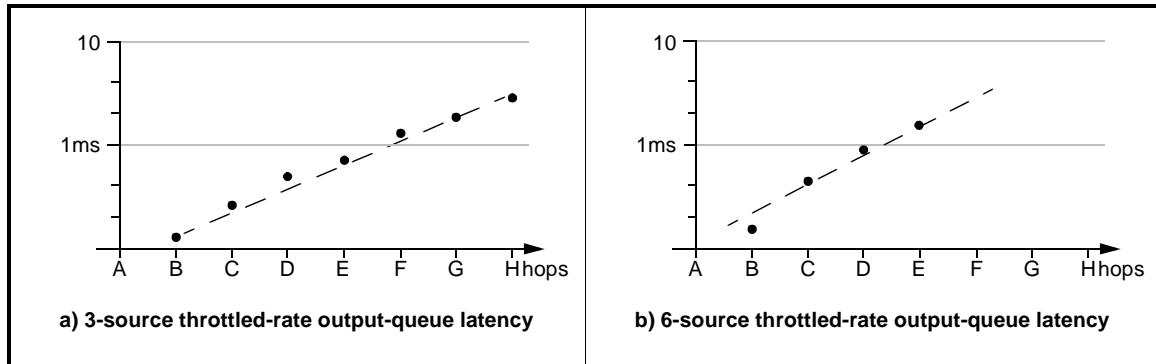


Figure F.21—Cumulative bunching latencies; classA throttled-rate output-queue bridge

Conclusion: On large topologies, the classA traffic latencies can accumulate beyond acceptable limits, even in the absence of conflicting lower-class traffic. Some form of receiver retiming may therefore be desired, even on higher speed links where the size of the MTU (in time) becomes much smaller than an assumed 8 kHz cycle time.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Annex G

(informative)

Denigrated alternatives

G.1 Stream frame formats

NOTE—The following streaming classA frame format options were considered but rejected. These options are retained for historical purposes and (if opinions change) possible reconsideration. For these reasons, the perceived advantages and disadvantages of each technique are listed.

G.1.1 VLAN routed frame formats (alternative 4)

Frames within a stream are no different than other Ethernet frames, with the exception of their distinct *da* (destination address) and *control* field values, as illustrated in Figure G.1.

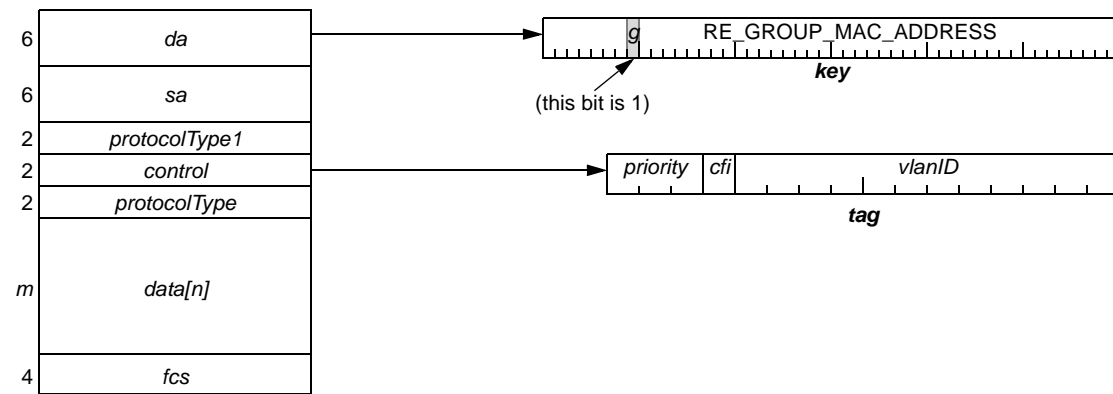


Figure G.1—classA frame formats

A single multicast address (labeled as *RE_GROUP_MAC_ADDRESS*) identifies the multicast time-sensitive nature of the frame. The following VLAN tag identifies the frame priority and provides a distinct *vlanID* identifier. The *vlanID* identifier is also the *streamID* identifier, allowing each stream to be independently selectively-switched through bridges.

The over-riding disadvantages of this design approach relates to its forwarding through bridges:

- a) Overloaded. This novel *vlanID* usage could conflict with existing bridge implementations.
- b) VLAN service. A method of generating distinct *vlanID* values would be required. (Some form of central server or distributed assignment algorithm would be required).

G.1.2 Broadcast routed frame formats (alternative5)

Frames within a stream are no different than other Ethernet frames, with the exception of their distinct fixed multicast *da* (destination address), as illustrated in Figure G.2.

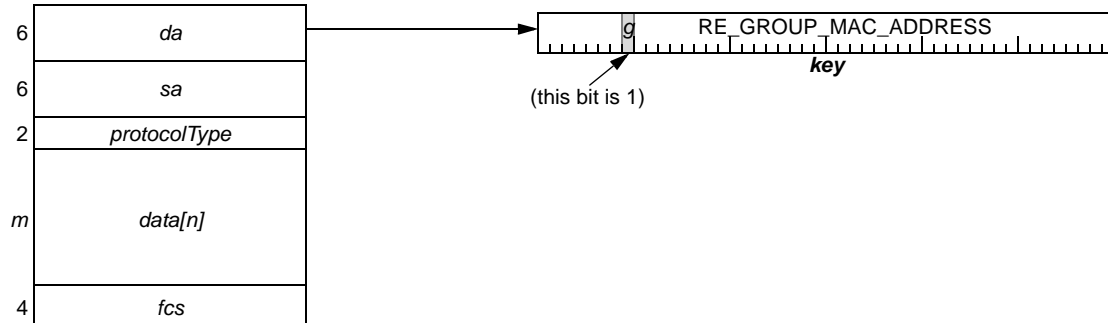


Figure G.2—ClassA frame formats

A single multicast address (labeled as RE_GROUP_MAC_ADDRESS) identifies the multicast time-sensitive nature of the frame.

The over-riding disadvantages of this design approach relates to its forwarding through bridges:

- a) Bandwidth. Bandwidth is wasted because frames are broadcast to all potential listeners, rather than only the subscribed listeners.
- b) Ambiguous. The *da* field is insufficient to identify the frame, mandating the presence of stream identifier information within the *data[]* payload.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

G.2 Subscription

G.2.1 Simple Reservation Protocol (SRP) overview

Subscription involves explicit negotiation for bandwidth resources, performed in a distributed fashion, flowing over the paths of intended communication. The RE subscription protocols are called Simple Reservation Protocols (SRP), due to their simplicity as compared to the Resource Reservation Protocol (RSVP). SRP shares many of the baseline RSVP features, including the following:

- a) SRP is simple, i.e. reservations apply to unidirectional data flows.
- b) SRP is receiver-oriented, i.e., the receiver of a classA stream initiates and maintains the resource reservation used for that stream.
- c) SRP maintains “soft” state in bridges, providing graceful support for dynamic membership changes and automatic adaptations to changes in network topology.
- d) SRP is not a routing protocol, but depends on transparent bridging and STP routing protocols.

SRP simplicity is derived from its restricted layer-2 ambitions, as follows.

- a) SRP is symmetric, i.e. the listener-to-talker path is the inverse of the talker-to-listener path.
- b) SRP does not provide for transcoding; any stream is fully characterized by its streamID and bandwidth.

