# Region-Based Rate Control and Bit Allocation for Wireless Video Transmission

Yu Sun, *Member, IEEE*, Ishfaq Ahmad, *Senior Member, IEEE*, Dongdong Li, *Student Member, IEEE*, and Ya-Qin Zhang, *Fellow, IEEE*

*Abstract*—In this paper, we propose a joint source-channel region-based rate control algorithm for real-time video transmissions over wireless systems. During the video transmission, the channel throughput available to the video encoder in the wireless systems is inherently variable, due to the retransmission of the error packets using the automatic repeat request (ARQ) error control. The variable data rate of the wireless system is characterized by the packet-level Gilbert two-state Markov Model, the parameters of which are extracted from the statistical properties of the channel information obtained from the wireless channel simulator. The proposed algorithm adopts a fast but effective block-based segmentation method to extract the regions of interest. Unlike traditional bit allocation methods used in the region/content-based rate control, the algorithm exploits the most effective criteria "coding qualities" as quantitative factors to directly control bit allocation among different regions so as to achieve better visual quality in the regions of interest. The computational complexity of the algorithm is low making it suitable for real-time applications. Compared with the MPEG-4 rate control algorithm, our algorithm can effectively enhance the perceptual quality for the regions of interest and significantly reduce the number of frame skipping; thereby, improve the smoothness of the video.

*Index Terms*—Channel model, moving region detection, region-based rate control, wireless systems.

## I. INTRODUCTION

**W**ITH the increasing bandwidth in the wireless systems and rapidly growing demand for video communications, wireless video transmission has received much attention during the last few years [1]. Due to the limited bandwidth of the wireless channels, video signals need to be highly compressed by efficient video coding standards [1] such as H.263 [2] and MPEG-4 [3]. However, real-time video transmission is very sensitive to burst errors caused by the time varying signal strength received from the wireless channels. Even one bit error might cause severe degradation in video quality. Therefore, it is obligatory for the video encoder to protect the video data

from the channel errors by using error control techniques such as forward error correction (FEC) and the automatic repeat request (ARQ) [1].

ARQ schemes have been used as an efficient mechanism to control packet errors in the real-time video transmission in the presence of a feedback channel [4]. Using ARQ schemes, the channel throughput depends on the channel condition. When the channel condition is good, the channel bandwidth is fully exploited for transmitting the signal. When the channel condition is bad, a lot of retransmission occurs and the channel throughput goes down [5]. From the video transmission point-of-view, the channel becomes a variable bit-rate channel with the throughputs depending on the channel conditions [5]. This requires rate control (RC) schemes to dynamically adjust the output bitrate of the video encoder to meet to the variable channel throughput. However, RC schemes recommended by current compression standards, such as MPEG-4 and H.263, are optimized for constant bit-rate channels, and they cannot adapt themselves in time to the variation of the channel bandwidth. Rate control for wireless video transmission is a challenging task due to the limited channel throughput and time-varying characteristics of wireless channels. Further, a RC algorithm has to jointly decide encoding parameters and estimate current channel conditions in order to optimize its encoding performance.

In very low bit-rate video coding, the quality of encoded frames always suffers from serious degradation due to the limited channel throughput. In order to improve the coding efficiency and keep good subjective qualities of frames, region-based coding is usually exploited to code the regions of interest (ROI) more accurately than the rest of the video content. This would reduce the amount of information to be transmitted, while keeping good subjective quality of images. The object-oriented video compression standard MPEG-4 has been established based on this idea. But object segmentation is generally difficult to apply for real-time video coding and transmission due to its high computational complexity [6]. In addition, object-based coding needs to send the shape information along with the regular video data. Thus, the amount of coding bits for the shape information has to be small at the low bit-rate encoding, but that may cause the object to become coarse. As a result, the number of objects must be limited [7].

Some research works on the region/content-based rate-control have been reported [6], [8]–[10]. The Lagrange multiplier method is employed for rate control in region-based coding [6], although the complexity of Lagrange multiplier was reduced largely, it is still a major concern in real time video applications over wireless. Research works in [8]–[10] adopted a heuristic

Y. Sun is with the Department of Computer Science, University of Central Arkansas, Conway, AR 72035 USA (e-mail: yusun@uca.edu).

I. Ahmad is with the Department of Computer Science & Engineering, University of Texas at Arlington, Arlington, TX 76019 USA (e-mail: iahmad@cse.uta.edu).

D. Li is with the Department of Electrical Engineering, University of Texas at Arlington, Arlington, TX 76019 USA (e-mail: ldd@ieee.org).

Y.-Q. Zhang is with the Mobile and Embedded Devices Division, Microsoft Corporation, Redmond, WA 98052 USA (e-mail: yzhang@microsoft.com).
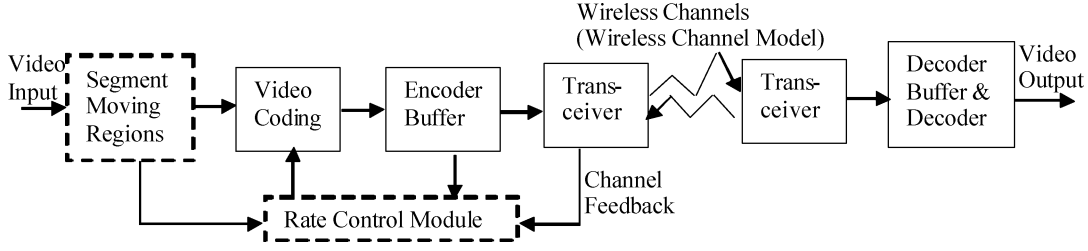
Fig. 1.   Diagram of the proposed region-based RC algorithm over wireless channels.

approach to decide the quantization parameters (QPs) for different regions in a frame, for example: an initial QP for a frame or macroblock (MB) is chosen first. Then the initial QP is directly decreased by a factor to obtain a finer quantizer for ROI, which results in more bits used in coding ROI; or it is increased by a factor to acquire a coarser quantizer for non-ROI and thus fewer bits are used in coding non-ROI. These factors are heuristically set to constants, and the contents of regions are not taken into consideration. The direction of the QP adjustments is correct, but these algorithms lack a quantitative method to perform bit allocation among different regions, this may cause improper QPs, unreasonable bits used for different regions, and the sub-optimal perceptual quality, which is crucial in the low-bit rate video coding.

In this paper, we propose a region-based rate control algorithm for transmitting video in real-time over wireless channels. The Rate Control Module, represented by the dashed line in Fig. 1, functions by interacting with other modules of the encoder. We adopt a wireless channel model to characterize wireless channels that are inherently time-varying. A fast, yet effective, segmentation method is used to detect ROI in real time based on digital image processing techniques. Region shapes are not necessary to be transmitted. The novelty of the proposed algorithm is that it adopts the most effective criteria "coding qualities" as quantitative factors to directly control bit allocations among different regions, and employs the priority concept to further adjust bit allocation under different channel conditions, aiming at improving the visual quality for ROI and reducing frame skipping to keep the motion smoothness under time-varying channels.

