
IEEE 802.11
Wireless Access Method and Physical Specification

Title: Priority in CSMA/CA to support distributed Time-Bounded Services.

Prepared by:

Wim Diepstraten
WCND-Utrecht
AT&T-GIS (NCR)
Nieuwegein The Netherlands
Tel: (31)-3402-76482
Fax: (31)-3402-39125
Email: Wim.Diepstraten@utrecht.ncr.com

Abstract: This paper introduces priority mechanisms in CSMA/CA as can be utilized to support Distributed Time Bounded Services. Possible priority mechanisms are already referred to in Doc. IEEE P802.11-93/190, but they were not yet explained and specified.

Applicability to the "Foundation MAC":

Because priority mechanisms were already referred to in Doc. 190, the priority mechanism subject can be seen as further detail of Doc 190. Application of this mechanism for the purpose of providing Distributed Time Bounded Service is to be regarded as new functionality.

Introduction

The concept of Distributed Time Bounded Service was first introduced in 802.11 by Kerry Lynn in the January meeting [2]. This was the result of work that was done within the ETSI RES10 HIPERLAN committee.

The DFWMAC supports Time Bounded Services by using an optional PCF on top of the basic DCF function. Its characteristics are that it is a connection oriented service, that is based on bandwidth reservation controlled in the AP by the PCF.

As discussed in doc 190, a major disadvantage of a PCF approach is that it does not support multiple overlapping PCF's on the same channel. Timing coordination between the PCF's would be required (across the DS) to prevent overlap situations.

This does limit the practical application of Time Bounded Services because full isolation between different overlapping BSS's will be difficult to achieve in a large installation.

This is true for the 2.4 GHz ISM band, the committee is currently focusing on, but it will be even more true for high speed operation in future bands like the 1.9 GHz (Async part) or possibly the 5.2 GHz band. Spreading is not needed in these bands, but when a factor 5-10 higher PHY speeds are developed (needing more bandwidth), then the number

of available channels will again be limited, so that overlap in the same band can not be prevented.

The latter is already demonstrated by HIPERLAN, where 4-6 high speed (10-20 Mbps) channels are expected to be used in the available 100-150 MHz bandwidth. This is not sufficient to obtain enough channel isolation in all installations.

Time Bounded Service requirement evaluation

The main characteristic for Time Bounded Service is that the (end to end) transfer delay should be bounded within a given limit. In addition for certain applications, the delay variance should be limited.

An example of this is Voice, where on the sender end, a time regular stream of bytes are generated by the voice encoder, which can be packed into a regular stream of frames at a given framing period (superframe).