G.2.2 Soft reservation state

SRP takes a “soft state” approach to managing the reservation state in bridges. SRP soft state is created and periodically refreshed by listener generated RequestRefresh messages; this state is deleted if no matching RequestRefresh messages arrive before the expiration of a “cleanup timeout” interval. Listener’s may also force state deletions by generating an explicit RequestLeave message.

RequestRefresh messages are idempotent. When a route changes, the next RequestRefresh message will initialize the path state to the new route, and future RequestRefresh messages will establish state there. The state on the now-unused segment of the route will be deleted after a timeout interval. Thus, whether a RequestRefresh message is “new” or a “refresh” is determined separately by each station, depending upon the existence of state at that station.

SRP soft state is also deleted in the continued absence of associated classA traffic; this state is deleted if no matching classA traffic arrives before the expiration of a “cleanup timeout” interval. Thus, talker stations or agents may force reservation-state deletions by stopping their transmissions of classA traffic.

SRP sends its messages as layer-2 datagrams with no reliability enhancement. Periodic transmissions by listener stations and agents is expected to handle the occasional loss of an SRP message.

In the steady state, state is refreshed on a hop-by-hop basis to allow merging. Propagation of a change stops when and if it reaches a point where merging causes no resulting state change. This minimizes the SRP control traffic and is essential for scaling to large audiences.

G.2.3 Subscription bandwidth constraints

The SRP subscription protocols limit cumulative bandwidth allocations to a fixed percentage less than the capacity of the link, much like IEEE 1394 limits isochronous traffic to less than the capacity of its bus. This guarantees that high priority management information can be transmitted across the link. For RE systems,

classA traffic is limited to 75% of the capacity of any RE link. Enforcement of such a limit is done in multiple ways:

- a) Admissions controls (described in previous subclauses) reject any RequestRefresh message that (when combined with previously accepted request) would consume more than 75% of link bandwidth.
- b) Transmit queue hardware of RE stations (including bridges) discards classA content that (if transmitted) would cause classA traffic to exceed 75% of the transmit link capacity.

Method (b) is desired to recovery from unexpected transient conditions (typically topology changes) that result in admission control violations, and is also useful for managing misbehaving devices

G.2.4 Bridge-resident agents

Subscription facilities establish multicast paths from a talker to one or more listeners. Streams of time-sensitive data can then flow over these established paths, as illustrated by the dark arrow paths in Figure G.3-a. Maintaining these established paths involves active participation of agents within the end-point talker, local listener, local talker, and end-point listener entities, as illustrated in Figure G.3-b.

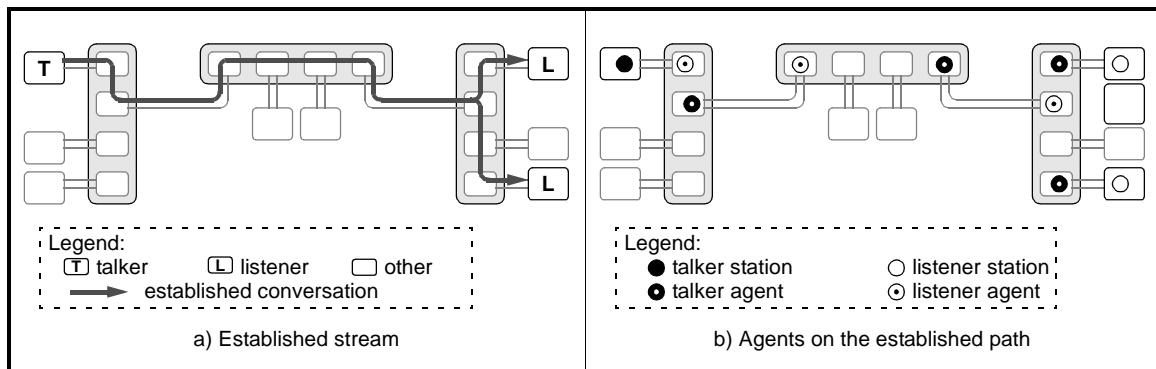


Figure G.3—Agents on an established path

The talker stations/agents are responsible for maintaining an account consisting of {streamID, bandwidth} pairs, one for each of their distinct flows. Requests for additional link bandwidth are checked against these accounts and rejected if the cumulative bandwidth would exceed 75% of the link capacity. The talker agents are also responsible for sustaining streams of classA data; their absence can result in disconnections of the attached listener agent.

The listener agents are responsible for periodically refreshing their adjacent talker agents, to confirm their continued presence. A persistent absence of refreshes causes the adjacent talker agent to disconnect its stream transmissions and (if appropriate) to inform other station-local agents.

For each established stream within a bridge, the listener agent remains active while all but the last downstream flows are disconnected. The upstream station receives its disconnect notice only after the last of the downstream flows has disconnected.

The listener agent's messages that establish and maintain the path are the same. This reduces design complexity and (most importantly) automatically re-routes stream flows after topology changes.

G.2.5 Controller entities

Subscription when a relative-intelligent controller discovers the need to establish a classA path between talker and listener entities. For example, user interactions with a television (called the controller) may cause streams flowing between the content source (called the talker) and speakers (the listeners), as illustrated in Figure G.4.

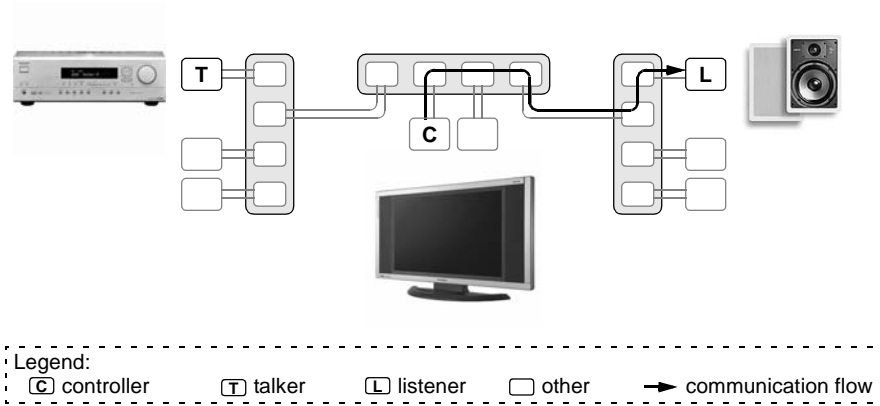


Figure G.4—Controller activation

A controller can potentially simplify the listener by reducing the need to providing user interface and device-discovery capabilities. However, a controller could also reside within talker and/or listener components. However, actions between controllers and talker/listener stations are beyond the scope of this working paper.

G.2.6 Pinging the talker

After being activated by a talker, listeners are expected to ping the talkers before initiating subscription operations, as illustrated in Figure G.5. The purpose of the ping is to ensure that bridges have learned listener and talker addresses, allowing frames to be sequentially passed from the listener-to-talkers.

