

## 9. MAC fairness

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

**References:**  
None.

**Definitions:**

**Abbreviations:**

**Revision History:**

Draft 0.1, February 2002

Draft 0.2, April 2002

Draft 0.3, June 2002

Initial draft document for RPR WG review.

Revised according to WG comments for TF review.

Revised according to WG comments for TF review.

### 9.1 Overview

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

Missing things in the overview section include:

- Functional block diagram and contents of clause
- Notational conventions

Detailed state diagrams are needed that describe:  
the events that cause fairness control messages to be sent;  
the actions when a fairness control message is received;  
congestion and congestion subsided, etc.

The PICS needs to be added.

Need to make sure that reference model clause is in sync with this one.

#### 9.1.1 Scope

This clause defines the fairness algorithm for RPR MACs how the MAC uses the algorithm to enforce fairness among stations on the ring. The RPR fairness algorithm (RPR-FA) is a local fairness mechanism that enforces fairness among all the stations on a ring, even during times when the ring is heavily utilized or congested. Only fairness eligible traffic (the class-B traffic in excess of the provisioned amount, and all class-C traffic) is subject to the fairness algorithm.

#### 9.1.2 Goals and objectives

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

We could not agree on a definition for stability, and hence the corresponding objective has been removed. If we find an acceptable way to describe stability, then we can add it later. For now, note that contributions are invited on a suitable definition (and corresponding objective) for stability. One possible definition is:  
**Stability**—The protocol should not oscillate regardless of the type of traffic presented to the network.

The fairness protocol has the following objectives:

- Source-based weighted fairness**—On any given link on the ringlet, the available bandwidth is allocated to each station in proportion to its relative weight. For example, if every station has an equal weight, then the available bandwidth on the link should be shared equally by all stations. On the

other hand, if one station has a higher weight, the bandwidth allocated to that station should be in proportion to the station's weight divided by the sum of the weights of all the stations.

- b) **Reclamation of unused committed bandwidth**—The fairness protocol should be able to reclaim unused bandwidth, including that which is provisioned.
- c) **Support for single choke and multi choke capable clients**—The fairness protocol should be able to support clients whether they utilize a single choke point or multiple choke points (such as those clients performing virtual destination queueing (VDQ)).
- d) **Fast response time**—Because data traffic tends to be bursty, in order to ensure maximum ring bandwidth utilization and to ensure that the protocol is responsive to instantaneous changes in traffic load, it must have a fast response time.
- e) **High bandwidth utilization on the ring**—The protocol should be able to achieve very high levels of bandwidth utilization even under heavy loads approaching 100% of the ring capacity.
- f) **Scalability**—The protocol should be scalable and should be able to function predictably for all ringlet speeds and ring diameters allowed by this standard.

### 9.1.3 Relationship to other clauses

The RPR-FA is implemented within a control entity called the Fairness Control Unit (FCU) located in the MAC Control Sublayer, as described in Clause 5.

The FCU uses the following variables defined in Clause 6: add\_rate, add\_rate\_congested, fw\_rate, fw\_rate\_congested.

The FCU provides the following variables for use in Clause 6: allowed\_rate, allowed\_rate\_congested, TTL\_to\_congestion. The usage of these variables by Clause 6 is described in 9.6 .

## 9.2 Acronyms

This clause contains the following acronyms:

FA	fairness algorithm
FCM	fairness control message
FCU	fairness control unit
FE	fairness eligible
LR	line rate
MC	multi choke
RTT	round trip time
SC	single choke

## 9.3 Variables and terminology used

This clause contains the following definitions and formula variables:

**9.3.1 active\_stations:** Calculated. The number of stations that were active (as measured by having at least 1 transited fairness eligible frame) on the outgoing link of a ringlet during the last aging\_interval.

**9.3.2 add\_rate:** Calculated. A running byte count of all FE marked packets transmitted to the ringlet from the add queues of the local station during the previous AGEcoef aging\_intervals (including the current interval), and smoothed over the previous AGEcoef aging\_intervals. This rate is enforced by the Sd shaper to be less than or equal to the allowed\_rate (+ 1 MTU - 1 byte).

**9.3.3 add\_rate\_congested:** Calculated. A running byte count of all FE marked packets transmitted to the ringlet from the add queues of local station during the previous AGEcoef aging\_intervals (including the

current interval), and smoothed over the previous AGECOEF aging\_intervals, that have a destination past the congestion\_point. This rate is enforced by the Sc shaper to be less than or equal to the allowed\_rate\_congested (+ 1 MTU - 1 byte).

**9.3.4 add\_rate\_congested\_ok:** Calculated. A boolean value indicating if the current add\_rate\_congested is within acceptable limits. This value is an input to the sendC signal defined in Clause 6.

**9.3.5 add\_rate\_ok:** Calculated. A boolean value indicating if the current add\_rate is within acceptable limits. This value is an input to the sendC signal defined in Clause 6.

**9.3.6 advertised\_fair\_rate:** Calculated. The normalized rate (in bytes) advertised in an FCM giving the fair number of bytes that can be sent during the next aging\_interval.

**9.3.7 advertisement\_interval:** Configured. The interval at which FCMs are sent. This time is less than the aging\_interval. The default value for this constant is a rate that uses .125% of the link bandwidth (e.g. 40 usec for a LINK\_RATE of 2.5 Gbps).

**9.3.8 AGECOEF:** Configured. The coefficient used for aging the running add\_rate, add\_rate\_congested, fw\_rate, and fw\_rate\_congested. The default value for this constant is 4. Values other than factors of 2 are not required to be supported.

**9.3.9 aging\_interval:** Configured. The interval at which aging and low pass filtering functions are performed. The value to use for each supported link speed is shown in Table 9.1.

**Table 9.1—Rate Coefficient**

Link Speed	aging_interval
OC-12 and higher rates	100 usec
OC-3	400 usec

**9.3.10 aggressive:** Configured. Aggressive at maximizing link utilization. Local station fairness (during transitions). Trades off stability for utilization (during transitions).

**9.3.11 allowed\_rate:** Calculated. The rate at which the local station is allowed to transmit FE marked packets to the ringlet, specified as the number of bytes per AGECOEF aging\_intervals. During any aging\_interval, the MAC may send up to allowed\_rate minus add\_rate FE bytes. This value is an input to the Sd rate shaper/limiter defined in Clause 6.

**9.3.12 allowed\_rate\_congested:** Calculated. The rate at which the local station is allowed to transmit FE marked packets to the ringlet that have a destination past the congestion\_point, specified as the number of bytes per AGECOEF aging\_intervals. During any aging\_interval, the MAC may send up to allowed\_rate\_congested minus add\_rate\_congested FE bytes beyond the congestion\_point. This value is the input to the Sc rate shaper/limiter defined in Clause 6.

**9.3.13 congested:** Calculated: The indication of whether the local station considers itself to be congested in trying to add traffic.

**9.3.14 congestion\_point:** Calculated. The most congested station on the ringlet reporting a congestion condition (as indicated by the SA and TTL in the most recent SC-FCM).

1 **9.3.15 conservative:** Configured. Conservative at maximizing link utilization. Congested span fairness  
2 (during transitions). Trades off utilization for stability (during transitions).  
3

4 **9.3.16 fairness eligible:** Definition. A quality of a packet that indicates (as specified in 8) if it is subject to  
5 the fairness algorithm. Packets that are subject to fairness are the class-B packets added in excess of the sta-  
6 tion's provisioned amount, and all class-C traffic. See also 6.  
7

8 **9.3.17 FULL\_RATE:** Constant. A special value for control\_rate, indicating no congestion. The value for  
9 this constant is an all "1"s value.  
10

11 **9.3.18 full\_threshold:** Configured. The STQ threshold that indicates that the STQ of a dual-queue MAC is  
12 almost full. Same as NEEDY in Clause 6. The default value for full\_threshold is STQ - 1 MTU.  
13

14 **9.3.19 fw\_rate:** Calculated. A running byte count of all FE marked packets transmitted to the ringlet from  
15 the transit queues of the local station during the previous AGECOEF aging\_intervals (including the current  
16 interval), and smoothed over the previous AGECOEF aging\_intervals.  
17

18 **9.3.20 fw\_rate\_congested:** Calculated. A running byte count of all FE marked packets transmitted to the  
19 ringlet from the transit queues of the local station during the previous AGECOEF aging\_intervals (including  
20 the current interval), and smoothed over the previous AGECOEF aging\_intervals, that have a destination  
21 past the congestion\_point.  
22

23 **9.3.21 high\_threshold:** Constant. The threshold that indicates that a MAC is experiencing some congestion.  
24 When the STQ depth or link utilization exceed this value, the outgoing link is considered congested.  
25

26 **9.3.22 LINK\_RATE:** Configured. Conceptually, the actual maximum transmission rate for the outbound  
27 link. For fairness formulas, it is interpreted as the maximum number of bytes that can be transmitted in  
28 AGECOEF aging\_intervals. This rate is less than AGECOEF \* aging\_interval worth of bytes.  
29

30 **9.3.23 local\_fair\_rate:** Calculated. The rate at which the local station may transmit FE marked packets,  
31 regardless of downstream congestion, specified as the number of bytes per AGECOEF aging\_intervals.  
32