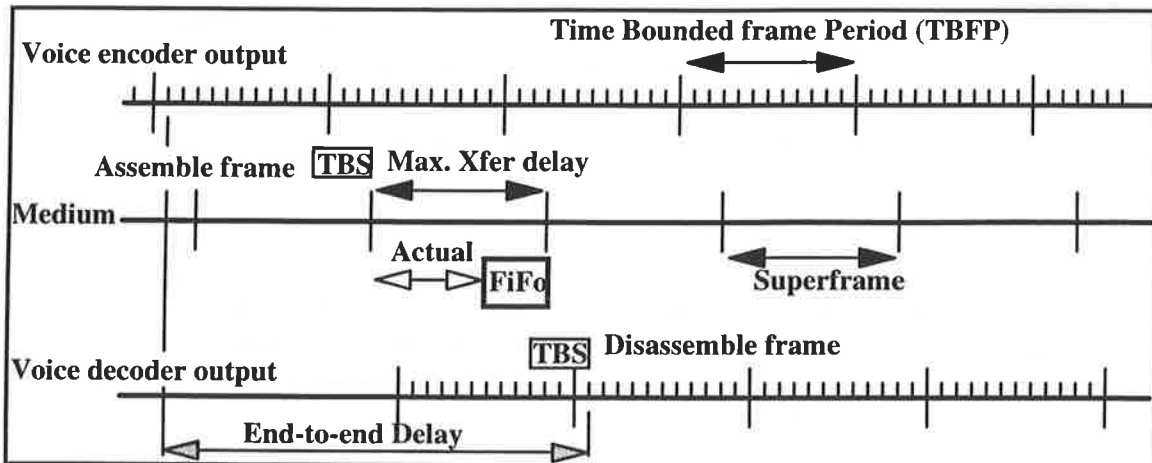


Fig 1: Voice TBS application

The regular stream of frames are offered to the MAC for delivery in a time bounded fashion. The actual transfer delay however can vary due to contention on the medium. At the receiver end the time regularity should be restored, to feed the voice decoder with a regular stream of bytes. This timing can be restored at the receiver, by utilising a FiFo or a simple buffer that can bridge the delay variance.

Please note that in the illustrated example, the frame assembly and disassembly and the timing restore FiFo is something above the MAC.

Also the medium can be several segments with some of them being wireless. Whether or not re-timing is done on a per segment basis or end-to-end only is something to be determined by I think the Link layer immediately above the MAC.

The functionality that needs to be provided by the MAC is a delivery service with a more strictly bounded delay. What needs to be done when the delay is larger then the assumed limit will be application dependent. For instance a voice application may decide to drop the frame, while an industrial application may have to do something different with it.

Question is what kind of parameters needs to be provided by the MAC to support this. One such parameter could for instance be an indication of the transfer delay that the frame encountered on that segment.

Distributed Time Bounded Service (DTBS).

Distributed Time Bounded Service can be characterised as a delivery service that does not rely on any bandwidth reservation mechanism, and (therefore) does not require a PCF. It therefore does not have the overlap limitations of a PCF.

The service uses the same distributed access mechanism as the Asynchronous Service, but with a higher access priority. Priority access can decrease the average delay and has a delay probability distribution (as function of the load) such that the probability that a certain maximum delay will be exceeded is acceptable low for the service. One of the factors that needs to be controlled is the total DTBS load that is offered to the medium at any one time.

DTBS can be provided when the DCF does support priority access. A total of two access priority levels seems sufficient to support the service.

Priority access in CSMA/CA

There are several mechanisms with which a relative access priority can be created from one traffic stream compared to an other. In general they are referred to on page 25 of the DFWMAC document, but they are not specified as such.

The following mechanism can be used to create relative priority levels in CSMA/CA. Using the basic CSMA/CA access scheme, the following parameters can be changed between the different priority levels:

- DIFS This is the minimum time that the medium should be idle before the medium can be accessed, and before a possible backoff delay will elapse.
- CW This is the window from which the backoff delay is calculated, and will control the average delay between two subsequent frames when a station will defer for a busy medium.
- CW increase This will determine the relative priority of the retransmissions.

Different DIFS periods:

Two different IFS periods can be used, such that the high priority traffic has a shorter IFS then the low priority traffic $HPIFS < LPIFS$. The effect of this is that low priority traffic will defer earlier. In addition the backoff counter of the high priority traffic will decrease earlier then for the low priority traffic. The net effect is that the probability that a high priority frame gets through is higher then the probability that a low priority frame achieves access.

Note that the backoff counter does not decrease when the medium is busy.

When the difference between the HPIFS and the LPIFS is larger than the high priority contention window, then there will be an absolute priority difference between the two, such that when there is contention, then the high priority traffic will always win.

Different Contention windows.

Different contention windows do translate in a relative priority difference when $HP_{cw} < LP_{cw}$. This will cause that the average backoff delay that is selected during the defer will be longer for the low priority traffic. Please note that the longer a frame is deferring, the higher will be its relative priority compared to the previous access try. It should be noted that these parameters do also have effect on the collision probability. The minimum CW is a tradeoff between collision probability and effect on throughput and delay distribution as function of priority.

CW increase difference.

