

Improving Node Behavior in a QoS Control Environment by Means of Load-dependent Resource Redistributions in LANs[†]

Bernd E. Wolfinger¹, Jürgen Wolf¹, Gwendal Le Grand²

¹*Dept. of Computer Science, TKRN, University of Hamburg*

Vogt-Kölln-Str. 30, 22527 Hamburg, Germany

Phone: +49 40 42883 2424; Fax: +49 40 42883 2345

{wolfinger|jwolf}@informatik.uni-hamburg.de

²*GET/ENST, INFRES, 46 rue Barrault, 75634 Paris, Cedex, France*

SUMMARY

An important means to guarantee an acceptable quality of service in networks with real-time communication requirements is the reservation of resources at connection setup time. However, such reserved resources, e.g. transmission bandwidth, may be unused as a consequence of the variations in the actual resource demands. Therefore, a more efficient resource utilization is possible if communicating stations or end-users dynamically hand over some of the free resources temporarily to the other communication partners, e.g. of a “broadcast network”.

This paper concentrates on two fundamental problems of such a demand-based sharing of resources: on the one hand, estimation of the current resource requirement on the basis of load measurements is investigated and, on the other hand, we elaborate efficient algorithms for resource sharing respecting real-time requirements. The algorithms proposed for load estimation and for resource sharing are evaluated analytically with respect to their efficiency for worst-case, average-case and realistic load scenarios. Our approach suggested for resource and traffic management allows one to achieve significantly better utilization of network resources. Copyright © 2000 John Wiley & Sons, Ltd.

KEY WORDS: Quality of Service (QoS), Resource Management, Traffic Management, Real-Time Communications, Performance Evaluation

1. INTRODUCTION

The trend towards distributed multimedia communications via local-area networks has led to the requirement to transmit continuous media streams, such as audio and video streams, with sufficiently good quality. Basically, two common possibilities exist in order to guarantee a certain QoS as required by the users of a network: either, one can use *resource reservation* or, alternatively, one could rely on *priorization combined with admission control*, cf. [1]. In the context of the Internet, the resource reservation approach is reflected e.g. by the IntServ proposal [2] and the priorization approach is inherent to the DiffServ proposal [3].

[†]A shorter preliminary version of this paper has appeared in the Proceedings of the Int. Symp. on Performance Evaluation of Computer and Telecommunication Systems (SPECTS), Montreal, July 2003, 361-371.

In this paper, which is a revised and extended version of [4], we focus on the resource reservation approach and its usage in local broadcast networks, especially networks with wireless or common bus communication infrastructure. In particular, we study ways of how communication resources which are already reserved and allocated to communicating end-users or end-systems can be dynamically handed over to other end-users or end-systems in case that they are not needed by their original "owner" for some time. Our interests center around methods which make resources given back sufficiently quickly, if they are needed again by their owner, in order to make sure that real-time requirements are not neglected.

The approach we suggest assumes that the problem of reserving and statically allocating resources to communicating entities has already been solved [5]. We also suppose that the resource reservation has been established for a rather long-term time interval, e.g. duration of an audio/video stream, typically in the order of minutes. In our approach, each owner of resources determines whether its communication load, i.e. the amount of time critical data waiting for transmission at the owner, justifies the continued reservation of all the resources as allocated to the owner. If a sufficiently large amount of resources is observed to be temporarily free, the owner of the resources will pass some of these resources to its "neighbors", i.e. the other end-systems of the broadcast LAN. In order to respect real-time requirements for its data to be sent each owner will continue to observe the arrivals to its transmission queue. If its local load is increasing again, the owner informs its neighbors that they are no longer allowed to make use of the resources which were offered for public access at an earlier instant. The freeing of resources and their reclaiming could lead to oscillations in the resource redistribution. Therefore, our approach for making resource allocation, depending on the actual communication load, is based on load estimators, which produce some kind of smoothed estimates for the load as it is generated locally over time. Moreover, we introduce systems of load thresholds which, only when being crossed by the load estimates of an owner of resources, leads to a broadcast of messages for freeing or reclaiming resources.

The remainder of the paper is organized as follows: Section 2 gives a survey of the state of the art in load estimation, scheduling resource allocations and dynamic redistribution of resource reservations in the context of QoS management. Section 3 presents and analyzes the behavior of a set of functions for estimating the actually generated load at a LAN station, based on geometric or on arithmetic weighting of measured utilization samples. In Section 4 we introduce several threshold systems which differ in their responsiveness to observed load changes and in their overhead, in terms of number of messages, for sending resource redistribution messages. Section 5 investigates, by means of mathematical models, how load estimators can be combined with threshold systems by a single station if it has to respect real-time requirements for sending data. The main lessons learned are summarized in Section 6.

2. RELATED WORK

Reservation of resources is an important means to achieve quality of service guarantees in computer and communication networks. In continuous media communications with real-time requirements, such as audio/video communications, it is typical to reserve the resources (as they are expected to be required later on) during connection setup of a stream, cf. IntServ [2]. Some real-time protocols and services, such as CMTP/CMTS [6], allow that resources may be re-allocated by means of setting up a new stream dynamically with different, in

particular less strict, resource reservations. Other approaches assume that variations of load within single stations are smoothed out to a great extent by means of multiplexing a larger number of streams and doing the reservation for the complete overlay of those streams [7]. Up to now, only few proposals exist which - as in our approach - advocate for a dynamic freeing and recalling of bandwidth between stations, cf. [8] and [9]; but these earlier proposals do not provide a redistribution of the unused prioritized bandwidth to the other nodes in the LAN. Thus, these solutions satisfy strict real-time bounds for transmission delays, but at the cost of a worst case reservation. Bandwidth reservation has already been studied in the literature. The various contributions include centralized and distributed approaches. Within the Integrated Services over Specific Link Layer (ISSLL) group at IETF, Yavatkar, et al. [10] proposes a centralized architecture for subnet bandwidth (SB) management using a centralized algorithm. A single SB manager (SBM) acts as a designated SBM and is responsible for admission control over the resource reservation requests originating from the source hosts. However, this approach depends highly on the Resource Reservation Protocol (RSVP) [11] and does not deal explicitly with best effort traffic related issues. Ghanwani, et al. [12] discuss methods for supporting RSVP in LAN environments. It introduces a bandwidth allocator (BA) responsible for performing admission control and maintaining resource allocation states in the subnet. End systems request services such as bandwidth reservation that are processed by the BA. BAP [13] represents a method to manage the dynamic bandwidth allocation; it can be used in arbitrary shared transport media. Bandwidth is allocated in discrete portions, that can be allocated and de-allocated at any point during the connection. The Controlled Load Ethernet Protocol (CLEP) described in [14] is an implementation of the Controlled Load service over Ethernet as defined by Wroclawski in [15]. It provides the client data flow with a quality of service approximating the QoS this flow would receive on an unloaded network. This distributed service is obtained by incorporating an access controller on the outgoing interfaces of the nodes.

Estimation of load, as a fundamental problem of our approach, has been an important topic, e.g., in adaptive routing and in load-balancing since more than a decade [16, 17]. Besides the idea of basing load estimation directly on the utilizations observed [18] it has also been suggested to observe arrival times, thus estimating arrival rate, and service-time requirements, thus estimating mean service-times, directly. However, in scenarios relevant to our paper service-time requirements of data transfer requests typically would not vary to a sufficient extent in order to justify the additional expenditure by roughly duplicating the measurement efforts (estimation of two parameters instead of one). Still, as an alternative to observing utilization one could observe the current backlog and base resource redistribution on this estimate.

Regarding the *weighting functions* needed for calculating load estimates, geometric weighting [19] seems to be the dominating approach. This can be explained by the easy and straightforward way to execute the calculations for geometric weighting and also by the fact it has been used successfully in the context of TCP [20]. In *load sampling* besides the periodic sampling which we suggest for simplicity, a more complex measurement procedure based on so-called random sampling has been suggested just recently [21]. Measurements, as required in our paper, refer to observations of primary load rather than of secondary load in the sense of [22].