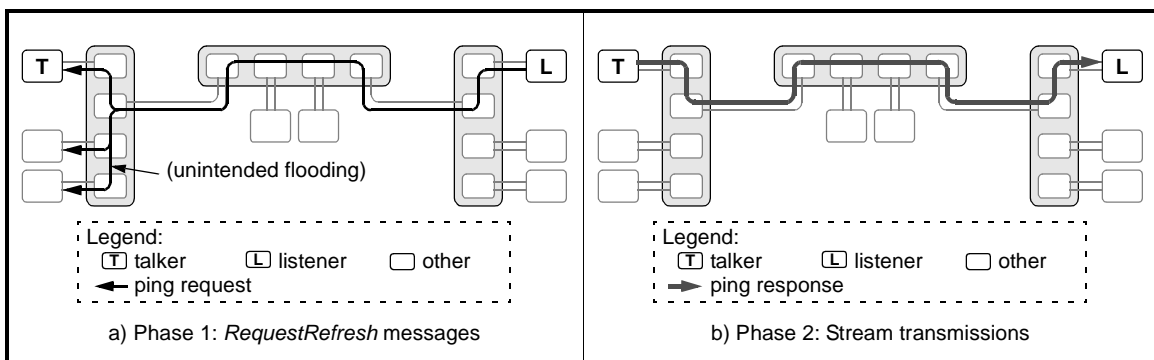


Figure G.5—Pinging the talker

G.2.7 Path creation

Establishing a conversation between a listener and a talker involves sending a RequestRefresh message from the listener towards the talker, illustrated by the dark arrow paths in Figure G.6-a. If available bandwidths are sufficient, the talker starts its stream transmissions, as illustrated by the gray arrow paths in Figure G.6-b.

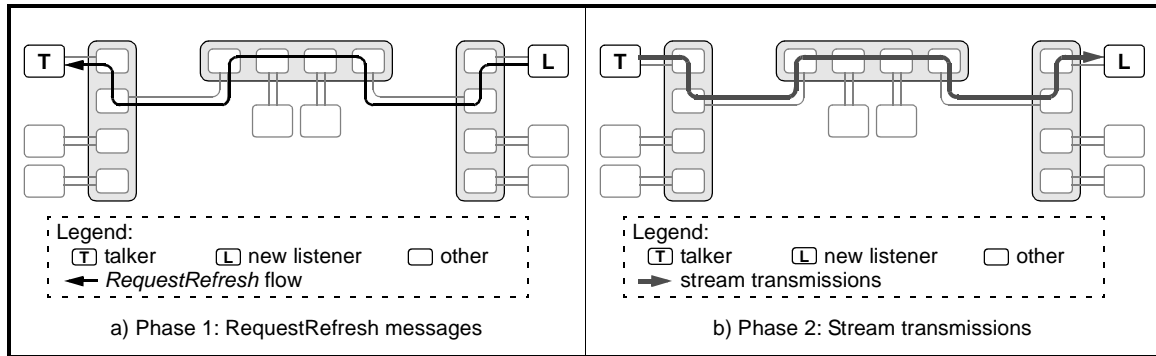


Figure G.6—Path creation

In rare circumstances, some talker addresses may not have been learned and the RequestRefresh message will terminate with a returned ResponseError message. The listener has the option of repeating the RequestRefresh after performing a ping (see G.2.6), which validates the talker presence and activates bridge learning.

Another timeout is associated with the absence of periodic RequestRefresh messages. In the continued absence of these expected messages, the listener is assumed to be absent or deactivated. Based on this assumption, the associated talker (station or agent) resources are released.

G.2.8 Side-path extensions

A second listener joins an established conversation by sending a RequestRefresh message towards the talker, as illustrated by the dark-arrow path in Figure G.7-a. When an established connection is discovered, the switch (not the talker) returns stream transmissions, as illustrated by the dark-gray path in Figure G.7-b.

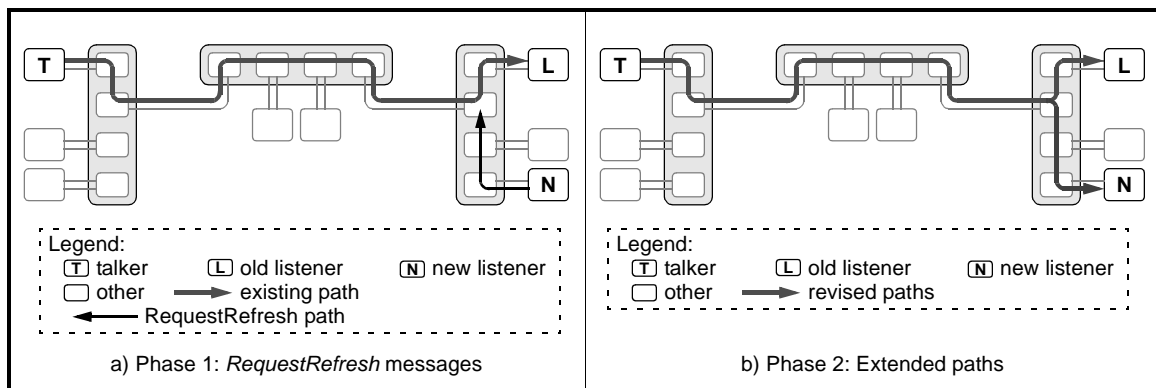


Figure G.7—Side-path extensions

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Each talker agent maintains separate state, so that classA traffic can be multicast to the applicable stations, rather than flooded downstream. The distinct markers also allow the switch to detect when the last listener disconnects, so that its previously shared upstream span can be released appropriately.

G.2.9 Side-path release

A retiring listener normally leaves an established conversation, by sending a RequestLeave message towards the talker. That message propagates to the nearest merging bridge connection, as illustrated by the dark-arrow path in Figure G.8-a. When an established/merged connection is discovered, the switch (not the talker) stops the stream transmissions, as illustrated by the disappearance of a side path in Figure G.8-b.

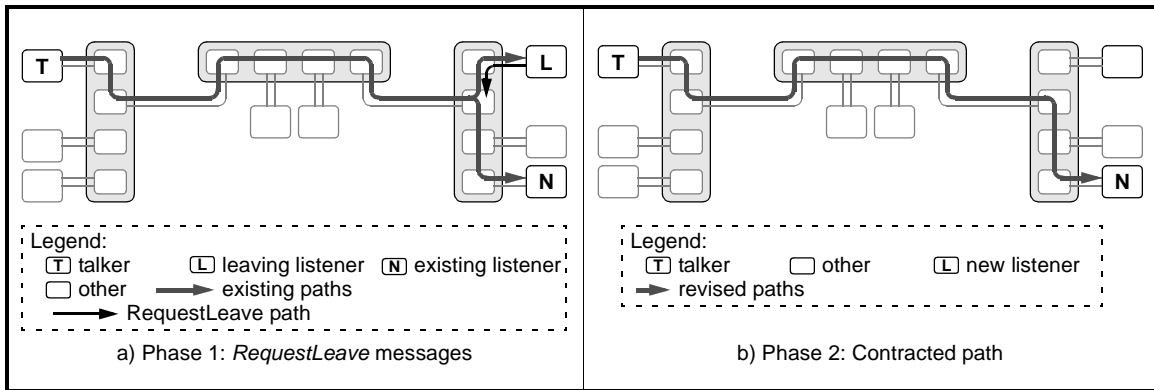


Figure G.8—Side-path demolition

G.2.10 Released path

The final listener bandwidth release involves sending a RequestLeave message towards the talker. In this case, that message propagates to the talker, as illustrated by the dark-arrow path in Figure G.9-a. The stream transmissions then stop, as illustrated in Figure G.9-b.