33 **9.3.24 local\_station:** Definition. The station in which the local MAC resides.  
34

35 **9.3.25 low\_threshold:** Constant. The threshold that indicates that a MAC may soon start experiencing con-  
36 gestion. When the STQ depth or link utilization exceed this value, the outgoing link is considered imminent  
37 to being congested.  
38

39 **9.3.26 lp\_add\_rate:** Calculated. A low pass filtered value of add\_rate.  
40

41 **9.3.27 lp\_add\_rate\_congested:** Calculated. A low pass filtered value of add\_rate\_congested.  
42

43 **9.3.28 lp\_fw\_rate:** Calculated. A low pass filtered value of fw\_rate.  
44

45 **9.3.29 lp\_fw\_rate\_congested:** Calculated. A low pass filtered value of fw\_rate\_congested.  
46

47 **9.3.30 lp\_nr\_xmit\_rate:** Calculated. A low pass filtered value of nr\_xmit\_rate.  
48

49 **9.3.31 LPCOEF:** Configured. The coefficient used for low pass filtering the add\_rate and fw\_rate to result  
50 in lp\_add\_rate and lp\_fw\_rate. The default value for this constant is 64. Values other than factors of 2 are not  
51 required to be supported.  
52

53 **9.3.32 MAX\_ALLOWED\_RATE:** Configured. The maximum value for allowed\_rate. The default value  
54 for this constant is LINK\_RATE.

**9.3.33 multi choke:** Definition. The ability to work with multi choke points (of different values) at the same time, e.g. as in a VDQ implementation. The primary intended use of the multi choke fairness control message (MC-FCM).

**9.3.34 norm\_local\_fair\_rate:** Calculated. The (NORMCOEF) normalized value of local\_fair\_rate. It is the value for advertised\_fair\_rate in SC-FCMs when the local station is more congested than downstream stations; and always the value for advertised\_fair\_rate in MC-FCMs.

**9.3.35 norm\_lp\_fw\_rate\_congested:** Calculated. The (NORMCOEF) normalized value of lp\_fw\_rate\_congested.

**9.3.36 NORMCOEF:** Configured. The coefficient used for normalizing advertised fair rates to a common rate. The value for this constant is the product of AGECOEF, RATECOEF, and WEIGHT.

**9.3.37 nr\_xmit\_rate:** Calculated. A running byte count of all non reserved (non subclass-A0 marked) packets transmitted to the ringlet (from all add queues and transit queues) of the local station during the previous AGECOEF aging\_intervals (including the current interval), and smoothed over the previous AGECOEF aging\_intervals.

**9.3.38 num\_stations:** Calculated. The number of stations that are known to be on the ringlet. Reported via rprIfNodesOnRing.

**9.3.39 RAMPCOEF:** Configured. The coefficient used for ramping up or down the allowed\_rate\_congested or local\_fair\_rate. The default value for this constant is 64. Values other than factors of 2 are not required to be supported.

**9.3.40 RATECOEF:** Configured. The coefficient used for normalizing advertised fair rates to a common link rate. The RATECOEF is used to ensure that the control value does not overflow 16-bits. The value to use for each supported link speed is shown in Table 9.2.

**Table 9.2—Rate Coefficient**

Link Speed	Rate Coefficient
Up to and including 2.5 Gbps	1
10 Gbps	4
40 Gbps	16

**9.3.41 rcvd\_fair\_rate:** Calculated. The normalized rate (in bytes) received in an FCM, giving the fair number of bytes that can be sent during the next aging\_interval.

**9.3.42 reserved\_rate:** Calculated. The total class A traffic that is reserved (i.e. the total of all sub class A0 traffic) on the outgoing link for all stations, including the local station, specified as the number of bytes per AGECOEF aging\_intervals.

**9.3.43 single choke:** Definition. The ability to work with only one choke point at one time. The primary intended use of the single choke fairness control message (SC-FCM).

**9.3.44 TTL\_to\_congestion:** Calculated. The number of hops to the congestion point.

**9.3.45 unreserved\_rate:** Calculated. The rate available for the outbound link that is not reserved, specified as the number of bytes per AGECOEF aging\_intervals.

**9.3.46 WEIGHT:** Configured. The local weight by which allowed\_rate, allowed\_rate\_congested, and local\_fair\_rate are calculated. The default value for this constant is 1.

## 9.4 MAC fairness operation

The fairness algorithm implemented within the FCU consists of the following functions:

- a) Determining when the congestion threshold is crossed and when the congestion has subsided;
- b) Determining the fair rate for advertisement;
- c) Determining the station's allowed rate;
- d) Sourcing and consuming fairness control messages;
- e) Communicating the allowed rate to the data path shapers for controlling access to the medium;
- f) Providing the information contained in MC-FCMs to the client.

Each station is assigned a weight, which allows the user to allocate more ring bandwidth to certain stations as compared with other stations. This is referred to as the *weighted fairness* property of RPR-FA.

In RPR-FA, a station advertises a fair rate to upstream stations via the ringlet opposite of the ringlet upon which the algorithm is running. The fair rate is run through a low pass filter function, with the result known as the local\_fair\_rate. The local\_fair\_rate is normalized by the local station weight (WEIGHT), the aging coefficient (AGECOEF), and the rate coefficient (RATECOEF) to yield the norm\_local\_fair\_rate and the advertised\_fair\_rate. The low-pass filter stabilizes the feedback, the division by WEIGHT normalizes the transmitted value to a weight of 1.0, the division by AGECOEF normalizes the transmitted value to one aging\_interval, and the division by RATECOEF normalizes the transmitted value to a link speed of 2.5 Gbps. When the upstream stations receive an advertised fair rate, they will adjust their transmit rates for fairness eligible traffic so as not to exceed the advertised value (adjusted by their respective weights).

Propagation of the advertised value to other stations on the ring is done using fairness control messages. The format of the fairness control messages is described in 9.7. There are 2 types of fairness control messages—Single Choke and Multi Choke.

Single Choke messages are propagated hop-by-hop around the opposite ringlet and are processed by the FCU. They are sent every advertisement\_interval.

A wrapped ring shall be treated as a folded ring from the point of view of SC-FCMs. A station with a wrapped attachment point receiving a SC-FCM shall change the ringlet\_id and wrap the SC-FCM. SC-FCMs are stripped only when received with SA equal to the local MAC address and with the correct ringlet\_id.

Single Choke fairness control messages contain the SA of the most congested station. If a station experiences congestion, it will include a non-FULL\_RATE value for the fair rate in its Single Choke fairness control messages. A station that receives a Single Choke message with a SA of its MAC address shall treat the message as having a control value of FULL\_RATE regardless of the actual control value.

A station that in the congested state shall advertise the minimum of its local\_fair\_rate and the last received fair rate (known as the rcvd\_fair\_rate) in its Single Choke message.

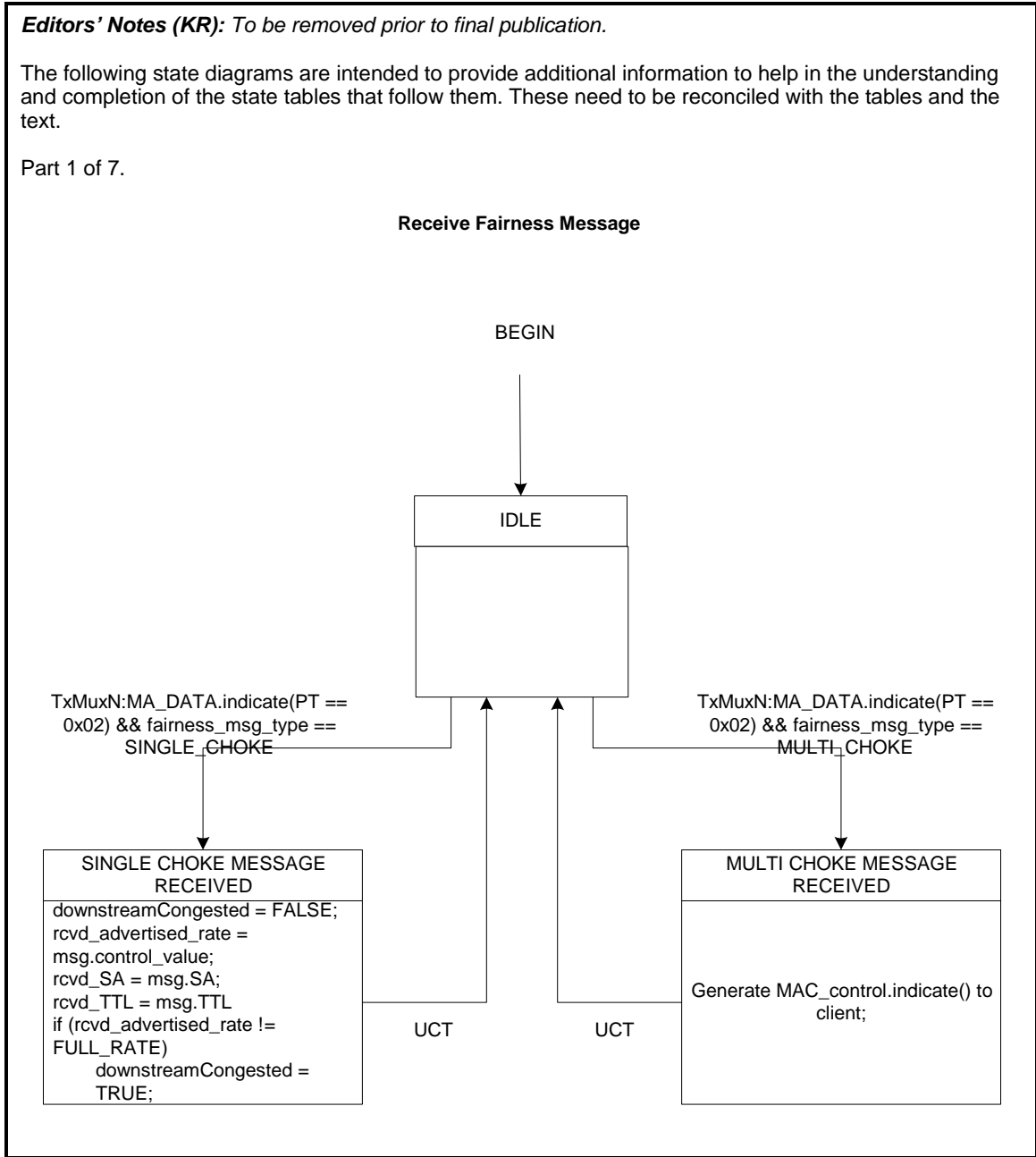
A station that is not congested and that receives a Single Choke message containing a non-FULL\_RATE rcvd\_fair\_rate shall either propagate the rcvd\_fair\_rate to its upstream neighbor (leaving the SA the same), or it will send a value of FULL\_RATE and set the SA to its own. Which of these is sent is determined by the