The third mechanism that we can play with is the contention window increase policy for every retransmission. An exponential CW increase is used for Asynchronous traffic to assure stability under high loads. This is undesirable for Time Bounded traffic, because a retransmission means that again a certain access delay will be needed to recover from the failed transmission. Therefore the policy is that the HP_{cw} should not increase. It would be better when we could decrease the HP_{cw} for every retransmission, because that would give the retransmission an even higher priority compared to both the low priority traffic and the initial high priority traffic.

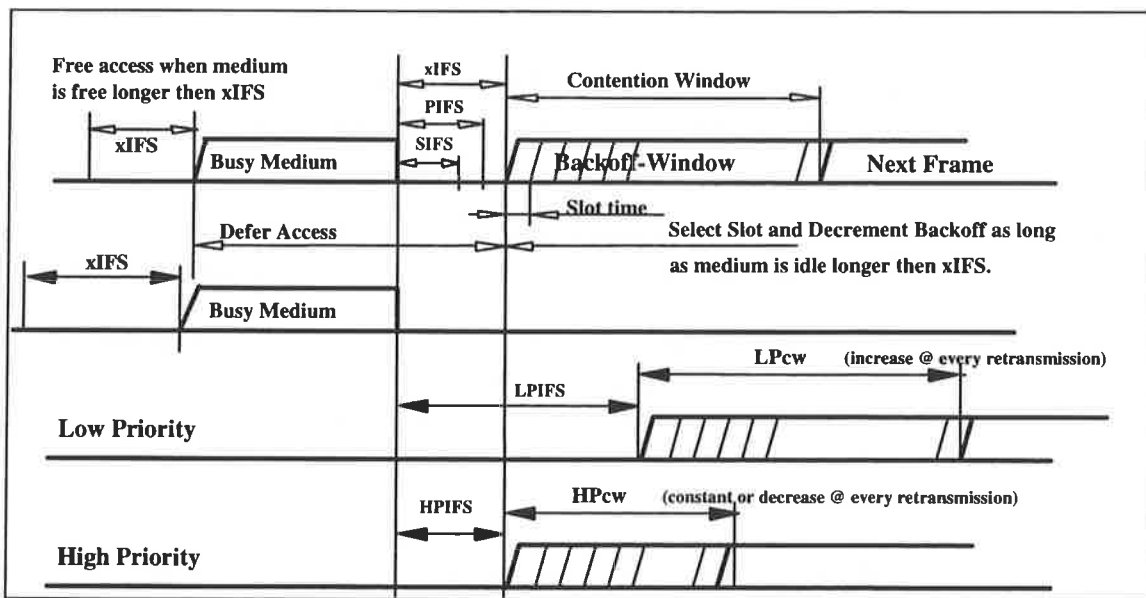


Fig 2: CSMA/CA with relative access priority

The CW increase algorithm is of prime importance for the establishment of priority in an access scheme where the backoff delay is decremented independent of the "medium Busy" status. It is of less importance for the current Foundation approach.

This mechanism assumes that the DTBS load is limited by some other means to a level that provides sufficient additional capacity to allow sharing the medium with "Bursty" asynchronous traffic.

Simulation conditions:

The following are the simulation conditions that were run:

The nominal parameters used during the simulations are:

- HPIFS = 2 slots
- LPIFS is controlled between HPIFS+0, HPIFS+16 and HPIFS+32
- LPcw is controlled between 32 and 64
- HPcw is 32 with either a constant or a decreasing HPcw on retransmission.

The PHY parameters used during the simulation are a 23 usec slot time and 150 usec PHY training time. A 2 Mbps channel speed is used.

The load consist of a 60%/40% frame length distribution of Short frames of 64 Bytes and Long frames of 576 Bytes.

The configuration is normally 3 high priority stations together with 6 low priority stations.

The figures show primarily delay and delay distribution curves as function of the offered load, for each of the priority levels.

