computer
communications

# Supporting MPEG video VBR traffic in wireless networks

Y. Iraqi[a,b,*], R. Boutaba[b]

[a]*University of Montreal, DIRO, C.P. 6128, Succ. A, Montreal, Quebec, Canada H3C 3J7*
[b]*Department of Computer Science, University of Waterloo, 417 605 Davenport Road, Waterloo, Ontario, Canada N2L 6H8*

## Abstract

In this paper, we analyse the statistical data sets of several MPEG encoded videos and propose a model of the elements of a scene. We also propose and evaluate a novel dynamic bandwidth allocation scheme for MPEG video sources suitable for wireless networks. The proposed model permits the characterisation of the elements of the stream scenes. The proposed scheme is dynamic and pro-active. It automatically adjusts the amount of reserved resources, while guarantying the required QoS. It exploits the structure of the MPEG video stream and allocates bandwidth on a scene basis. The dynamic bandwidth allocation algorithm which has been presented is evaluated using simulation and actual MPEG video data. The performance evaluations showed a major improvement in bandwidth utilisation as compared to other proposed schemes. © 2001 Elsevier Science B.V. All rights reserved.

*Keywords*: Dynamic bandwidth allocation; Wireless networks; MPEG; Scene; Quality of Service; Cell loss ratio; Variable bit rate; Statistical properties

## 1. Introduction

In a wireless network, bandwidth is perhaps the most precious and limited resource of the whole communication system. Therefore, it is extremely important to use this resource in the most efficient way.

One of the main differences between wireless cellular and wired networks is the inherent mobility of the terminals. Handoff is a time-critical feature in wireless mobile communications that has to be addressed to provide seamless multimedia communications under changing radio resource conditions. Handoff ensures the continuity of a call, while the dedicated radio resource changes within one cell or during cell crossing. Handoff has a significant impact on system capacity and performance. Effective and reliable handoff is highly desirable from the subscriber's point of view. The handoff is already a key process in current systems and it is foreseen to gain increasing importance in third and fourth generation cellular systems as cells radius is going to decrease and the number of users is expected to grow dramatically.

A prevalent underlying theme is the techniques used to control the handoff of users as they move between shrinking cells, at greater speeds, and with stricter requirements on both the QoS delivered to the user and the operational costs associated with a connection. The wireless network must provide the requested level of service even if the user moves to an adjacent cell. A handoff could fail due to insufficient bandwidth in the new cell, and in such case, the connection is dropped. The Call Dropping Probability is a very important connection level QoS parameter. Also users already in the system should have higher priority over new users. This is because, from a user point of view, receiving a busy signal is more bearable than having a forced termination. All these considerations make bandwidth management in wireless networks a very complex task.

Video applications produce large amount of data. As a result, video is transmitted in compressed format to reduce the generated data rates. Among the used compression techniques, MPEG is the standard that has recently gained a considerable attention. The MPEG coding scheme is widely used for any type of video applications.

Compressed video sources produce a Variable Bit Rate (VBR) with a considerable degree of burstiness. To guarantee Quality of Service (QoS) for such VBR applications when used over a wireless link, specific resource management solutions must be considered.

Whenever a mobile terminal connects to a base station, the base station will allocate bandwidth to this mobile terminal. This bandwidth will remain constant throughout the duration of the connection.

Resource allocation could be performed according to the

---

* Corresponding author. Department of Computer Science, University of Waterloo, 417 605 Davenport Road, Waterloo, Ontario, Canada. Tel.: +1-519-883-7342; fax: +1-519-885-1208.

*E-mail addresses:* iraqi@iro.umontreal.ca (Y. Iraqi), rboutaba@bbcr.uwaterloo.ca (R. Boutaba).

peak cell rate of the VBR sources. Such an approach leads to under utilisation of wireless resources due to the bursty nature of the sources. The wireless bandwidth will be wasted and the wireless network will experience high call blocking and forced termination probabilities. Resource allocation could be performed based on the sources' mean cell rates. In such approach, video sources will suffer from unacceptable losses and delays (especially those with hard real-time constraints).

These problems can be solved using a dynamic bandwidth allocation algorithm. In this paper, we propose a predictive resource allocation scheme that provides high wireless network utilisation by dynamically reserving only those resources that are needed. The proposed scheme is dynamic and pro-active, i.e. the amount of bandwidth to be reserved is determined "on-the-fly". It requires some communication between the mobile terminal and the base station, but the amount of extra information generated by the mechanism is acceptable in comparison to the capacity gain obtained.

The proposed algorithm exploits the structure of the MPEG video stream and allocates bandwidth on a scene basis. This will result in a high bandwidth gain, which will affect the overall network performance.

However, this can be done only if we have a characterisation of the elements of a scene. In this work, we will also analyse the statistical data sets of several MPEG encoded videos and will propose a model of the elements of a scene.

The proposed scheme aims at facilitating bandwidth management. Using the proposed bandwidth allocation algorithm, VBR traffic can be supported while decreasing its inherent burstiness. With this scheme users need only to change their bandwidth requirements at scenes boundaries, which can take place, as indicated by simulations (Section 9), in time intervals in the order of seconds to minutes. Let us take the following example. Assuming a scene duration average of 10 s. If a mobile terminal is travelling at the speed of 50 km/h, in 10 s the mobile travels about 138 m. If the cells were 10 m in diameter, the mobile would have crossed about 13 cells without changing its bandwidth requirements. This is very important in bandwidth/handoff management. This aspect makes our scheme better suited to wireless cellular networks.

Also information gathered by the base station in terms of the amount of bandwidth a user is using and for how long will help neighbouring cells to make clear-sighted decision about their bandwidth management. Sharing resource management information between neighbouring cells proves to be very important and achieves better performance than schemes considering only local resource information [20].

Although the proposed scheme in this paper can be used in the context of wired networks, we concentrate on its application in the wireless context for the reasons stated above.

The paper is organised as follows. Section 2 formulates the problem and introduces the proposed approach. Section 3 introduces the scene concept. Section 4 presents evaluations of the Cell Loss Ratio (CLR). Sections 5 and 6 present a study of the distribution of the elements of a scene. In Section 7, the corresponding analysis and obtained results are presented. In Section 8, the dynamic bandwidth allocation algorithm is described. Section 9 presents simulations and discusses the performance results. Conclusions and future directions are presented in Section 10.

## 2. Problem statement and proposed approach

Consider a wireless network system that is able to support mobile terminals running applications that require a varying range of bandwidth resources. The wireless network users expect good QoS from the system, for example low call dropping and packet loss probabilities.

Whenever a mobile terminal connects to a base station, the base station will allocate bandwidth to this mobile terminal. This bandwidth will remain constant throughout the duration of the connection. In Ref. [1] the authors suggest the use of different amount of bandwidth depending on user requirements. For example, a voice call user will use a single bandwidth unit (BU) while a video mobile terminal will require several BUs, where a BU is the minimum quota of bandwidth resources that can be assigned to any mobile user.

The above approach is a good solution for CBR sources, however it is clearly inadequate for VBR sources. The bit rate of this kind of sources varies over time and they have most of the time a bursty nature. Compressed video sources are known to produce a VBR with a high degree of burstiness, which needs specific resource management solutions, especially for guaranteed QoS networks.