following. If the low pass filtered `fw_rate_congested` is less than the `allowed_rate_congested` divided by the local station weight, then a `FULL_RATE` value is propagated to the upstream neighbor instead of the `rcvd_fair_rate`. Otherwise, there cannot be an upstream station that is the cause of congestion. Thus there is no need to propagate the Single Choke fairness control message indicating congestion upstream of this station.

The value of the rate that is allowed for the fairness eligible traffic (`allowed_rate_congested`) that can be added by the station is derived from the received advertised rate in the Single Choke messages and the local station normalization factor (`NORMCOEF`). If the number of hops to the destination of a data packet from the MAC client is less than the number of hops to the station that generated the Single Choke message, then the MAC client can take advantage of the available bandwidth on the ringlet; otherwise, if the destination of a data packet would allow the packet to go beyond the station that generated the fairness control message, the client will at least receive its fair share of the bandwidth from the most congested span that it contends for.

Multi Choke messages are broadcast on the ringlet opposite to the ringlet upon which the fairness algorithm is running and contain the SA of the station that originated the message. They are sent every 10 advertisement\_intervals. If a station experiences congestion, it will include a non-`FULL_RATE` value for the fair rate in its Multi Choke fairness control messages; otherwise it will include a `FULL_RATE` value for the fair rate. Multi Choke messages are not required to be processed by the FCU, but the information is passed to the MAC client by the FCU and is used by the MAC client as described below.

NOTE—The client can also take advantage of Virtual Destination Queueing (VDQ) by utilizing the multi-choke concept of RPR-FA. VDQ combined with RPR-FA can increase ring utilization. The multi-choke concept deals with the case where a station wants to send traffic to a destination that is closer than a congested link. As an example, consider the case where station 1 wants to send traffic to station 2, and the link between stations 2 and 3 is congested. RPR-FA will allow station 1 to send as much traffic as it wants to station 2, and will only limit traffic to stations beyond the congested link to the fair rate. In a multi-choke implementation of the RPR-FA, each client will track advertised fair rates for congested stations. A station is allowed to send unlimited traffic to any station between itself and the first congested station (choke point). It can send traffic to stations between the first and second choke point based on the first choke point's advertised fair rate. In general, a station can send traffic to a particular destination if it has satisfied the fair rate conditions for all choke points between itself and the destination. The maximum possible number of choke points is equal to the number of stations on the ring.