Simulation interpretation:

It should be noted that the simulation model is such that a new frame is not generated until the previous frame is transmitted by the MAC. The load is controlled by changing the average time interval that a new frame is generated per station.

The effect of this is that the high priority to low priority throughput ratio will increase when the load is increased, simply because the transfer delay of the high priority traffic is decreased, which results in an earlier generation of the next frame that is offered to the MAC.

A separate simulation is run that assumes a fixed (maximum) DTBS load with a variable low priority load. This is considered a more realistic situation when we consider a continuous time regular stream of frames of for instance a voice application.

A simulation is run using a Superframe of 30 msec, and a fixed high priority load equivalent to 3 ADPCM coded voice channels of 32 Kbps each (simplex). The results do not show "Total Throughput" performance, but in general the throughput is inversely proportional to the average delay.

Delay characteristics:

The effect of priority can be shown when looking at the delay characteristics. The characteristic behaviour of this access mechanism is such that there is no increase of the transfer delay at low loads (see fig 3) for low priority only operation. For higher loads some effect is seen that will generally result in a slightly lower throughput. Fig 4 shows that the priority difference will increase with increasing load, such that the transfer delay of the high priority traffic will increase much slower with increased load than the transfer delay of the low priority traffic. Fig 4 shows double traces at high and low priority. These are of two separate runs with different HPcw increase policies. It shows a constant HPcw and a decreasing HPcw on every retransmission.

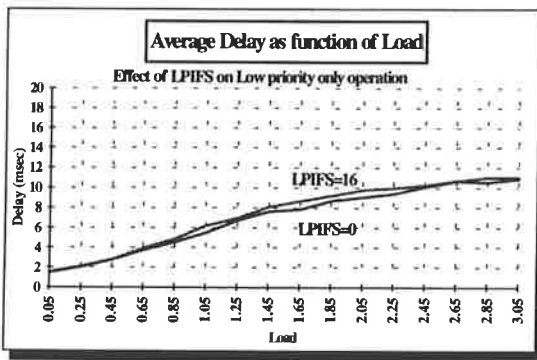


Fig 3: Effect of LPIFS on low priority only.

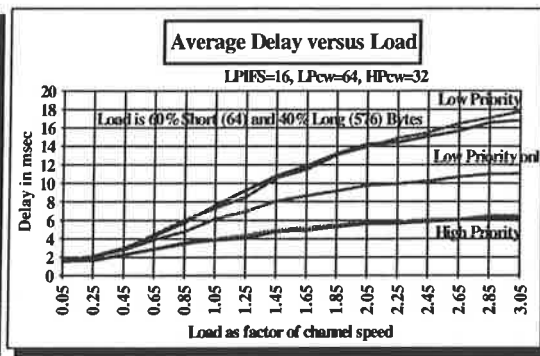


Fig 4: Effect of priority.

Figures 5 and 6 show the delay distribution curves as function of the load for the high and low priority traffic respectively. This is one of the most important curves to judge the Time Bounded maximum transfer delay performance. The X-axis shows the delay in units of 1000 Bytes, where 5 is equivalent to 20 msec in a 2 Mbps environment. It should be understood that the curves show the delay distribution at a constant given load. This means that for instance the actual probability that the delay exceeds 20 msec is the percentage in the figure multiplied with the probability for that (over)load.

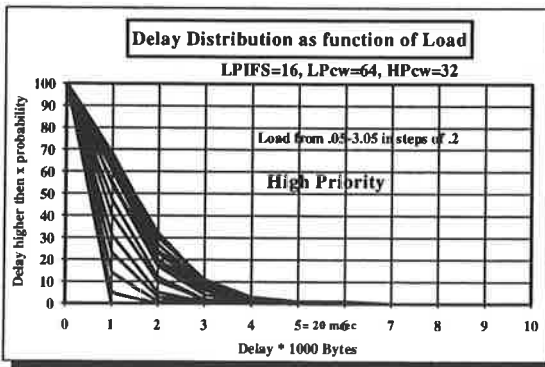


Fig 5: High priority delay variance.