If resource allocation is performed according to the peak cell rate of the VBR source, the network will be most of the
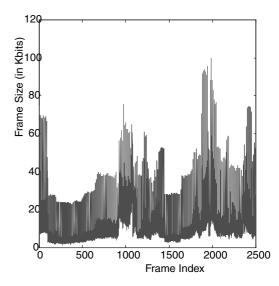


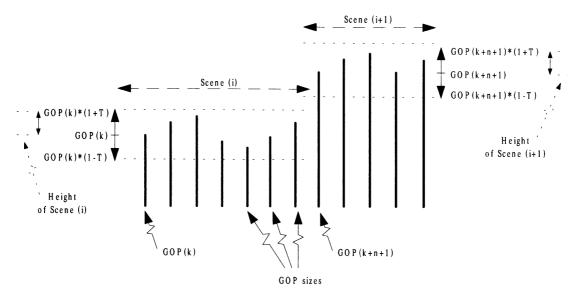Fig. 1. Segment of the frame size sequence for Bond trace.

Fig. 2. Scene duration.

time highly under-utilised when the peak-to-average rate ratios are high. The wireless bandwidth will be wasted and the wireless network will experience high call blocking and forced terminations. On the other hand, if resource allocation is performed based on the source mean cell rate, it is expected for the source to suffer from unacceptable losses and delays (especially for video sources imposing hard real-time constraints).

Several studies investigated the characterisation of MPEG traffic [2–4]. In Ref. [5], the authors concluded that the traffic probability density function (p.d.f.) of some VBR MPEG sources could be modelled by a Gamma distribution. In Ref. [6], the authors have shown that the Gamma and the Log-Normal distributions are good fits for the p.d.f. of $I$, $P$ and $B$ subsets of MPEG sequences. The Log-Normal distribution can also be used as a model for some VBR MPEG sources.

In the case that the model of the VBR source is known, we can calculate the required capacity[1] to have a certain CLR. Even with this approach, big frames are more likely to be affected by a cell loss than small ones, which will affect the visual QoS. For example, consider the VBR source depicted in Fig. 1 offered to a bufferless switch (we consider only hard real-time services) on a wireless link of capacity $C$. Frames with numbers around 2000 will experience a very high cell loss which will be noticed by the user.

Instead of allocating the wireless bandwidth for the lifetime of the connection (as in traditional wireless systems using FDM or CDMA channel access schemes) we will allocate capacity dynamically for each scene. Here a scene represents a group of successive GOPs with close sizes. This capacity will remain constant and will not change until the beginning of another scene. This allocation scheme allows, as we shall see in Section 9, a better bandwidth management. It will lead to an increase in the number of users that can be supported by a mobile wireless network cell without affecting the QoS of the connections.

However, this can be done only if we have a characterisation of the elements of a scene. In the following, we will analyse the statistical data sets of several MPEG encoded videos and will propose a model of the elements of a scene.

## 3. The scene concept

Video applications produce large amount of data. As a result, video is transmitted in compressed format to reduce the generated data rates. Among the used compression techniques, MPEG is the standard that has recently gained a considerable attention. The MPEG coding scheme is widely used for any type of video applications.

An MPEG encoder generates three types of compressed frames: Intra-coded ($I$), Predictive ($P$), and Bi-directional ($B$) frames. An $I$ frame is encoded independently of other frames based on the Discrete Cosine Transform (DCT) and entropy coding. A $P$ frame uses a similar coding algorithm to $I$ frames, but with the addition of motion compensation with respect to the previous $I$ or $P$ frame, and is used as a reference point for the next $P$ frame. A $B$ frame is an interpolated frame that requires both a past and a future reference frames ($I$ or $P$).

Typically, $I$ frames require more bits than $P$ frames. $B$ frames have the lowest bandwidth requirement. After coding, the frames are arranged in a deterministic periodic sequence, for example "*IBBPBB*" or "*IBBPBBPBBPBB*", which is called Group of Pictures (GOP).

From Fig. 1, it is observed that an MPEG trace consists of several segments such that the sizes of $I$ frame in each

---

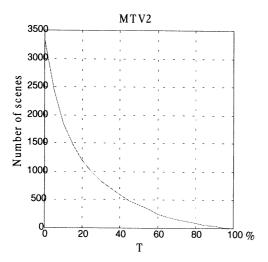[1] The terms 'capacity' and 'bandwidth' are used interchangeably throughout the paper.

Fig. 3. Number of scenes for MTV2 trace.

segment are close in value. In Refs. [7,8], such segments were referred to as scenes. In this paper we consider scenes with respect to GOP sizes. The goal behind this choice is two folds. First, it permits to not distinguish between frame types (*I*, *P* or *B*). Second, it will allow for a uniform characterisation of the scene elements.

To model the length of a scene, the authors in Refs. [7,8] have proposed a method that computes scene duration using the fact that a "sufficient" difference between the sizes of two consecutive *I* frames is a strong indication of the start of a new scene. However, this approach requires the availability of the VBR trace. It takes into account only *I* frames and do not permit a uniform characterisation of all frame types (*I*, *P* and *B*) within a scene.

In this work, we consider two requirements that will lead us to a new algorithm for determining the scene duration in an MPEG stream: First, the proposed algorithm must work on-the-fly, which means that the decision of determining the scene boundaries must only take into consideration the past GOPs. One advantage of such algorithm is the ability to handle MPEG streams for which we do not have a trace. The second requirement concerns the size of the first GOP in each scene, which has to be as close as possible to the mean GOP size of the scene. This will allow us to have a characterisation of the elements of a scene knowing only the size of the first GOP of that scene.

With respect to the above two requirements we compute scene duration differently (see Fig. 2). Let $\{GOP(j) : j = 1, 2, \ldots\}$ be the GOP sequence in an MPEG stream. This sequence consists of the sizes of consecutive GOPs in a given MPEG trace. Suppose that the current scene is the $i$th scene that started with the $k$th GOP. The $(n + k + 1)$th GOP of the sequence indicates the start of the $(i + 1)$th scene if

$$|GOP(n + k + 1) - GOP(k)| \geq T * GOP(k) \qquad (1)$$

where $T$ is a thresholds ($T \geq 0$). $n + 1$ in this case represents

the length of the $i$th scene. Notice that the length of a scene is measured by the number of consecutive GOPs in that scene.

With this definition of scene, all the GOP sizes within a scene $i$ are located between First_GOP$(i) * (1 - T)$ and First_GOP$(i) * (1 + T)$. Where First_GOP$(i)$ is the size of the first GOP in scene $i$.

Clearly, the value of $T$ impacts the shape of the scene length distribution. It determines the amount of correlation between successive scenes; the larger these value, the less correlated the scenes. The value of the $T$ parameter impacts also the number of scenes in a particular trace. Larger values of $T$ produce smaller number of scenes. For example, for the MTV2 trace, a value of $T = 20\%$ i.e. (0.2) produces 1200 scenes while a value of $T = 80\%$ produces 100 scenes (see Fig. 3).

The traces used in our study and simulations were provided by Rose[2] [6]. Rose's movies were taken from VCR tapes, and were digitised at rate of 25 frames/s using a Sun Video card. The movies were compressed using MPEG [9,10] Berkeley's software encoder [11]. Each MPEG video consists of 40,000 frames, which is equivalent to approximately half an hour.