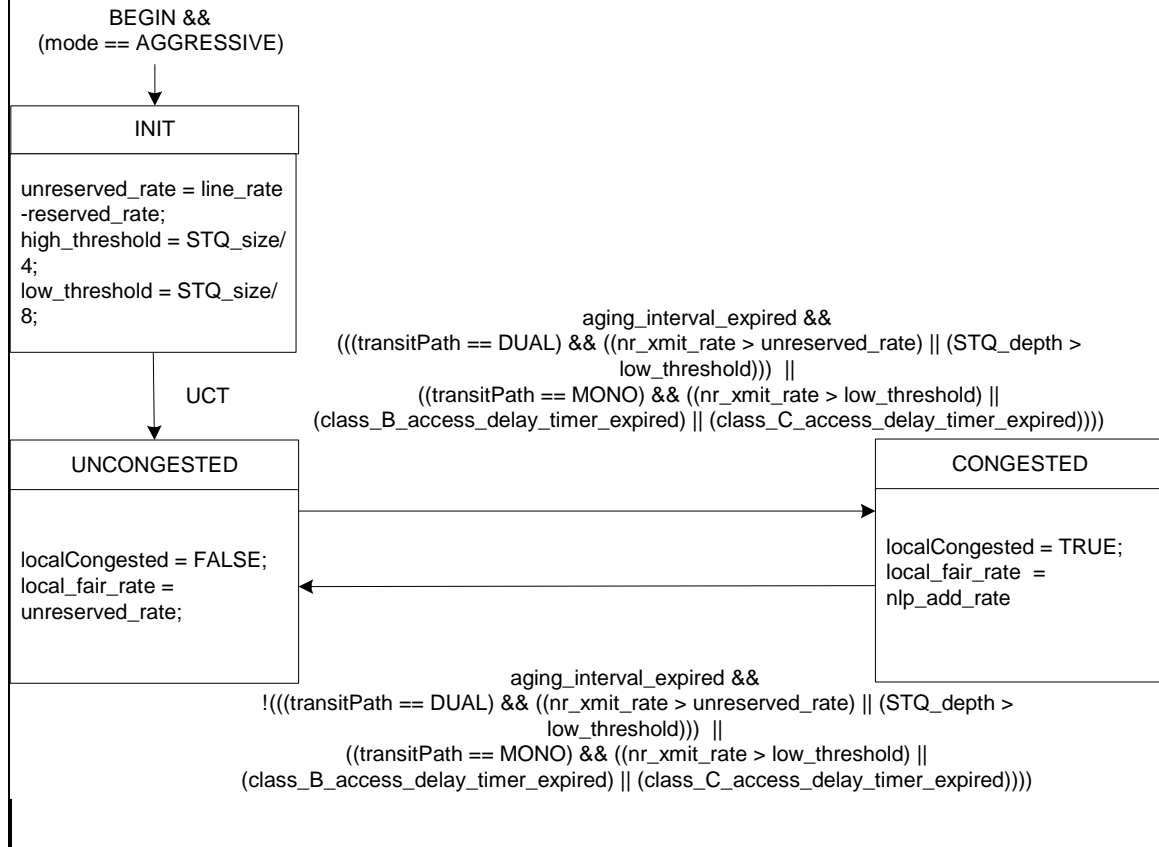


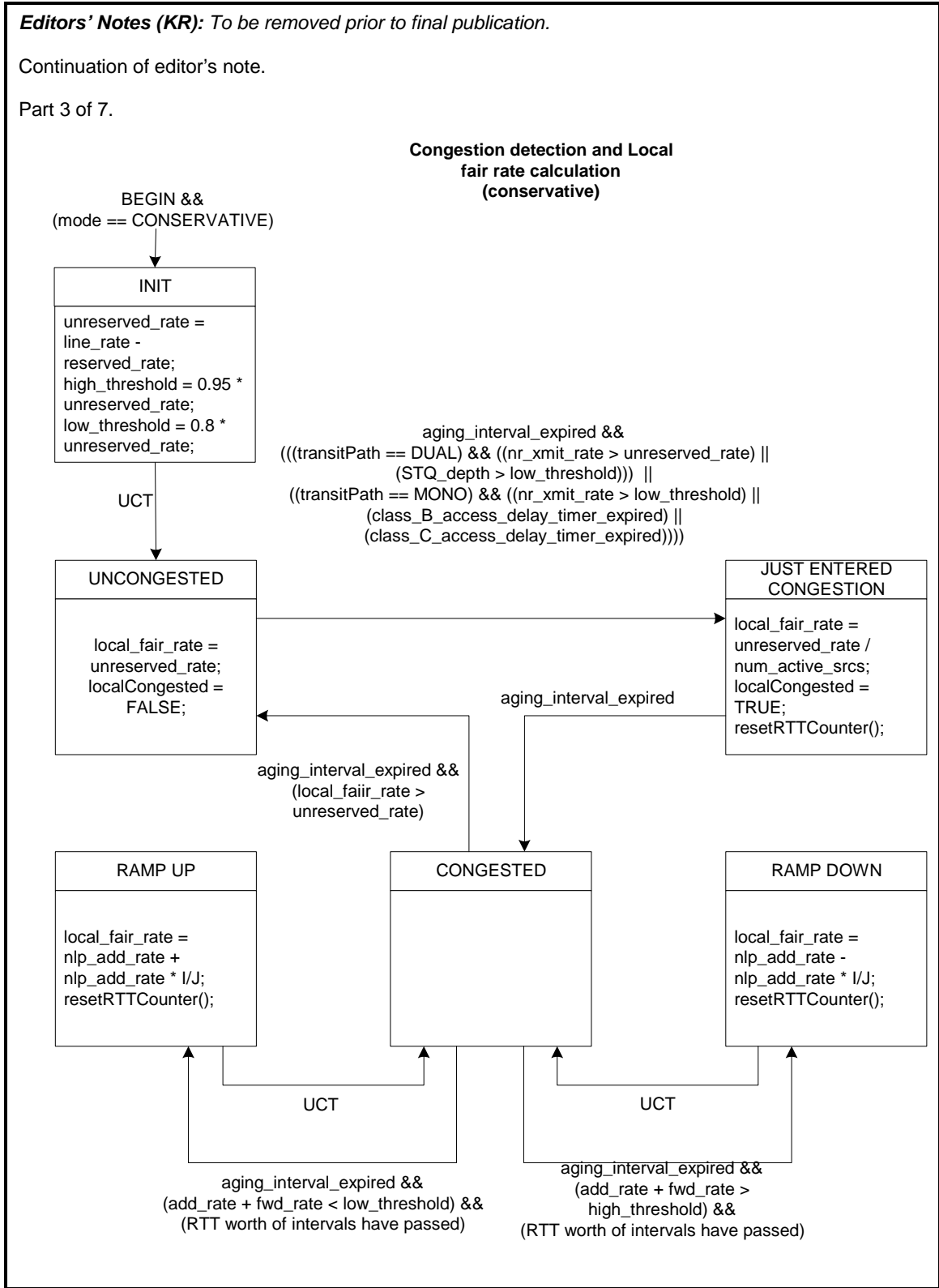
**Editors' Notes (KR):** To be removed prior to final publication.

Continuation of editor's note.

Part 2 of 7.

### Congestion detection and local fair rate calculation (aggressive)



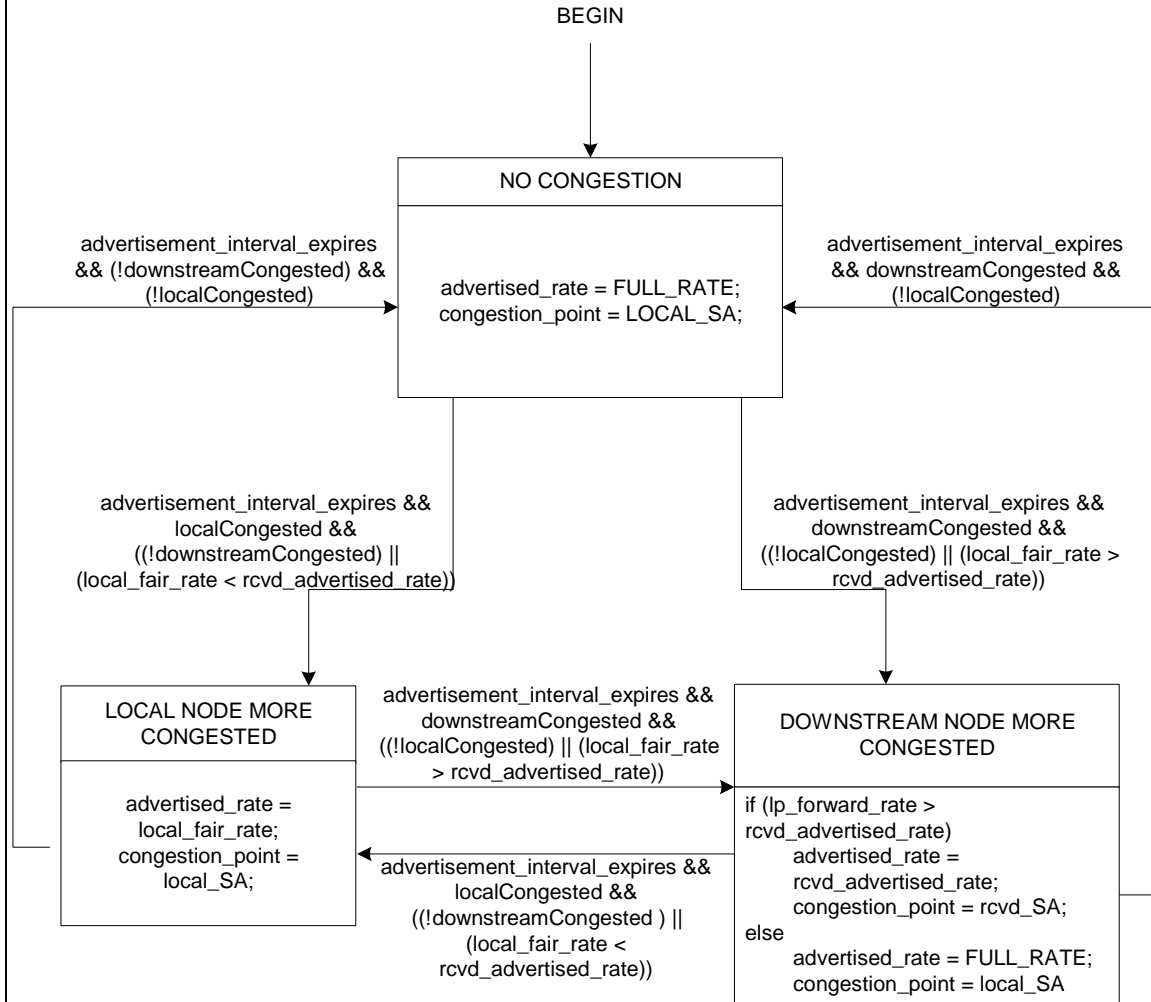


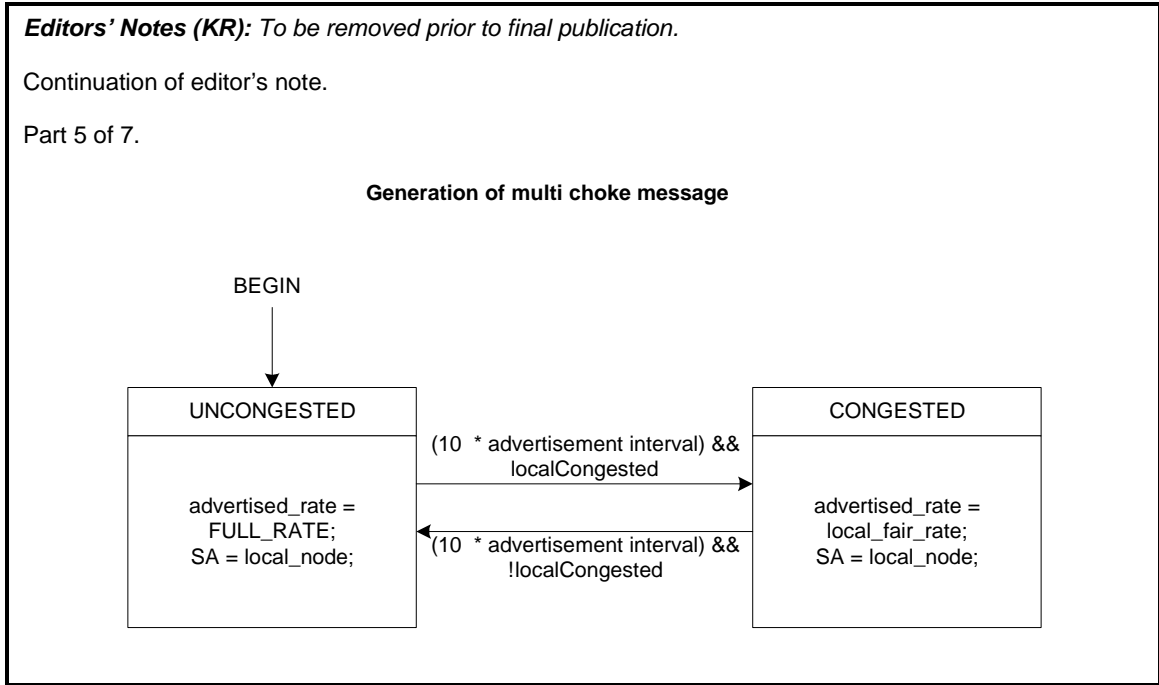
**Editors' Notes (KR):** To be removed prior to final publication.

Continuation of editor's note.

Part 4 of 7.

