

# A Simulation Study on a High-Speed Slotted-Ring Access Mechanism with Dynamically Adaptive Slot Sizes

H.L. Pasch, I.G. Niemegeers  
Tele-Informatics and Open Systems Group  
University of Twente  
Enschede, Netherlands

## Abstract

*Two access mechanisms complementary in performance are the token ring and the slotted ring. While the token ring outperforms the slotted ring for long messages, the latter performs significantly better for short messages. As we show in this paper, the factor that causes this difference is the number of tokens on the ring. We propose a new network design in which the number of tokens has been made adaptive, resulting in a network that can be made to behave like both the token ring and the slotted ring, or anything in between, i.e. it performs optimally for any given message length. We present first simulation results that show the performance gain that is achieved by using this principle. Further, we introduce a new priority mechanism in which an upper bound for the access delay of high priority messages is guaranteed.*

## 1 Introduction

Over the last 20 years a lot of research has been conducted on new access mechanisms for LANs, and more recently High Speed LANs and MANs. One of the main concerns has been the performance, in particular the transfer delay. Most of these networks are designed to perform at their best for a particular type of load (usually data traffic), but their performance decreases under loads that have strongly different characteristics. Because of the uncertainty about the characteristics and the mix of future (broadband) traffic (see e.g. [1] and [2]) it is important to search for networks which can adapt to the traffic characteristics.

Two access mechanisms that are to a certain extent complementary when it comes to the influence of the message length on the performance, are the token ring and the slotted ring. It has been shown in [3] and [4] that for various types of slotted and token rings, the

latter perform significantly better for long messages, while the former are at their best for short messages<sup>1</sup>. Ideally a network should perform well for both types of load, i.e. it should combine the advantages of both slotted and token rings in a single network. In this paper we propose a design which goes in this direction.

In Section 2 of this paper we analyse the difference between the slotted and the token ring, focusing on the performance. In Section 3, we introduce the basic concepts of a new network called Universal Channel Network (UCN), which is aimed at HSLANs and MANs. We argue that this network combines the advantages of the slotted and the token ring, by making the mechanism that causes the difference in performance adaptive. We also present simulation results that show the gain that is achieved by using this mechanism. In Section 4 we discuss the priority mechanism of UCN. The priority mechanism guarantees an upperbound to the access delay of the messages with a high priority. In Section 5 we draw some conclusions and indicate some open issues regarding the UCN concept.

## 2 Comparing slotted rings and token rings

Since there are various types of token ring and slotted ring principles let us first discuss which specific types are of interest to us. Next we compare the basic differences between the two types of access mechanisms and show that this the cause of the difference in performance. In the sequel we will repeatedly use the term message with the meaning of a Data Link layer Service Data Unit (D-SDU).

---

1. The Orwell ring has a different behaviour than other slotted rings. Its behaviour is not considered here.

## 2.1 Token rings

Three types of token rings are usually distinguished (see [5]): single message, single token and multiple-token rings. We consider only the latter, since the other types are not suitable for high-speed networks. In a multiple-token ring, e.g. FDDI, a sending station releases the token as soon as it stops sending. This can result in multiple tokens circulating on the ring, at most one of which is a free token. We further assume that the service discipline is gated, i.e. when a station has acquired the token, it sends all messages that were queued up at the moment the token was acquired. Compared with the exhaustive service discipline, which leads to the highest performance, the performance of the gated service discipline is slightly worse. However, unfairness between stations, as it could occur with the exhaustive service discipline, does not occur with the gated service discipline.

## 2.2 Slotted rings

There are also various types of slotted-ring networks. An overview and performance comparison of these has been presented in [3]. We confine ourselves to rings exclusively using channel slots, source release of slots and without a limit on the number of slots a station can use simultaneously. The channel-slot mechanism was introduced in the Cambridge Fast Ring [6]. A channel slot is characterised by the fact that it can be reused by the (source) station which just released it. In the other slotted-ring access mechanisms, slots have to be passed on to the next downstream station after being released by the source.

Compared to normal slots, channel slots offer a higher throughput, which is advantageous for large bulk transfers. For the transfer of short messages, this advantage disappears. A drawback of channel slots is the risk of hogging, which requires a channel slot management system. In our proposal, this drawback has been eliminated by means of an interrupt mechanism, and by using a gated service discipline. An additional advantage of channel slots (as we intend to use them in our new network) is that unlike in the case of normal slots, MAC PCI such as source and destination address need to be put in a slot only once. This is because after the first time a station has used a channel slot, there is an implicit association between source and destination. This implies a reduction of the overhead.