The rest of this paper is organized as follows: Section II presents the segmentation method. Section III describes how to establish the channel model in wireless systems to estimate the variable channel-bandwidth. Section IV describes the proposed region-based RC algorithm. Section V includes the simulation results demonstrating the performance of the algorithm. Section VI concludes the paper with final observations.

## II. FAST SEGMENTATION OF MOVING FOREGROUND AND STATIC BACKGROUND

Since the human visual system (HVS) is more sensitive to the moving regions [6], it is reasonable to sacrifice the perceptual quality of still regions while enhancing that of moving regions. For instance, in the head-and-shoulder types of video sequences, people tends to focus on the face and give more emphasis on important facial features, such as the mouth and eyes that are usually most intensively observed. Therefore, it is worthwhile

to allocate more bits for coding these regions of the scene more accurately at the expense of coarser coding of less important regions, which is the underlying technique exploited in the region-based coding [11]. In region-based coding, image segmentation must be employed to identify the locations and shapes of ROI within the video scene. While a large number of research works on how to efficiently segment ROI, including color segmentation, texture segmentation, and motion segmentation, the computational complexity of segmentation approaches must be significantly reduced in order to meet the requirement of low delay in real-time video applications.

In this paper, using the MB as a fundamental unit, we adopt a fast method to detect moving regions from the coding frame in real time, the moving regions are classified as the foreground (ROI) while the still regions are regarded as the background (non-ROI).

In order to detect changes between two successive frames, each frame is smoothed by a simple $3 \times 3$ low-pass filter to reduce the high frequency noise [6], and then the difference between the current frame and previous frame is computed. We can calculate the difference mask $DM_t$ [5] as

$$DM_t(i,j) = \begin{cases} 1, & |I_t(i,j) - I_{t-1}(i,j)| > Thr_t \\ 0, & else \end{cases} \quad (1)$$

where $(i,j)$ represents the coordinate of the processed pixel, $I_t$ and $I_{t-1}$ represent the current and previous frames respectively, and $Thr_t$ is a threshold for $DM_t$ given by

$$Thr_t = \frac{1}{M \cdot N} \sum_{i=1}^{M} \sum_{j=1}^{N} |I_t(i,j) - I_{t-1}(i,j)| \quad (2)$$

where $M$ and $N$ denote the numbers of the row and column in a frame respectively. All pixels in $DM_t(i,j)$ with value 1 are regarded as the moving pixels, and the moving ratio $r_t(B_m)$ of the $m^{th}$ macroblock $B_m$ is defined as

$$r_t(B_m) = \sum_{(i,j) \in B_m} \frac{DM_t(i,j)}{L} \quad (3)$$

where $L$ is the pixel number in $B_m$. Then, the threshold $MB\_Thr_t$, determining if a $MB$ is a moving or still $MB$, is set to

$$MB\_Thr_t = k \cdot \left( \sum_{m=1}^{N_{MB}} \frac{r_t(B_m)}{N_{MB}} \right). \quad (4)$$

If $r_t(B_m)$ is larger than $MB\_Thr_t$, $B_m$ is selected as the moving $MB$; otherwise, $B_m$ is the still $MB$. $N_{MB}$ in (4) is the number of MBs in a frame, $k$ is a constant factor which

decides the number of moving MBs in a frame, a larger $k$ results in a smaller number of moving MBs, while a smaller $k$ gives rise to a larger number of moving MBs. In the simulation, $k$ is chosen empirically to be 1.4. In practice, the value of $k$ could be increased or decreased if we wanted a larger or smaller number of moving MBs in a frame.

Finally, we sort the moving MBs according to their moving ratios in a decreasing order. If the moving ratio $r_t(B_m)$ of an isolated moving MB, whose neighboring MBs are all still MBs, is in the last 40% of the sorted list, it is merged to the still region since it is not a heavily moving MB among all of the moving MBs. Similarly, an isolated still MB is also assigned to be a moving MB when its moving ratio is in the first 40% of the sorted list for still MBs. By using this segmentation method, we are able to keep tracking of moving regions temporally; this may improve the flexibility of the tracking ability since these moving regions might not be restricted to one specific object such as a speaker in MPEG-4 concept.

## III. WIRELESS CHANNEL MODEL

Wireless communication channels suffer from fading, which causes the received signal envelope to drop below a certain threshold for a period of time [12]–[17]. When the envelope of the received signal is below a certain value, the possibility that bit errors occur increases. In the fade duration, the bit errors embody property of burstness, which was initially modeled by the hidden Markov model in [18]–[20] and is called the Gilbert two-state Markov model. In the Gilbert model, the channel is characterized by two states, corresponding to the good (bit error free) and bad (bit error) states, and the states transfer between each other at a certain probability.

The Gilbert model on the bit level was extended to characterize the wireless channel at the packet level in [5], [21], [22]. Again, two states are taken to represent the good and bad channel states, but in terms of the packet error. Since the bits are grouped into packets, the bit-error burstness is alleviated in terms of packet error and the Gilbert model is more accurate in the predication of the packet error rate than the bit error rate. By denoting $S_0$ as the packet error-free state and $S_1$ as the packet error state, the two-state Gilbert model is shown in Fig. 2. The channel state transition probability matrix is given by $P$ with the element $P_{ij} = P[S(k) = S_j|S(k-1) = S_i], i, j \in \{0, 1\}$, which is the transfer probability from the state $S_i$ at time $t_{k-1}$ (i.e. $S(k-1)$) to the state $S_j$ at time $t_k$ (i.e., $S(k)$) [23]–[25]. Without loss of generality, assuming that the channel state $S(k-m)$ at time $t_{k-m}$ is known, the probability that the channel in state $S_j$ at time $t_k$ is denoted as $\pi_j(k|S(k-m))$. Therefore, the state probability vector at time $t_k$ is written as $\pi(k|S(k-m)) = [\pi_0(k|S(k-m)), \pi_1(k|S(k-m))]$. Since the channel state at time $t_{k-m}$ is known to be $S(k-m)$, the state probability vector at time $t_{k-m}$ is initialized by

$$\pi_j[k-m|S(k-m)] = \begin{cases} 1, & S(k-m) = S_j \\ 0, & \text{otherwise} \end{cases} \quad j \in \{0, 1\}. \tag{5}$$

Thus, the state probability vector at time $t_k$ can be derived from the state probability vector at time $t_{k-m}$ by using

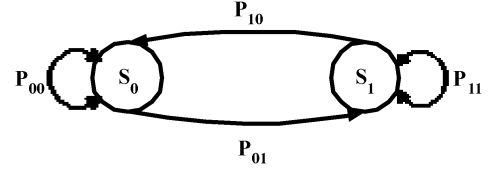$$\pi(k|S(k-m)) = \pi(k-m|S(k-m))P^m. \tag{6}$$



Fig. 2. Two-state Markov channel model.