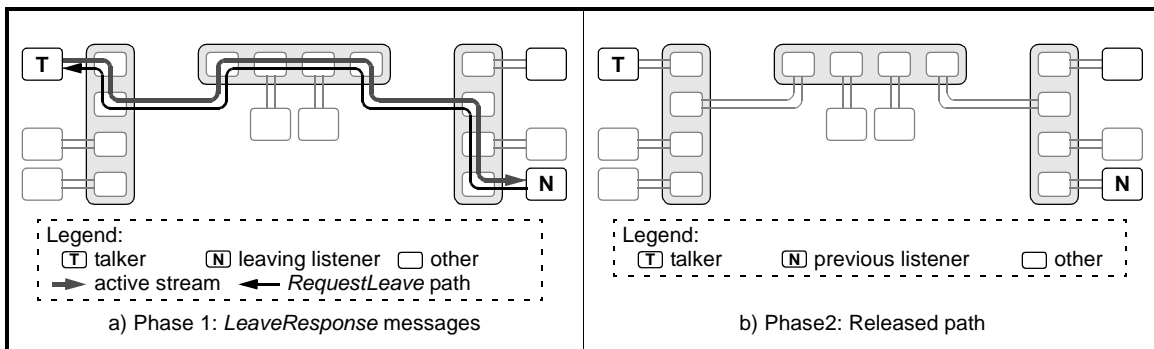


Figure G.9—Released path

G.2.11 Errors and timeouts

G.2.11.1 Subscription failures

A RequestRefresh message can encounter an error while flowing from the listener towards the talker, illustrated by the dark arrow paths in Figure G.10-a. When such errors occur, a ResponseError message is normally returned to the listener, as illustrated by the gray arrow paths in Figure G.10-b.

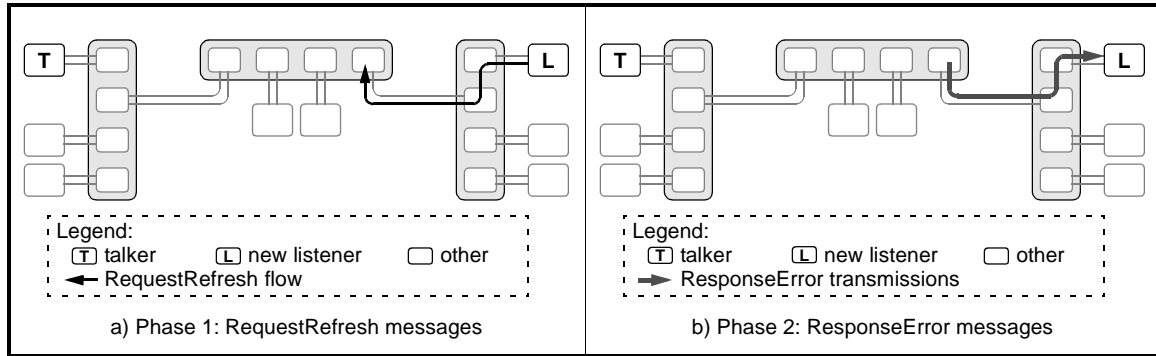


Figure G.10—Error responses

Errors may be associated with a variety of errors including (but not limited to) the following:

- a) Insufficient resources. Necessary resources are available within the bridge:
 - 1) Insufficient bandwidth is available on the link from the talker agent to its adjacent listener.
 - 2) Insufficient path-related resources are available in the bridge's talker agent.
 - 3) Insufficient path-related resources are available in the bridge's upstream listener agent.
 - 4) Insufficient link or memory bandwidth is available with the bridge.
- b) Unlearned address. The route from the bridge to the talker is unknown.
(To avoid complexities and inefficiencies, RequestRefresh messages are never flooded.)

G.2.11.2 Listener-presence timeouts

Listener agents and stations are responsible for refreshing their local talkers, to demonstrate their continued presence. In the absence of these refresh messages, the talkers assume the listener is absent and teardown the inactive path (or inactive branch from the path).

Thus, sustaining the active paths of Figure G.11-a requires periodic refresh messages on each hop, as illustrated in Figure G.11-b. The refresh messages and associated timeouts are performed independently on each span. The messages that establish the path (see G.2.7 and G.2.8) are the same as these listener-initiated messages that sustain the established path.

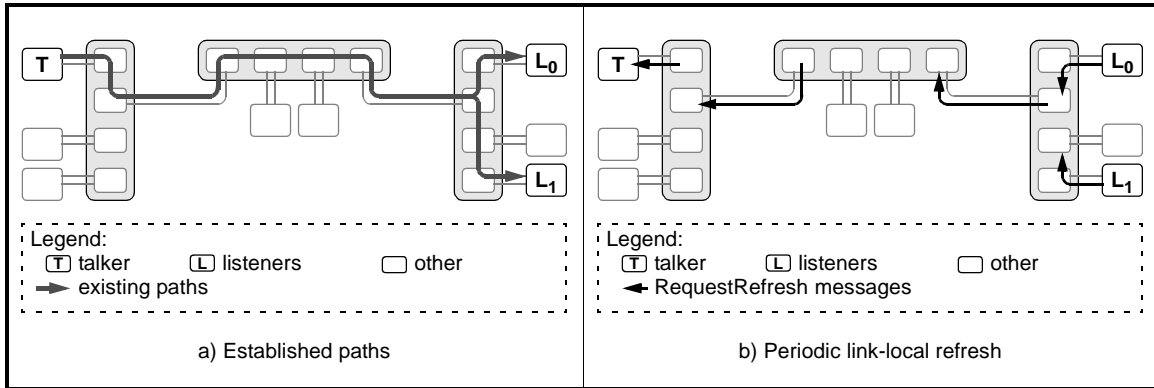


Figure G.11—Side-path demolition

G.2.11.3 Talker-presence timeouts

Talker agents and stations are responsible for updating their local listeners, to demonstrate their continued presence. In the absence of these updates, the listeners assume the talker is absent and teardown the inactive path (or inactive branch from the path).

Thus, sustaining the active paths of Figure G.11-a requires periodic transmissions of classA traffic on each hop (not illustrated). The associated timeouts are performed independently on each span. The frames that transfer classA data are the same as these talker-initiated frames that sustain the established path.

Annex H

(informative)

Frequently asked questions (FAQs)

H.1 Unfiltered email sequences

H.1.1 Bandwidth allocation

Question (AM): Is bandwidth allocation really necessary to meet RE requirements? Over-provisioning and best-effort (with class of service) may be adequate. You can get a lot of data through a conventional gigabit switch with very low latencies. The RE traffic can be given a higher priority and so not be held up by less urgent traffic.

Answer (MJT): I think admission control is needed. In an unmanaged layer 2 environment there is no way to *guarantee* that the streaming QoS parameters can be met ... you can only say *probably*. With GigE and a fully bridge-based environment with class of service you can get to a pretty good *probably*, but you can't get to the *it will always work* QoS that the wonderful BER of Ethernet promises. On the other hand, a simple admission control system and simple pacing mechanism can get you there, even with an FE-only network.

H.1.2 Best effort

Question (AM): With access control what happens if access is denied? My assumption is that a user connecting to a RE network would prefer best-effort service to no service at all if there is no spare bandwidth to be allocated. If you decide you need to support best-effort as a fallback then you need buffers in your end stations and the reason for using time slots goes away.