## 2.3 Qualitative considerations

The station which holds the token in a token ring has the exclusive right to use the ring. Similarly in a slotted ring, the station which changed the busy flag of a slot from 'idle' to 'busy' has the exclusive right to use the slot. This right is relinquished by putting the (free) token back onto the ring or by changing the busy flag from 'busy' back to 'idle', respectively. Therefore we will also in case of the slotted ring speak of the token of a slot when we mean the exclusive right to use the slot.

The fundamental difference between the two systems is that in the token ring there is a single token that controls the access to the entire bandwidth, whereas in the slotted ring, there are a number of tokens, each controlling access to part of the bandwidth (the transfer capacity of a single slot recurring periodically).

This difference is also reflected in the performance models of the two systems: they can be modelled by single and multiple cyclic server queueing models respectively [7]. Let us now explain how this difference also causes the difference in performance characteristics between the token ring and the slotted ring<sup>2</sup>.

## 2.4 Performance considerations

In [3] an extensive performance evaluation has been made of a number of high-speed ring protocols. Among the protocols that have been compared are the multiple token ring with exhaustive service and a slotted ring with a combination of normal and channel slots and in which a station can use several slots simultaneously. The influence of the average message length on the mean total sojourn time<sup>3</sup> that a message suffers has been investigated. The performance evaluation shows that the token ring outperforms the slotted ring when the average message length is large. For a short average message length, however, the slotted ring outperforms the token ring. In other words: for a given message length the performance depends, amongst other parameters, on the number of tokens. Therefore, for each type of load there is a particular number of tokens that provides the best performance.

This property can be explained by efficiency arguments. For long messages the token-passing mechanism becomes more and more efficient, since the ratio of the overhead per message and the data field

---

2. There is also another major difference between a token ring and a slotted ring. We discuss this difference later on.

3. With sojourn time we mean the time interval from the arrival of the first bit of a message at a (source) station until the last bit of the message has been received by the destination.

(SDU) is small and decreases as the message length increases. For slotted rings this is not the case. Long messages are segmented into mini-packets, which fit in a slot, and which have a constant amount of overhead. As the message length increases, the overhead grows quasi linearly. The net result is that for a given effective network utilization (MAC SDU traffic load), the gross utilization of the transmission medium will be larger for slotted rings than for token rings, leading to lower expected sojourn time for the latter.

For short messages on the other hand the opposite occurs. Here two effects cause the token passing mechanism to become less and less efficient as the average message size decreases. First there is the token passing overhead due to preamble (assuming an asynchronous transmission scheme) and PCI, which is larger than the overhead per mini-packet in a slotted ring. Typical values are 200 bits for token passing vs. 48 bits for a high-speed slotted ring [7]. The second effect occurs at low and medium loads. When the message transmission time becomes short compared to the ring latency, the token latency, i.e. the time a message at the head of a station queue has to wait until a free token arrives, starts playing an important role. A token ring, with its single free token is here at a disadvantage, compared to a slotted ring with multiple slots, where a number of free tokens are circulating. This becomes particularly important when, for a given type of traffic load, transmission speeds are increased.

The conclusion one can draw is that for long messages the token-passing mechanism is more efficient and tends to offer better delay performance, while for short messages the slotted ring is better. An important design goal of UCN is that its behaviour can be varied in discrete steps between that of a ring with multiple slots and that of an efficient multiple token passing ring. In UCN the number of slots, i.e. the number of tokens can be adapted to the traffic characteristics. In particular for a given average message size the number of slots or the slot size can be optimized with respect to the expected sojourn time. In [8] it is shown that such an optimum exists.

### 3 The UCN access mechanism

In this section we discuss specific features of the UCN ring: the slot layout (Section 3.1), and the grouping mechanism which adapts the slot size by joining adjacent basic slots (Section 3.2). We further argue in Section 3.3 that the UCN ring can be made to behave like a multiple token ring, with the corresponding efficiency. In Section 3.4 we present

some simulation results that show the gain that is achieved by adapting the group sizes to the load.