In order to *reduce the fluctuations of load* at the source besides using methods of traffic shaping, e.g. based on rate-based flow control such as [23], one could also apply techniques for smoothing the traffic generated at the source [24]. Our use of *threshold-systems* to reduce

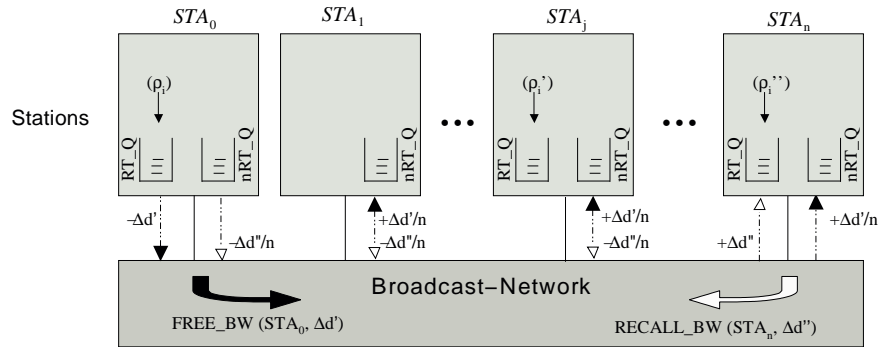


Figure 1. Redistribution of bandwidth: the depicted example scenario illustrates the freeing of $\Delta d'$ by station STA_0 and recalling of $\Delta d''$ by STA_n

the number of control messages to be exchanged has been inspired again by similar techniques which have been successfully applied in practice in adaptive routing protocols [25].

As an alternative to provisioning deterministic *QoS guarantees*, as we strive for in our approach, stochastic QoS guarantees are possible [5]. Then, evidently, no longer a reservation for the worst case would be required but QoS requirements could become violated temporarily, unless applications have the capability to adapt to network changes [26]. For applications tolerating transient fluctuations in the QoS, different approaches for resource management exist, in particular, in the context of wireless networks [27, 28].

3. ESTIMATION OF ACTUAL NODE UTILIZATIONS BASED ON WEIGHT FUNCTIONS

In the following we present the basic proceeding which we suppose throughout this paper for resource reservation and redistribution between the stations of a broadcast network.

3.1. Assumptions and underlying model for resource redistribution

- Resource reservation: We assume:
 - $n + 1$ stations STA_0, \dots, STA_n communicating via a local broadcast network (Ethernet, WLAN, ... [29])
 - communication load of the stations consisting of time-critical transmissions and non-time-critical data transfers, cf. the real-time (RT_Q) and non-real-time queues (nRT_Q) in Figure 1
 - in each station with real-time communication requirements there exists at least one owner of resources, where reserved data rate (allocated “bandwidth”) is the resource considered by way of example; scenarios could be that the owner reflects the complete station or an owner could model a single end-user
 - resources are supposed to be reserved before starting a new stream with real-time requirements and are available to its owner until it releases connection after end

of the stream; the amount of bandwidth reserved is such that the owner's QoS requirements can be guaranteed.

- Temporary resource redistribution:

We assume that an owner of resources can pass resources to other stations temporarily, whenever these resources are not required locally for some time. We introduce a message of type

- FREE_BW ($STA, \Delta d$) with which a station STA hands over a bandwidth (respectively data rate) of Δd bit/s to other stations, and
- RECALL_BW ($STA, \Delta d$) to recall bandwidth Δd .

With receipt of a RECALL_BW message the recalled bandwidth has to be released immediately. The question of how to distribute free bandwidth among stations – though it represents a very important topic, too – is out of the scope of this paper; an elementary solution could be, that all neighboring stations get equal shares. A more advanced approach to distribute free bandwidth efficiently is proposed in [8].

- Approach to estimate the level of load resulting from time-critical transmission requests: We assume periodic measurements available at instants $t_i := t_0 + i \cdot \Delta t$, for $i \geq 1$, where the measured samples ρ_i characterize the load of time-critical transmission requests which have been generated by the source during the interval $T_i := [t_{i-1}, t_i)$. Evidently, $\rho_i := d_i / (r \cdot \Delta t)$, where d_i denotes the amount of time-critical data generated during T_i and r denotes the data rate as it was allocated to its owner at connection setup time. This sequence of samples (ρ_i) may be used at each instant t_i to calculate an estimate $\hat{\rho}(t_i)$ of the actual level of load.

The traffic load generated by the communicating users of a station must be estimated in order to be able to adjust the allocated capacity to the actual requirement. If the capacity, in terms of transmission bandwidth, allocated is too high, only a poor level of utilization for the network can be achieved. Some flows may have more reserved capacity than what they really use. This capacity is lost for all the nodes of the network. On the other hand, if the reserved capacity is too low in comparison to the flow activity, a situation of congestion occurs and the general QoS of the flow will be degraded. It is important to find out good estimators to allow the algorithm controlling the bandwidth to react in an appropriate manner. Evidently, different estimators can be chosen for different streams and different demands on the expected QoS.

The estimator must take into account some kinds of fluctuations in order to enable adequate node reactions. Major changes in data rates for a complex source of traffic must be detected quickly. This is an important condition for flows with a required level of QoS. If the reaction is too slow, a QoS degradation is viewed by these flows, whereas a minor change in data rates within single streams may be ignored in order to get a certain form of smoothness over time. This constraint is to avoid an overreaction with respect to the observed fluctuation of offered load within a station.

3.2. Geometric and arithmetic weighting of utilization estimates

The methods proposed and investigated by us for load estimation are based on *geometric* and *arithmetic* weighting.

Here, *geometric weighting* is defined as follows:

$$\hat{\rho}(t_i) = \alpha \rho_i + (1 - \alpha) \cdot \hat{\rho}(t_{i-1}), \quad i \geq 1$$

where $\alpha \in (0, 1]$ and $\hat{\rho}(t_0)$ to be initialized, e.g. $\hat{\rho}(t_0) := 0$. The α factor indicates how fast the estimator changes with a strongly different new sample. The geometric weighting is also known as the *exponential weighted moving average* (EWMA) [19]. As an abbreviation we denote by G_α geometric weighting with parameter $\alpha \in (0, 1]$.

The *arithmetic weighting* estimator is defined as:

$$\hat{\rho}(t_i) = C_0 \cdot \sum_{j=0}^{w-1} \frac{w-j}{w} \rho_{i-j}, \quad i \geq 1,$$

where $w \in \{1, 2, \dots\}$ and ρ_k to be initialized for all $2-w \leq k \leq 0$. $C_0 = \frac{2}{w+1}$ denotes a normalization constant. This estimator is equivalent to a sliding window with window size w and weights increasing in a linear way when the samples are more recent. Arithmetic weighting uses $\frac{k}{w}$ as weights, where k is the position in the window of size w with the most recent sample taking the position $k = w$. The size of w indicates how many past samples including the present sample must be kept to estimate the current load. As an abbreviation, A_w denotes arithmetic weighting with a window size of $w \in \mathbb{N}$.

3.3. Impact of weight functions on worst case backlog of load and on delays

In this section we discuss the impact of the weight functions and their parameterization on estimation of $\hat{\rho}(t)$. In particular, we analyze how much backlog (BL) is possible if the “worst case” occurs and determine how much “delay” is induced at most by the BL. It has been demonstrated previously [30] that all measured data (offered load) is taken into account once and only once by both of our estimators.

We constitute that the throughput available for each interval T_i (with boundaries $[t_{i-1}, t_i)$) is identical to $\hat{\rho}(t_{i-1}) \cdot r \cdot \Delta t$, except that lowering throughput from one interval to the next is not admitted if backlog exists in the sending queue. Here we disregard the time required to claim back transmission capacity, as we assume that control messages are exchanged with high priority and, moreover, in local networks we can neglect propagation delays. As mentioned previously, within the scope of these studies we presume that the amount of data d_i generated in each sample interval does not exceed $r \cdot \Delta t$. Hence ρ_i does not exceed 1. Such behavior could be realized, e.g., by a “leaky bucket” component, which could establish some rate-based flow control [23].