## 4. Cell loss ratio evaluation

Consider a hard real-time VBR service where the stream produced by the VBR source is directed to a bufferless switch on a wireless link of capacity $C$. Let this VBR source bit-rate at time $t$ be $Rt$. The CLR can be estimated by the fluid approximation [12] as follows:

$$CLR = \frac{E\{(Rt - C)^+\}}{E\{Rt\}} \qquad (2)$$

Where $E\{\cdot\}$ represents the expectation operator and $\mathbf{X}^+$ is defined as $X^+ = X$ if $\mathbf{X} > 0$ and $X^+ = 0$ if $\mathbf{X} < 0$. If the p.d.f. of the source rate is defined by $f(u)$, then Eq. (2) can be written as:

$$CLR = \frac{\int_C^\infty (u - C)f(u)\,du}{\int_0^\infty uf(u)\,du} \qquad (3)$$

If the MPEG source can be modelled by a Gamma distribution with a cumulative distribution function (c.d.f.) $F_{(a,b)}$ (where $a$ and $b$ are the parameters of the gamma distribution) as suggested in Ref. [5], Eq. (3) will be:

$$CLR = \frac{a * b * (1 - F_{(a+1,b)}(C)) - C * (1 - F_{(a,b)}(C))}{a * b} \qquad (4)$$

For a normal distribution $\aleph(\mu, \sigma)$ with a p.d.f. $f$ and a c.d.f. $F$

---

[2] The traces can be obtained from the ftp site ftp-info3.informatik.uni-wuerzburg.de in the directory /pub/MPEG/.
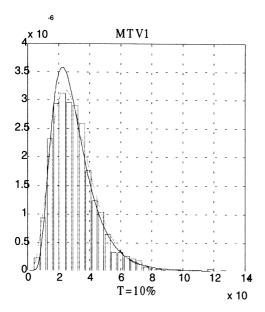
Fig. 4. The histogram of $\{GOP_S(i)\}$ for MTV1 trace with $T = 10\%$.

and using Eq. (3) we obtain:

$$CLR = \frac{\mu * (1 - F(C)) + \sigma 2 * f(C) - C * (1 - F(C))}{\mu} \quad (5)$$

Using Eq. (5) we can calculate the required capacity $C$ for a pre-specified CLR using binary search.

## 5. GOP size distribution within a scene

Based on our definition of a scene (see Eq.(1)), the sizes of GOPs within a scene could be considered close to each other. The GOP sizes fluctuate around an average value that represents the level of activity of the scene. We recall that $GOP(n)$ is the size of the $n$th GOP in the MPEG stream. We suppose that $GOP(n)$ is the sum of two independent random variables:

$$GOP(n) = GOP^*(n) + GOPD(n) \quad (6)$$

where $GOP^*(n)$ reflects the level of activity of the scene, while $GOPD(n)$ represents the fluctuation of the $n$th GOP around $GOP^*(n)$. This assumption will simplify greatly the model.

The quantity $GOP^*(n)$ is constant for all GOPs in the same scene, and it varies from one scene to another. Hence, for the $i$th scene that started with the $k$th GOP, we have

$$GOP^*(k) = GOP^*(k + 1) = ... = GOP^*(k + N_i - 1)$$

$$= GOP_S(i) \quad (7)$$

$N_i$ represents the length of the $i$th scene. Notice that the length of a scene is measured by the number of consecutive GOPs in that scene.

Suppose $\{GOP_S(i): i = 1, 2, ...\}$ is a sequence of independent and identically distributed (i.i.d.) random variables with common p.d.f. $f$. The histogram of $\{GOP_S(i)\}$ for the MTV1 trace is shown in Fig. 4. The shape of this histogram suggests the use of a Log-Normal or a Gamma fit. Solid line represents the Log-Normal fit for the empirical data, dash–dot line represents the Gamma fit.

Now, consider $GOPD(n)$ the quantity representing the variation of the size of the $n$th GOP around $GOP^*(n)$. For a given trace, we compute the empirical sequence $\{GOPD(n)\}$, where $GOPD(n) = GOP(n) - GOP_S(i_n)$ ($i_n$ is the scene index to which the $n$th GOP belongs). Note that $\{GOPD(n)\}$ was assumed as independent of $\{GOP_S(i)\}$, and thus invariant with respect to scene changes, and depends only on the $T$ parameter.

The histogram of $\{GOPD(n)\}$ for the MTV1 trace is shown in Fig. 5. The shape of this histogram suggests the use of a Normal fit.

Since $\{GOPD(n)\}$ is the variation of the GOP sizes around the mean size value of the scene to which each GOP belongs, the mean value of $\{GOPD(n)\}$ is equal to zero. Therefore, we set $f_{GOPD} = \aleph(0, \sigma_{GOPD}^2)$.

Thus the GOP sizes within a scene $i$ can be modelled by a normal distribution $\aleph(\mu, \sigma^2)$ with mean $\mu = GOP_S(i)$ and variance $\sigma^2 = \sigma_{GOPD}^2$. $\mu$ varies from one scene to another
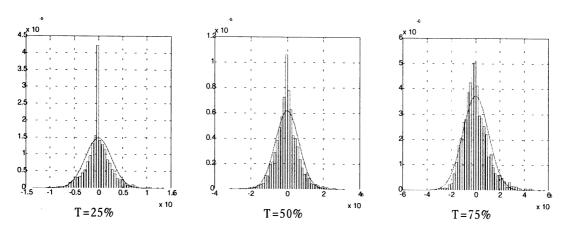


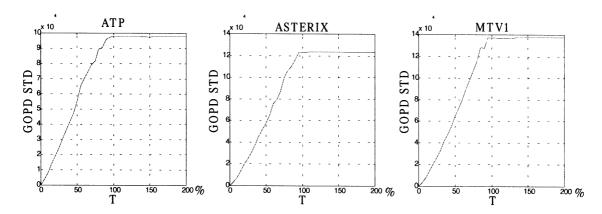Fig. 5. GOPD for $T = 25\%$, 50 and 75% for MTV1 trace.

Fig. 6. GOPD STD for ATP, ASTERIX and MTV1 traces.

(since it reflects the level of activity of the scene) while $\sigma$ is invariant to scene changes and depends only on the $T$ parameter.

## 6. The influence of the $T$ parameter

Our ultimate goal is to allocate bandwidth dynamically on a per scene basis. For this purpose, it is important to characterise the elements of the scene. As this characterisation depends on the $T$ parameter, the study of the influence of the $T$ parameter on the characterisation is of great importance.

Let us recall that $T$ is the parameter that determines the height of the scenes (see Fig. 2). The value of $T$ impacts the shape of the scene length distribution. It determines the amount of correlation between successive scenes; the larger these value, the less correlated the scenes. The value of the $T$ parameter impacts also the number of scenes in a particular trace. Larger values of $T$ produce smaller number of scenes (see Fig. 3).

Our goal is the characterisation of the GOPs within a scene. The idea is to see if from the study of MPEG traces we can approximate or even calculate the mean and the variance of the GOPs within a scene knowing only the size of the first GOP of that scene. This will allow a characterisation of the entire scene based on the size of the first GOP, which can be used to allocate bandwidth for the scene.