#### 3.1 Slot layout

UCN is a slotted ring (see Figure 1). Every slot (also referred to in this paper as a basic slot) consists of:

- a Token bit (T-bit): Access to a slot is controlled by the token bit. A free slot is characterised by a T-bit equal to 1. A station captures a free slot by setting the T-bit equal to 0. It can use the slot as long as it needs it (within the constraints of the gated service discipline). A slot is released by changing the T-bit back to 1. Several slots can be used simultaneously by a single station.

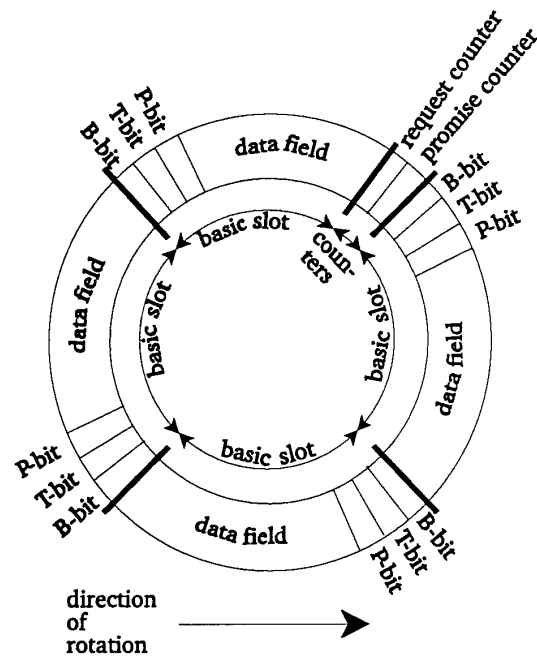


Figure 1

- a Contain\_PCI bit (P-bit): In UCN a station can use a slot repeatedly for transferring successive parts of a message. Only the first time a slot is used it has to contain the PCI since afterwards both sender and receiver know that this particular slot is used for transferring a certain message (Figure 2). The P-bit indicates that a slot contains the beginning of a message. This allows substantial savings in overhead compared to the slotted-ring mechanism used e.g. in the

Cambridge Fast Ring [6] and Cambridge Backbone Network [9].

- a Beginning\_of\_group bit (B-bit): In the initial situation every basic slot is separately accessible and therefore has its own token. Adjacent basic slots can be joined to form a group<sup>4</sup>, i.e. a single large slot with a single token. A group can consist of any number of basic slots, ranging from one to all the basic slots on the entire UCN ring. Basic slots keep their B-, T- and P-bits when they are joined. The B-bit is used for marking the border between adjacent groups. A B-bit set to 1 indicates the beginning of a new group, while a B-bit set to 0 indicates that a basic slot belongs to the same group as the previous one.
- a data field. This field contains segments of the message (Figure 2).

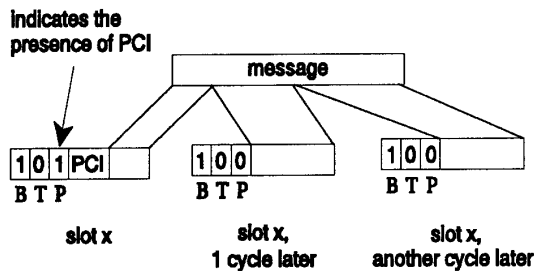


Figure 2

### 3.2 Grouping

As we saw, the way a system performs depends on the number of tokens, and therefore on the size of a group. Let us describe how a station can change the group size. The initial situation is illustrated in Figure 3a: two groups consisting of a single basic slot each. Group 1 and 2 are joined by changing the first (and in this case the only) B-bit of group 2 from 1 to 0 and by removing the token from group 2, because access to a group is controlled by a single token. The resulting situation is illustrated in Figure 3b.

The size of a group can be decreased by splitting a group into a number of smaller groups, each consisting of an integer number of basic slots. It is necessary to mark the beginning of each new group with a B-bit of 1, and to put a token into each of the new groups. This process of splitting is exactly the opposite of grouping.