Notice that the probability of packet error-free at time $t_k$ is $\pi_0(k|S(k-m))$.

## IV. REGION-BASED RATE CONTROL FOR WIRELESS VIDEO TRANSPORT

Designed primarily for real-time video transmission over wireless channels, the proposed algorithm is able to improve the perceptual quality of the foreground under the limited channel bandwidth. This section describes the principles and foundations of the proposed algorithm.

### A. Initial Target Bit Estimation

In real time video communications, it is impossible to calculate remaining bits and remaining frames, since the frame number of the total sequence is not known at the moment of the encoding. Thus, we set the initial target bits for a frame to the average bits available per frame by

$$T_t = \frac{R}{F_r} \tag{7}$$

where $R$ and $F_r$ are the channel and frame rates, respectively.

### B. Adjust Target Bits Based on the Buffer Fullness

To get more accurate target bit estimation, the initial bit target is further refined based on the buffer fullness. To keep the low end-to-end delay, we adopt a small buffer and set the default buffer size $B_s$ to $0.125R$, the target buffer fullness is the middle level of the buffer size. Here, we adopt our proposed Proportional-Integral-Differential (PID) buffer control technique [26]. The error signal $E_t$ at time $t$, which represents the deviation between the target buffer fullness $(B_s/2)$ and the current buffer fullness $B_{f,t}$, is defined as

$$E_t = \frac{B_s}{2} - B_{f,t}. \tag{8}$$

This error signal is sent to the PID controller as

$$PID_t = K_p \cdot \left( E_t + K_i \cdot \int_0^t E_\tau \cdot d\tau + K_d \cdot \frac{dE_t}{dt} \right) \tag{9}$$

where $K_p$, $K_i$ and $K_d$ are the Proportional, Integral and Differential control parameters, respectively, and are set empirically to 0.1, 0.25, and 0.3 respectively in the simulations. Then the total target bits $T_t$ can be further adjusted by

$$T_t := T_t + T_t \cdot PID_t. \tag{10}$$

To obtain a minimum visual quality for each frame, the lower bound of the target bits imposed to each frame is $R/4F_r$. And

to avoid buffer overflow, the maximum number of bits is given by $B_s - B_{f,t} + R/F_r$.

### C. Compute the Number of Retransmitted Bits per Frame-Interval

Based on the Gilbert channel model in Section III, the number of retransmitted bits in the presence of ARQ is derived in this section. In order to compensate for the retransmission, the ARQ retransmission bits, which will be deducted from the frame target bits, are estimated at a frame basis as follows: First, we assume the length of the encoded frame $r$ to be $L_r$, then $N_r = L_r/L_p$ packets have been transmitted in a packet size of $L_p$ bits. Second, the average error packet ratio of the previous $L_{rp}$ frames is calculated by

$$r_{avg}(L_{rp}) = \frac{\sum\limits_{r=1}^{L_{rp}} \sum\limits_{i=1}^{N_r} P_e(r,i)}{\sum\limits_{r=1}^{L_{rp}} N_r}$$

where $P_e(r,i)$ equals to 1 when the $i^{th}$ packet in the $r^{th}$ frame is in error, otherwise 0. If $r_{avg}(L_{rp})$ is less than a threshold value $W_{th}$, the current state of the wireless channel is said to be good. Otherwise, the channel is said to be in the bad state. Third, using the current channel state denoted as $S(0)$ and (6), we can estimate the channel state for the next $L_{rf}$ frames. The average number of packets per frame-interval to be transmitted can be calculated by $L_c = (R/F_r)/L_p$. Therefore, we need to generate channel states for the next $N_{rf} = L_{rf} L_c$ packets as follows: To obtain a smooth estimation of the channel state, the average probability of an error-free transmission of the $m^{th}$ packet given the current channel state $S(0)$ is derived based on (6) as

$$p_{avg}(m|S(0)) = \frac{1}{m} \sum_{i=1}^{m} \pi_0(i|S(0)). \tag{11}$$

Following a commonly accepted statistical approach [25], we generate a uniform distributed random variable $\alpha$ in the interval [0, 1] and compared it with the $P_{avg}(m|S(0))$ to decide the channel state of the $m^{th}$ packet. This approach can more dynamically characterize the channel state than that used in [5]. If $\alpha$ is greater than $P_{avg}(m|S(0))$, the channel of the $m^{th}$ packet is said to be bad and $P_e(m) = 1$; otherwise, $P_e(m) = 0$. Therefore, the total retransmitted bits $RTB_t$ in the current frame can be calculated by

$$RTB_t = \frac{1}{L_{rf}} \sum_{m=1}^{N_{rf}} P_e(m) \times L_p. \tag{12}$$

### D. Adjust Target Bits Based on the Number of Retransmitted Bits

During the retransmissions of error packets when the channel is in deep fades, the video data in the encoder buffer are not transmitted. Due to the reduced channel throughput, the encoder buffer fills up quickly which may cause the rate control algorithm to skip frames or significantly reduce the bits allocated to

each frame [5]. In order to prevent the buffer fill-up in the future, the retransmission bits should be taken into account when allocate the target bits to a frame [5]. Thus, the frame target is deducted by the retransmission bits as

$$T_t := T_t - RTB_t. \tag{13}$$

### E. Target Bits Distribution Between the Foreground and Background

In order to obtain better perceptual quality in ROI, a rate control scheme should allocate more bits to the foreground than the background, especially when the channel is in deep fades. To achieve this goal, the algorithm sets weights for the foreground and background to control bit allocation between them. The larger the weight is, the more target bits should be allocated to the corresponding region. To quantitatively and directly control bit allocation, the most effective criteria, foreground PSNR ($PSNR_{F,t}$) and background PSNR ($PSNR_{B,t}$), are employed in the weight adjustment. They are defined as

$$PSNR_{F,t} = 10 \cdot \log_{10} \frac{255^2}{MSE_{F,t}},$$

$$MSE_{F,t} = \frac{1}{N_{F,t}} \sum_{(i,j) \in F} (I_t(i,j) - I'_t(i,j)),$$

$$PSNR_{B,t} = 10 \cdot \log_{10} \frac{255^2}{MSE_{B,t}},$$

$$MSE_{B,t} = \frac{1}{N_{B,t}} \sum_{(i,j) \in B} (I_t(i,j) - I'_t(i,j)) \tag{14}$$

where $N_{F,t}$ and $N_{B,t}$ are the numbers of pixels belonging to the foreground and background respectively at time $t$, $I_t(i,j)$ and $I'_t(i,j)$ denote the original intensity and the reconstructed intensity of the pixel $(i,j)$, $MSE_{F,t}$ and $MSE_{B,t}$ are the mean square errors for the foreground and background, respectively.