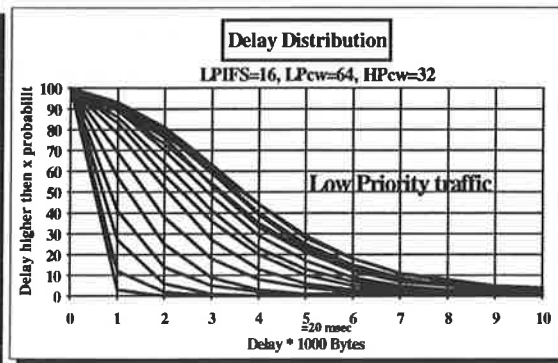


Fig 6: Low priority delay variance.

The figures 7-10 show similar curves for an alternative setting of the priority relevant parameters. the settings are such that the contention windows do not overlap, so that the high priority has absolute priority over the low priority.

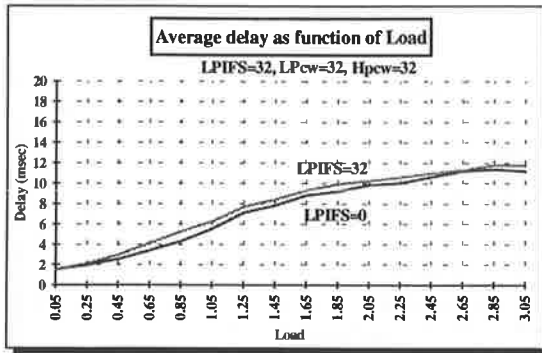


Fig 7: Effect of LPIFS on low priority only.

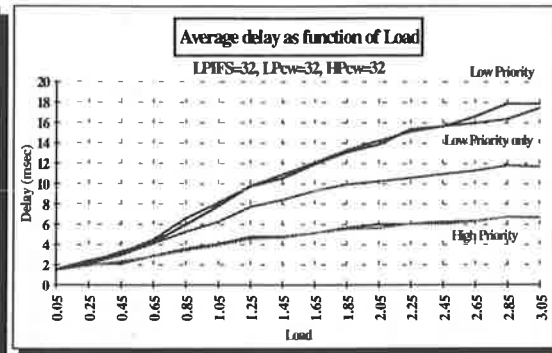


Fig 8: Effect of priority.

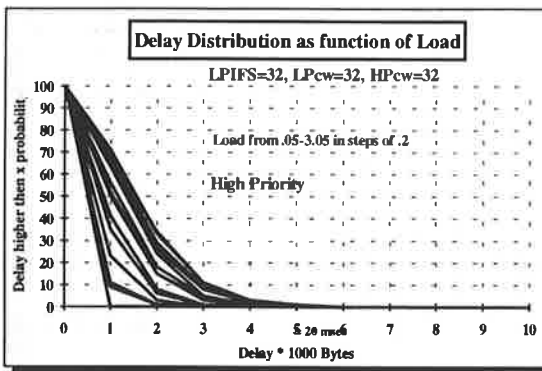


Fig 9: High priority delay variance.

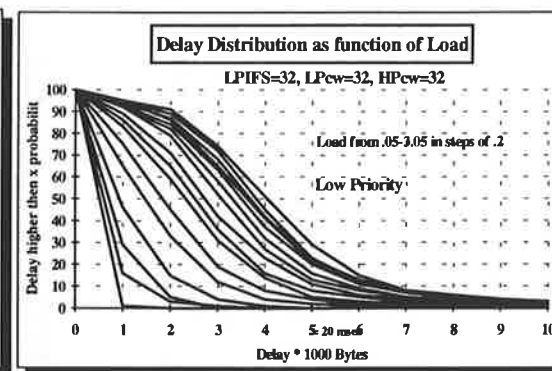


Fig 10: Low priority delay variance.

There is not much difference in the results between the two parameter combinations. The main visible difference is that the delay distribution of the low priority traffic is broader for the second case.