Answer (MJT): Your assumption is only correct if the service the consumer is subscribing to *is* a best-effort service. Right now, consumers expect that when they select a channel, or a CD, or a DVD they will get it *perfectly*. Cable companies get lots of calls if a stream is substandard for any reason. The general procedure to select a stream on a CE-oriented network would be something like:

- a) Hit the *directory* or *guide* button on your remote control
- b) Find the content you want (note that the content entries might be labeled with *not currently available* or *low quality only* or not even present depending on the state of the path to the source).
- c) Hit the *play* button.

Once the consumer hits that *play* button, the endpoints and network need to make a contract to deliver the content with the QoS expected by the consumer. So, in the case you describe where there is no guaranteed bandwidth available, you *may* present an alternative method (such as the *low quality* tag). This may be perfectly OK. If, on the other hand, the consumer wants to see the HD movie with full quality, they can yell at their kid to stop watching the movie that is causing the network link of interest to saturate.

H.2 Formulated responses

TBD

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54

Annex I

(informative)

Comment responses

<p>NOTE—This clause should be skipped on the first reading (reading starts at Clause 1). This clause is provided for communicating detailed responses to reviewer comments.</p>
--

I.1 Recent review-comment resolutions

TBD.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

Annex J

(informative)

C-code illustrations

NOTE—This annex is provided as a placeholder for illustrative C-code. Locating the C code in one location (as opposed to distributed throughout the working paper) is intended to simplify its review, extraction, compilation, and execution by critical reviewers. Also, placing this code in a distinct Annex allows the code to be conveniently formatted in 132-character landscape mode. This eliminates the need to truncate variable names and comments, so that the resulting code can be better understood by the reader.

This Annex provides code examples that illustrate the behavior of RE entities. The code in this Annex is purely for informational purposes, and should not be construed as mandating any particular implementation. In the event of a conflict between the contents of this Annex and another normative portion of this standard, the other normative portion shall take precedence.

The syntax used for the following code examples conforms to ANSI X3T9-1995.

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

```

// *****
//
//      1      2      3      4      5      6      7      8      9      0      1      2      3
//3456789012345678901234567890123456789012345678901234567890123456789012345678901234567890123456789012
//
#include <assert.h>
#include <stdio.h>

// unsigned char      uint8_t;          // 1-byte unsigned integer
// unsigned short     uint16_t;         // 2-byte unsigned integer
// unsigned int       uint32_t;         // 4-byte unsigned integer
// unsigned long long uint64_t;         // 8-byte unsigned integer

// signed char       int8_t;           // 1-byte signed integer
// signed short      int16_t;          // 2-byte signed integer
// signed int        int32_t;          // 4-byte signed integer
// signed long long  int64_t;          // 8-byte signed integer

#define OPTION_FAST 0                  // 1 if precedence-change is very-quick
#define OPTION_BASE 0                  // 1 if baseTimer hardware is provided
#define DIFF_SCALE ((double)4096 * ((uint64_t)1 << 31)) // Changes <200PPM to a 32-bit signed integer
#define EXTRACT_CORE(a, b) (((a) << 32) | ((b) >> 32)) // Extract seconds and fraction component
#define FULL_SCALE (0x7FFFFFFF)       // Biggest 32-bit positive integer
#define LIMIT(a, b, c) MAX(MIN((a), (b)), (c)) // Force base/bounds constraints
#define MAX(a, b) ((a) < (b) ? (b) : (a)) // Maximum value definition
#define MIN(a, b) ((a) > (b) ? (b) : (a)) // Minimum value definition
#define MINIMUM_WIDE(a, b) (CompareWide((a), (b)) < 0 ? (a) : (b))
#define ONES64 ~((uint64_t)0)         // 64-bit all-ones value
#define SCALE64 ((double)16 * (1 << 30) * (1 << 30)) // Floating-point equivalent of (1<<64)
#define BASE_PORT(siPtr) (siPtr->portPtr) // A pointer to the first station port
#define NEXT_PORT(piPtr) (piPtr->portPtr) // A pointer to the next station port
#define GRAND 2                       // An indication of grand-master mode
#define SLAVE 1                       // If not grand-master, slave mode

// The grand-master precedence check is based on concatenated fields, as follows:
//
//      MSB                                  LSB
//      |                                     |
//      |-----hi-----|-----lo-----|
//      |-----+-----+-----+-----+-----+-----+-----+-----+-----+
//      | 0000 systemTag eui64 00 hops portTag |
//      +---16---+---16---+-----64-----+---8---+---8---+---16---+
//
// If hops == ONES, this value is considered VOID and has the worse precedence
// Otherwise, the best precedence corresponds to the smallest of two tested values.
//
#if (CPU_TYPE == BIG)
typedef struct
{
    uint64_t hi; // more-significant portion
    uint64_t lo; // less-significant portion
} DoubleData;
typedef struct

```



```

{
    unsigned fill16:16;
    unsigned systemTag:16;
    unsigned uniqueHi:32;
    unsigned uniqueLo:32;
    unsigned fill08:8;
    unsigned hopsCount:8;
    unsigned portLevel:4;
    unsigned portNumber:12;
} DoubleInfo;
#else
typedef struct
{
    uint64_t lo;                // less-significant portion
    uint64_t hi;                // more-significant portion
} DoubleData;
typedef struct
{
    unsigned portNumber:12;
    unsigned portLevel:4;
    unsigned hopsCount:8;
    unsigned fill08:8;
    unsigned uniqueLo:32;
    unsigned uniqueHi:32;
    unsigned systemTag:16;
    unsigned fill16:16;
} DoubleInfo;
#endif

typedef union
{
    DoubleData data;           // As 64-bit data values
    DoubleInfo info;          // As data fields
} PrecedenceInfo;

typedef struct _PortInfo
{
    struct _PortInfo *portPtr; // Points to the next linked port
    unsigned portLevel:4;      // Relative priority number of ports
    unsigned portNumber:12;    // Port number
    DoubleData portPrecedence; // Incoming frame parameters

    uint8_t skipCount;         // Number of 10ms intervals
    uint32_t cableDelay;       // The cable delay, from local master
    uint32_t linkOffset;      // The cable difference, from local master
    uint64_t deltaTime;       // For inclusion in transmitted frames

    uint64_t latchRxFlexTime;  // Best if captured accurately by the PHY
    // Snapshot of flexTimer, on clockSync arrival,
    // available for this clockSync reception.
    uint64_t latchTxFlexTime;  // Snapshot of flexTimer, on clockSync departure,
    // available for next clockSync transmission

    uint64_t savedRxFlexTime;  // Previous latchRxFlexTime value
    uint64_t savedRxFlexData;  // Previous clockSync.lastFlexTime value
}

```

```

uint32_t latchRxBaseTime;
uint32_t latchTxBaseTime;
} PortInfo;

typedef struct
{
    PortInfo *portPtr;
    double nominalFrequency;
    int8_t clockDeviation;
    uint64_t eui64;
    unsigned systemLevel:4;
    unsigned systemNumber:12;

    unsigned selectCount;

    uint8_t skipCount;
    uint32_t myDiffRate;
    uint32_t diffRate;
    uint32_t linkOffset;
    uint64_t deltaTime;

    DoubleData thisPrecedence;
    DoubleData bestPrecedence;
    int16_t bestPort;

    uint64_t savedRxFlexTime;
    uint32_t savedRxBaseTime;

    uint64_t timeOfDay;
    uint64_t flexTimerHi;
    uint64_t flexTimerLo;
    uint64_t flexOffset;
    uint64_t flexRate;

    uint64_t baseTimer;
    uint64_t baseRate;

    uint32_t savedRxBaseTickTime;
    uint32_t savedRxBaseTickData;
} StationInfo;

typedef struct
{
    uint32_t da_hi;
    uint16_t da_lo;
    uint16_t sa_hi;
    uint32_t sa_lo;
    uint16_t protocolType;
    uint8_t subType;
    uint8_t syncCount;
    uint8_t hopsCount;
    uint8_t reserved;
    uint16_t systemTag;