In addition, since a simple and efficient way to express application requirements is to specify priorities between the foreground and background, priority is also employed in the weight adjustment to assist target bit allocation.

We adopt the background as a referential base, its weight $W_{B,t}$ is always 1.0. Let $W_{F,t}$ be the weight for the foreground at time $t$, its initial value is 1.0. $U_F$ is the priority of the foreground, $U_F > 0$ (dB) means a higher priority while $U_F < 0$ (dB) corresponds to a lower priority, it is specified based on the application requirements in practice. $(PSNR_{F,t-1} - U_F)$ for the foreground ($F$) at time $t-1$ is compared with the $PSNR_{B,t-1}$ for the background ($B$), if $(PSNR_{F,t-1} - U_F)$ is lower than $PSNR_{B,t-1}$, the algorithm improves the weight of the foreground, thus the foreground obtains more target bits and thus achieves a higher quality; otherwise, decreases $W_{F,t}$ to achieve a lower quality. The weight for the foreground is updated as follows:

$$W_{F,t} = W_{F,t-1} \cdot e^{\left(\frac{PSNR_{B,t-1} - PSNR_{F,t-1} + U_F}{\theta}\right)} \tag{15}$$

here, the tuning factor $\theta$ is selected to 4 empirically. To avoid heavily unbalanced bits allocation between the foreground and background, $W_{F,t}$ is bounded to a range from 1 to 10.

Then the normalized weight for the foreground $(NW_{F,t})$ and background $(NW_{B,t})$ can be obtained by

$$NW_{F,t} = \frac{W_{F,t}}{(W_{F,t} + W_{B,t})}, \quad NW_{B,t} = \frac{W_{B,t}}{(W_{F,t} + W_{B,t})}. \tag{16}$$

The coding complexities and perceptual importance of the foreground and background must be considered during bit allocation between them. We have chosen the normalized weight, size and variance as three factors in the target bit distribution. The variances of the foreground and background for the current frame, $VAR_{F,t}$ and $VAR_{B,t}$, are defined as

$$VAR_{F,t} = \frac{1}{N_{F,t}} \sum_{(i,j)\in F} \left(P_t(i,j) - \overline{P}_{F,t}\right)^2,$$

$$VAR_{B,t} = \frac{1}{N_{B,t}} \sum_{(i,j)\in B} \left(P_t(i,j) - \overline{P}_{B,t}\right)^2 \tag{17}$$

where $P_t(i,j)$ is the luminance value of the pixel $(i,j)$ in the motion-compensated residual frame, $\overline{P}_{F,t}$ and $\overline{P}_{B,t}$ are the arithmetic average pixel value of the foreground and background respectively. Therefore, as long as the target bits are given for a frame, the number of target bits for the foreground $(T_{F,t})$ and background $(T_{B,t})$ are allocated by

$$T_{F,t} = \frac{NW_{F,t} \cdot (NMB_{F,t} \cdot NVAR_{F,t})}{\sum\limits_{j=\{F,B\}} NW_{j,t} \cdot (NMB_{j,t} \cdot NVAR_{j,t})} \cdot T_t,$$

$$T_{B,t} = \frac{NW_{B,t} \cdot (NMB_{B,t} \cdot NVAR_{B,t})}{\sum\limits_{j=\{F,B\}} NW_{j,t} \cdot (NMB_{j,t} \cdot NVAR_{j,t})} \cdot T_t \tag{18}$$

where $NMB_{F,t}$ and $NMB_{B,t}$ are the number of MBs in the foreground and background respectively, normalized by the total number of MBs in a frame, $NVAR_{F,t}$ and $NVAR_{B,t}$ are the normalized variances of the foreground and background correspondingly, and can be obtained by $NVAR_{F,t} = VAR_{F,t}/(VAR_{F,t} + VAR_{B,t})$, $NVAR_{B,t} = VAR_{B,t}/(VAR_{F,t} + VAR_{B,t})$.

### F. Macroblock-Level Rate Control

Once the target bit budgets for the foreground and background are obtained, the target bits for coding the $i^{th}$ macroblock $(MB_i)$, $T_{MB,i}$, can be allocated by

`If` $(MB_i \in$ Foreground$)$

$$T_{\text{MB,i}} = \left(\frac{VAR_{MB,i}}{\sum\limits_{k=i}^{NMB_{F,t}} VAR_{MB,k}}\right) \cdot T_{F,t}^i,$$

`else if` $(MB_i \in$ Background$)$

$$T_{\text{MB,i}} = \left(\frac{VAR_{MB,i}}{\sum\limits_{k=i}^{NMB_{B,t}} VAR_{MB,k}}\right) \cdot T_{B,t}^i$$

where $VAR_{MB,i}$ is the variance of the motion-compensated residual $MB_i$, $T_{F,t}^i$ and $T_{B,t}^i$ are the remaining available target bits for the foreground and background respectively when coding $MB_i$, initially $T_{F,t}^i = T_{F,t}$ and $T_{B,t}^i = T_{B,t}$. The marcoblock-layer rate control of MPEG-4 [27] is then used to compute the QP for $MB_i$ and encode $MB_i$.

### G. Update Buffer Fullness

The buffer fullness is updated after encoding, the number of actual bits used for encoding the current frame $(A_t)$ and the retransmission bits are added to the previous buffer level, $B_{f,t-1}$, at the same time, the number of bits to be output from the buffer per encoding time, $R/F_r$, is decreased from the buffer fullness [5]:

$$B_{f,t} = B_{f,t-1} + A_t + RTB_t - \frac{R}{F_r}. \tag{19}$$

The encoder needs to examine the current buffer level before encoding the next frame, if the buffer occupancy exceeds 80% of the buffer size $(80\% \cdot 0.125 \cdot R = 0.1 \cdot R)$, the encoder skips the next frame, and the buffer fullness is subtracted by $R/F_r$.

## V. SIMULATION RESULTS

### A. Simulation Environment

Our algorithm is generic to any wireless channels by changing the transitional probability in Markov model. Specifically, we take a wide-band code division multiple access (WCDMA) system as an example. To evaluate the performance of the proposed region-based rate control strategy, we compare it with the MPEG-4 Q2 scheme [27] under four sets of channel conditions. These four sets of channel conditions are as follows.

— Channel condition set 1 (CCS1): walking speed (3 km/ hour) and a user data rate 32 Kbps (Downlink);
— Channel condition set 2 (CCS2): car speed (40 km/hour) and a user data rate 32 Kbps (Downlink);
— Channel condition set 3 (CCS3): walking speed (3 km/ hour) and a user data rate 64 Kbps (Downlink);
— Channel condition set 4 (CCS4): walking speed (3 km/ hour) and a user data rate 64 Kbps (Uplink).