The normalized “Backlog of workload” (BL) is defined in the following as the quotient of the amount of data which has to be stored temporarily due to missing transmission capacity and the amount of data that can be sent with full transmission capacity during one sample interval of length Δt . This means that $BL = 1$ corresponds to the amount of data $r \cdot \Delta t$.

(a) Analysis of “worst case” load scenario

A “worst case” scenario occurs, e.g., when

$$\rho_i = \begin{cases} 0 & \text{for } i = 0 \\ 1 & \text{for } i \geq 1 \end{cases} \quad (1)$$

- Using **geometric weighting**, the estimator at instant t_i , initialized with $\hat{\rho}(t_0) := 0$, for $i \geq 0$ is given by

$$\hat{\rho}(t_i) = \sum_{k=0}^i \alpha(1-\alpha)^k \rho_{i-k} \stackrel{(1)}{=} 1 - (1-\alpha)^i.$$

Hence, as normalized backlog $BL(i)$ which accumulates in just one sample interval, we obtain for each interval T_i with $i \geq 1$ the difference between actual level of load ρ_i and the estimated level of load $\hat{\rho}(t_{i-1})$ at the beginning of the interval:

$$BL(i) = \rho_i - \hat{\rho}(t_{i-1}) = (1-\alpha)^{i-1}$$

which gives us the following expression for the overall backlog BL_i having accumulated up to time instant t_i :

$$BL_i = \sum_{j=1}^i BL(j) = \frac{1 - (1-\alpha)^i}{\alpha} \quad \forall i \geq 1 \quad (2)$$

The total size of generated (offered) load to be buffered, expressed in multiples of $r \cdot \Delta t$, is therefore for all $\alpha \in (0, 1]$ limited by

$$BL_i \leq BL_\infty = \lim_{k \rightarrow \infty} \frac{1 - (1-\alpha)^k}{\alpha} = \frac{1}{\alpha} \quad \forall i \geq 1$$

Examples:

$\alpha = 1$: Since, at each instant t_i , ρ_i is the only value used for estimation ($\hat{\rho}(t_i) = \rho_i$), as a maximum the amount of data received during the last interval T_i has to be queued in RT_Q for future transmission. Therefore the amount of data to be buffered can not exceed $1 \cdot r \cdot \Delta t$, i.e. $BL_{max} = 1$.

$\alpha = 0.5$: The maximum amount of data accumulating in the worst case is $BL_{max} = 2$.

- Using **arithmetic weighting**, and taking into account Equation (1), the estimator at instant t_i is given by

$$\hat{\rho}(t_i) = \begin{cases} \frac{2}{w+1} \sum_{j=0}^{i-1} \frac{w-j}{w} = \frac{2wi+i-i^2}{(w+1)w} & 0 \leq i \leq w \\ 1 & i > w. \end{cases} \quad (3)$$

Again we obtain the maximum backlog for each interval T_i as difference between actual and estimated load

$$BL(i) = \begin{cases} 1 - \frac{(i-1) \cdot (2w-i+2)}{(w+1)w} & \text{for } 1 \leq i \leq w \\ 0 & \text{for } i > w \end{cases}$$

which leads to the maximum backlog at instant t_i with $1 \leq i \leq w$

$$BL_i = \sum_{j=1}^i \left(1 - \frac{(j-1) \cdot (2w-j+2)}{(w+1)w} \right) \quad (4)$$

and $BL_i = BL_w$ for $i > w$. The total size of generated (offered) load to be buffered is limited by BL_w , i.e.

$$\begin{aligned} BL_i &\leq BL_w = \sum_{j=1}^w \left(1 - \frac{(j-1) \cdot (2w-j+2)}{(w+1)w}\right) \\ &= \frac{w+2}{3} \quad \forall w \geq 1, i \geq 1 \end{aligned}$$

Examples:

$w = 1$: With $w = 1$ the maximum backlog is $BL_{max} = 1$ as the amount of data received within T_i has to be queued in the worst case, but is completely served during T_{i+1} .

$w = 4$: Maximum backlog results to $BL_{max} = 2$.

(b) **Analysis of maximum backlog delay**

Now we determine how much backlog delay at most is induced by the data being stored temporarily. The backlog delay is the time elapsing between receiving data and being ready to send it and its maximum for each interval T_i is denoted with τ_i . The data is processed with throughput $\hat{\rho}(t_{i-1}) \cdot r \cdot \Delta t$ within each interval T_i following the FIFO-principle. As mentioned previously, we constitute in our study that the available throughput is not lowered while further backlog is present. Nevertheless, the generated load might decrease at any instant whereby the remaining backlogged data would be sent with a data rate less than the maximum, possibly resulting in higher delay. Hence, considering the highest increase of generated load, which is described by Equation (1), it follows that the maximum delay is upper bounded by

$$\tau_{max} = \max_{i \geq 1} \tau_i = \max_{i \geq 1} \frac{BL_i \cdot r \cdot \Delta t}{\hat{\rho}(t_i) \cdot r}.$$

- Using **geometric weighting**, the maximum backlog delay is given by

$$\begin{aligned} \tau_{max} &= \max_{i \geq 1} \frac{BL_i \cdot r \cdot \Delta t}{\hat{\rho}(t_i) \cdot r} \\ &\stackrel{(2)}{=} \max_{i \geq 1} \frac{\hat{\rho}(t_i) \cdot r \cdot \Delta t}{\hat{\rho}(t_i) \cdot r} = \frac{1}{\alpha} \cdot \Delta t. \end{aligned} \quad (5)$$

- Using **arithmetic weighting**, the maximum backlog delay is given by

$$\begin{aligned} \tau_{max} &= \max_{i \geq 1} \frac{BL_i \cdot r \cdot \Delta t}{\hat{\rho}(t_i) \cdot r} \\ &\stackrel{(4)(3)}{=} \max_{1 \leq i \leq w} \frac{\frac{1}{3}i^2 - (w+1)i + w^2 + 2w + \frac{2}{3}}{-i + 2w + 1} \cdot \Delta t. \end{aligned}$$

Since the above fraction is monotonically decreasing (with $1 \leq i \leq w, w \geq 1$) it follows

$$\tau_{max} = \frac{BL_1}{\hat{\rho}(t_1)} \cdot \Delta t = \frac{w+1}{2} \cdot \Delta t. \quad (6)$$

Estimator	$G_{0.1}$	$G_{0.3}$	A_{10}	A_{20}
max. Delay	10	3.3	5.5	10.5

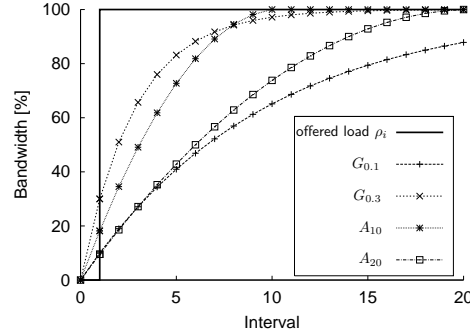
Table I. Upper bounds for τ_{max} given as multiples of Δt .

Figure 2. Estimators in the worst case

From the equations above an upper bound for τ_{max} can be computed, which allows us to determine the acceptable range of values which can be chosen for α or for w , if a given delay bound τ_{QoS} and a sampling interval size Δt have to be respected. Table I shows the maximum backlog delay for different estimators. As we will show in Section 5.2 upper bounds for τ_{max} for estimators combined with threshold systems can be computed using a recursive approach leading to significantly better bounds than the closed-form expression as derived in [30].

3.4. Responsiveness of different load estimators

Besides requiring appropriate throughput, applications with real-time requirements need guarantees for upper limits of delay. In order to keep these limits, the parameters, in particular Δt and α or w respectively, have to be chosen carefully.