```

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

```

uint64_t  uniqueID;           // Identifier for grand-master election
uint64_t  lastFlexTime;      // flexTimer on last clockSync transmission
uint64_t  deltaTime;        // Time difference on opposing link
uint64_t  offsetTime;       // Cumulative grand-master offset differences
uint32_t  diffRate;         // Cumulate grand-master rate differences
uint32_t  lastBaseTime;     // baseTimer on last clockSync transmission
uint32_t  fcs;              // Frame check sequence
} ClockSyncFrame;

uint32_t  BaseTimerChange(uint64_t, uint64_t, double);
void      ClockSyncArrived(StationInfo *, PortInfo *);
void      ClockSyncDeparted(StationInfo *, PortInfo *);
void      ClockSyncReceive(StationInfo *, PortInfo *, ClockSyncFrame *, uint8_t);
void      ClockSyncTransmit(StationInfo *, PortInfo *, ClockSyncFrame *);
int       CompareWide(DoubleData, DoubleData);
DoubleData PrecedenceMerge(uint16_t, uint64_t, uint8_t, uint8_t, uint16_t);
void      TimerTick(StationInfo *);
int       UpdatePrecedence(StationInfo *, PortInfo *);

// Called with:
// stationInfoPtr -- the station information context
void
StationSetup(StationInfo *stationInfoPtr)
{
    PortInfo *portPtr;
    StationInfo *siPtr = stationInfoPtr;
    uint16_t systemTag;
    uint16_t count;

    assert(siPtr != NULL);
    siPtr->baseRate = SCALE64 / siPtr->nominalFrequency;
    siPtr->systemLevel = 0X8;
    systemTag = ((uint16_t)(siPtr->systemLevel) << 12) | siPtr->systemNumber;
    siPtr->thisPrecedence =
        PrecedenceMerge(systemTag, siPtr->eui64, 0, 0, 0);
    count = 0;
    for (portPtr = BASE_PORT(siPtr); portPtr != NULL; portPtr = NEXT_PORT(portPtr)) {
        portPtr->portLevel = 0X8;
        portPtr->portNumber = count;
        count += 1;
    }
}

// Called with:
// stationInfoPtr -- the station information context
// portInfoPtr -- the port information context
// clockSyncPtr -- the contents of a clockSync frame
void
ClockSyncReceive(StationInfo *stationInfoPtr, PortInfo *portInfoPtr, ClockSyncFrame *clockSyncPtr, uint8_t rateAdjust)
{
    PortInfo *piPtr = portInfoPtr, *portPtr;
    PrecedenceInfo precedence;
    StationInfo *siPtr = stationInfoPtr;
    ClockSyncFrame *csPtr = clockSyncPtr;

```

```

uint32_t measuredDelta, receivedDelta, diffRate;
uint64_t rxDelta, txDelta, clockDelta, cableDelay;
double tempRate;
int8_t grand, slave;

assert(siPtr != NULL && piPtr != NULL);
if (csPtr != NULL && csPtr->hopsCount != 0xFF) // Compute the precedence,
    piPtr->portPrecedence = PrecedenceMerge(csPtr->systemTag, // in the absence of timeouts
    csPtr->uniqueID, csPtr->hopsCount, piPtr->portLevel, piPtr->portNumber);
else // Compute the precedence,
    piPtr->portPrecedence.hi = piPtr->portPrecedence.lo = ~((uint64_t)0); // in the presence of timeouts

if (OPTION_FAST && UpdatePrecedence(siPtr, piPtr)) {
    for (portPtr = BASE_PORT(siPtr); portPtr != NULL; portPtr = NEXT_PORT(portPtr))
        siPtr->bestPrecedence = MINIMUM_WIDE(piPtr->portPrecedence, siPtr->bestPrecedence);
    siPtr->selectCount += 1;
}
if (csPtr == NULL)
    return;

rxDelta = piPtr->savedRxFlexTime - csPtr->lastFlexTime; // Measured receive-link delay
txDelta = csPtr->deltaTime; // Reported transmit-link delay
clockDelta = (txDelta - rxDelta)/2; // Local timer differences
cableDelay = (txDelta + rxDelta)/2; // Cable transmission delay

precedence.data = siPtr->bestPrecedence;
grand = (precedence.info.hopsCount == 0) ? GRAND : 0; // Grand-master properties
slave = (precedence.info.portNumber == piPtr->portNumber) ? SLAVE : 0; // Slave port identification
switch(grand | slave) {
case GRAND: // Grand-master properties
    case GRAND | SLAVE: // override slave-port ID
        siPtr->diffRate = 0; // Grand-master reference
        siPtr->flexRate = siPtr->baseRate; // runs at the base rate
        break;
case SLAVE:
    if (rateAdjust) { // Low-rate adjustments
        measuredDelta = (siPtr->savedRxBaseTime - siPtr->savedRxBaseTickTime); // Clock-slave difference
        receivedDelta = (csPtr->lastBaseTime - siPtr->savedRxBaseTickData); // Clock-master difference
        siPtr->savedRxBaseTickTime = siPtr->savedRxBaseTime; // Previous saved value
        siPtr->savedRxBaseTickData = csPtr->lastBaseTime; // Previous saved value
        tempRate = DIFF_SCALE * ((double)(receivedDelta - measuredDelta)/receivedDelta); // Local rate difference
        siPtr->myDiffRate = LIMIT(tempRate, FULL_SCALE, -FULL_SCALE); // Rate difference limits
    }
    siPtr->diffRate = diffRate =
        LIMIT(siPtr->myDiffRate + csPtr->diffRate, FULL_SCALE, -FULL_SCALE); // Rate-range limitation
    siPtr->flexOffset = csPtr->offsetTime + clockDelta + siPtr->linkOffset; // Offset compensation
    siPtr->flexRate = siPtr->baseRate + siPtr->baseRate * (diffRate / DIFF_SCALE); // Rate compensation
    if (OPTION_BASE)
        siPtr->savedRxBaseTime = piPtr->latchRxBaseTime;
    else
        siPtr->savedRxBaseTime += // Receiver's baseTimer snapshot
            BaseTimerChange(siPtr->savedRxFlexTime, piPtr->latchRxFlexTime, diffRate);
    break;
default:
    break;
}