4. The term 'group' is introduced (instead of using the term 'slot') to prevent confusion with the term 'basic slot'.

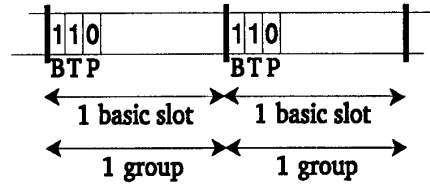


Figure 3a

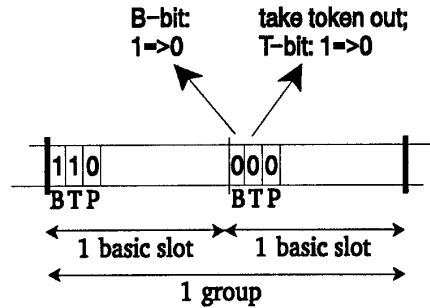


Figure 3b

Adapting the group size to the general characteristics of the traffic is only done by a management station. Grouping could be based on monitoring the traffic on the ring or on a-priori knowledge about the network load as a function of the time of the day, e.g. mainly short-message telephone traffic between 9.30 and 11.00 am, and large file transfers between 4 and 5 pm.

### 3.3 Early token release

We have explained how the number of slots and hence also the slot size can be adapted in the UCN ring. By decreasing the number of slots to one we claim that the system can be made to approach the behaviour of a multiple token ring quite closely, provided the size of the basic slots is sufficiently small compared to the average message length. Let us give the qualitative arguments which lead to this conclusion.

In both a slotted ring containing a single slot, and a multiple-token ring, access to the entire bandwidth is controlled by a single token. However there are still significant differences:

- If, in a multiple-token ring, a station has a backlog of several messages, the transmission of the next message is started immediately after the transmission of the previous message. Furthermore the token is released immediately after the transmission of the last message. Waste of transmission capacity is limited to the passing

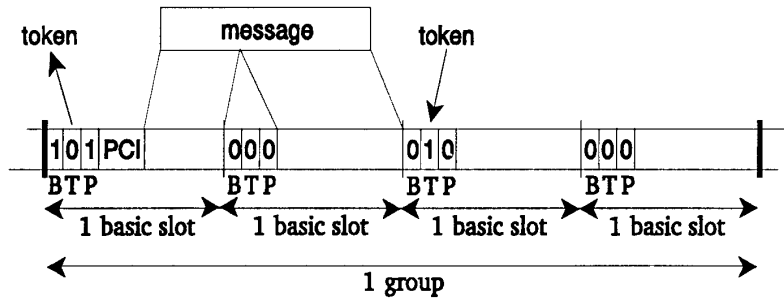


Figure 4a

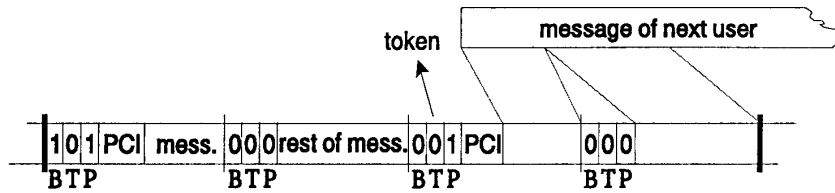


Figure 4b

of the token from station to station, i.e. it is determined by the ring latency.

- In a slotted ring a message has to start at the beginning of a slot. Therefore when a message does not exactly fit in an integer number of slots, waste occurs due to internal fragmentation. Furthermore, since a token is passed to a station at the beginning of a slot, the time between the end of the last message a station has to transmit and the beginning of the next slot is wasted. This type of waste generally increases with the slot size for a given message-length distribution. It can therefore be substantial in case of a ring containing a single, large slot.

Hence, compared to the multiple-token ring, the slotted ring has some inefficiencies which tend to become worse when the slot size increases. In UCN this problem has been significantly reduced because large slots, the groups, have a substructure of basic slots. A message can in principle start at the beginning of each basic slot within a group, and not just at the beginning of a group. So, if a message ends and leaves a number of basic slots within a group unused, they are not wasted, but used for the transfer of the next message. Fragmentation only occurs within a basic slot and does not depend on the group size.

The same holds for the passing of the token. While in ordinary slotted rings the token is put at the beginning of a slot, in UCN it can be put at the beginning of each basic slot of a group, therefore minimizing the waste. An example will clarify this. Consider a group consisting of four basic slots (Figure 4a) and a message with a length of two basic slots. Assume further that in our initial situation the token is in the first basic slot of the group. The station claims the token, changes the P-bit of the first basic slot to 1 and puts its message into the slot. After the message has been transmitted, the token of the group is released by putting it into the next basic slot (basic slot number 3) of the group. The next backlogged station will take the token out of this basic slot, change the P-bit to 1 and start transmitting its message (Figure 4b).