For large Δt values the efforts for measuring and calculating the actual load are reduced. Furthermore the measurements are smoothed in a stronger way. The less smoothed the measurements are the more data rate has to be reserved in advance to be able to guarantee small delays and $\rho_i \leq 1$. In contrast, in order that the system can react within an appropriate response time, Δt has to be small. As can be seen from Equations (5) and (6) an upper bound for Δt is defined depending on the delay limit from the application and on the deployed weighting estimator and its parameter α or w respectively. The choice of the system's parameters depends – of course – on the real-time flow. Concerning video streams, in particular with 25 frames per second, $\Delta t = 40ms$ seems to be a good choice since each measurement interval consists of the data for one frame while the responsiveness is well enough for typical video applications.

Figure 2 compares the reactivity of four estimators, two estimators for either arithmetic and geometric weighting, again assuming the worst case load scenario being the one showing the

highest possible fluctuation after the estimator has taken its lowest value (cf. Equation (1)).

We observe that functions A_{10} and $G_{0.3}$ are more reactive than functions A_{20} and $G_{0.1}$. This results from the fact that parameterization with α values near to 1 using geometric weighting or small w values using arithmetic weighting leads to the advantage, that an increase of load is taken into account rather quickly and thus the past (“history”) is “faded out” fast as well. But with this also quite high variations in function $\hat{\rho}(t)$ come along.

We also observe that $G_{0.3}$ is more reactive than A_{10} during the first samples (it yields the highest estimates during 8 samples) but $G_{0.3}$ takes much longer than A_{10} to estimate the traffic as being approximately 100% (4 more samples). This means that $G_{0.3}$ reacts faster to traffic fluctuations, but it needs more time to converge.

4. SYSTEM OF LOAD THRESHOLDS TO CONTROL DISTRIBUTION OF BANDWIDTH

In this section we initially introduce a system of load thresholds, which determines when to free or claim back bandwidth based on current load estimation as discussed in the preceding section. The main tasks of such a system are on the one hand to assure real-time requirements to be respected and secondly to avert oscillations caused by small variations of the estimator’s value. After giving the definition, including the introduction of reasonable constraints, we present three sample parameterizations used to illustrate as well their representation by state models as their analysis concerning potential increase of network utilization.

Definition 1. We define an *n*-state threshold system as tuple $TS(S, \vartheta)$. The *n*-tuple $S = (S_1, S_2, \dots, S_n) \subset [0, 1]^n$, with $S_i < S_j$ for all $1 \leq i < j \leq n$, denotes the set of states while the state-transitions are defined by the $n \times (n - 1)$ -matrix

$$\vartheta = \begin{pmatrix} \vartheta_{1,2} & \cdots & \vartheta_{1,n} \\ \vdots & \ddots & \vdots \\ \vartheta_{n,2} & \cdots & \vartheta_{n,n} \end{pmatrix} \subset (0, 1]^n \times (0, 1]^{n-1}$$

with $i < j \Rightarrow \vartheta_{k,i} < \vartheta_{k,j} \quad \forall k$. With given state S_i and estimated load $\hat{\rho}$ the next state is given by

$$r(S_i, \hat{\rho}) := S \max_{j \in \{1, \dots, n\}} \{\vartheta_{i,j} \leq \hat{\rho}, 1\}$$

Interpretation of ϑ :

- Up-thresholds: For $i < j$, each entry $\vartheta_{i,j}$, determines the up-threshold which has to be crossed upwards (or at least reached) by the estimator $\hat{\rho}$, in order to change the system’s state from S_i to S_j .
- Down-thresholds: For $i \geq j$, each entry $\vartheta_{i,j}$, determines the down-threshold which has to be crossed downwards by the estimator $\hat{\rho}$, in order to change the system’s state from S_i to S_{j-1} .

If neither any up-threshold nor any down-threshold is crossed, the system remains in the current state. Transitions are passed through periodically after each multiple of Δt . Figure 3

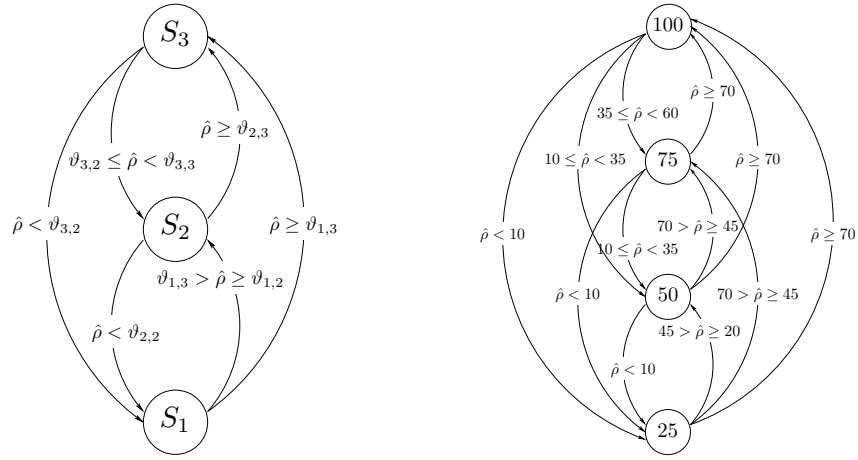


Figure 3. A general 3-state threshold system (left picture) and M2, an example 4-state threshold system; states and transitions of M2 are labeled in percent.

shows a general 3-state (left picture) and a specific, i.e. actually parameterized, 4-state threshold system by way of example.

Definition 1 is kept little restrictive in order not to limit its universality. Nevertheless, concerning our study we assume reasonable constraints for the rest of the paper:

C1 $S_n = 1$ (or 100%), otherwise full capacity can not be reached

C2 down-thresholds *in* S_i are strictly less than up-thresholds *to* S_i ,
formally $\forall i, k \quad 1 \leq k < i \Rightarrow \vartheta_{i,i} < \vartheta_{k,i}$

C3 up-thresholds *in* S_i are greater than up-thresholds *to* S_i ,
formally $1 \leq k < i \Rightarrow \vartheta_{i,i+1} > \vartheta_{k,i}$

C4 down-thresholds *in* S_i are less than down-thresholds *to* S_i ,
formally $n \geq k \geq i + 1 \Rightarrow \vartheta_{i,i} < \vartheta_{k,i+1}$

Nota bene: The control messages FREE_BW and RECALL_BW are sent if a transition $S_i \rightarrow S_j$ occurs with $i > j$ or $i < j$ respectively.

4.1. Potential increase of overall network utilization: worst case study

Let us consider resource reservation and static allocation to communicating entities. Hence, if just a small part of the reserved data rate is used, the remainder is lost for the overall network – in the worst case this sacrifice can reach 100% of the initially reserved bandwidth, obviously, if no real-time data is transmitted at all. Unlike static allocation, our approach is focused on methods to reduce such sacrifice even if quality of service guarantees are provided. So it is interesting to investigate, for different boundary conditions, how much capacity is made available using the presented threshold system. In section 4.1 we want to calculate the amount

of bandwidth which is made available to other stations even in the case that the stream with real-time requirements shows the most unfavourable behaviour possible at all.

So, let us consider how much network capacity can, in the worst case, be reserved but remain unused by a station though a system for sharing bandwidth using our proposal was employed. This sacrificed capacity, in multiples of the total reserved data rate r , can not exceed the maximum difference between the estimated load $\hat{\rho}$ and its associated state. As mentioned in Section 3.3 the entire load measured in each sample interval is considered adequately in the long run by the presented estimators. Therefore, it is acceptable to consider the estimated load $\hat{\rho}$ instead of the actual load in order to assess the utilized and the provided capacity respectively.