```

```

    }
    piPtr->cableDelay = cableDelay; // Local cable-delay knowledge
    piPtr->savedRxFlexTime = piPtr->latchRxFlexTime; // Saved reference time
    piPtr->deltaTime = rxDelta; // Saved for retransmission
}
// Called with:
// stationInfoPtr -- the station information context
// portInfoPtr -- the port information context
int
UpdatePrecedence(StationInfo *stationInfoPtr, PortInfo *portInfoPtr)
{
    PortInfo *piPtr = portInfoPtr;
    DoubleData pastPrecedence;
    PrecedenceInfo precedence;
    StationInfo *siPtr = stationInfoPtr;

    assert(siPtr != NULL && piPtr != NULL); // Set grand-master precedence
    precedence.data = pastPrecedence = siPtr->bestPrecedence; // Compare precedence values
    if (piPtr->portNumber == precedence.info.portNumber) // If this is the best port,
        siPtr->bestPrecedence = MINIMUM_WIDE(piPtr->portPrecedence, siPtr->thisPrecedence); // update baseline precedence
    else // If this is not the best port,
        siPtr->bestPrecedence = MINIMUM_WIDE(piPtr->portPrecedence, siPtr->bestPrecedence); // update overall precedence
    return(CompareWide(siPtr->bestPrecedence, pastPrecedence) != 0); // A precedence-change result
}

// Called with:
// stationInfoPtr -- the station information context
// portInfoPtr -- the port information context
// clockSyncPtr -- the contents of a clockSync frame
void
ClockSyncTransmit(StationInfo *stationInfoPtr, PortInfo *portInfoPtr, ClockSyncFrame *clockSyncPtr)
{
    ClockSyncFrame *csPtr = clockSyncPtr;
    PortInfo *piPtr = portInfoPtr;
    PrecedenceInfo precedence;
    StationInfo *siPtr = stationInfoPtr;

    assert(siPtr != NULL && piPtr != NULL && csPtr != NULL);
    if (UpdatePrecedence(siPtr, piPtr)) // If precedence has changed,
        siPtr->selectCount += 1; // start fast transmissions

    // An absent baseTimer is emulated by properly scaling time differences,
    // measured from the last recorded received-clockSync event.
    // - baseTime value was computed
    // - a different normDiffRate value has taken effect
    if (!OPTION BASE)
        piPtr->latchTxBaseTime = siPtr->savedRxBaseTime + // Derived from latchTxFlexRate
            BaseTimerChange(siPtr->savedRxFlexTime, piPtr->latchTxFlexTime, siPtr->diffRate);

    precedence.data = siPtr->bestPrecedence;
    csPtr->hopsCount = precedence.info.hopsCount; // Increment hop-count value
    csPtr->systemTag = precedence.info.systemTag; // Supply systemTag values
    csPtr->uniqueID = ((uint64_t)(precedence.info.uniqueHi) << 32) | precedence.info.uniqueLo; // Unique number tie-breaker
}

```

```

csPtr->lastFlexTime = piPtr->latchTxFlexTime; // Send last timer value 1
csPtr->deltaTime = piPtr->deltaTime; // Send received-link delay 2
csPtr->lastBaseTime = piPtr->latchTxBaseTime; // Send last baseTimer value 3
csPtr->offsetTime = siPtr->flexOffset; // This station's cumulative offset 4
csPtr->diffRate = siPtr->diffRate; // Send current diffRate value 5
} 6
// Called when a clockSync frame is received, to latch timer values. 7
// Latches timers are available when ClockSyncReceive() is called. 8
// 9
// Called with: 10
// stationInfoPtr -- the station information context 11
// portInfoPtr -- the port information context 12
void 13
ClockSyncArrived(StationInfo *stationInfoPtr, PortInfo *portInfoPtr) 14
{ 15
    PortInfo *piPtr = portInfoPtr; 16
    StationInfo *siPtr = stationInfoPtr; 17
    assert(siPtr != NULL); 18
    piPtr->latchRxFlexTime = EXTRACT_CORE(siPtr->flexTimerHi, siPtr->flexTimerLo); // Latch seconds:fraction fields 19
    if (OPTION_BASE) // If a baseTimer is present, 20
        piPtr->latchRxBaseTime = siPtr->baseTimer >> 32; // latch its fraction field 21
} 22
// Called when a clockSync frame is transmitted, to latch timer values. 23
// Latches timers are available for the next ClockSyncTransmit() call. 24
// 25
// Called with: 26
// stationInfoPtr -- the station information context 27
// portInfoPtr -- the port information context 28
void 29
ClockSyncDeparted(StationInfo *stationInfoPtr, PortInfo *portInfoPtr) 30
{ 31
    PortInfo *piPtr = portInfoPtr; 32
    StationInfo *siPtr = stationInfoPtr; 33
    assert(siPtr != NULL); 34
    piPtr->latchTxFlexTime = EXTRACT_CORE(siPtr->flexTimerHi, siPtr->flexTimerLo); // Latch seconds:fraction fields 35
    if (OPTION_BASE) // If a baseTimer is present, 36
        piPtr->latchTxBaseTime = siPtr->baseTimer >> 32; // latch its fraction field 37
} 38
// Called at a high clock rate (less than 20 ns) to update flexTimer and baseTimer (if present). 39
// This routine is intended to illustrate the computations involved in updating hardware timers; 40
// this code is _not_ expected to be incorporated into firmware. 41
// 42
// 43
// 0000 0000 0000 0000(hex) |-----+-----+-----+-----+-----+-----+-----+-----+-----+
// | fraction | subfraction | flexRate 44
// |-----32-----| |-----32-----| 45
// |-----v-----| |-----v-----| 46
// |-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+ 47
// ( flexAdd:128 ) 48
// |-----+-----+-----+-----+-----+-----+-----+-----+-----+ 49

```

```

//
//
// +-----+-----+-----+-----+-----+ flexTimer
// | superSeconds  seconds  fraction  subfraction |
// +-----32-----32-----32-----32-----+
//
//      ^ |
//      | v
//      +-----+
//      | offsetAdd:64 |
//      +-----+
//      ^ |
//      | v
// +-----+-----+ flexOffset
// | seconds  fraction |
// +-----32-----32-----+
//
//
// +-----+-----+ baseRate
// | fraction  subfraction |
// +-----32-----32-----+
//
//      |
//      v
// +-----+
// | baseAdd:64 |
// +-----+
//      ^ |
//      | v
// +-----+-----+ baseTimer
// | fraction  subfraction |
// +-----32-----32-----+
//
//
// Called with:
// stationInfoPtr -- the station information context
void
TimerTick(StationInfo *stationInfoPtr)
{
    StationInfo *siPtr = stationInfoPtr;
    int64_t pastTimerLo;

    assert(siPtr != NULL);
    pastTimerLo = siPtr->flexTimerLo;
    siPtr->flexTimerLo += siPtr->flexRate; // Saved to detect overflows
    siPtr->flexTimerHi += (pastTimerLo > siPtr->flexTimerLo) ? 1 : 0; // Addition of the subfractions
    siPtr->timeOfDay = EXTRACT_CORE(siPtr->flexTimerHi, siPtr->flexTimerLo) + siPtr->flexOffset; // Propagate carry into seconds
    siPtr->timeOfDay = EXTRACT_CORE(siPtr->flexTimerHi, siPtr->flexTimerLo) + siPtr->flexOffset; // Compensate for offset drifts
    if (OPTION_BASE)
        siPtr->baseTimer += siPtr->baseRate;
}

// Called with:
// stationInfoPtr -- the station information context
// currentTime -- the current flexTimer value
// diffRate -- the scaled rate difference
uint32_t