We saw that in UCN the amount of waste depends on the basic slot size, and not on the group size. We now discuss why the size of a basic slot can be chosen very small in order to minimize waste and the difference between the behaviour of UCN and a multiple token ring with exhaustive service.

In most slotted systems the slot size is a trade-off between two types of waste:

- PCI waste: since normally every slot needs PCI, an increase in the slot size means that a message

can be transferred using less slots and therefore with less overhead;

- Fragmentation waste: a message does not normally fill up an integer number of slots; on the average the last slot will only be filled half; therefore, the larger the slots the larger the fragmentation waste.

In UCN, however, PCI is not sent once per basic slot, but once per group, regardless of how often the group is used. Therefore the PCI overhead is independent of the basic slot size which can be decreased, without incurring an efficiency penalty, in order to minimize the fragmentation waste. This further minimizes the difference between the behaviour of UCN and the multiple token ring with exhaustive service.

### 3.4 Simulation results

The design goal of UCN was a network that can be made to behave like the slotted ring or the token ring or anything in between (depending on what offers the lowest sojourn time). In this section we evaluate by means of simulation if this goal has been achieved. In our simulations we compare token ring and slotted ring (as described in the Sections 2.1 and 2.2) with UCN. Three different situations have been simulated : one in which the token ring performs very well, one in which the slotted ring performs very well, and one situation more or less 'in between' of the other two.

In our first simulation we compare the three networks using a load consisting of long, exponentially distributed messages. As can be seen in Figure 5, for this load the token ring performs substantially better than the slotted ring (i.e. a ring with a single server performs better than a ring with many servers)<sup>5</sup>. In this figure it can further be seen that the performance of UCN approximates that of the token ring very closely.

In our second simulation we used a load consisting of short, fixed-length messages that fit exactly into a single slot. For this load the slotted ring outperforms the token ring, as can be seen in Figure 6 (i.e. a ring with many servers outperforms a ring with only a single server). As can further be seen, the sojourn time of UCN approaches that of the slotted ring. Only for very high load the sojourn time increases faster than in the case of a slotted ring.

The third simulation concerns short messages with exponentially distributed lengths. For this type of load

5. The curve of the slotted ring seems flat since most of the sojourn time consists of transfer delay, which is in case of channel slots independent of the load. Besides, a logarithmic scale has been used, which further flattens the curve.

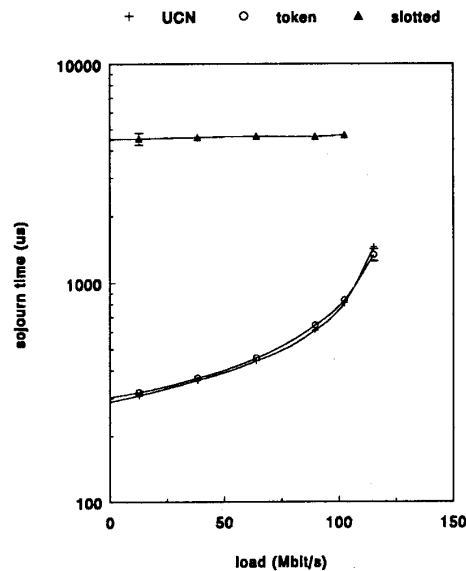


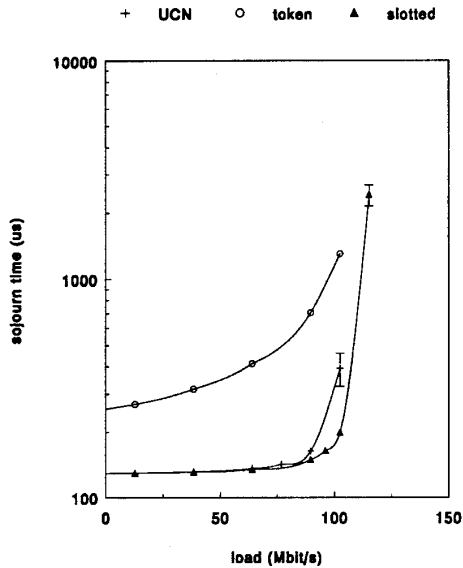
Figure 5: The sojourn time as a function of the load on the ring, for the token ring, the slotted ring, and UCN. The load consists of exponentially distributed messages with a mean length of 800 bytes (excluding the 6-byte header). (Ring-configuration parameters: medium bandwidth 136 Mbit/s, 40 stations connected to the ring, 250  $\mu$ s ring latency, exponential interarrival times of messages).

we expect a ring with 'few' servers to offer a higher performance than a ring with only a single server (i.e. the token ring) or a ring with many servers (i.e. the slotted ring). In Figure 7 it can be seen that this expectation is true: UCN (configured with few servers) offers a lower sojourn time than both the token ring and the slotted ring.