Thus, if we apply an arbitrary n -state threshold system, the bandwidth σ_{max} being lost in the worst case is given by

$$\sigma_{max} = r \cdot \max \left(S_1, \max_{2 \leq i \leq n} (S_i - \vartheta_{i,i}) \right)$$

or, equivalently, a utilization of $1 - \sigma_{max}/r$ can be achieved in the network even in the worst case. To give an example, we obtain $\sigma_{max} = r \cdot (S_2 - \vartheta_{2,2}) = r \cdot 40\%$ for the 4-state threshold system (M2) of Figure 3.

4.2. Potential increase of overall network utilization: average case study

With the worst case assessment above, we can derive hard boundaries. In this section our objective is to estimate the potential increase of overall network utilization, in form of upper bounds of sacrifice, for an average case. We do not assume further restrictions on the sending behavior of the communicating stations. Therefore, for this evaluation we presume $\hat{\rho}(t_i)$ as being uniformly distributed in the interval $[0,1]$.

To analyze the overall sacrifice, we determine sacrificed capacity σ for ranges of possible values for $\hat{\rho}(t_i)$. Since the station's state is not necessarily determined uniquely, we derive two values, namely σ_o and σ_p , respectively, which correspond to the optimistic (or pessimistic) assumptions that the station's state is the most favorable (or most unfavorable) one all along.

First, we consider the optimistic assumption. The most favorable state, regarding used capacity, is the one with the highest value. Each state is the most favorable in the range between its highest down-threshold and the highest down-threshold of the next higher state or 1 respectively. Below this range, the state would have changed in consequence of crossing a threshold and above it another (more favorable) could be possible. Since we presume $\hat{\rho}(t_i)$ as being uniformly distributed in the interval $[0,1]$, for each state we calculate its sacrifice by the difference between its stated transmission capacity and the mean of the described range, multiplied by its share. Thus, presuming C1-C4, with optimistic assumptions we derive

$$\sigma_o = \sum_{i=1}^n (\vartheta_{i,i+1} - \vartheta_{i-1,i}) \cdot \left(S_i - \frac{\vartheta_{i,i+1} + \vartheta_{i-1,i}}{2} \right)$$

with $\vartheta_{0,1} := 0$ and $\vartheta_{n,n+1} := 1$. Analogically, with pessimistic assumptions we derive

$$\sigma_p = \sum_{i=1}^n (\vartheta_{i+1,i+1} - \vartheta_{i,i}) \cdot \left(S_i - \frac{\vartheta_{i+1,i+1} + \vartheta_{i,i}}{2} \right)$$

with $\vartheta_{1,1} := 0$ and $\vartheta_{n+1,n+1} := 1$.

Sample parameterizations are given below with their individual assessments, cf. Table II. As we will see in Section 5, they meet well reality.

4.3. Sample parameterizations of threshold systems

We present and discuss three parameterized threshold systems which illustrate as well their representation by state models as their analysis concerning potential increase of network utilization. We define three models $M1$, $M2$, and $M3$ as 2-state, 4-state, and 9-state threshold system respectively. $M2$ is shown in the right picture of Figure 3 by way of example.

$$M1 = TS \left((0.5, 1); \begin{pmatrix} 0.30 \\ 0.45 \end{pmatrix} \right) \quad M2 = TS \left((0.25, 0.5, 0.75, 1); \begin{pmatrix} 0.20 & 0.45 & 0.70 \\ 0.10 & 0.45 & 0.70 \\ 0.10 & 0.35 & 0.70 \\ 0.10 & 0.35 & 0.60 \end{pmatrix} \right)$$

$$M3 = TS \left((0.2, 0.3, \dots, 0.9, 1); \begin{pmatrix} 0.15 & 0.25 & 0.35 & 0.45 & 0.55 & 0.65 & 0.75 & 0.85 \\ 0.10 & 0.25 & 0.35 & 0.45 & 0.55 & 0.65 & 0.75 & 0.85 \\ 0.10 & 0.20 & 0.35 & 0.45 & 0.55 & 0.65 & 0.75 & 0.85 \\ \vdots & & & & & & & \\ 0.10 & 0.20 & 0.30 & 0.40 & 0.50 & 0.60 & 0.70 & 0.80 \end{pmatrix} \right)$$

These models exemplify different parameterizations for the decisions how much bandwidth to be allocated to a specific node depending on its momentary traffic estimation. By taking care of recalling bandwidth before the level of utilization is reached, real-time requirements can be respected. The three sample models were assorted considering the following principles: In this study we analyze equidistant upper and lower thresholds. Since perception of load variations is delayed as a consequence of the discrete sample interval size, the actual available bandwidth should always exceed the estimate, e.g. by at least 5%. Furthermore, to effectively avoid oscillations it is essential that upper and lower threshold have sufficient distance. We use a minimum difference of 10%. As mentioned previously, we presume that any host will distribute transmission capacity only in the case that no local backlog exists.

The fundamental differences between the state models are their granularities. State model $M1$ consists of two states, state model $M2$ consists of four states, and model $M3$ is composed of nine states. Intuitively, more states will be more efficient in terms of local loss, but globally, more states will cause more signaling and this will not necessarily be more efficient from a network's point of view. Evidently, different stations can use different threshold systems as well as different load estimators depending on their application. Both granularity and levels of thresholds can be customized to individual needs.

Using the proposed models $M1$, $M2$, or $M3$, respectively, at least 30%, 60%, or 80% of the a-priori reserved capacity is utilized or placed at other stations' disposal. These values regarding the worst case result from the difference between reserved data rate and maximum sacrificed capacity.

So, the minimum utilized or freed capacity can be enhanced by increasing the number of states of the model. Evidently, communication overhead used for signaling increases as a consequence. We neglect this overhead, depending on individual implementations, for all of the analysis within Section 4. As results from Section 5 will show, this simplification can be

Model	σ_o	σ_p	used/lent out
M1	27.5%	35%	65-72.5%
M2	16.25%	23.75%	76.25-83.75%
M3	10%	14%	86-90%

Table II. Average case: Estimated sacrifice σ_o and σ_p as well as estimated overall utilization of bandwidth

acceptable in real load scenarios, particularly if the user data rate is quite high as for typical video applications.

It is worth mentioning that already a very simple 4-state model such as Model M2 is, in any case, able to make sure that at least 60% of the available real-time bandwidth is utilized by the communicating stations – assuming that the offered load, e.g. in terms of best effort traffic, is sufficiently high.

Concluding, Table II shows the average case assessments under optimistic and pessimistic assumptions as well as the estimated utilization of the reserved data rate for the proposed threshold systems. The estimated values let us expect significant enhancements for real-time communications within broadcast networks. In Section 5 we will check those potential enhancements under realistic load scenarios – and show that the derived assessments meet well reality.

5. NODE BEHAVIOR BASED ON LOAD ESTIMATOR AND THRESHOLD ALLOCATION SYSTEM

5.1. Investigation to study the responsiveness of capacity allocation

The main goal of this section is to evaluate how well applications' real-time constraints can be guaranteed by using our proposed threshold system in combination with different load estimation functions. In particular, we evaluate the time until transmission capacity is called back from other stations when it is needed again by its owner. We verify that resources are called back sufficiently quickly in order to ensure that real-time requirements are met – that is as a rule within 100ms, but depends on the relevant application.

First, we investigate the worst case, that is a situation when the offered load increases to its maximum after a node has redistributed all of its capacity to the other nodes in the network, cf. Equation (1). In Section 3.4 we discussed the responsiveness of different estimators and now we evaluate how fast the up-thresholds are crossed. Our concern is to get back full capacity as soon as possible. Therefore, the reactivity of each model corresponds to the number of intervals needed to get back full capacity.

Using $G_{0.3}$, full capacity is reached within only 2 sampling intervals using state model M1. That is faster than in the case of state model M2 which needs 4 sampling intervals and faster than in the case of state model M3 (6 sampling intervals). This is due to the fact that an estimator $\hat{\rho}$ greater than 45% in model M1, 70% in model M2, and greater than 85% in model M3 is needed to get back full capacity. In other words, a high value of the estimator is needed to reach an allocated value of 100% in the case of model M3. Moreover, $G_{0.3}$'s slope decreases

significantly when 80% of the capacity is reached; therefore, in this case model M2 is more reactive than model M3.