```

```

BaseTimerChange(uint64_t lastTime, uint64_t nextTime, double diffRate)
{
    uint64_t delta;

    delta = nextTime - lastTime;
    delta -= delta * (diffRate / DIFF_SCALE);
    return(delta);
}

// Merge multiple fields into an 128-bit integer, for comparisons
// Called with:
//   systemTag    -- the 16-bit most-significant precedence subfield
//   uniqueID     -- the 64-bit unique identifier (EUI-64)
//   hopsCount    -- the hop-count distance from the grand master
//   portTag      -- the tag associated with the port
DoubleData
PrecedenceMerge(uint16_t systemTag, uint64_t uniqueID, uint8_t hopsCount, uint8_t portLevel, uint16_t portNumber)
{
    PrecedenceInfo result;

    result.info.fill16 = 0;
    result.info.systemTag = systemTag;
    result.info.uniqueHi = (uniqueID >> 32);
    result.info.uniqueLo = uniqueID;
    result.info.fill08 = 0;
    result.info.hopsCount = hopsCount;
    result.info.portLevel = portLevel;
    result.info.portNumber = portNumber;
    return(result.data);
}

// Performs a comparison of 128-bit preceision unsigned values
// Called with:
//   a -- the first of two 128-bit values
//   b -- the final of two 128-bit values
int
CompareWide(DoubleData a, DoubleData b)
{
    if (a.hi != b.hi)
        return(a.hi > b.hi ? 1 : -1);
    if (a.lo != b.lo)
        return(b.lo > b.lo ? 1 : -1);
    return(0);
}

```


Index**C**

classA frame	
<i>da</i>	59
<i>sa</i>	59
<i>protocolType</i>	59
<i>serviceDataUnit</i>	59
<i>fcs</i>	59
clockSync frame	
<i>da</i>	60
<i>sa</i>	60
<i>protocolType</i>	60
<i>subType</i>	60
<i>hopsCount</i>	60
<i>syncCount</i>	60
<i>cycleCount</i>	61
<i>systemTag</i>	61
<i>systemLevel</i>	61
<i>systemNumber</i>	61
<i>uniqueID</i>	61
<i>oui</i>	62
<i>extension</i>	62
<i>ouiDependent</i>	62
<i>lastFlexTime</i>	61
<i>seconds</i>	62
<i>fraction</i>	62
<i>deltaTime</i>	61
<i>seconds</i>	62
<i>fraction</i>	62
<i>offsetTime</i>	61
<i>seconds</i>	62
<i>fraction</i>	62
<i>diffRate</i>	61
<i>lastBaseTime</i>	61
<i>fcs</i>	61

cycleCount
 See clockSync frame

D

<i>da</i>	
<i>See</i> classA frame	
<i>See</i> clockSync frame	
<i>See</i> RequestRefresh frame	
<i>deltaTime</i>	
<i>See</i> clockSync frame	
<i>diffRate</i>	
<i>See</i> clockSync frame	

E

<i>extension</i>	
<i>See</i> clockSync frame	

F

<i>fcs</i>	
<i>See</i> classA frame	
<i>See</i> clockSync frame	
<i>See</i> RequestRefresh frame	
<i>fraction</i>	
<i>See</i> clockSync frame	
<i>See</i> time field	

H

<i>hopsCount</i>	
<i>See</i> clockSync frame	

I

<i>info</i>	
<i>See</i> RequestRefresh frame	
<i>info</i> field	
<i>reserved</i>	64
<i>talkerID</i>	64
<i>plugID</i>	64
<i>maxCycles</i>	64
<i>multicastID</i>	64
<i>maxBw</i>	64
<i>reserved</i>	64

L

<i>lastBaseTime</i>	
<i>See</i> clockSync frame	
<i>lastFlexTime</i>	
<i>See</i> clockSync frame	

M

<i>maxBw</i>	
<i>See</i> <i>info</i> field	
<i>See</i> RequestLeave frame	
<i>See</i> RequestRefresh frame	
<i>See</i> ResponseError frame	
<i>maxCycles</i>	
<i>See</i> <i>info</i> field	
<i>See</i> RequestLeave frame	
<i>See</i> RequestRefresh frame	
<i>See</i> ResponseError frame	
<i>mcastID</i>	
<i>See</i> RequestLeave frame	
<i>See</i> ResponseError frame	
<i>mcastSrc</i>	
<i>See</i> RequestRefresh frame	
<i>multicastID</i>	
<i>See</i> <i>info</i> field	

O

<i>offsetTime</i>	
<i>See</i> clockSync frame	

1	<i>oui</i>	<i>reserved</i>	64
2	<i>See</i> clockSync frame		
3	<i>ouiDependent</i>		
4	<i>See</i> clockSync frame		
5			
6	P		
7	<i>pad</i>		
8	<i>See</i> RequestRefresh frame		
9	<i>plugID</i>		
10	<i>See</i> <i>info</i> field		
11	<i>See</i> RequestLeave frame		
12	<i>See</i> RequestRefresh frame		
13	<i>See</i> ResponseError frame		
14	<i>protocolType</i>		
15	<i>See</i> classA frame		
16	<i>See</i> clockSync frame		
17	<i>See</i> RequestRefresh frame		
18			
19	R		
20	RequestLeave frame		
21	<i>info</i>		
22	<i>mcastID</i>	64	
23	<i>talkerID</i>	64	
24	<i>plugID</i>	64	
25	<i>maxCycles</i>	64	
26	<i>maxBw</i>	64	
27	<i>reserved</i>	64	
28	RequestRefresh frame		
29	<i>da</i>	63	
30	<i>sa</i>	63	
31	<i>protocolType</i>	63	
32	<i>subType</i>	63	
33	<i>count</i>	63	
34	<i>info</i>	63	
35	<i>mcastID</i>	64	
36	<i>talkerID</i>	64	
37	<i>plugID</i>	64	
38	<i>maxCycles</i>	64	
39	<i>maxBw</i>	64	
40	<i>reserved</i>	64	
41	<i>pad</i>	63	
42	<i>fcs</i>	63	
43	<i>reserved</i>		
44	<i>See</i> <i>info</i> field		
45	<i>See</i> RequestLeave frame		
46	<i>See</i> RequestRefresh frame		
47	<i>See</i> ResponseError frame		
48	ResponseError frame		
49	<i>info</i>		
50	<i>mcastID</i>	64	
51	<i>talkerID</i>	64	
52	<i>plugID</i>	64	
53	<i>maxCycles</i>	64	
54	<i>maxBw</i>	64	

S	
<i>sa</i>	
<i>See</i> classA frame	
<i>See</i> clockSync frame	
<i>See</i> RequestRefresh frame	
<i>seconds</i>	
<i>See</i> clockSync frame	
<i>See</i> time field	
<i>serviceDataUnit</i>	
<i>See</i> classA frame	
<i>subType</i>	
<i>See</i> clockSync frame	
<i>See</i> RequestRefresh frame	
<i>syncCount</i>	
<i>See</i> clockSync frame	
<i>systemLevel</i>	
<i>See</i> clockSync frame	
<i>systemNumber</i>	
<i>See</i> clockSync frame	
<i>systemTag</i>	
<i>See</i> clockSync frame	
T	
<i>talkerID</i>	
<i>See</i> <i>info</i> field	
<i>See</i> RequestLeave frame	
<i>See</i> RequestRefresh frame	
<i>See</i> ResponseError frame	
time field	
<i>seconds</i>	62
<i>fraction</i>	62
U	
<i>uniqueID</i>	
<i>See</i> clockSync frame	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54