Other effects like the collision probability and throughput, should be taken into account to determine the optimum parameterisation.

The basic assumption is that the Time Bounded load will be limited to a certain percentage of the total capacity, which needs to be shared with bursty Asynchronous traffic. So an overload situation will consist of the same constant DTBS load with an Asynchronous (low priority) overload. The simulations were however run with simultaneous increase of both the low and high priority load.

Fixed repetitive load simulations:

As discussed under the "simulation interpretation" section a simulation is run with a fixed time repetitive high priority load, as function of the low priority load. This represents a

more realistic situation that you would have to support for instance voice. Figures 11-16 show the results for the two different access parameter settings. Please note the different scale in figure 13 and 14. These simulations clearly show that the DTBS service based on the available CSMA/CA priority mechanisms is very feasible.

More exhaustive simulations are needed to determine the maximum DTBS load limit that should be applied.

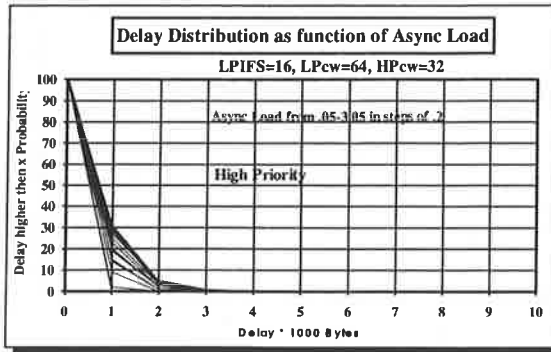


Fig 11: High priority delay variance.

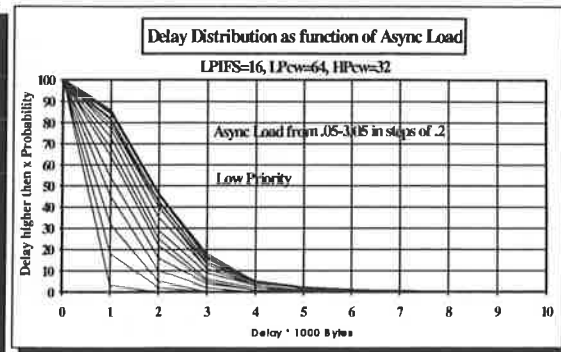


Fig 12: Low priority delay variance.

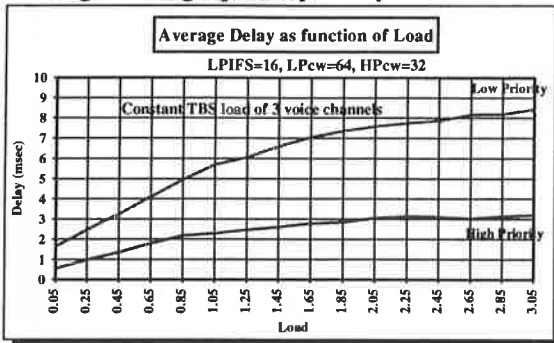


Fig 13: Priority with constant TBS load.