### Generation of single choke message

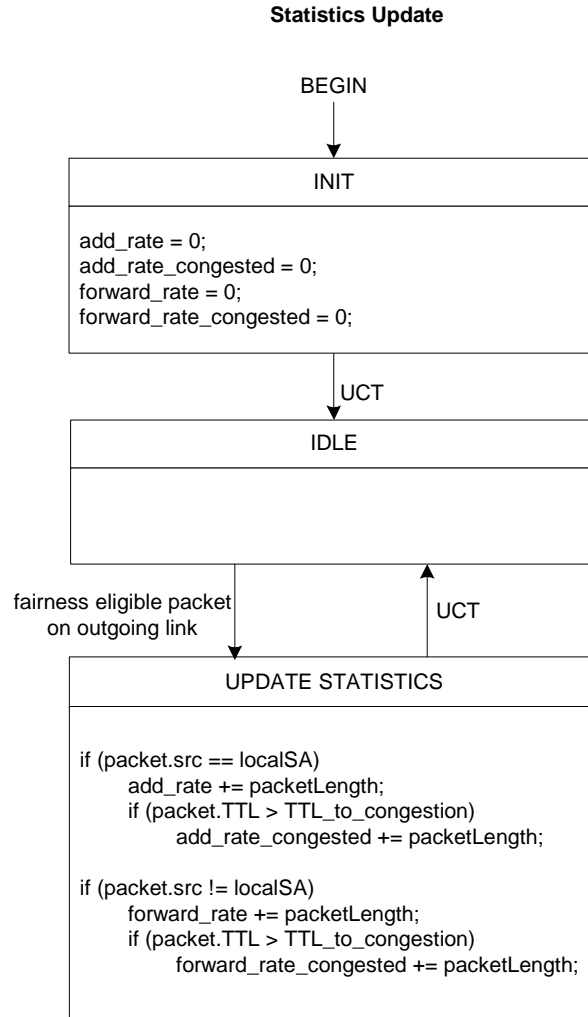




**Editors' Notes (KR):** To be removed prior to final publication.

Continuation of editor's note.

Part 7 of 7.



9.4.1 FCM Processing

9.4.1.1 FCM Receive

For every aging\_interval, all variables used to determine congestion condition, congestion location, and rate control are recomputed. Those values that depend upon values received in SC-FCMs use the most recently received SC-FCM. The flow chart and state machine for processing received FCMs are as follows:

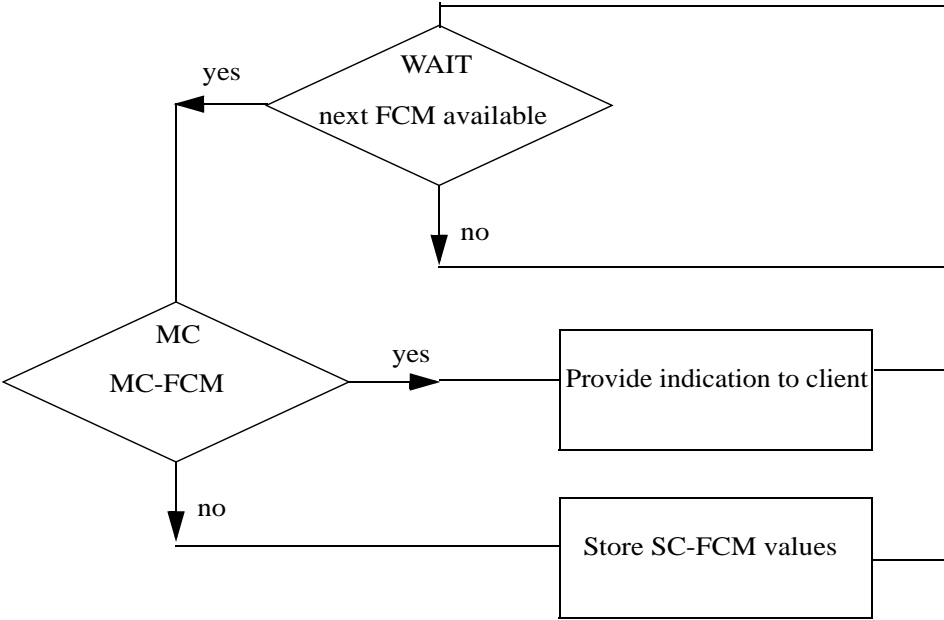


Figure 9—1— Fairness Control Message Reception

Table 9.3—FCM Reception States

Last state		Row	Next state	
state	condition		action	state
MC	MC-FCM	9.0.1	MA_CONTROL.indication(msg);	WAIT
	—	9.0.2	rcvd_fair_rate = msg.control_value; rcvd_SA = msg.SA; rcvd_TTL = msg.TTL;	
WAIT	next FCM available	9.0.3	—	MC
	—	9.0.4	—	WAIT

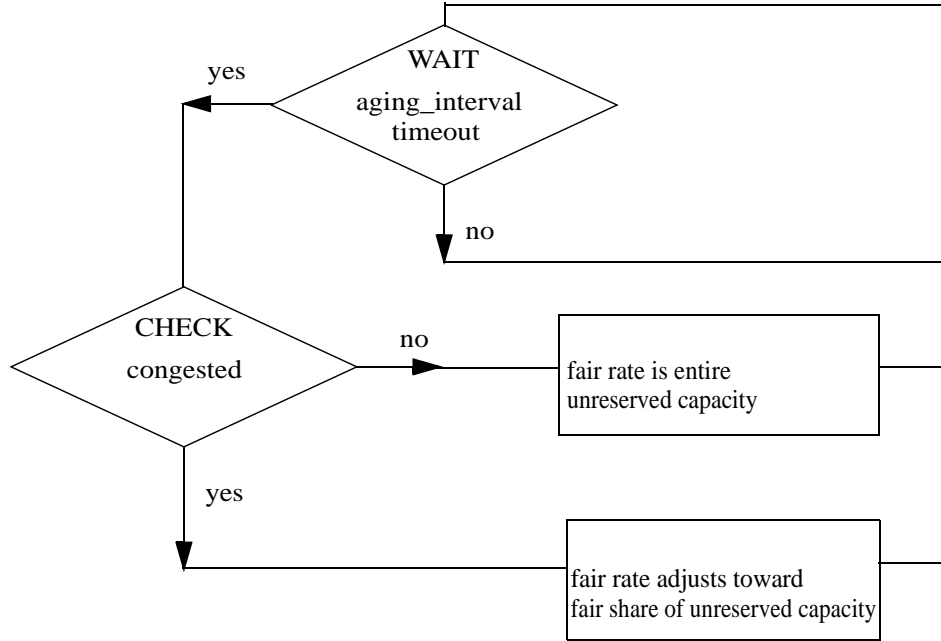
**Row 9.0.1:** The FCM is a MC-FCM. Make its relevant values available to the client.  
**Row 9.0.2:** The FCM is a SC-FCM. Store its relevant values for use after the next aging\_interval expires.

**Row 9.0.3:** When a new FCM is received, start processing it.

**Row 9.0.4:** Wait for another FCM.

#### 9.4.1.2 Local fair rate calculation

Using the most recently received SC-FCM, the local\_fair\_rate is recomputed every aging\_interval.



**Figure 9—2—Local Fair Rate Calculation**

**Table 9.4—Local Fair Rate Calculation States**

Last state		Row	Next state	
state	condition		action	state
CHECK	!congested	9.0.1	local_fair_rate = unreserved_rate;	WAIT
	aggressive	9.0.2	local_fair_rate = lp_add_rate;	
	just_entered_congestion	9.0.3	set local_fair_rate to weight adjusted equal share of unreserved_rate	
	FE transmit rate below low_threshold	9.0.4	ramp up local_fair_rate	
	FE transmit rate above high_threshold	9.0.5	ramp down local_fair_rate	
	—	9.0.6	—	
WAIT	aging_interval timeout	9.0.7	recalculate the world;	CHECK
	—	9.0.8	—	WAIT

**Row 9.0.1:** Outgoing link is not congested. Local station can try to use entire (unreserved) capacity.

**Row 9.0.2:** Aggressively try to use current add rate.

**Row 9.0.3:** When first becoming congested, try to use an equal share of the unreserved capacity, adjusted for station weight.

**Row 9.0.4:** FE transmit rate is too low. Ramp up fair rate towards unreserved capacity.

**Row 9.0.5:** FE transmit rate is too high. Ramp down fair rate towards 0.

**Row 9.0.6:** FE transmit rate is not too low and not too high. Don't change amount of outgoing link currently trying to use.

**Row 9.0.7:** When the next aging\_interval expires, recalculate all aging\_interval dependent variables, except for local\_fair\_rate and norm\_local\_fair\_rate.

**Editors' Notes (JL):** To be removed prior to final publication.

Should 9.8.3 be referenced, or should the calculations/formulas be stated in this section? Also, is it helpful to describe the calculation of local\_fair\_rate as state machine instead of just a formula?