We compute the value $GOP(n) - GOP^*(n)$ for every $n$ as the difference between each GOP size and the mean size of the scene to which the GOP belongs. Let us recall that

$$GOPD(n) = GOP(n) - GOP^*(n) \qquad \text{for all } n$$

GOPD depends on the value of the $T$ parameter. For values of $T$ below 200%, the shape of the histogram of {GOPD} (Fig. 5) suggests the use of a normal fit. Values above 200% are not important for us since the entire MPEG stream will not have many scenes, and hence, will not profit from the scene-based bandwidth allocation (SBA). Indeed, for $T$

above 200%, the majority of studied streams are composed of a single scene.

Fig. 5 depicts the histogram for empirical sequences {GOPD($n$)} for different values of $T$, namely 25, 50 and 75% for the MTV1 trace. Solid lines represent the Normal fits for the empirical data.

The normal distribution $f_{GOPD}$ is fully characterised by the mean $\mu_{GOPD}$ and the variance $(\sigma_{GOPD})^2$ of the empirical sequence {GOPD}. Since {GOPD($n$)} depends on the $T$ parameter, the values $\mu_{GOPD}$ and $\sigma_{GOPD}$ depend also on the $T$ parameter.

As said before, since GOP$^*$($n$) represents the mean GOP size of the scene to which GOP($n$) belongs, $\mu_{GOPD}$ is equal to zero for all $T$. We compute $\sigma_{GOPD}$ for different values of $T$ (from $T = 0\%$ to 200% step 1%). The results are presented in Fig. 6.

Fig. 6 shows the variation of $\sigma_{GOPD}$ as a function of $T$ for the ATP, ASTERIX and MTV1 traces.

When $T$ is near 0 ($<10\%$), the scenes contains only few GOPs (we suppose that the GOP sizes are not constant since we consider VBR traffic), $\sigma_{GOPD}$ is also near 0. $\sigma_{GOPD}$ will increase with the value of $T$ until some maximum value and then it remains constant (see Fig. 6). The constant value, $\sigma_{GOP}$, is the standard deviation (STD) of the entire stream. Indeed, when $T$ is very high the entire stream is considered as one scene and the STD of {GOPD} is nothing but $\sigma_{GOP}$.

Since there is some value $T_M$ from which there is only one scene in the entire stream, we can set $\sigma_{GOPD} = \sigma_{GOP}$ for all $T \geq T_M$.

It is worth noting that the value of $\sigma_{GOPD}$ increases almost linearly from zero to $\sigma_{GOP}$. Also, $\sigma_{GOPD}$ begins to stabilise, sometimes, before $T$ reaches $T_M$. This occurs for values of $T$ for which the entire stream contains only few GOPs (2 or 3). Let the first value of $T$ for which this stabilisation occurs be $T_m$ (we have $T_m \leq T_M$).

We can then model the value of $\sigma_{GOPD}$ by

$$\sigma_{GOPD} = T * \sigma_{GOP}/T_m \quad \text{for } T \leq T_m \text{ and}$$

$$\sigma_{GOPD} = \sigma_{GOP} \qquad\qquad \text{for } T \geq T_m$$

(8)

Table 1
$T_M$ and $T_m$ values for some traces

| Film | $T_M$ (%) | $T_m$ (%) | Number of scenes for $T = T_m$ |
|---|---|---|---|
| News1 | 84.90 | 84.90 | 1 |
| MrBean | 88.49 | 88.49 | 1 |
| Lambs | 89.76 | 89.76 | 1 |
| MTV1 | 130.13 | 94 | 2 |
| Simpsons | 158.90 | 97 | 2 |
| ATP | 160.20 | 97 | 2 |
| Soccer | 160.37 | 91 | 2 |
| Bond | 298.66 | 88 | 3 |
| Terminator | 301.19 | 90 | 2 |
| Dino | 304.29 | 92 | 3 |
| Race | 404.59 | 91 | 3 |
| Movie | 404.76 | 97 | 3 |
| MTV2 | 460.89 | 96 | 3 |
| Asterix | 855.45 | 98 | 3 |

From the study of MPEG traces, we notice that the value of $T_m$ is always around 90%. If the maximum GOP size, $GOP_{MAX}$, and the minimum GOP size, $GOP_{min}$, of the entire GOP stream are known, we can calculate the value of $T_M$ as follows:

$$T_M = MAX[(1 - GOP_{min}/GOP(1)),$$
$$(GOP_{MAX}/GOP(1) - 1)] \tag{9}$$

Where $MAX[X,Y]$ is the maximum operator, and $GOP(1)$ is the size of the first GOP of the stream.

Table 1 shows the values of $T_M$ and $T_m$ for different MPEG streams. Note that values of $T_m$ are only approximations found while plotting the variation of $\sigma_{GOPD}$ as a function of $T$. Table 1 shows also the number of scenes for the MPEG traces with $T = T_m$.
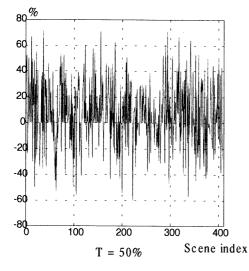


Fig. 7. PERC with $T = 50\%$ and $\mu^*(s) = First\_GOP(s)$ for the MTV1 trace.

## 7. The scene model

Our ultimate goal is to allocate bandwidth dynamically on a per scene basis knowing the first GOP of the scene only. For that purpose, it is important to characterise the elements of the scene.

It is worth noting that we consider the variance and the average of the films as known values for two reasons: Either we have these films beforehand, hence we can calculate these values, or use an encoder that allows us to specify a desired variance. A number of works [12–16] have been done on designing rate control (or rate shaping) mechanisms to enforce the encoder to respect some predefined characteristics like a certain mean rate and target variance.

As shown earlier, the GOP sizes within a scene can be modelled by a normal distribution. For MPEG streams, for which we have already traces, $\mu$ is accurately computed based on those traces. However, for MPEG streams for which we do not have traces beforehand, $\mu$ must be approximated. We must then find $\mu^*$ and $\sigma^*$ that approximate or overestimate $\mu$ and $\sigma$ for each scene. We can use these values to allocate bandwidth for the scene, and if $\mu^* \geq \mu$ and $\sigma^* \geq \sigma$ for each scene then user requirements will not be violated.

The STD $\sigma$ of the GOPs within a scene can be calculated using Eq. (8). To check if we can replace the mean GOP size with the size of the first GOP without violating user requirements (a specified CLR), we compute the difference {GOPDS}. Where $GOPDS(s) = (\mu^*(s) - \mu(s))$ for all scene $s$. This corresponds to check if $\mu^*(s) = First\_GOP(s)$ is viable for each scene $s$.