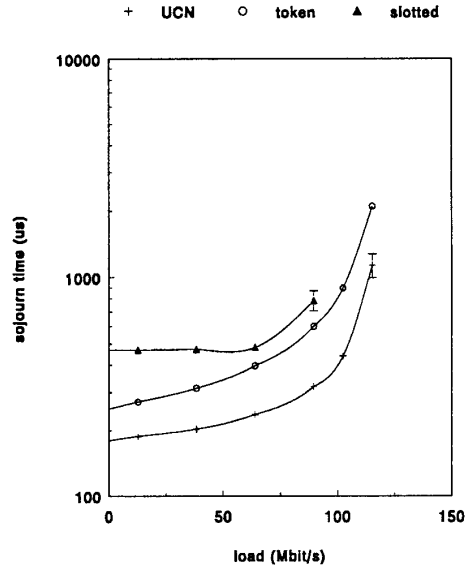
From this it can be concluded that UCN meets the expectation in that its performance is better than or at least equal to that of the slotted ring and the token ring.

## 4 The priority mechanism

UCN also has a priority mechanism. The aim of this priority mechanism is to provide a guaranteed maximum access delay to traffic with high delay requirements. The mechanism provides two priorities: a high one (with a guaranteed maximum access delay) and a low one (without such a guarantee). The basic



**Figure 6:** The sojourn time as a function of the load on the ring. The load consists of messages with a constant length of 44 bytes (excluding the 6-byte header).



**Figure 7:** The sojourn time as a function of the load on the ring for messages with exponentially distributed lengths. The mean length is 80 bytes (excluding the 6-byte header).

idea behind the priority mechanism is to interrupt the transmission of low-priority messages when high-priority messages are in imminent danger of missing their delay deadlines. Let us now discuss the operation of the priority mechanism in more detail.

The priority mechanism is based on two counters: a request counter  $C_{req}$  and a promise counter  $C_{prom}$ .  $C_{req}$  indicates the number of groups that stations have not been able to claim and that are still needed for the transmission of high-priority messages.  $C_{prom}$  indicates the number of groups that stations have promised to release by interrupting their transmission of low-priority messages. Both counters are part of the PCI continuously broadcast on the medium (see Figure 1).

Further, each station has a so-called High-Priority Requests counter ( $C_{hpreq}$ ). This counter indicates the number of free groups a station has to let pass in order to satisfy all outstanding high-priority requests from stations that fear to miss their delay deadlines. Only when the value of this counter is zero, a station is allowed to claim free groups that pass.  $C_{hpreq}$  is decreased by 1 each time a free group passes. Each time  $C_{prom}$  passes, stations increase  $C_{hpreq}$  with the value of  $C_{prom}$ .

When a low-priority message is generated at a station, the station is allowed to claim a free group as long as (or as soon as)  $C_{hpreq}$  is zero. When a high-priority message is generated at a station, that station tries to claim a free group in order to serve the message. When it is not able to do so within a certain amount of time (either because no free group passes, or because the value of  $C_{hpreq}$  exceeds zero), it increases  $C_{req}$  by 1.

All stations on the ring continuously monitor  $C_{req}$  and  $C_{prom}$ . Since  $C_{req}$  indicates the number of groups that are needed for high-priority messages, and  $C_{prom}$  indicates the number of groups that have already been promised,  $C_{req} > C_{prom}$  indicates that a number of groups (equal to  $C_{req} - C_{prom}$ ) are still needed for the transmission of high-priority messages. When noticing this, a station checks the priority of the message it is currently transmitting. If this message has a low priority, the station:

- increases  $C_{prom}$  by 1;
- suspends the transfer of the low-priority message;
- and releases the group.

As soon as stations are allowed again to claim groups for transmitting low-priority messages (i.e. as soon as

$C_{hpreq}$  is zero again), the station will claim a group in order to resume the suspended transmission of the low-priority message.