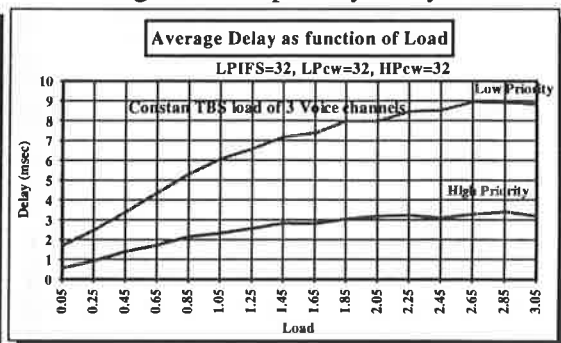


Fig 10: Priority with constant TBS load

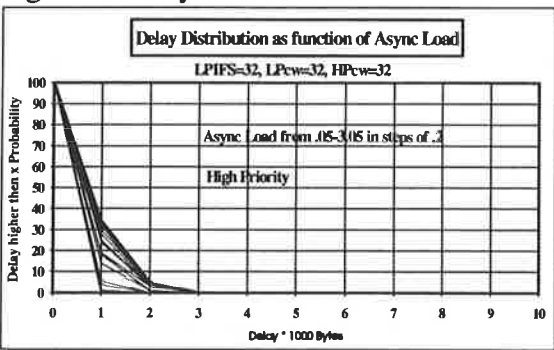


Fig 15: High priority delay variance.

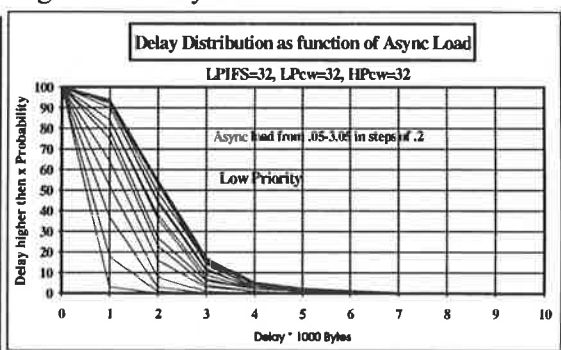


Fig 16: Low priority delay variance.

Other applications of priority:

As discussed in document 190, different priority access can also be used to give busy stations a higher access priority. This will however have effect on the overall fairness between stations.

A main difference however is the Access Point. In general it can be assumed that most of the traffic will be passing through the AP. Moreover there are a lot of applications where most data traffic is outbound coming from a wired network based server. This would clearly mean that most of the traffic within a BSS is generated by an AP.

It would therefore be advantageous for the total application throughput when an AP would have a higher (Asynchronous) access priority than stations.

The following three levels could be considered:

High priority:	HPIFS=2, HPcw=32	For DTBS.
Medium priority:	MPIFS=2+16, MPcw=32	For Async in AP.
Low priority:	LPIFS=2+16, LPcw=64	For Async in station.

Further the priority levels can also be applied for some of the MAC management frames, like for instance a Beacon.

Other priority mechanisms:

In addition mechanisms can be implemented within the MAC such that queuing delays are optimized for the higher priority service. This technique can be applied more intelligently when the MAC has some knowledge of the transfer delay that is acceptable for that segment. This could be in the form of a "Time to Live" parameter, which can be maintained by the MAC by subtracting the actual transfer delay.

Data service specification:

The question is what MAC interface should be provided for the distributed service. For the Asynchronous service this should map onto the LLC to MAC interface specification as specified within 802. This interface does currently have the capability to specify the priority. This would be one way of utilising the priority access service provided by the MAC.

An other approach is that a separate interface is provided for Time Bounded Services. Such an interface will need to be defined to a yet not existing TBS-LLC.

Document 190 assumes that this interface does contain a "Quality of Service" (QoS) parameter to control the TBS service.

This QoS parameter could be a "Maximum transfer delay" kind of thing. This is to be translated by the MAC into an access priority (and could be used for a implementation specific queuing priority). As discussed in the previous section, it is desirable that the MAC does maintain a "Time to Live" parameter. This can be the QoS parameter which decrements when the transfer delay elapses, so that the destination (or intermediate hop) can determine whether to transfer the frame further to the final destination or application, or to drop it when the maximum time is exceeded.