From a networking point of view, the gain obtained using model M1 is even greater since each threshold crossing is accompanied by signaling information and thus more signaling is needed if model M2 or model M3 is used.

5.2. Maximum delay induced due to threshold systems

In Section 3 we analytically derived upper bounds for the maximum backlog delay which resulted from the estimators' underestimations of the actual load and from the time lag between load appearance and its measurement. In the following we discuss the impact of the threshold systems combined with weight functions as introduced previously. In particular, we upper-bound the maximum delay τ_i that can occur for data received during interval T_i .

For this we determine that instant $t' \in T_i$ at which data is delayed longest. The delay is computed then as difference between the instant t'' , that is the instant when all data waiting in the sending queue at instant t' has just been sent, and t' .

The backlog delay within T_i is maximal for the instant $t' := t_{i-1} + \frac{\rho_i \cdot r \cdot \Delta t}{r} = t_{i-1} + \rho_i \cdot \Delta t$. Assuming the worst case that data is generated with the maximum possible data rate r during the first part of the interval, t' is the instant when the last data unit within T_i has been added to the sending buffer. In order to determine t'' , with S_i denoting the portion of r allocated during T_i by the threshold system, let B_i denote the amount of backlogged data buffered after interval T_i :

$$B_i = \max(0, B_{i-1} + (\rho_i - S_i) \cdot r \cdot \Delta t) \quad \text{for } i > 0$$

Algorithm 1 DELAY(x, i)

Require: Amount of data x and index i of interval; S_i denotes portion of r allocated within T_i .

Ensure: Time needed to send amount of data x starting at interval T_i ;

```

 $x_r \leftarrow x - S_i \cdot r \cdot \Delta t$  /* data not sent during  $T_i$  */
if  $x_r > 0$  then
    return  $1 + \text{DELAY}(x_r, i + 1)$  /* start sending remainder in  $T_{i+1}$  */
else
    return  $\frac{x}{S_i \cdot r}$  /* time to send  $x$  in  $T_i$  */
end if

```

Furthermore, we introduce DELAY(x, i), described by Algorithm 1, computing the maximum time needed to send the amount of data x where sending starts at the beginning of interval T_i . With

$$t'' = t_{i-1} + \text{DELAY}(B_{i-1} + \rho_i \cdot r \cdot \Delta t, i)$$

it follows

$$\tau_i \leq t'' - t' = \text{DELAY}(B_{i-1} + \rho_i \cdot r \cdot \Delta t, i) - \rho_i \cdot \Delta t.$$

Applying the load scenario given by Equation (1), the maximum delays shown by Table III can be computed. It can be seen that even hard real-time requirements can be solved by selecting adequate estimators and models. To give an example, by using a geometric estimator with $\alpha = 0.3$ combined with Model M1, the maximum backlog delay does not exceed Δt .

	$G_{0.1}$	$G_{0.3}$	A_{10}	A_{20}
Model M1	3	1	1.5	3
Model M2	5.25	1.75	2.5	4.75
Model M3	6.7	2.2	3.1	5.5

Table III. Maximum backlog delay in multiples of Δt in the load scenario given by Equation (1) for different threshold systems.

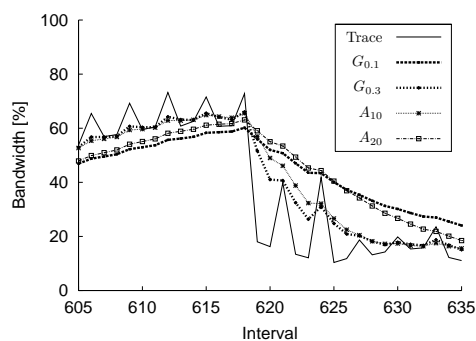


Figure 4. Jurassic Park/Mr Bean trace – different estimators

5.3. Dynamic behavior of utilization for realistic load situations within a station

We evaluate now our scheme in a realistic situation. We studied the behavior of the dynamic resource management for several different real-time streams represented by real traces. As representatives for real-time streams with in fact variable bit rates, which we consider in this paper, we investigated video streams. We studied traces generated from MPEG-2 and MPEG-4 streams as well as video sequences coded according to the ITU Standard H.263. For the latter some of the resource redistribution's possible impacts on the perceived video quality are pointed out in [31]. In this section we single out one typical example.

The input MPEG-4 trace chosen is a mix of the two MPEG-4 traces obtained from [32]. The traces are that of the movies “Jurassic Park” and “Mr. Bean” which we superimpose. We study the reactivity of each transition model using arithmetic and geometric estimators with different weights with this typical stream of total size 573.8 MB, lasting 59.37 minutes. The size of the sample intervals is $\Delta t = 40$ ms (corresponding to 25 frames per sec.) with maximum measured data within one interval of $d_{max} = 178.41$ kbit which induces an a-priori reserved capacity of $r = 4.4$ Mbit/s in order to fulfill the precondition $\rho_i \leq 1 \forall i$.

Figure 4 shows a clipping from the trace and its corresponding estimator values. We have selected sampling intervals during which many fluctuations can be observed in the input trace in order to see how reactive the models are. As mentioned previously, in all cases, $G_{0.3}$ is the most reactive of the traffic estimators compared. The bandwidth allocated by each of the models M1-M3 with estimators $G_{0.3}$ and A_{10} are represented in Figure 5. This figure describes how the different state models realize the bandwidth allocation in the case of real

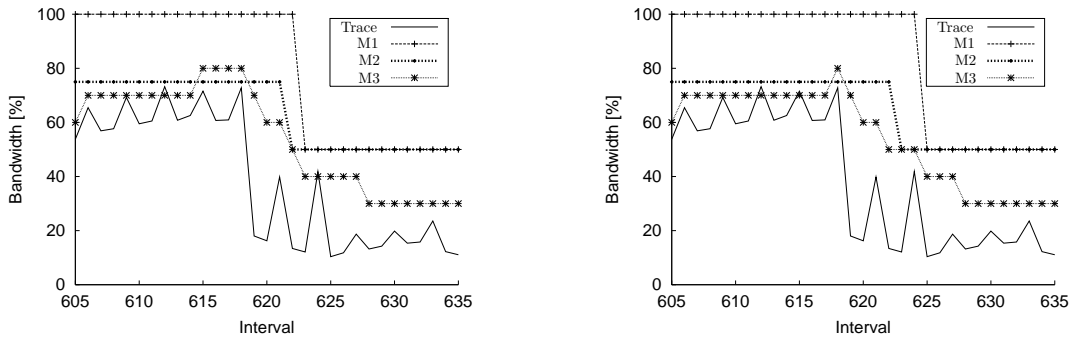


Figure 5. Jurrassic Parc/Mr Bean trace – $G_{0.3}$ (left picture) and A_{10} combined with models M1-M3

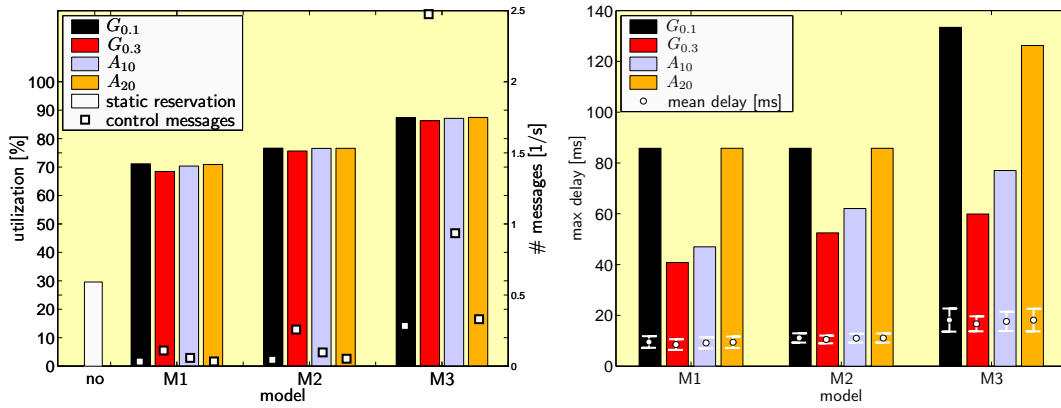


Figure 6. Video mix sent with $r = 4.4 \text{ Mbit/s}$ using different combinations of estimators and threshold systems: left diagram: overall bandwidth available for usage within the network and mean message rate for signaling purposes; right diagram: maximum and mean delay (including 95% confidence intervals)

traffic fluctuations.