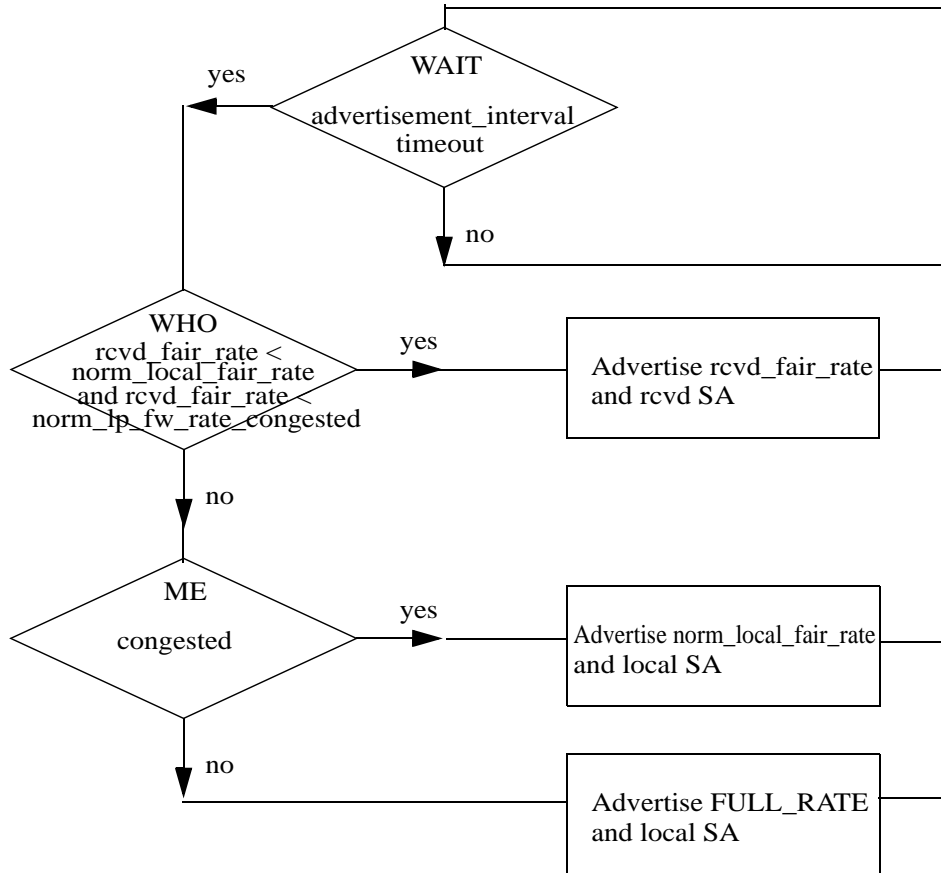
**Row 9.0.8:** Wait for the aging\_interval to expire.



## 9.4.2 FCM Transmit

### 9.4.2.1 SC-FCM Transmit

For every advertisement\_interval, the advertised\_fair\_rate and congested SA are recomputed (as described in 9.8.4 ) and advertised in a SC-FCM.



**Figure 9—3—Single Choke Fairness Control Message Advertisement**

**Table 9.5—SC-FCM Advertisement States**

Last state		Row	Next state	
state	condition		action	state
WHO	rcvd_fair_rate < norm_local_fair_rate && rcvd_fair_rate < norm_lp_fw_rate_congested	9.0.1	advertised_fair_rate = rcvd_fair_rate; congestion_point = rcvd_SA;	WAIT
	—	9.0.2	—	ME
ME	congested	9.0.3	advertised_fair_rate = norm_local_fair_rate; congestion_point = local_SA;	WAIT
	—	9.0.4	advertised_fair_rate = FULL_RATE; congestion_point = local_SA;	
WAIT	advertisement_interval timeout	9.0.5	—	WHO
	—	9.0.6	—	WAIT

**Row 9.0.1:** A downstream station is more congested than the local station.

**Row 9.0.2:** No downstream station is more congested than the local station

**Row 9.0.3:** The local station is congested.

**Row 9.0.4:** The local station (and therefore all stations) is not congested.

**Row 9.0.5:** When the next advertisement\_interval expires, transmit another SC-FCM.

**Row 9.0.6:** Wait for the advertisement\_interval to expire.

9.4.2.2 MC-FCM Transmit

For every 10 advertisement\_intervals, the local\_fair\_rate is advertised in a MC-FCM.

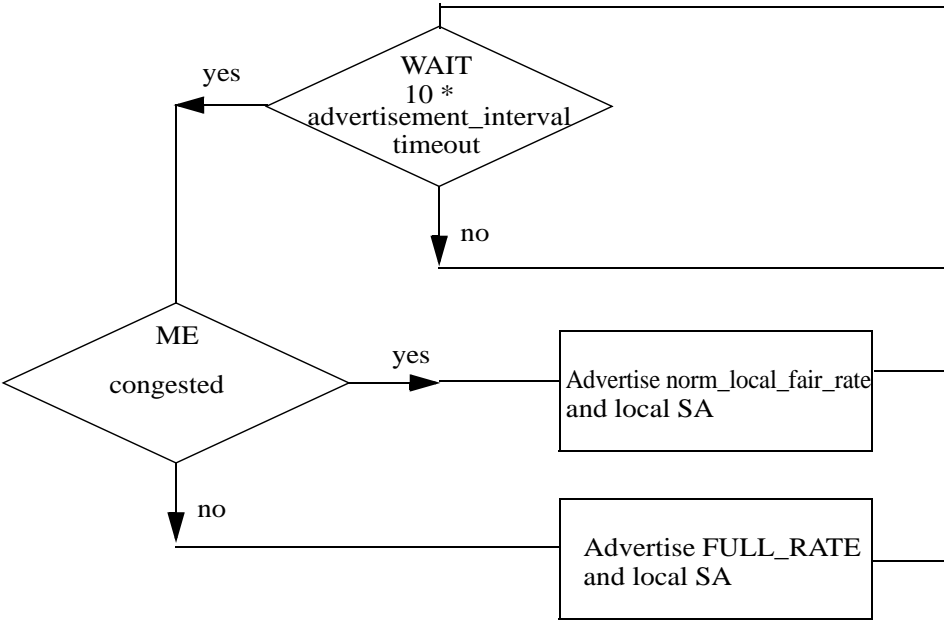


Figure 9—4—Multi Choke Fairness Control Message Advertisement

Table 9.6—MC-FCM Advertisement States

Last state		Row	Next state	
state	condition		action	state
ME	congested	9.0.1	advertised_fair_rate = norm_local_fair_rate; congestion_point = local_SA;	WAIT
	—	9.0.2	advertised_fair_rate = FULL_RATE; congestion_point = local_SA;	
WAIT	10 * advertisement_interval timeout	9.0.3	—	ME
	—	9.0.4	—	WAIT

- Row 9.0.1:** The local station is congested.
- Row 9.0.2:** The local station is not congested.
- Row 9.0.3:** When 10 advertisement\_intervals expire, transmit another MC-FCM.
- Row 9.0.4:** Wait for 10 advertisement\_intervals to expire.

## 9.5 Congestion detection

Congestion is declared by using the following criteria. The formula for determining congestion is specified in 9.8.3.4 .

Congestion is detected when any of the following conditions is true:

- a) The rate of outgoing non-reserved traffic ( $lp\_nr\_xmit\_rate$ ) is more than the  $unreserved\_rate$  (i.e.  $LINK\_RATE - reserved\_rate$ ), or the  $low\_threshold$ , if specified as a rate.
- b) The access delay timer for Class B add packets expires.
- c) The access delay timer for Class C add packets expires.
- d) The depth of the STQ exceeds the  $low\_threshold$  congestion depth threshold, or the rate of non-A0 traffic exceeds the  $low\_threshold$  congestion rate threshold.

Conditions b) and c) are optional for dual-queue MACs. Condition d) is not applicable for mono-queue MACs.

The above metrics are maintained independently for each ringlet. Any of them could indicate the onset of congestion for the ringlet. The absence of all that are maintained indicates the lack of congestion.

**Editors' Notes (JL):** *To be removed prior to final publication.*

Adding (optional) hysteresis to each of a - d need to be added.

The access delay timers are maintained by the MAC for each ringlet. There are separate access delay timers for Class B and Class C traffic. The timer corresponding to the class of packet is started at the time the packet becomes the head-of-line packet in the MAC. The timer is reset when the MAC successfully transmits that packet. If the MAC is unable to transmit the packet before the timer expires, the MAC enters a congested state. The access delay timers will have a range of 1-10 msec, with a default value of 5 msec.

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

*Is there any justification for the above numbers? It should probably depend at least on the link speed and the RTT.*

*For Class A traffic there is no access delay timer. The access delay timer expiring for Class A traffic means there is something wrong with respect to provisioning. Perhaps OAM&P should look into providing an alarm for this condition. However, this issue doesn't concern fairness.*

When a station determines that it is a congestion point, it will send a fairness control message containing a non-FULL\_RATE fair rate, based on a normalized value of its low pass filtered add rate.

### 9.5.1 Threshold settings

Three thresholds,  $full\_threshold$ ,  $high\_threshold$ , and  $low\_threshold$ , are used to specify congestion thresholds. For dual-queue MACs, they are based on occupancy of the STQ. For mono-queue MACs, they are based on the rate at which traffic is flowing through the MAC. The formulas for calculating these thresholds are provided in 9.8.1 . When the  $low\_threshold$  is crossed, the station is considered to be approaching a congested state, and advertises a reduced fair rate. When the  $high\_threshold$  is crossed, the station is considered to be in a congested state, and no more fairness eligible traffic is accepted from the MAC client for transmission on to the ringlet. For dual-queue MACs, if the occupancy of the STQ crosses  $full\_threshold$ , it is almost full, and the traffic from the STQ will be prioritized above all other traffic until the occupancy of the STQ drops below  $full\_threshold$  (by any amount).

**Editors' Notes (JL):** To be removed prior to final publication.

The remainder of this section deals with rationale, and should probably be moved to Annex I: Implementation Guidelines.