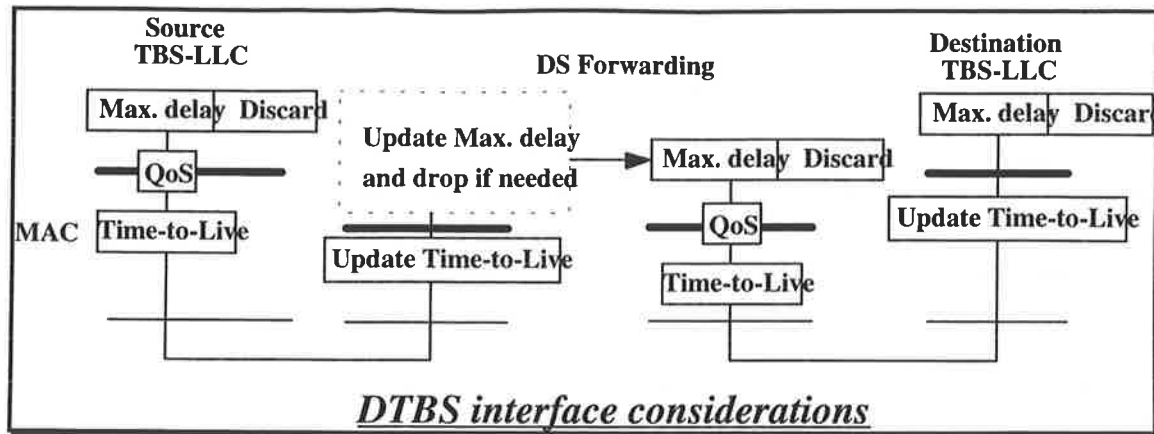


Fig 17: MAC to TBS-LLC interface

This is illustrated in figure 17.

The DTBS frame payload contains a "Max. Delay" and "Discard" parameter.

In the source station the "Max. Delay" parameter is provided to the MAC in the QoS field, but it is also part of the payload. The MAC can copy this in a "Time-to-Live" field, which is put into the MAC header (possible by means of the element method). This Time-to-Live field should be decreased as the time elapses until the frame is transmitted.

At the receiver end, which can be the AP, this field will then represent the "remaining Time-to-Live", which is delivered to the forwarding layer in the AP. Here this layer could make a decision whether to drop the frame when the delay is exhausted and the "Discard" parameter indicates that an overdue frame could better be dropped.

The "remaining Time-to-Live" can then become the QoS parameter of the next hop/segment.

Additional MAC support for DTBS:

Although bursty DTBS traffic can be handled within limits, as the simulations have shown, it may be desirable to control the DTBS load. This should be done by something above the MAC, who wants to establish a TBS connection. That layer would however need an indication of the total TBS load, before it would establish another TBS connection.

For this purpose the MAC may need to provide some indication to the layer above about the total high priority or TBS load on the network. A function like that may need to be specified to adequately support a service like that.

Conclusion:

This document clearly shows that priority access can be provided by the basic DCF, by parameterisation of the IFS, CW and CW increase parameters. The big advantage of this is that this can be applied to provide a Distributed Time Bounded service (DTBS) that can be supported in all environments. It does not have the limitation of the optional PCF based service, which is handicapped by PCF overlap limitations.

Very adequate delay distribution performance is demonstrated that can support for instance voice applications.

A separate provision in the PHY, of a single slot signal burst in the first slot after the DIFS as suggested by Kerry Lynn in his presentation of [2] does not seem needed. This is because the effect of the larger LPIFS on the low priority only traffic, compared to results with the original smaller DIFS is not significant, and seem to result only in a slightly lower throughput under overload conditions.

The DTBS MAC interface is been discussed, and some idea's of its format and parameters and support functions are given.

I hereby strongly propose to adopt the priority access functionality in the Foundation MAC as a non optional part of the DCF. A full specification of the DTBS service is advisable but this may be postponed to a later version of the standard.

References:

- [1] DFWMAC Distributed Foundation Wireless MAC Protocol", W. Diepstraten NCR-WCND-Utrecht, G. Ennis Symbol Technologies, P. Belanger Xircom; November 93, IEEE P802-93/190. See also P802.11-93/191, P802.11-93/192, P802.11-93/193.
- [2] A Distributed Time Bounded Service", Philip Rakity, Larry Taylor, Kerry Lynn Apple Computer; January 1994, IEEE P802.11-94/21