Figure 6 shows the network capacity put at other stations' disposal and the number of FREE/RECALL messages exchanged by the system (left chart) as well as the maximum and mean time the transmission of data was delayed for. Mentionable, we neglect the signaling communication overhead needed for controlling the sharing of bandwidth; the actual gained capacity depends on the details of individual implementations.

Previously, we have chosen the data rate r so that it is never exceeded by the stream. Since we can not always act on this presumption of excessive over-reservation, we start further experiments with the trace. The data rate is set to $r = 125 \text{ kbit}/40 \text{ ms}$, whereby the video mix uses 42.2% of the bandwidth (instead of 29.6%) which appears more pragmatical. Figure 7 shows the results of the run without prior smoothing, tolerating $\rho > 1$.

Contrary to the worst case, in all experiments transition model M3 achieves the minimal

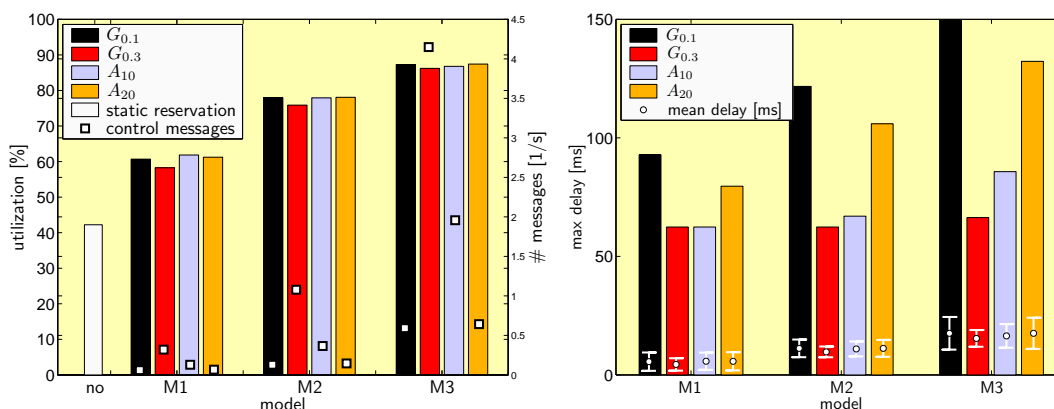


Figure 7. Video mix sent with $r = 3.1$ Mbit/s using different combinations of estimators and threshold systems: left diagram: overall bandwidth available for usage within the network and mean message rate for signaling purposes; right diagram: maximum and mean delay (including 95% confidence intervals)

capacity loss. This is due to the fact that real traces often produce variations which are less radical than the worst case, whereas more detailed models can better adapt the resource allocation to the actual requirements. However, the allocated bandwidth of M3 is often smaller than the value required by the input trace. M2 fits this requirement better since it is very rarely smaller than the bandwidth required by the trace. Therefore, a more detailed model achieves better performance in terms of lost bandwidth, but also introduces more delay if the trace fluctuations are rapid and also more communication overhead is generated.

It can be seen, that significant enhancements under realistic load scenarios are achieved by the threshold systems. It is also worth mentioning that the derived assessments from Section 4.2 meet well the realistic traces we tested, cf. Table II. Further experiments, in particular with smoothed traces, show similar results. Though offered load sometimes exceeds the presumed upper bound of $\rho_i \leq 1$, neither mean nor maximum delay differ significantly, although the number of exchanged messages rises.

An overall interpretation of the results obtained in our various series of experiments, covering traces obtained from a variety of different video codecs, leads to the following insights:

- As to be expected, the gained bandwidth G increases starting with model M1 via M2 towards M3 and this holds independently of the weighting function used ($G_{0.1}, G_{0.3}, A_{10}, A_{20}$). Evidently, the advantages of more refined threshold models become much more significant if one reduces the amount of resources (bandwidth) as reserved during connection-setup-time.
- Some direct dependency between the results regarding buffer occupancy B and additional backlog delay τ exists as a direct consequence of *Little's Law* [33].
- The mean backlog delay $\bar{\tau}$, in all cases, remains below 20 ms, which is encouraging and seems to be fully acceptable for a large variety of multimedia applications with real-time requirements. The maximum backlog delay τ_{max} , too, remains pleasingly small, at least for cases where the weight function $G_{0.3}$ or A_{10} is used.
- The number of control messages to be exchanged remain practicable even in the worst

case as the control messages have a rather small size. So we conclude that the overhead of control messages to be broadcasted seems to be perfectly acceptable, too. And this holds even for network configurations where a relatively large number of stations (e.g. 50 to 100) with real-time communication requirements is present in each one of possibly many interconnected broadcast networks.

In existing LANs the signalling could be directly based on UDP using the following formats by way of example: A signalling message could be encoded using 16 *byte* (it contains the protocol version, some flags, and the table entries (see below) used to run the algorithm). Therefore, once encapsulated in UDP, IP, and Ethernet, a message is 512 *bit* long. Assume each node is sending 2.5 messages per second when it sends high priority traffic (the signalling overhead represents 1280 *bit/s*) and one message every five seconds otherwise (102.4 *bit/s*). Assuming that each stream, e.g., reserves about 150 *kbit/s* and that nearly all of the network's bandwidth is reserved statically, leads to the total signalling overhead of about 0.8% for high priority traffic and less than 0.1% in case of low priority traffic, independent from the number of communicating nodes. Actually, in [34], the simulation scenarios chosen demonstrate an overhead of the order of 0.5% of the network capacity. Therefore, we conclude that our proposal scales well in broadcast LANs.

A possible implementation of our proposal is described in [34]. This implementation relies on the DiffServ paradigm and implements two service classes (BE and EF) that are classified according to the content of the DS field in the IP header. Since BE traffic must limit its rate to spare bandwidth for real-time EF flows, we use a dynamic shaper, tuned by our algorithm, at the outgoing interface of nodes. Real-time traffic is thus isolated by reserving allocated bandwidth resulting from the threshold system and by sharing the remainder among BE flows. Each node must be aware of real-time load in the network to reserve that quantity, and of current demands of each flow. This is implemented with the help of a table representing the usage of bandwidth in the overall network. An entry in the table comprises the node address, the bandwidth allocated to EF flows of this node, and the estimated BE rate. To keep tables synchronized, signalling messages are broadcasted within the LAN. A timer, representing the last broadcast from a corresponding host is associated to each entry. The entry is removed from the table if the aging interval expires. On receiving a signalling message, the entry is updated, if it exists, and the timer is rescheduled; otherwise, a new entry is created. A message broadcast occurs on the following events:

- variations in EF traffic, resulting in a state transition in the threshold allocation system,
- major changes in local rate in terms of BE traffic,
- necessity to refresh the state information and thus to avert being removed from tables of the other nodes.

6. CONCLUSIONS

In this paper, we have investigated how the efficiency of reservation strategies that execute some fixed resource reservations (in particular reservations for the complete duration of continuous media streams) can be increased. On the one hand, we have proposed two classes of parameterized weight functions for calculating estimates of the actual load as induced by

the connections for which resource reservations have been established for the present. On the other hand, we have introduced threshold systems in order to reduce the update traffic for exchanging control information the purpose of which is to hand over free bandwidth and to recall own resources between the stations of a broadcast network. We have demonstrated how weight functions and threshold systems can be combined.