A WCDMA channel simulator is used to generate the bit-error pattern under the four sets of channel conditions, which is setup using the corresponding channel condition. For a WCDMA system with a user data rate of 64 Kbps, the data packet size is 640 bits (without counting the addition 16 CRC bits) in a 10 ms interleaving block and the data packet size is 320 bits for a user data rate of 32 Kbps for both links.

After obtaining the bit-error pattern from the WCDMA simulator, the packet-error pattern is generated by grouping the bits into packets. The average packet burst-error length statistics, the average packet error-free run length statistics and the packet error rate can be obtained from the packet-error pattern. By using the fact that the average packet burst error length equals to $1/P_{10}$, the average packet error-free run length equals to $1/P_{01}$, and the packet error rate equals to $P_{01}/(P_{01} + P_{10})$, the transition probabilities of the Gilbert two-state Markov model are

TABLE I
PERFORMANCE COMPARISON FOR MPEG-4 Q2 AND THE PROPOSED ALGORITHM UNDER THE CCS1

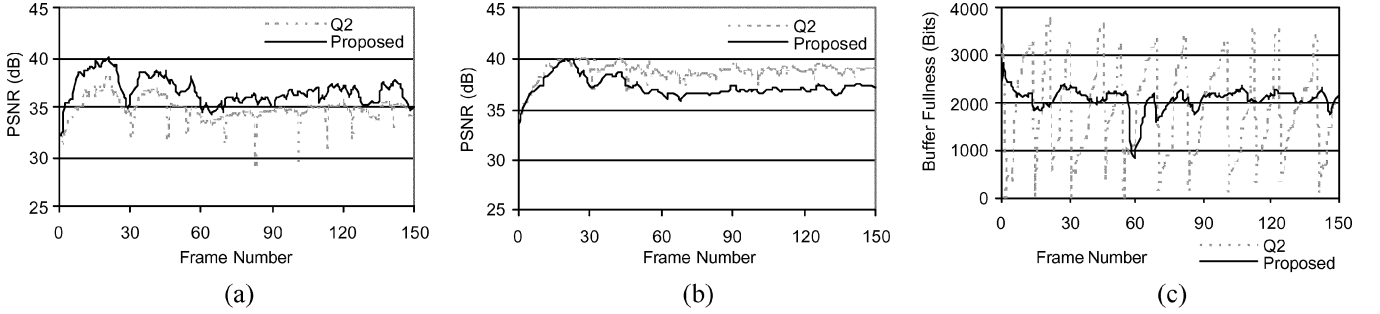| Video Sequence | Algorithm | Bit count ratio (%) | | # Skipped Frames | Average PSNR (dB) | | |
|---|---|---|---|---|---|---|---|
| | | Foreground | Background | | Foreground | Background | Overall |
| Mother_ Daughter | Q2 | 46% | 54% | 12 | 34.66 | 38.68 | 37.36 |
| | Proposed | 76% | 24% | 0 | 36.66 | 37.22 | 37.05 |
| Silent Voice | Q2 | 61% | 39% | 14 | 32.50 | 34.51 | 33.85 |
| | Proposed | 73% | 27% | 0 | 33.69 | 33.59 | 33.61 |
| Salesman | Q2 | 42% | 58% | 18 | 32.18 | 35.37 | 34.54 |
| | Proposed | 61% | 39% | 1 | 34.03 | 34.44 | 34.26 |



Fig. 3.    PSNR and buffer curves for the *Mother_Daughter* sequence under the WCDMA channel condition CCS1. (a) Foreground PSNR; (b) background PSNR; (c) buffer occupancy.

obtained from the simulation results of the WCDMA channel simulator as follows. For the CCS1, $P_{00} = 0.97538$, $P_{01} = 0.02462$, $P_{10} = 0.30367$, $P_{11} = 0.69633$, which corresponds to an average packet burst error length of 3.3 packets and a packet error rate of 0.075. For the CCS2, $P_{00} = 0.97928$, $P_{01} = 0.02072$, $P_{10} = 0.335592$, $P_{11} = 0.64409$, which corresponds to an average packet burst error length of 3.0 packets and a packet error rate of 0.058. For the CCS3, $P_{00} = 0.96024$, $P_{01} = 0.039759$, $P_{10} = 0.17154$, $P_{11} = 0.82846$, which corresponds to an average packet burst error length of 5.8 packets and a packet error rate of 0.19. For the CCS4, $P_{00} = 0.9566$, $P_{01} = 0.0434$, $P_{10} = 0.1538$, $P_{11} = 0.8462$, which corresponds to an average packet burst error length of 6.5 packets and a packet error rate of 0.22. All the above values obtained are similar to the results given in [24], Table III. It is worth noticing that when the packet size increases at a higher user data rate, the average packet burst error length and packet error rate also increase. This is because larger packet size is more susceptible to the sporadic bit errors. As mentioned earlier in Section III, with the increase of vehicle speed, the fade duration of the received signal is shorter; thereby, the average packet burst error length at a car speed is shorter than that at a walking speed. Furthermore, packet error threshold $W_{th}$ is empirically set to 0.2. Both the number of past $L_{rp}$ frames and the future $L_{rf}$ frames are empirically chosen to be equal to 2.

The coded video sequences are corrupted using the packet-error pattern files for these four sets of channel conditions. All test sequences are in QCIF format ($176 \times 144$ pixels/frame) and encoded at a target frame rate of 10 frames/s, the number of frames to be encoded is 150 frames. The first frame is intra-coded and the remaining frames are all intercoded ($P$ frames). To meet the requirement of the low-delay in video communications, a small buffer size $0.125 \cdot R$ is used in both the MPEG-4 Q2 and our algorithms. Frame skipping occurs when the number
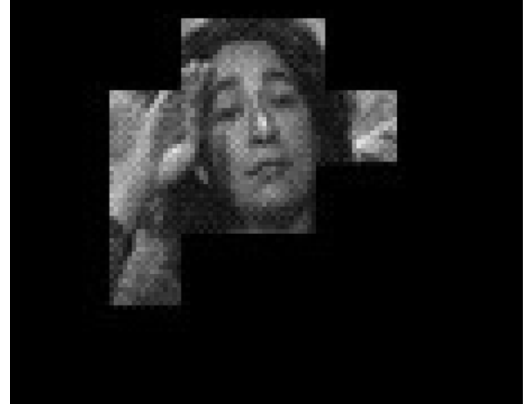


Fig. 4.    Segmented result of the 23rd frame of the "silent voice" sequence (original frame), 200% enlarged local area.

of bits in the buffer is above 80% of the buffer size $(0.1 \cdot R)$, the same buffer level for frame skipping as in TMN8 [28] when the frame rate is 10 frames/s. When the channel is in the good state, we dynamically set $U_f$ to 0 to obtain balanced coding qualities between the foreground and background; when the channel state is bad, in order to ensure the foreground's visual quality, we set $U_F$ to 2 to improve the foreground's priority. All the erroneous packets are assumed to be retransmitted successfully in a single attempt. Since a skipped frame is represented in the decoded sequence by repeating the previously coded frame according to MPEG-4 core experiments, the PSNR of a skipped frame is computed by using the previous encoded frame [29].