We compute the percentage that represents $GOPDS(s)$ to the height of scene $s$. The height of a scene $s$ is defined by $First\_GOP(s) * T$ (see Fig. 2). This percentage is determined as follows:

$$PERC(s) = [GOPDS(s)/(First\_GOP(s)^*T)]$$

If $PERC(s)$ is positive, then $First\_GOP(s)$ is bigger than the mean GOP size in the scene $s$. So if we take $\mu^*(s) = First\_GOP(s)$ for that scene, the requested CLR will be respected. If on the other hand, $PERC(s)$ is negative then $First\_GOP(s)$ is smaller than the mean GOP size within the scene $s$ and the requested CLR will not be respected. $PERC(s)$ shows also how far is $\mu^*(s)$ from $\mu(s)$ for each scene $s$.

As shown in Fig. 7, PERC is not always positive. Similar results were found for all values of $T$ below 100% and for all other MPEG traces. This means that for many scenes the size of the first GOP of the scene is smaller that the mean GOP size within the scene. Thus, if $\mu^*(s) = First\_GOP(s)$ is used to evaluate the bandwidth required to satisfy a certain CLR, user requirements will not be guaranteed.

This problem can be solved by correcting the value of $\mu^*$. We propose to add the standard deviation $\sigma_{GOP}$ of the entire stream, obtaining consequently $\mu^*(s) = First\_GOP(s) + \sigma_{GOP}$ for each scene $s$.
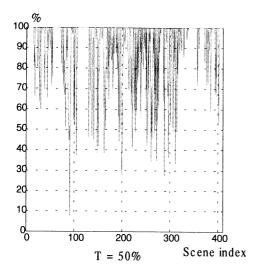
Fig. 8. PERC with the new value of $\mu^*$ for the MTV1 trace.

As all GOP sizes within a scene $s$ are below Frist_GOP$(s)^*(I + T)$, we set $\mu^*(s) = \text{MIN}$ [First_GOP$(s) + \sigma_{\text{GOP}}$, First_GOP$(s) * (1 + T)$] for each scene $s$. MIN[$X,Y$] here represents the minimum operator.

The new percentage that represents $(\mu^*(s) - \mu(s))$ to the height of scene $s$ is now computed. Fig. 8 shows that this new percentage is always positive. Similar results are obtained for all values of $T$ below 100% and for all other studied MPEG traces.

The new value of $\mu^*$ can be used to characterise the elements of each scene without violating user requirements.

In our simulations we have allocated bandwidth to each scene $s$ supposing a normal distribution of the GOPs within the scene as stated in Section 5. We have taken $\mu^*(s) =$ First_GOP$(s) + \sigma_{\text{GOP}}$ and $\sigma^*$ as in Eq. (8). The performed simulations show that the required CLR is always respected.

Fig. 9 shows the obtained CLR for the ATP trace and for different required CLRs. We remark that the required CLR is always respected. For values of CLR below $10^{-4}$, the obtained CLR is always below $10^{-10}$. Similar results are obtained for all the studied MPEG traces.

We have computed the obtained CLR for 16 MPEG streams for different values of $T$ (from $T = 10\%$ to 100% step 10%) and 10 different values of CLR (from $10^{-10}$ to $10^{-1}$ multiplying by 10 in each step). We always have the same and even better CLR (for a CLR below $10^{-3}$).

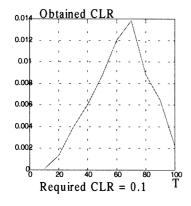## 8. The dynamic bandwidth allocation algorithm

### 8.1. Assumptions

As illustrated in Fig. 10, we consider a mobile terminal that wants to send or receive an MPEG video stream over a wireless link. We assume an underlying mechanism that allows us to allocate bandwidth dynamically throughout the duration of a connection between a base station and a mobile terminal. An example of such mechanism is the ETSI UMTS Terrestrial Radio Access (UTRA) [17]. The proposed algorithm in this paper is not intended to implement such mechanism but to compute the required bandwidth for each scene.
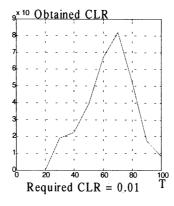
In addition, we assume that when the required capacity, if known, is allocated to the mobile terminal without any delay. In this study we do not take into account the delay between the demand and the acquisition of the capacity. We also assume a wireless system with a capacity that can accommodate the MPEG video traces considered in this study. An example is the fourth-generation mobile wireless networks capable of offering capacities up to 150 Mb/s to fully mobile users in various environments [19].

It is worth noting that in this work the problem of call admission is not addressed and that the capacity gain is obtained supposing that the required capacity is always available. The call admission algorithm proposed in Ref. [20] can be used or some other mechanism that will satisfy the assumption stated earlier.

### 8.2. The allocation algorithm

Our bandwidth allocation algorithm is based on the idea that the GOP sizes within a scene are close. We propose to allocate bandwidth requirement for each scene depending on the GOP sizes mean and variance within the scene. As
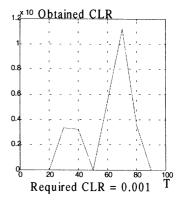
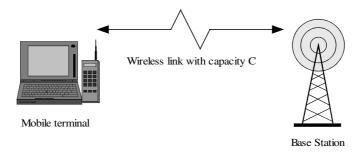

Fig. 9. obtained CLR for ATP trace.

Fig. 10. Mobile terminal wireless access.

shown later in Section 9, the proposed approach requires less capacity than the traditional scheme while guaranteeing the same and even better user QoS requirements.

Based on our definition of a scene, Eq. (1), the sizes of GOPs within a scene could be considered close. The GOP sizes fluctuate around an average value that represents the level of activity of the scene. In Section 5 we showed that the GOP sizes within a scene can be modelled by a normal distribution with mean $\mu$ and variance $\sigma^2$ ($\aleph(\mu, \sigma)$). $\mu$ varies from one scene to another while $\sigma$ is invariant to scene changes and depends only on the $T$ parameter.

We have also shown that $\sigma$ can be approximated by $T^*\sigma_{\text{GOP}}$ when $T$ is below 100%, where $\sigma^2_{\text{GOP}}$ is the variance of the film, and that $\mu$ for some scene $i$ can be replaced by First_GOP($i$) + $\sigma_{\text{GOP}}$ to calculate the capacity required for a specified CLR without violating user requirements, where First_GOP($i$) is the size of the first GOP in scene $i$.

We can then use (First_GOP($i$) + $\sigma_{\text{GOP}}, T * \sigma_{\text{GOP}}$) as an approximation for ($\mu, \sigma$) for the scene $i$ to compute the required capacity and as confirmed by simulation (see Section 7) the predetermined CLR is still respected.

The following algorithm (SBA) can be used to allocate bandwidth dynamically for each scene and can be used either in the mobile terminal or in the mobile station depending on the transfer direction.

The algorithm begins by allocating the required capacity (for a pre-specified CLR) for the first scene depending on scene's first GOP size. Then, it checks the following GOP sizes to detect the beginning of a new scene using Eq. (1). Moreover, depending on the size of the first GOP of the new scene it allocates a different amount of bandwidth using Eq. (5). This capacity will remain constant until the beginning of another scene, and can be, for example, used by neighbouring base stations to reserve bandwidth for the mobile terminal in case it emigrates to another cell.