The setting of high\_threshold depends on the following factors:

- a) If set too low, it will trigger congestion reports too frequently which in turn will adversely affect the link utilization.
- b) If set higher, the difference between the high\_threshold and full\_threshold is reduced. Because this reduces the time between congestion onset and the STQ filling up, it has the following effect:
  - 1) For a given level of reclaimable Class A bandwidth (i.e. Class A1 bandwidth), the link distances must be shorter.
  - 2) For a given link distance, a smaller level of reclaimable Class A traffic can be supported. (See Clause 6.x where this calculation is provided.)

The more severe consequences of a) vs. b) are the reason that a threshold of less than 50% is chosen. This recommendation applies only when a STQ of many MTUs is used to support high amounts of reclaimable Class A bandwidth.

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

As requested at the May meeting, Vasan Karighattam will have a contribution for an analytical method to compute the STQ size.

The primary transit queue needs to hold at least 2 MTUs.

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

As requested at the May meeting, Necdet Uzun and David James will have a justification for this threshold.

NOTE—high\_threshold should be set to about 25% of the total buffer available. full\_threshold should be set to at least 1 MTU less than the total buffer size. The recommendation for the setting of the threshold values is done based on delivering the best possible end-to-end delay for the low priority traffic without penalizing the Class A traffic. Lower values will result in higher end-to-end delays for low priority data packets. If either Class A, B, or C traffic is extremely bursty, then a lower threshold value should be considered. A full\_threshold of at least 1 MTU less than the total buffer size is needed to ensure that the station has sufficient buffer space for a packet that may arrive on the transit path. If the Class A traffic has a bursty nature, a more conservative (i.e. lower) value of high\_threshold is recommended in order to avoid overflow of the STQ.

## 9.6 Interaction with data path

RPR-FA uses the add\_rate, add\_rate\_congested, fw\_rate, and fw\_rate\_congested values provided by the data path, as specified in Clause 6.

RPR-FA provides the add\_rate\_ok, add\_rate\_congested\_ok, allowed\_rate, allowed\_rate\_congested, and TTL\_to\_congestion values to the data path rate control function to police the fairness eligible traffic that is added by the station to the ringlet. The add\_rate\_ok and add\_rate\_ok\_congested variables are used to enable or disable the sendC signal and to determine if the TTL\_to\_congestion value should be applied. The TTL\_to\_congestion variable stores the value of the distance to the present or more recent choke point. The allowed\_rate\_congested variable stores the value of the fair rate for fairness eligible traffic passing through the congestion\_point, as given by the most recently received Single Choke fairness control message. The

allowed\_rate variable stores the value of the rate for all fairness eligible traffic. The data path maintains rate shapers to limit the add\_rate (Sd) and add\_rate\_congested (Sc) values below the provided allowed\_rate and allowed\_rate\_congested values, respectively, and to shape their output at or below the provided rates to smooth bursty transmissions.

**Editors' Notes (JL):** To be removed prior to final publication.

This following paragraph belongs in clause 6.

The RPR-FA dynamic shaper consists of a token bucket. The token bucket has a default maximum depth of MTU bytes. The tokens are generated at a rate equal to the allowed\_rate\_congested. For every byte that is transmitted, the number of tokens is decremented by one. When a packet is waiting for access to the ring at the head of the queue, it will be granted ring access only if its rate conforms to the fair rate as computed by RPR-FA and there is at least one byte in the token bucket. It is therefore possible for the number of tokens to become negative after subtracting the number of bytes in the transmitted packet. The maximum number of tokens that can accumulate in the token bucket is equal to the size of an MTU.

## 9.7 Fairness control messages

### 9.7.1 Generation of fairness control messages

Single Choke and Multi Choke fairness control messages shall be sent periodically to propagate fair rate information to upstream stations. Single Choke messages are sent every advertisement\_interval. The recommended advertisement\_interval is between the transmission delay for a single MTU and one-half of the aging\_interval. Multi Choke messages are sent every 10 advertisement\_intervals.

### 9.7.2 Impact of lost fairness control messages

The allowed\_rate\_congested maintains its current value until the next Single Choke fairness control message is received. Therefore, the loss of a Single Choke fairness control message would result in one of the following. If the lost message contained a new value of FULL\_RATE, it would prevent the station from ramping up its allowed rate for the ringlet. If the value contained in the lost message was a new non-FULL\_RATE value, the MAC will not be able to react by setting its allowed rate to the received advertised rate, which could have been higher or lower than the current allowed rate. If the lost message did not contain a new value from the previous advertisement, there will be no effect on the fairness algorithm.

Refer to Clause 10 for protection implications of not receiving consecutive Single Choke fairness control messages.

The Multi Choke messages are not processed by the MAC, and therefore loss of Multi Choke messages doesn't impact behavior of the MAC.

### 9.7.3 Validation of fairness control messages

**Editors' Notes (JL):** To be removed prior to final publication.

This following paragraph probably belongs in clause 6 and should be generalized for all control messages.

If the source of the rcvd\_fair\_rate is the same station that received it then the rcvd\_fair\_rate will be disregarded. When comparing the rcvd\_fair\_rate source address, the ring identifier of the fairness control mes-

sage must match the receiver’s ring identifier in order to qualify as a valid compare. The exception is if the receive station is in the wrap state, in which case the fairness control message’s ring identifier is ignored.

9.7.4 Packet format

The fairness control message format is as shown in Figure 9.5. The fairness protocol determines the number of hops to the station that originated the fairness control message by subtracting the TTL in the received message from 256. Therefore all fairness control messages shall be sourced with a TTL of 255. And all single choke fairness control messages shall reset the TTL to 255 upon changing the value of the SA to the local SA.

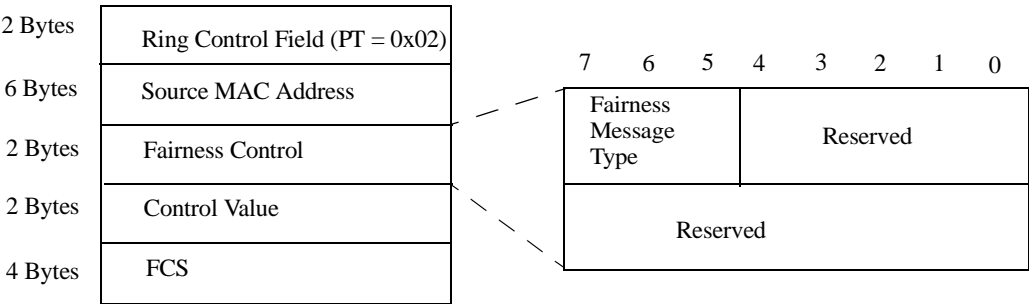


Figure 9.5—fairness control message Format

9.7.4.1 Fairness Control Message Type (3 bits)

This field specifies the type of fairness control message. Table 9.7 shows the values of the Fairness Control Message Type. Single Choke fairness control messages are used to implement RPR-FA and provide information about the most congested span on the ringlet. Multi Choke fairness control messages are needed to support a multi-choke implementation.

Table 9.7—Fairness Control Message Type values

Value (binary)	Type of fairness control message	How is it used
000	Single Choke	Generated by the FCU every advertisement_interval. SA is the MAC address of the most congested station in the fairness domain

**Table 9.7—Fairness Control Message Type values**

Value (binary)	Type of fairness control message	How is it used
001	Multi Choke	Generated by the FCU every 10 advertisement_intervals. The SA is the SA of the MAC, and the messages are broadcast. The received fair rate is passed to the client.
010 to 111	Reserved	For future use.

**9.7.4.2 Reserved field (13 bits)**

These bits are ignored on receipt and set to zero on transmit for both Single Choke and Multi Choke fairness control messages.

**9.7.4.3 Control value (16 bits)**

This field carries the fair rate encoded as a 16-bit quantity. A fair rate of all ones indicates a value of FULL\_RATE which corresponds to advertising a fair rate equal to the line rate. The fair rate, when not FULL\_RATE, indicates the number of bytes that a station may add to the ringlet during the next aging\_interval. The reported fair rate is a rate that has been normalized by the normalization coefficient (NORMCOEF), to normalize the number of aging\_intervals incorporated into the rate (AGECOEF), the local station weight (WEIGHT), and the rate coefficient (RATECOEF).

**9.8 Fairness algorithm formulas**

**Editors' Notes (JL):** To be removed prior to final publication.

A large amount of the following section probably belongs in clause I. Some of the following will probably remain as equations inlined in the text to augment descriptions that need a more precise definition.

This state machines in this clause use the following formulas:

**9.8.1 Initialization**

At the beginning of the state machine, initialize:

- allowed\_rate
- full\_threshold
- high\_threshold
- low\_threshold
- unreserved\_rate

```
allowed_rate = MAX_ALLOWED_RATE;
```



```

unreserved_rate = LINK_RATE - reserved_rate;

if (dual_queue_MAC)
    full_threshold = STQ_SIZE - MTU_SIZE; // can subtract more than MTU_SIZE

if (dual_queue_MAC && aggressive)
    high_threshold = STQ_SIZE / 4;
else
    high_threshold = .95 * unreserved_rate;

if (dual_queue_MAC && aggressive)
    low_threshold = STQ_SIZE / 8;
else
    low_threshold = .8 * unreserved_rate;

```

### 9.8.2 Per byte

**Editors' Notes (JL):** To be removed prior to final publication.

This section (per byte operations) should be moved to Clause 6.