Usage of mathematical performance predictions allowed us to prove that the algorithms we suggest are able to respect real-time bounds for maximum delays in the order of 50 to 100 ms (as they are quite typical in a large variety of audio/video applications) and that, at the same time, they achieve a high overall utilization of resources even under worst case load fluctuations. Additional studies have provided the encouraging result that in practice, i.e. for realistic load scenarios, the efficiency to be expected for our proposal to execute dynamic and load-adaptive resource reservations will even be considerably higher than the "worst case bounds".

Future investigations are planned by us to get a still better understanding regarding the dependencies between the parameterization as used for specifying the load estimators and the granularity of the threshold system being chosen. As a final goal we strive to solve the optimum configuring of adaptive reservation systems, for given load scenarios and given QoS requirements, in a largely automatic way and at the same time limiting the overhead resulting from the exchange of control information supporting the dynamic redistribution of resources.

REFERENCES

1. A. Jamalipour and J. Kim. Measurement-Based Admission Control Scheme with Priority and Service Classes for Application in Wireless IP Networks. *Int. Journal of Communication Systems*, 16(6):535–551, August 2003.
2. R. Braden, D. Clark, and S. Shenker. Integrated Services in the Internet Architecture: An Overview. RFC 1633, IETF, June 1994.
3. S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss. An Architecture for Differentiated Services. RFC 2475, IETF, Dec. 1998.
4. B. E. Wolfinger, J. Wolf, and G. Le Grand. Improving Node Behavior in a QoS Control Environment for Local Broadcast Networks. In *Proc. of the Int. Symposium on Performance Evaluation of Computer and Telecommunication Systems (SPECTS)*, volume 35, pages 361–371, Montreal, July 2003.
5. H. J. Chao and X. Guo. *Quality of Service Control in High-Speed Networks*. J. Wiley, 2002.
6. M. Moran and B. Wolfinger. Design of a Continuous Media Transport Service and Protocol. Technical Report TR-92-019, Int. Computer Science Institute, Berkeley, CA, April 1992.
7. L. Massoulié and J. Roberts. Bandwidth Sharing: Objectives and Algorithms. *IEEE/ACM Trans. on Networking*, 10(3):320–328, 2002.
8. P. Anelli and G. Le Grand. Differentiated Services over Shared Media. In *IWQoS*, pages 288–293, Karlsruhe, 2001.
9. G. Le Grand. *Qualité de service dans des environnements Internet mobile*. PhD thesis, Université P. et M. Curie, LIP 6, Paris, July 2001.
10. R. Yavatkar, D. Hoffman, Y. Bernet, F. Baker, and M. Speer. SBM (Subnet Bandwidth Manager): A Protocol for RSVP-based Admission Control over IEEE 802-Style Networks. RFC 2814, IETF, May 2000.
11. R. Braden, L. Zhang, S. Berson, S. Herzog, and S. Jamin. Resource Reservation Protocol (RSVP) - Verison 1 Functional Specification. RFC 2205, IETF, Sept. 1997.
12. A. Ghanwani, W. Pace, V. Srinivasan, A. Smith, and M. Seaman. A Framework for Integrated Services Over Shared and Switched IEEE 802 LAN Technologies. RFC 2816, IETF, May 2000.
13. C. Richards and K. Smith. The PPP Bandwidth Allocation Protocol (BAP), The PPP Bandwidth Allocation Control Protocol (BACP). RFC 2125, IETF, Mar. 1997.
14. M. Bouyer and E. Horlait. Bandwidth Management and Reservation over Shared Media. In *Proc. of SFBSID'97*, Fortaleza, Brasil, Nov. 1997.
15. J. Wroclawski. Specification of the Controlled-Load Network Element Service. RFC 2211, IETF, Sept. 1997.
16. C. R. Kalmanek, H. Kanakia, and S. Keshav. Rate Controlled Servers for Very High-Speed Networks. In *Proc. of the Conference on Global Communications (GLOBECOM)*, pages 12–20, 1990.

17. A. Shalkh, J. Rexford, and K. G. Shin. Load-Sensitive Routing of Long-Lived IP Flows. In *Proc. of ACM SIGCOMM '99*, pages 215–226, Cambridge, USA, 1999.
18. D. D. Clark and W. Fang. Explicit Allocation of Best-Effort Packet Delivery Service. *IEEE/ACM Trans. on Networking*, 6(4):362–373, August 1998.
19. E. P. Rathgeb. Modelling and Performance Comparison of Policing Mechanisms for ATM Networks. *IEEE Journal On Selected Areas In Comm.*, 9(3):325–334, 1991.
20. W. R. Stevens. *TCP/IP Illustrated: The Protocols*, volume 1. Addison-Wesley, 1994.
21. B.-Y. Choi, J. Park, and Z.-L. Zhang. Adaptive Random Sampling for Load Change Detection. *Performance Evaluation Review*, 30(1):272–273, 2002.
22. B. E. Wolfinger, M. Zaddach, K. D. Heidtmann, and G. Bai. Analytical Modeling of Primary and Secondary Load as Induced by Video Applications using UDP/IP. *Computer Commun.*, 25(11/12):1094–1102, 2002.
23. S. Radhakrishnan, S. V. Raghavan, and A. K. Agrawala. A Flexible Traffic Shaper for High Speed Networks: Design and Comparative Study with Leaky Bucket. *Computer Networks and ISDN Systems*, 28(4):453–469, 1996.
24. G. Bai and B. E. Wolfinger. Possibilities and Limitations in Smoothing MPEG-coded Video Streams: A Measurement-based Investigation. In *ITG/GI-Fachtagung MMB'97*. VDE-Verlag, 1997.
25. M. Schwartz and T. Stern. Routing Techniques Used in Computer Communication Networks. *IEEE Trans. on Communications*, 28(4):539–552, April 1980.
26. C. Bouras and A. Gkamas. Multimedia Transmission with Adaptive QoS based on Real-Time Protocols. *Int. Journal of Communication Systems*, 16(3):225–248, April 2003.
27. M. El-Kadi, S. Olariu, and H. Abdel-Wahab. A Rate-Based Borrowing Scheme for QoS Provisioning in Multimedia Wireless Networks. *IEEE Trans. on Parallel and Distributed Systems*, 13(2):156–166, February 2002.
28. A. Malla, M. El-Kadi, S. Olariu, and P. Todorova. A Fair Resource Allocation Protocol for Multimedia Wireless Networks. *IEEE Trans. on Parallel and Distributed Systems*, 14(1):63–71, January 2003.
29. A. S. Tanenbaum. *Computer Networks*. Prentice-Hall, 4th edition, 2003.
30. J. Wolf, B. E. Wolfinger, G. Le Grand, and P. Anelli. Leistungsbewertung von Algorithmen zur dynamischen Ressourcenverwaltung in lokalen “Broadcast”-Netzen. In *Proc. of the GI/ITG-Conf. “Kommunikation in Verteilten Systemen” (KiVS)*. Springer Verlag, 2003.
31. J. Wolf. Network Resource Management for Real-Time Streams within a Multimedia Document Server Architecture. In *Proc. of 49. Int. Wissenschaftliches Kolloquium (IWK)*, volume 2, pages 297–302, Ilmenau, September 2004.
32. A. Wolisz. MPEG-4 and H.263 Video Traces for Network Performance Evaluation, Available under <http://www-tnk.ee.tu-berlin.de/research/trace/trace.html>, 2000.
33. J. D. C. Little. A Proof of the Queueing Formula $L = \lambda W$. *Operations Research*, 9, 1961.
34. F. Harivelo, G. Le Grand, P. Anelli, J. Wolf, and B. E. Wolfinger. Expedited Forwarding for WiFi. In *Proc. of ISWCS'04, 1st Int. Symp. on Wireless Communication Systems*, Mauritius, September 2004.