The algorithm pseudo-code:

*Initialisation:*
    Set the value for $T$
    Set the value for CLR
    $\sigma^2_{\text{GOP}}$ = the stream variance
    $N = 1$
    First_GOP = GOP($N$)   // the first GOP size
    $\mu$ = First_GOP + $\sigma_{\text{GOP}}$

    $\sigma = T * \sigma_{\text{GOP}}$
    Allocate the capacity for $\aleph(\mu, \sigma)$   // using Eq. (5)
    Send GOP($N$)
    $N = N + 1$
*Loop:*
    While not end of stream do
        $S$ = GOP($N$)   // the $n$th GOP size
        If$|S - $ First_GOP$| > T * $ First_GOP then
        // another scene starts
            First_GOP = $S$
            $\mu$ = First_GOP + $\sigma_{\text{GOP}}$
            Allocate the capacity for $\aleph(\mu, \sigma)$   // using Eq. (5)
        End if
        Send GOP($N$)
        $N = N + 1$
    End while

## 9. Simulations and results

In this section we will show how our dynamic bandwidth allocation algorithm surpasses the traditional scheme.

In this work, we will compare the performance of our algorithm in terms of the used capacity with the ones calculated supposing a Gamma distribution MPEG source and a Log-Normal distribution MPEG source. Only the results obtained supposing a Gamma distribution MPEG source are presented. Similar results were found while using a Log-Normal distribution as a model for the MPEG sources.

The capacity gain is calculated as follows:

$$\text{Gain} = 100 * \left(1 - \frac{\text{Dynamic\_Cap}}{\text{Static\_Cap}}\right) \qquad (10)$$

Where Dynamic_Cap is the total capacity allocated by our dynamic algorithm and Static_Cap is the total capacity allocated by the static approach.

To clarify the meaning of the Gain variable, let us take the following example: assume that we have allocated bandwidth to an MPEG stream for 30 min. Knowing that the considered MPEG source can be modelled by a gamma distribution source and using Eq. (4), the static scheme allocates the capacity $C$ for this film (to have a certain CLR). Suppose also that our dynamic allocation algorithm has
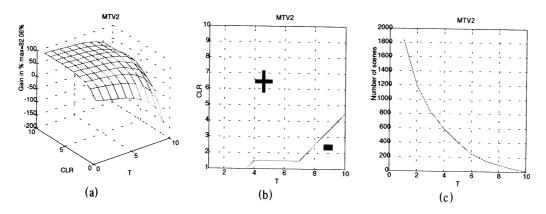
Fig. 11. Simulation results for MTV2 trace.

identified three scenes that last 5, 15 and 10 min, respectively, and allocates the capacity $C1$, $C2$ and $C3$ for the three scenes, respectively. We then have the following capacities:

Static_Cap $= 30 * C$

and

Dynamic_Cap $= 5 * C1 + 15 * C2 + 10 * C3$

and the Gain represents the percentage of the Static_Cap that the dynamic approach did not use. If for example, Gain is equal to 70% then if the static approach uses a particular amount $B$ of bandwidth to have a certain CLR, our algorithm uses only 30% of $B$ to have the same CLR.

We calculate the capacity gain for different values of $T$ (from $T = 10\%$ to 100% step 10%) and different values of CLR (from $10^{-10}$ to 0.1 step 0.1) and for different MPEG video streams.

Although we use traces that are already available, our algorithm remains valid for online traces. This is the case because we did not use any information already available on the used traces. We use only the variance that we consider as a known value.

In the following figures (Figs. 11–15):

- for the axis labelled '$T$', a value of $s$ means that $T = 10 * s\%$
- for the axis labelled 'CLR', a value of $s$ means that the required CLR $= 10^{-s}$ (e.g. the value 5 means that the CLR is $10^{-5}$)

Fig. 11(a) shows the capacity gain obtained for the MTV2 trace in comparison with the static allocation scheme with a Gamma distribution as a model for the MPEG source (the capacity needed for the static allocation is obtained using Eq. (4)). Similar results were found while using a Log-Normal distribution source model. We can notice a significant gain in the area where the required CLR is below $10^{-5}$ and where $T$ is less than 70%.

From Fig. 11(a), we can notice that for $T = 10\%$ and a required $CLR = 10^{-10}$ the capacity gain is 82.06%. This means that if the static approach uses a particular amount $B$ of bandwidth to have a $CLR = 10^{-10}$, our algorithm uses only 17.94% of $B$ to have the same CLR.

Fig. 11(b) shows when the gain obtained is positive and when it is negative. The curve represents the boundary where the two compared schemes have the same performance (i.e. the two schemes require the same amount of bandwidth). The area containing the plus sign represents
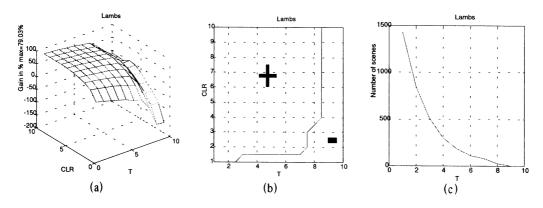


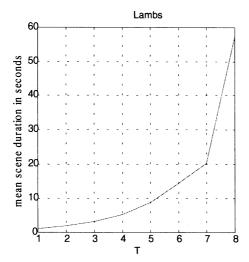Fig. 12. Simulation results for Lambs trace.

Fig. 13. Lambs mean scene duration.

the area where the Gain is positive, which means that the dynamic scheme requires less bandwidth than the static one. The area containing the minus sign represents the area where the Gain is negative, i.e. where the dynamic scheme requires more bandwidth than the static one.

From Fig. 11(b), we can notice that for high values of $T$ ($T > 70\%$) and high values of CLR ($\approx 0.1$), our algorithm use more capacity than the static approach. But this area ($T > 70\%$ and CLR $\geq 0.1$) is not very important since in practice the required CLR is usually below $10^{-5}$.

Fig. 11(c) shows the number of scenes depending on the value of the $T$ parameter. A high number of scenes mean many capacity request messages between the base station and the mobile terminal. But in comparison to the capacity gain obtained by our algorithm, a little overhead is acceptable.

For example, in a wireless ATM context, the total ATM cells for the MTV2 trace is 2080076 cells. But for a capacity gain of 82.06% (obtained for $T = 10\%$ and a CLR equal to $10^{-10}$), an overhead of 1800 ATM cells (corresponding to the number of scenes obtained for $T = 10\%$) is acceptable. We suppose that the capacity requests can be transmitted

using ATM Operation And Management (OAM) [18] cells. We consider one OAM cell per request if the transfer direction is from the base station towards the mobile terminal, and two OAM cells otherwise (request-confirmation).

Fig. 12(a) shows the capacity gain obtained for the Lambs trace in comparison with the static allocation scheme. For this MPEG stream the maximum capacity gain is obtained for $T = 10\%$ and a cell loss ratio CLR $= 10^{-10}$. From Fig. 12(a) and (c), we can notice that for an overhead of 1480 messages the capacity gain is 79.03%. This means that if the static approach uses a particular amount $B$ of bandwidth to have a $CLR = 10^{-10}$, our algorithm uses only 20.97% of $B$ to have the same CLR with an overhead of 1480 additional messages.

As illustrated by Fig. 12(b), our dynamic bandwidth allocation algorithm is always better than the static scheme for practical values of CLR ($\leq 10^{-3}$).