### B. Simulation Results

Table I shows the encoding performance by using MPEG-4 Q2 and the proposed algorithm under the channel condition CCS1. Because the bit-rates are decreaseed in the background,

Fig. 5.    Subjective results of the 23rd frame of the "Silent Voice" on the WCDMA channel condition CCS1 (32 Kbps), 200% enlarged local area. (a) MPEG-4 Q2 (reconstructed frame); (b) proposed (reconstructed frame).

TABLE  II
PERFORMANCE COMPARISON FOR MPEG-4 Q2 AND THE PROPOSED ALGORITHM UNDER THE CCS2

| Video Sequence | Algorithm | Bit count ratio (%) | | # Skipped Frames | Average PSNR (dB) | | |
|---|---|---|---|---|---|---|---|
| | | Foreground | Background | | Foreground | Background | Overall |
| Mother_Daughter | Q2 | 46% | 54% | 10 | 34.65 | 38.66 | 37.34 |
| | Proposed | 76% | 24% | 0 | 36.76 | 37.24 | 37.09 |
| Silent Voice | Q2 | 61% | 39% | 12 | 32.48 | 34.60 | 33.89 |
| | Proposed | 75% | 25% | 0 | 33.82 | 33.65 | 33.68 |
| Salesman | Q2 | 41% | 59% | 18 | 32.17 | 35.34 | 34.49 |
| | Proposed | 59% | 41% | 1 | 34.10 | 34.56 | 34.37 |

the overall average PSNR values of the proposed algorithm are a little lower than those of MPEG-4 Q2. To evaluate the picture quality more properly in the region-based RC algorithm, the foreground and background PSNRs are used here. One can see that the average foreground PSNRs of the proposed algorithm in Table I are much higher than those of MPEG-4 Q2, while the average background PSNRs of our algorithm are lower. Therefore, the proposed scheme can apparently enhance the visual quality in the foreground, at the expense of the quality degradation on the background. Meanwhile, our algorithm can reduce frame skipping effectively when compared with MPEG-4 Q2.

The bit count ratio of the foreground or background, which is defined as the individual used bit count normalized by the total used bit count of the whole sequence, is also shown in Table I. Using the proposed algorithm, the bit count ratio of the foreground increases a lot. This indicates our bit allocation method is effective. As a result, the visual quality of the foreground has been enhanced.

Fig. 3 shows the foreground PSNR, background PSNR and buffer curves for the *Mother_Daughter* sequence on the WCDMA channel condition CCS1. PSNR curves in Fig. 3(a) and (b) show that the foreground PSNRs of our algorithm are higher than those of MPEG-4 Q2, but our background PSNRs are lower than those of MPEG-4 Q2. They further exhibit that our algorithm obtains smoother foreground and background qualities among frames. In addition, our buffer fullness curve in Fig. 3(c) is more stable, as it is around the target buffer fullness ($B_s/2$) with a small fluctuation.

The segmented result in Fig. 4 shows our segmentation method is effective in detecting moving regions. The subjective

improvement achieved by our algorithm is depicted in Fig. 5, which shows the 200% enlarged local area of the 23rd frame of the "Silent Voice". This indicates our algorithm achieves its objective in enhancing the perceptual quality of the moving regions, the foreground including face and hand parts is clearer than that of MPEG-4 Q2. Although it may introduce some degradation on the background, it is almost invisible to human perception. Hence, the subjective quality is improved when compared with the overall objective quality (PSNR).

The WCDMA channel condition CCS2 is taken to investigate the performance of our algorithm for rapid movement of the vehicle. The results are shown in Table II. Since the average packet error-rate under this channel condition is 5.8%, close to 7.5% of the first channel condition CCS1, we obtain the similar results. These results show our algorithm is also stable under this kind of channel condition.

Traditionally, due to the limited bandwidth, video sequences transmitted under wireless channels normally have low motion, such as head and shoulder scenarios. Since higher bandwidths are available by the current WCDMA technique, such as 64 Kbps, this may support to transmit some sequences with medium or fast motion. Here, we also study the effect of transmitting some medium or fast motion sequences under the WCDMA 64 Kbps downlink and uplink wireless channels (CCS3 and CCS4). By using our algorithm, the quality of the foreground for "Stefan" has been improved about 2.09 dB (Table III), while the quality of its background has been degraded when compared with the MPEG-4 Q2. The results in Table III and Table IV demonstrate again that, for medium or fast motion sequences under higher bandwidth channels, the

TABLE III
PERFORMANCE COMPARISON FOR MPEG-4 Q2 AND THE PROPOSED ALGORITHM UNDER THE CCS3

| Video Sequence | Algorithm | Bit count ratio (%) | | # Skipped Frames | Average PSNR (dB) | | |
|---|---|---|---|---|---|---|---|
| | | Foreground | Background | | Foreground | Background | Overall |
| Stefan | Q2 | 57% | 43% | 29 | 23.22 | 28.60 | 25.70 |
| | Proposed | 84% | 16% | 0 | 25.66 | 26.92 | 26.33 |
| Foreman | Q2 | 43% | 57% | 26 | 31.76 | 34.30 | 33.25 |
| | Proposed | 66% | 34% | 0 | 34.30 | 33.48 | 33.65 |
| Mobile | Q2 | 28% | 72% | 29 | 23.14 | 25.07 | 24.60 |
| | Proposed | 52% | 48% | 0 | 25.05 | 24.27 | 24.41 |

TABLE IV
PERFORMANCE COMPARISON FOR MPEG-4 Q2 AND THE PROPOSED ALGORITHM UNDER THE CCS4

| Video Sequence | Algorithm | Bit count ratio (%) | | # Skipped Frames | Average PSNR (dB) | | |
|---|---|---|---|---|---|---|---|
| | | Foreground | Background | | Foreground | Background | Overall |
| Stefan | Q2 | 57% | 43% | 30 | 23.11 | 28.48 | 25.59 |
| | Proposed | 84% | 16% | 0 | 25.53 | 26.83 | 26.21 |
| Foreman | Q2 | 42% | 58% | 31 | 31.55 | 34.06 | 33.06 |
| | Proposed | 66% | 34% | 0 | 34.01 | 33.27 | 33.40 |
| Mobile | Q2 | 28% | 72% | 31 | 23.12 | 25.06 | 24.58 |
| | Proposed | 54% | 46% | 0 | 25.02 | 24.17 | 24.29 |

proposed algorithm effectively reduces the number of frame skipping and improves the perceived quality for the foreground. Furthermore, the overall PSNRs of some sequences of the proposed algorithm are a little higher than those of MPEG-4 Q2. This may be due to the large number of frame skipping in MPEG-4 Q2, which causes overall PSNR degradation.