After having increased  $C_{req}$ , the station with the high-priority message copies the value of  $C_{hpreq}$  to the count-down counter  $C_{cd}$ . Each time a free group passes,  $C_{cd}$  is decreased by 1. When the value of  $C_{cd}$  has become zero, the station is allowed to claim the next free group that passes. In this way, a global first-come first-serve service discipline for high-priority messages is achieved. One ring rotation after incrementing  $C_{req}$  the station decreases  $C_{req}$  again, as well as  $C_{prom}$ .

The priority mechanism is currently being tested by means of simulation.

## 5 Conclusions

The token ring and the slotted ring are to a certain extent complementary with respect to the influence of the message length on the performance: the token ring outperforms the slotted ring for long variable-length messages, while the slotted ring outperforms the token ring for high-priority messages that exactly fit into a single slot. In this paper we argued that the advantages of both networks can be combined in a slotted ring in which the slot size can be adapted to the characteristics of the load. We presented the basic concepts of a network with this property, the Universal Channel Network (UCN), and we showed the performance gain that is achieved by using this concept. We further introduced a priority mechanism in which an upper bound for the access delay is guaranteed for high-priority messages.

In this paper it was explained why it is advantageous to have an adaptive slot size and how the slot size can be made adaptive. The question of the adaptive control of the slot size has not been answered yet. Our next research effort to develop the concept further is aimed at determining how the 'state' of the network can be measured effectively, and to develop a simple algorithm to determine the optimal slot size for a given state.

The research on UCN is part of a research project at the Tele-Informatics and Open Systems group at the University of Twente, the Netherlands. In the context of this project we investigate local area broadband (B-ISDN) communication infrastructures.

## References

- [1] Lambarelli L., Luvison L., Roffinella D., Sposini M., 'Service Integration in Wideband Local Area Networks: Problems and System Solutions', Proceedings Intl. Tirrema Workshop on Digital Communications, Tirrema, 1985, North Holland, 1985.
- [2] Newmann S., 'The Communications Highway of the Future', IEEE Communications Magazine, October 1988, pp.45-50.
- [3] Zafirovic-Vukotic M., Niemegeers I.G., Valk D.S., 'Performance analysis of slotted ring protocols in HSLAN's', IEEE Journal on selected areas in communications, Vol. 6, No. 6, July 1988, pp 1011-1024.
- [4] Kamal A.E., Hamacher V.C., 'Utilizing bandwidth sharing in the slotted ring', IEEE Transactions on Computers, Vol. 39, No. 3, March 1990, pp. 289-299.
- [5] Bux W., 'Local-area subnetworks: a performance comparison', IEEE Transactions on Communications, Vol. 29, No. 10, October 1981, pp. 1465-1473.
- [6] Temple S., 'The Design of the Cambridge Ring', in 'Ring Technology Local Area Networks', Dallas I.N. and Spratt E.B. editors, Elsevier Science Publications, 1984, pp.74-88.
- [7] Zafirovic-Vukotic M., 'Performance modelling and evaluation of high-speed serial interconnection structures', Ph.D. thesis, University of Twente, Department of Computer Science, December 1988.
- [8] Pasch H.L., Niemegeers I.G., 'A High-Speed Slotted Ring Access Mechanism with Dynamically Adaptive Slot Sizes', Memoranda Informatica 90-36, June 1990, University of Twente, Enschede, Netherlands.
- [9] Greaves D., Hopper A., 'The Cambridge Backbone Network', Proc. EFOC/LAN '88, Amsterdam, June 1988, pp. 399-402.
- [10] Conway R.W., Maxwell W.L., Miller L.W., 'Theory of Scheduling', Addison-Wesley Publishing Company, 1967.
- [11] Morris R.J.T., Wang Y.T., 'Some Results for Multi-queue Systems with Multiple Cyclic Servers', in 'Performance of Computer Communication Systems', H. Rudin and W.Bux editors, Elsevier Science Publishers, 1984, pp.245-258.
- [12] Pasch H.L., Niemegeers I.G., 'A High-Speed Slotted Ring Access Mechanism with Dynamically Adaptive Slot Sizes', EFOC/LAN 91, London, June 19-21, 1991, LAN proceedings, pp 42-47 & pp 248-253.
- [13] Pasch H.L., Niemegeers I.G., 'Simulation results on UCN' to appear as Memorandum Informatica, University of Twente, Netherlands