Fig. 13 depicts the Lambs mean scene duration for different values of $T$. Higher values of $T$ lead to a small number of scenes and hence to a high mean scene duration. It is interesting to notice that for $T = 50\%$ the mean scene duration is around 10 s. This means that there is no overhead within this period. With our dynamic algorithm we obtain a 67.65% capacity gain for a CLR $= 10^{-10}$.

Fig. 14(a) shows the capacity gain obtained for the MrBean trace in comparison with the static allocation scheme. For this MPEG stream the maximum capacity gain is again obtained for $T = 10\%$ and a cell loss ratio CLR $= 10^{-10}$. Similar results were found while using a Log-Normal distribution source model.

Here again (see Fig. 14(b)), our dynamic algorithm surpasses the static scheme for practical values of CLR.

For this stream, the number of scenes and consequently the number of overhead messages is lower than the number of scenes for the two traces seen before (MTV2 and Lambs). This can be explained by the fact that MrBean trace has long segments that have a close GOP size values (see Fig. 15).

The following table (Table 2) shows the obtained results for other MPEG video streams.

We have applied our algorithm to a total of 16 MPEG streams for different values of $T$ (from $T = 10\%$ to 100%
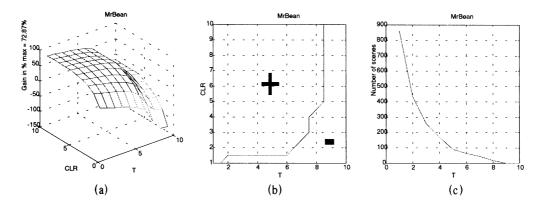


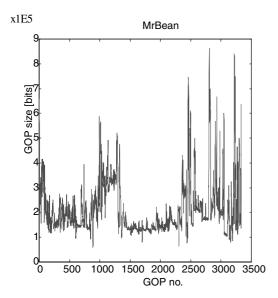Fig. 14. Simulation results for MrBean trace.

Fig. 15. MrBean GOP sizes.

step 10%) and different values of CLR (from $10^{-10}$ to $10^{-1}$ multiplying by 10 in each step). We always have a better capacity use than the static approach while guaranteeing the same and even better CLR for the region (CLR $< 10^{-3}$).

### 9.1. Investigating the statistical multiplexing gain

In this section, we investigate the statistical multiplexing gain (SMG) achievable using our SBA scheme. More specifically, we compare the SMG of SBA with that of another scenario (Fig. 16).

The first scenario (a) multiplexes $n$ streams without any restriction on a server with rate $C$ and buffer size $nB$. This gives the maximum achievable SMG for the given sources.

The second scenario (b) represents the SBA approach. The total service rate is $C$ in the two scenarios.

In this section we present simulations results comparing the performance in the two scenarios of Fig. 16. The stream we have used is the MPEG encoded trace of the Bond movie. The $n$ sources are randomly shifted versions of this trace. The dynamic allocation scheme used in the experiments is the on-line bandwidth allocation algorithm described in Section 8.2 with $T = 50\%$.

To assess the SMG for the two scenarios, we have determined the channel service rate per stream $C/n$, as a function

of $n$, needed to guarantee a desired bit loss probability. In scenario (a), bits are lost due to buffer overflow. In scenario (b), bits are lost due to failure in renegotiating for a higher rate (in which case we assume the source has to temporarily settle for whatever bandwidth remaining in the link until more bandwidth becomes available). Determining $C$ is straightforward for scenario (a). For scenario (b), we find for each $n$ the minimum $C$ that guarantees the desired loss probability: for each $n$, we do a binary search on $C$; for each step in the search, we do many simulations, where each simulation has a randomised phasing of the sources, and we compute the average fraction of bits lost as an estimate of the loss probability. At each step, we repeat the simulations until the sample standard deviation of the estimate is less than 20% of the estimate. Results for $10^{-5}$ loss probability requirement are depicted in Fig. 17.

Our scheme achieves slightly less SMG than the unrestricted case. Nevertheless, our scheme is able to extract most of the SMG, especially for a large number of multiplexed streams. When the number of sources is above 30, the normalised channel service rate per stream needed to guarantee $10^{-5}$ bit loss probability is less than 1.2, which means that for the two schemes each source requires less than 1.2 times its average rate.

To summarise, our solution has the following advantages:

- It allocates much less capacity than the constant allocation scheme while guarantying the required QoS.
- For a required CLR below $10^{-5}$ we notice to have a much lower CLR. This will lead to a good visual quality.
- For 50% in capacity gain, we have only around 100 to 200 messages overhead.
- With respect to video transmission on WATM, in average, we have an overhead in the order of 1 cell for each 1300 cell. We suppose that the capacity requests can be transmitted using OAM cells.
- With this method, we have a distributed CLR over all the film scenes while with the constant allocation method the CLR is not distributed. This means that big GOPs are not disadvantaged by comparison with small ones. This will improve the visual quality of the film and hence solves the problem stated in Section 2.
- It allows a higher number of users per cell since it uses less bandwidth.
- For scenes with a low level of activity (with GOP sizes

Table 2
Some results for other films

| Film | Gain for $T = 10\%$ (%) | Number of scenes for $T = 10\%$ | Gain for $T = 50\%$ (%) | Number of scenes for $T = 50\%$ |
|---|---|---|---|---|
| MTV1 | 74.61 | 1962 | 61.56 | 407 |
| Asterix | 72.59 | 1745 | 58.86 | 302 |
| Dino | 72.44 | 1361 | 58.33 | 191 |
| News1 | 71.09 | 667 | 56.29 | 76 |
| Simpsons | 70.65 | 1799 | 54.11 | 302 |
| Race | 69.49 | 1411 | 53.84 | 211 |
| Bond | 66.33 | 1154 | 47.95 | 120 |

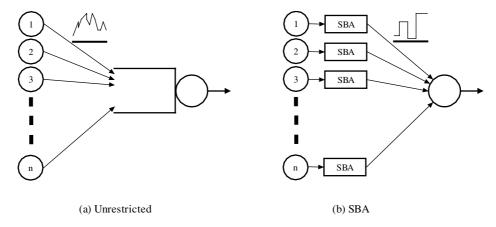(a) Unrestricted                                    (b) SBA

Fig. 16. The scenarios to assess statistical multiplexing gain of our scheme.

lower than the peak cell rate). The leftover bandwidth can be used by other users. We believe that this will reduce the network call blocking and forced termination probabilities.

- It requires no complex computations.
- It can be easily added to the base stations and mobile terminals.
- It is easy to implement.

These properties make our dynamic algorithm well suited for practical application. However, in our scheme we have assumed that whenever a user requests an amount of bandwidth, the request is accommodated by the base station. When the base station is overloaded, fairness of bandwidth assignment among the terminals becomes an important issue. Therefore, there must be some sort of strategies to

warrant the fairness among the multiple users in this situation. When distributing the available bandwidth among the users to achieve fairness, the base station can use information on user's bandwidth requirements, call holding times, call dropping and packet loss probabilities. For example, in case of system overload, the base station can distribute the loss evenly over the users according to their respective bandwidth requirements and QoS tolerance. More sophisticated fairness schemes can be used and this can be a subject of future research. For example the Generalised Processor Sharing scheme [21,22] or one of its new adaptations [23–25] can be used to allocate bandwidth fairly among all users.