The uplink WCDMA channel condition CCS4 is slightly worse than that of CCS3. Correspondingly, its simulation results (Table IV) are slightly worse when compared with those obtained for CCS3 (Table III). However, our proposed algorithm is still effective under CCS4, that is, it reduces the number of frame skipping and improves the perceptual quality of the foreground. Therefore, we can conclude that the proposed algorithm is effective for both the uplink and downlink of WCDMA systems and at various packet error rates.

### C. Complexity Analysis

Generally speaking, the rate control algorithm is regularly employed in video coding. Except an additional procedure to segment ROI (regions of interest) from a frame, the computational complexity of our proposed rate control algorithm is almost the same as that of MPEG-4s algorithm, which can be ignored when compared with that of the whole video coding process. The complexity of this procedure is: $O(n) + O(m \log m)$, where $n$ is the number of pixels in a frame ($176 \times 144$) and $m$ is the number of macroblocks in a frame (99).

To illustrate the real computing expenses, we have done some experiments to record the CPU time for both our segmentation method and our whole rate control algorithm. Table V shows the percentage CPU time spent on ROI segmentation and the whole rate control algorithm respectively.

From the Table V, we can see that the CPU time of segmentation is only 2.3% to 2.43% of the total encoding time, and the total expenses of our whole rate control algorithm including segmentation are from 2.6% to 2.84% of the entire computing

TABLE V
THE PERCENTAGE CPU TIME FOR SEGMENTATION AND THE WHOLE RATE CONTROL ALGORITHM UNDER THE CCS1

| Sequence (150 frames) | CPU Time Percentage (%) | |
|---|---|---|
| | (Segmentation Time) / (Total Encoding Time) | (Rate Control) / (Total Encoding Time) |
| Salesman | 2.3% | 2.60% |
| Silent Voice | 2.29% | 2.63% |
| Mother_Daughter | 2.43% | 2.84% |

time, which are very light burden and are quite affordable to the CPU computing power in wireless video applications.

### VI. CONCLUSIONS

We have proposed a joint source-channel region-based rate control algorithm for real time video transmissions over the variable data rate wireless channels. The packet-level Gilbert two-state Markov model is used to characterize wireless channels. The proposed algorithm adopts a fast segmentation method to extract the regions of interest in real time, then it exploits the most effective criteria "coding qualities" to directly and quantitatively control bit allocation among different regions, which is more advantageous than traditional bit allocation methods used in the region/content-based rate control. The algorithm has low computational complexity and implementation cost and therefore is suitable for real-time wireless applications. From the experimental results, it was observed that, when applied the proposed algorithm, obviously more target bits are allocated to ROI at the expense of fewer bits distributed to non-ROI. This gives rise to the improved subjective performance of ROI while visual quality on non-ROI is degraded. By virtue of using our proposed PID buffer technique and taking channel conditions into consideration during bit allocation, our algorithm significantly reduce the number of frame skipping; thereby, improve the motion continuity in wireless video transmissions.

Regarding future work directions, we will continue our research on developing segmentation methods to real time detect the regions of interest more efficiently and accurately; exploring joint power-source-channel optimized approaches which simultaneously control the source rate, transmission power, and error resilience, so as to minimize the consumption of the total power while maintain the desired quality for robust video communications over wireless networks; and improving bit allocation strategies, etc.

## REFERENCES

[1] Z. He, J. Cai, and C. W. Chen, "Joint source channel rate-distortion analysis for adaptive mode selection and rate control in wireless video coding," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 12, no. 6, pp. 511–523, Jun. 1999.

[2] Video Coding for Low Bit Rate Communications, Jan. 1998.

[3] T. Sikora, "The MPEG-4 video standard verification model," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 7, pp. 19–31, Feb. 1997.

[4] M. Khansari, A. Jalalali, E. Dubois, and P. Mermelstein, "Low bit-rate video transmission over fading channels for wireless microcellular systems," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 6, no. 1, pp. 1–11, Feb. 1996.

[5] S. Aramvith, I.-M. Pao, and M. Sun, "A rate-control scheme for video transport over wireless channels," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 11, no. 5, pp. 569–580, May 2001.

[6] H. Song and C.-C. Jay Kuo, "A region-based H.263+ codec and its rate control for low VBR video," *IEEE Trans. Multimedia*, vol. 6, no. 3, pp. 489–500, Jun. 2004.

[7] Y. Yokoyama, Y. Miyamoto, and M. Ohta, "Very low bit rate video coding using arbitrarily shaped region-based motion compensation," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 5, no. 6, pp. 500–507, Dec. 1995.

[8] C.-H. Lin and J.-L. Wu, "Content-based rate control scheme for very low bit-rate video coding," *IEEE Trans. Consumer Electron.*, vol. 43, no. 2, pp. 123–133, May 1997.

[9] S. Aramvith, H. Kortrakulkij, D. Tancharoen, and S. Jitapankul, "Joing source-channel coding using simplified block-based segmentation and content-based rate-control for wireless video transport," in *Proc. Int. Conf. Information Technology: Coding and Computing (ITCC) 2002*, Las Vegas, Apr. 2002, pp. 71–76.

[10] C.-W. Lin, Y.-J. Chang, and Y.-C. Chen, "A low-complexity face-assisted coding scheme for low-bit-rate video telephony," *IEICE Trans. Inform. Syst.*, vol. E86-D, no. 1, pp. 101–108, Jan. 2003.

[11] A. H. Sadka, *Compressed Video Communications*. New York: Wiley, 2002.

[12] T. S. Rappaport, *Wireless Communications-Principles and Practice*, 1st ed. Englewood Cliffs, NJ: Prentice-Hall, 1996.

[13] J. D. Parsons, *The Mobile Radio Propagation Channel*, 2nd ed. New York: Wiley, 2002.

[14] W. C. Jakes, Ed., *Microwave Mobile Communications*. Piscataway, NJ: IEEE Press, 1974.

[15] M. D. Yacoub, J. E. V. Bautistu, and L. G. de Rezende Guedes, "On higher order statistics of the nakagami-m distribution," *IEEE Trans. Veh. Technol.*, vol. 48, no. 3, pp. 790–794, May 1999.

[16] C.-D. Iskander and P. T. Mathiopoulos, "Analytical level crossing rates and average fade durations for diversity techniques in nakagami fading channels," *IEEE Trans. Commun.*, vol. 50, no. 8, pp. 1301–1309, Aug. 2002.

[17] M. D. Yacoub, C. R. C. M. da Silva, and J. E. V. Bautista, "Second order statistics for diversity-combining techniques in nakagami-fading channels," *IEEE Trans. Veh. Technol.*, vol. 50, no. 6, pp. 1464–1470, Nov. 2001.

[18] P. Sweeney, *Error Control Coding-From Theory to Practice*. New York: Wiley, 2002.