For each byte transmitted, update:

- add\_rate
- add\_rate\_congested
- fw\_rate
- fw\_rate\_congested
- nr\_xmit\_rate

```

if (FE byte added by local station)
    add_rate = add_rate + 1;
if (FE byte added by local station && sent beyond congestion_point)
    add_rate_congested = add_rate_congested + 1;
if (FE byte forwarded by local station)
    fw_rate = fw_rate + 1;
if (FE byte forwarded by local station && sent beyond congestion_point)
    fw_rate_congested = fw_rate_congested + 1;
if (non-A0 byte transmitted by local station)
    nr_xmit_rate = nr_xmit_rate + 1;

```

### 9.8.3 Per aging\_interval

For each aging\_interval, based upon the last received FCM, update:

- active\_stations
- add\_rate
- add\_rate\_ok
- add\_rate\_congested
- add\_rate\_congested\_ok
- fw\_rate
- fw\_rate\_congested
- nr\_xmit\_rate

```

— congested 1
— allowed_rate 2
— allowed_rate_congested 3
— congestion_point 4
— TTL_to_congestion 5
— local_fair_rate 6
— norm_local_fair_rate 7
— lp_add_rate 8
— lp_add_rate_congested 9
— lp_fw_rate 10
— lp_fw_rate_congested 11
— lp_nr_xmit_rate 12
— norm_lp_fw_rate_congested 13

```

### 9.8.3.1 active\_stations

```

for (active_stations = 0; station = 0; station < num_stations; station++)
    if (sent_FE_during_aging_interval(station))
        active_stations = active_stations + 1;

```

### 9.8.3.2 lp\_add\_rate, lp\_add\_rate\_congested, lp\_fw\_rate, lp\_fw\_rate\_congested, lp\_nr\_xmit\_rate, norm\_lp\_fw\_rate\_congested

The following low pass filterings must be done before the first pass agings specified in 9.8.3.3 .

```

lp_add_rate = ((LPCOE-1) * lp_add_rate + add_rate) / LPCOE;
lp_add_rate_congested = ((LPCOE-1) * lp_add_rate_congested +
    add_rate_congested) / LPCOE;
lp_fw_rate = ((LPCOE-1) * lp_fw_rate + fw_rate) / LPCOE;
lp_fw_rate_congested = ((LPCOE-1) * lp_fw_rate_congested +
    fw_rate_congested) / LPCOE;
lp_nr_xmit_rate = ((LPCOE-1) * lp_nr_xmit_rate + nr_xmit_rate) / LPCOE;
norm_lp_fw_rate_congested = lp_fw_rate_congested / NORMCOEF;

// alternatively, each can be simplified as shown for lp_add_rate
// lp_add_rate = lp_add_rate -
//     lp_add_rate >> (log2(LPCOE)) + add_rate >> (log2(LPCOE));

```

### 9.8.3.3 add\_rate, add\_rate\_congested, fw\_rate, fw\_rate\_congested, nr\_xmit\_rate

The following first pass agings must be done after the low pass filterings specified in 9.8.3.2 .

```

add_rate = (add_rate * (AGECOEF - 1)) / AGECOEF;
add_rate_congested = (add_rate_congested * (AGECOEF - 1)) / AGECOEF;
fw_rate = (fw_rate * (AGECOEF - 1)) / AGECOEF;
fw_rate_congested = (fw_rate_congested * (AGECOEF-1)) / AGECOEF;
nr_xmit_rate = (nr_xmit_rate * (AGECOEF - 1)) / AGECOEF;

```

**9.8.3.4 congested**

**Editors' Notes (JL):** To be removed prior to final publication.

This does not currently take into account any hysteresis. This needs to be added. The hysteresis should allow delta to be set to a range of values including 0 (no hysteresis).

```

if (dual_queue_MAC)
{
    if (lp_nr_xmit_rate > unreserved_rate)
    || (STQ_depth > low_threshold))
        congested = TRUE;
    else if (aggressive) // clearing of congested for conservative mode
        congested = FALSE; // is handled in calculation for local_fair_rate
}
else // mono_queue_MAC
{
    if (lp_nr_xmit_rate > unreserved_rate)
    || (lp_nr_xmit_rate > low_threshold)
    || (access delay timer for add class B expired)
    || (access delay timer for add class C expired))
        congested = TRUE;
    else if (aggressive) // clearing of congested for conservative mode
        congested = FALSE; // is handled in calculation for local_fair_rate
}

```

**9.8.3.5 allowed\_rate**

```

if (aggressive)
    allowed_rate = MAX_ALLOWED_RATE;
else // conservative
{
    if (congested) // local station is congested
        allowed_rate = local_fair_rate;
    else // local station is not congested
        allowed_rate = allowed_rate + (MAX_ALLOWED_RATE - allowed_rate) / RAMPCOEFF;
}

```

**9.8.3.6 allowed\_rate\_congested, congestion\_point, TTL\_to\_congestion**

```

if (rcvd_fair_rate != FULL_RATE)
{
    allowed_rate_congested = rcvd_fair_rate * NORMCOEFF;
    congestion_point = rcvd_SA;
    TTL_to_congestion = 256 - rcvd_TTL;
}
else // no congestion downstream, ramp up
{
    allowed_rate_congested = allowed_rate_congested +
        (MAX_ALLOWED_RATE - allowed_rate_congested) / RAMPCOEFF;
    congestion_point = local_SA;
    TTL_to_congestion = TTL_to_congestion;
    // TTL_to_congestion should be left at the last congestion point so that
    // the allowed_rate is not all of a sudden presented to the former
    // congestion point. One could add the following statement:
    // if (allowed_rate_congested == allowed_rate)
    //     TTL_to_congestion = 255;
    // But it is unnecessary since it has no effect upon the rate at which

```

```

    // the station may transmit.
}

```

### 9.8.3.7 local\_fair\_rate, norm\_local\_fair\_rate

```

if (aggressive)
{
    if (congested)
        local_fair_rate = lp_add_rate;
    else // not congested
        local_fair_rate = unreserved_rate;
}
else // conservative
{
    if (congested)
    {
        if (just_entered_congested) // first time when congested
            local_fair_rate = (unreserved_rate / active_stations) * WEIGHT;
            // optionally, instead of active_stations, which uses a weight of 1
            // for each of the other stations, one could use sum(ActiveWeights)
        else // was congested last time
        {
            if ((add_rate + fw_rate < low_threshold)
                && (RTT worth of aging_intervals passed since last update))
            {
                local_fair_rate = min(unreserved_rate,
                    local_fair_rate +
                    (unreserved_rate - local_fair_rate) / RAMPCOEFF;
                if (local_fair_rate >= unreserved_rate)
                    congested = FALSE;
            }
            else if ((add_rate + fw_rate > high_threshold)
                && (RTT worth of aging_intervals passed since last update))
                local_fair_rate = local_fair_rate - local_fair_rate / RAMPCOEFF;
            else
                local_fair_rate = local_fair_rate; // no change
        }
    }
    else // not congested
        local_fair_rate = unreserved_rate;
}

norm_local_fair_rate = local_fair_rate / NORMCOEFF;

```

### 9.8.3.8 add\_rate\_ok, add\_rate\_congested\_ok

```

add_rate_ok =
    (add_rate < allowed_rate) \\ current allowance is not exceeded
    && (lp_nr_xmit_rate < unreserved_rate) \\ space for reserved traffic
    && ((STQ_depth = 0) \\ upstream stations are not in need
        || ((fw_rate > add_rate) \\ upstream stations are not starved
            && (STQ_depth < high_threshold))); \\ node is not fully congested

```

NOTE—An optional optimization is to change (fw\_rate > add\_rate) to (fw\_rate > add\_rate/WEIGHT). This is not needed under almost all circumstances. The one exception is if the sum of weights of upstream stations is less than the weight of the local station. If this optimization is chosen, the implementation may require that WEIGHT values be factors of 2.

```
add_rate_congested_ok = add_rate_ok
                        && (add_rate_congestion < allowed_rate_congestion);
```

### 9.8.3.9 sendC

**Editors' Notes (JL):** To be removed prior to final publication.

This section (sendC) should be moved to Clause 6.

```
if (!add_rate_ok)
    sendC = 0;
else if (!add_rate_congested_ok)
    sendC = TTL_to_congestion;
else
    // see comment in 9.8.3.6 on TTL_to_congestion
    sendC = TTL_to_congestion; // calculation with no downstream congestion
```

### 9.8.4 Per advertisement\_interval

For each advertisement\_interval, update:

— advertised\_fair\_rate

```
if ((rcvd_fair_rate < norm_local_fair_rate) // downstream is more congested,
    && (rcvd_fair_rate < norm_lp_fw_rate_congested)) // upstream is contributing
    advertised_fair_rate = rcvd_fair_rate; // to downstream congestion
else if (congested) // local station is congested (and more than downstream)
    advertised_fair_rate = norm_local_fair_rate;
else // no downstream or local congestion
    advertised_fair_rate = FULL_RATE;
```

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