## 10. Conclusion

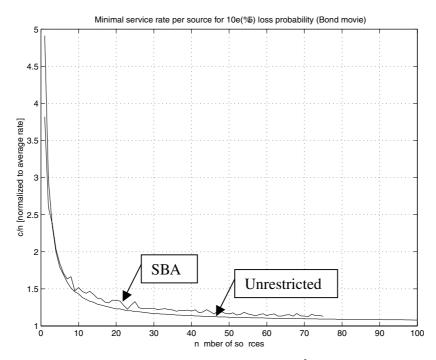In this paper, we analysed the statistical data sets of



Fig. 17. Statistical multiplexing gain achievable for $10^{-5}$ loss probability.

several MPEG encoded videos and proposed a model of the elements of the stream scenes. The proposed model permits the characterisation of the elements of the stream scenes, which can be used to allocate bandwidth dynamically for each scene and will result in a high capacity gain while guaranteeing the same and even better QoS.

In this paper we also proposed a dynamic bandwidth allocation scheme that can significantly improve bandwidth utilisation in wireless networks. It automatically adjusts the amount of reserved resources, while guarantying the required QoS.

The proposed algorithm exploits the structure of the MPEG video stream and allocates bandwidth on a scene basis. This will result in a high bandwidth gain, which will affect the overall network performance. The proposed scheme is dynamic and pro-active. It requires some communication between the mobile terminal and the base station, but the amount of extra information generated by the mechanism is acceptable in comparison to the capacity gain obtained. It is worth noting that the proposed scheme can be applied to any type of video coding as long as the scene is defined in the same notion. In this case a statistical study of the new encoded streams is needed to derive a bandwidth allocation algorithm similar to the one we proposed here.

Taking $\mu^*(s) = \text{First\_GOP}(s) + \sigma_{\text{GOP}}$ as an approximation of $\mu$ for each scene $s$ is a good approximation of the mean GOP size within the scene $s$. Indeed, this results in a lower CLR. However, this approximation is not the best one because it usually overestimates the mean $\mu$. Future work will involve studying a better estimation of $\mu$, which, while guaranteeing the CLR requirement, will further increase the capacity gain by allocating less bandwidth. Such study is particularly important for MPEG streams for which we do not have traces beforehand. For MPEG streams, for which we have already traces, $\mu$ is accurately computed based on those traces.

Future work will also involve studying the impact of the delay on the performance of the proposed algorithm as well as aspect related to call admission control. Studying the choice of the $T$ parameter is also an important research avenue as well as fairness among the different users when the system is overloaded.

## References

[1] D.A. Levine, I.F. Akyildiz, M. Naghshineh, The shadow cluster concept for the resource allocation and call admission in ATM-based wireless networks, IEEE/ACM Transaction on Networking 5 (1) (1997).

[2] C. Huang, M. Devetsikiotis, I. Lambadaris, A. Kaye, Modelling and simulation of self-similar variable bit rate compressed video: a unified approach, Proceedings of SIGCOMM'95, 1995.

[3] B. Jabbari, F. Yegengolu, Y. Kuo, S. Zafar, Y.-Q. Zhang, Statistical characterisation and block-based modelling of motion-adaptive coded video, IEEE Transactions on Circuits and Systems for Video Technology 3 (3) (1993) 199–207.

[4] P. Pancha, M. El Zakri, MPEG coding for variable bit rate video transmission, IEEE Communications Magazine 32 (5) (1994) 54–66.

[5] D.P. Heyman, T.V. Lakshman, Source models for VBR broadcast-video traffic, IEEE/ACM Transactions on Networking 4 (1) (1996).

[6] O. Rose, Statistical properties of MPEG video traffic and their impact on traffic modelling in ATM systems, 20th Conference on Local Computer Networks (LCN'95), Minneapolis, MN, October 1995.

[7] A.A. Lazar, G. Pacifici, D.E. Pendarakis, Modelling video sources for real-time scheduling, Technical Report 324-93-03, Columbia University, Department of Electrical Engineering and Center for Telecommunications Research, April 1993.

[8] Marwan Krunz, Satish K. Tripathi, Modelling bit rate variations in MPEG sources, University of Maryland, Institute for Advanced Computer Studies, Department of Computer Science, University of Maryland, December 1995.

[9] ISO/IEC International Standard 11172-2, Coding of Moving Pictures and Associated Audio for Digital Storage Media up to 1.5 Mbits/s Part2, Video, 1993.

[10] D. Le Gall, MPEG: a video compression standard for multimedia applications, Communications of the ACM 34 (4) (1991) 46–58.

[11] K.L. Gong, Berkeley MPEG-1 video encoder, user's guide, Technical Report, Computer Science Division-EECS, 1994.

[12] M. Hamdi, W. Roberts, Rate control for VBR video coders in broadband networks, IEEE Journal on Selected Areas in Communications August (1997).

[13] Ayman A.M. Ibrahim, Statistical rate control for efficient admission control of MPEG-2 VBR video sources, ATM'98, Fairfax, USA, May 1998.

[14] Ayman A.M. Ibrahim, M. Hamdi, Statistical rate control for real-time video multiplexing, Proceedings of the Eighth International Workshop on Packet Video (AVSPN'97), Aberdeen, September 1997.

[15] A.R. Reibman, B.G. Haskell, Constraints on variable bit rate video for ATM networks, IEEE Transactions on Circuits and systems for Video Technology 2 (4) (1992) 361–372.

[16] M.R. Pickering, J.F. Arnold, A perceptually efficient VBR rate control algorithm, IEEE Transactions on Image Processing 3 (5) (1996) 527–532.

[17] http://www.etsi.org/smg/utra/utra.htm.

[18] ITU-T I.610, B-ISDN Operation and Maintenance Principles Functions, Geneva, March 1993.

[19] M. Progler, C. Evci, M. Umehira, Air interface access schemes for broadband mobile systems, IEEE Communications Magazine September (1999).

[20] Y. Iraqi, R. Boutaba, A novel distributed call admission control for wireless mobile multimedia networks, Third ACM International Workshop on Wireless Mobile Multimedia (WoWMoM-2000), Boston, 11 August 2000.

[21] A.K. Parekh, R.G. Gallager, A. generalised, processor sharing approach to flow control in integrated services networks: the single-node case, IEEE/ACM Transactions on Networking 1 (1993) 344–357.

[22] A.K. Parekh, R.G. Gallager, A. generalised, processor sharing approach to flow control in integrated services networks: the multiple node case, IEEE/ACM Transactions on Networking 2 (1994) 137–150.

[23] R. Szabo, P. Barta, F. Nemeth, J. Biro, C.-G. Perntz, Call admission control in generalised processor sharing schedulers using non-rate proportional weighting of sessions, INFOCOM 2000 4 (2000) 1243–1252.

[24] A. Elwalid, D. Mitra, Design of generalised processor sharing schedulers with statistically multiplex heterogeneous QoS classes, INFOCOM'99 3 (1999) 1220–1230.

[25] K. Kumaran, G.E. Margave, D. Mitra, K.R. Stanley, Novel techniques for the design and control of generalised processor sharing schedulers for multiple QoS classes, INFOCOM 2000 4 (2000) 932–941.