[19] E. N. Gilbert, "Capacity of a burst-noise channel," *Bell Syst. Tech. J.*, vol. 39, pp. 1253–1265, Sept. 1960.

[20] W. Turin and R. van Nobelen, "Hidden markov modeling of flat fading channels," *IEEE J. Select. Areas Commun.*, vol. 16, no. 9, pp. 1809–1817, Dec. 1998.

[21] H. S. Wang, "On verifying the first-order Markovian assumption for a Rayleigh fading channel model," in *Proc. ICUPC'94*, pp. 160–164.

[22] M. Zorzi, R. R. Rao, and L. Milstein, "On the accuracy of a first-order Markov model for data transmission on fading channels," in *Proc. ICUPC'95*, pp. 160–164.

[23] UE radio transmission and reception (FDD), in 3rd Generation Partnership Project (3GPP) Technical Specification Group (TSG) Radio Access Networks (RAN) WG4.

[24] C. Y. Hsu, A. Ortega, and M. Khansari, "Rate control for robust video transmission over burst-error wireless channels," *IEEE J. Select. Areas Commun.*, vol. 17, no. 5, pp. 756–773, May 1999.

[25] L. Garcia, *Probability and Random Processes for Electrical Engineering*, 2nd ed. Norwell, MA: Addison-Wesley, 1994.

[26] Y. Sun and I. Ahmad, "A robust and adaptive rate control algorithm for object-based video coding," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 14, no. 10, pp. 1167–1182, Oct. 2004.

[27] H.-J. Lee, T. Chiang, and Y.-Q. Zhang, "Scalable rate control for MPEG-4 video," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 10, pp. 878–894, Sep. 2000.

[28] J. Ribas-Corbera and S. Lei, "Rate control in DCT video coding for low-delay communications," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 9, pp. 172–185, Feb. 1999.

[29] W. Li, J.-R. Ohm, M. van der Schaar, H. Jiang, and S. Li, MPEG-4 Video Verification Model V18.0, Pisa, Italy, Jan. 2000.

**Yu Sun** (S'04–M'05) received the B.S. and M.S. degrees in computer science from the University of Electronic Science and Technology of China, Chengdu, China, in 1996, and the Ph.D. degree in computer science and engineering from The University of Texas at Arlington in 2004.

From 1996 to 1998, she was a Lecturer in the Department of Computer Science, Sichuan Normal University, China. Since August 2004, she has been an Assistant Professor in the Department of Computer Science, University of Central Arkansas, Conway. Her main research interests include video compression, multimedia communication, and image processing.

**Ishfaq Ahmad** (S'88–M'92–SM'03) received the B.Sc. degree in electrical engineering from the University of Engineering and Technology, Lahore, Pakistan, in 1985, and the M.S. degree in computer engineering and the Ph.D. degree in computer science from Syracuse University, Syracuse, NY, in 1987 and 1992, respectively.

He is currently a Full Professor of computer science and engineering in the CSE Department, University of Texas at Arlington. Prior to joining UT Arlington, he was an Associate Professor in the Computer Science Department, Hong Kong University of Science and Technology (HKUST), where he was also the director of the Multimedia Technology Research Center, an officially recognized research center that he conceived and built from scratch. The center was funded by various agencies of the Government of the Hong Kong Special Administrative Region as well as local and international industries. With more than 40 personnel including faculty members, postdoctoral fellows, full-time staff, and graduate students, the center engaged in numerous R & D projects with academia and industry from Hong Kong, China, and the U.S. Particular areas of focus in the center are video (and related audio) compression technologies, video telephone and conferencing system. The center commercialized several of its technologies to its industrial partners worldwide. His recent research focus has been on developing parallel programming tools, scheduling and mapping algorithms for scalable architectures, heterogeneous computing systems, distributed multimedia systems, video compression techniques, and web management. His research work in these areas is published in over 125 technical papers in refereed journals and conferences.

Dr. Ahmad has participated in the organization of several international conferences and is an associate editor of *Cluster Computing*, *Journal of Parallel and Distributed Computing*, IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY, *IEEE Concurrency*, and *IEEE Distributed Systems Online*. He received best paper awards at Supercomputing '90 (New York), Supercomputing '91 (Albuquerque, NM), and 2001 International Conference on Parallel Processing (Spain).

**Dongdong Li** (S'03) received the B.Sc., M.Sc. in Electrical Engineering from Beijing University of Posts and Telecom, Beijing, China, in 1997 and 2000, respectively. He is currently a Ph.D. candidate in Electrical Engineering at the University of Texas at Arlington.

Since 2004, he has been with Cerion, Inc., Frisco, TX, as a System Engineer focusing on GSM/GPRS/EDGE/UMTS traffic modeling and network optimization for AT&T/Cingular Wireless. His current research interests include CDMA system optimization and performance evaluation, channel modeling, diversity receiver performance evaluation, and multimedia video transmission in wireless systems.

**Ya-Qin Zhang** (S'87–M'90–SM'93–F'97) received the B.S. and M.S. degrees in electrical engineering in 1983 and 1985, respectively, from the University of Science and Technology of China (USTC) and the Ph.D. degree in electrical engineering from George Washington University, Washington DC, in 1989. He received executive business training from Harvard University, Cambridge, MA.

He is Corporate Vice President of Microsoft Corporation's Mobile and Embedded Devices Division, Seattle, and was the Managing Director of Microsoft Research Asia, Beijing, China, which he joined in January 1999. Before that, he was Director of Multimedia Technology Laboratory, Sarnoff Corporation, Princeton, NJ (formerly David Sarnoff Research Center, and RCA Laboratories). He has been engaged in research and commercialization of MPEG2/DTV, MPEG4/VLBR, and multimedia information technologies. He was with GTE Laboratories, Inc., Waltham, MA, and Contel Technology Center, Chantilly, VA, from 1989 to 1994. He has more than 70 U.S. patents—granted or pending—and has authored or contributed to more than a dozen books and 300 influential technical papers and journal articles. Many of the technologies he and his team developed have become the basis for start-up ventures, commercial products, and international standards.

Dr. Zhang served as the Editor-in-Chief for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY from July 1997 to July 1999. He was the Chairman of Visual Signal Processing and Communications Technical Committee of IEEE Circuits and Systems. He serves on the editorial boards of seven other professional journals and over a dozen conference committees. He has been a key contributor to the ISO/MPEG and ITU standardization efforts in digital video and multimedia. He received numerous awards, including several industry technical achievement awards and IEEE awards such as Jubilee Golden Medal. He was awarded as the "Research Engineer of the Year" in 1998 by the New Jersey Engineering Council for his "leadership and invention in communications technology, which has enabled dramatic advances in digital video compression and manipulation for broadcast and interactive television and networking applications." He received the Prestigious National Award as "The Outstanding Young Electrical Engineer of 1998," given annually to one electrical engineer in the U.S.