

Wavelet-Based Video Compression Using Long-Term Memory Motion-Compensated Prediction and Context-Based Adaptive Arithmetic Coding

Detlev Marpe¹, Thomas Wiegand¹, and Hans L. Cycon²

¹ Image Processing Department,
Heinrich-Hertz-Institute (HHI) for Communication Technology, Einsteinufer 37,
10587 Berlin, Germany
{marpe,wiegand}@hhi.de

² University of Applied Sciences (FHTW Berlin), Allee der Kosmonauten 20–22,
10315 Berlin, Germany
hcycon@fhtw-berlin.de

Abstract. In this paper, we present a novel design of a wavelet-based video coding algorithm within a conventional hybrid framework of temporal motion-compensated prediction and transform coding. Our proposed algorithm involves the incorporation of multi-frame motion compensation as an effective means of improving the quality of the temporal prediction. In addition, we follow the rate-distortion optimizing strategy of using a Lagrangian cost function to discriminate between different decisions in the video encoding process. Finally, we demonstrate that context-based adaptive arithmetic coding is a key element for fast adaptation and high coding efficiency. The combination of overlapped block motion compensation and frame-based transform coding enables blocking-artifact free and hence subjectively more pleasing video. In comparison with a highly optimized MPEG-4 (Version 2) coder, our proposed scheme provides significant performance gains in objective quality of 2.0–3.5 dB PSNR.

1 Introduction

Multi-frame prediction [11] and variable block size motion compensation in a rate-distortion optimized motion estimation and mode selection process [12, 10] are powerful tools to improve the coding efficiency of today's video coding standards. In this paper, we present the design of a video coder, dubbed *DVC*, which demonstrates how most elements of the state-of-the-art in video coding as currently implemented in the test model long-term [2] (TML8) of the ITU-T H.26L standardization project can be successfully integrated in a blocking-artifact free video coding environment. In addition, we provide a solution for an efficient macroblock based intra coding mode within a frame-based residual coding method, which is extremely beneficial for improving the subjective quality as well as the error robustness.

We further explain how appropriately designed entropy coding tools, which have already been introduced in some of our previous publications [6, 7] and which, in some modified form [5], are now part of TML8, help to improve the efficiency of a wavelet-based residual coder.

In our experiments, we compared our proposed wavelet-based DVC coder against an improved MPEG-4 coder [10], where both codecs were operated using a fixed frame rate, fixed quantization step sizes and a search range of ± 32 pels. We obtained coding results for various sequences showing that our proposed video coding system yields a coding gain of 2.0–3.5 dB PSNR relative to MPEG-4. Correspondingly, the visual quality provided by the DVC coder compared to that of the block-based coding approach of MPEG-4 is much improved, especially at very low bit rates.

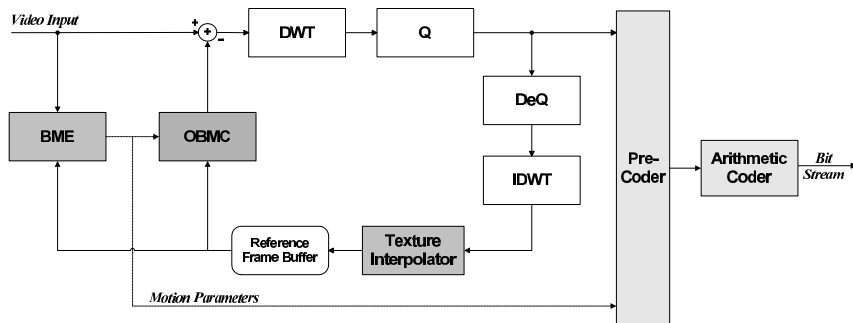


Fig. 1. Block diagram of the proposed coding scheme

2 Overview of the DVC Scheme

Fig. 1 shows a block diagram of the proposed DVC coder. As a hybrid system, it consists of a temporal predictive loop along with a spatial transform coder. Temporal prediction is performed by using a *block motion estimation* (BME) and an *overlapped block motion compensation* (OBMC), such that the reference of each predicted block can be obtained from a long-term *reference frame memory*. Coding of the motion compensated *P*-frames as well as of the initial intra (*I*) frame is performed by first applying a *discrete wavelet transform* (DWT) to an entire frame. Uniform scalar *quantization* (*Q*) with a central dead-zone around zero similar to that designed for H.263 is then used to map the dynamic range of the wavelet coefficients to a reduced alphabet of decision levels. Prior to the final *arithmetic coding* stage, the *pre-coder* further exploits redundancies of the quantized wavelet coefficients in a 3-stage process of partitioning, aggregation and conditional coding.

Table 1. Macroblock partition modes

Mode	Block Size	Partition
1	16×16	Leave MB as a whole
2	16×8	Split MB into 2 sub-blocks
3	8×16	Split MB into 2 sub-blocks
4	8×8	Split MB into 4 sub-blocks

3 Motion-Compensated Prediction

3.1 Motion Model

As already stated above, the motion model we used is very similar to that of the H.26L TML8 design [2]. In essence it relies on a simple model of block displacements with variable block sizes. Given a partition of a frame into macroblocks (MB) of size 16×16 pels, each macroblock can be further sub-divided into smaller blocks, where each sub-block has its own displacement vector. Our model supports 4 different partition modes, as shown in Table 1.

Each macroblock may use a different reference picture out of a long-term frame memory. In addition to the predictive modes represented by the 4 different MB partition modes in Table 1, we allow for an additional macroblock-based intra coding mode in P-frames. This local intra mode is realized by computing the DC for each 8×8 sub-block of each spectral component (Y,U,V) in a macroblock and by embedding the DC-corrected sub-blocks into the residual frame in a way, which is further described in the following section.

3.2 Motion Estimation and Compensation

Block motion estimation is performed by an exhaustive search over all integer pel positions within a pre-defined search window around the motion vector predictor, which is obtained from previously estimated sub-blocks in the same way as in TML8 [2]. In a number of subsequent steps, the best integer pel motion vector is refined to the final $\frac{1}{4}$ -pel accuracy by searching in a 3×3 sub-pel window around the refined candidate vector. All search positions are evaluated by using a Lagrangian cost function, which involves a rate and distortion term coupled by a Lagrangian multiplier. For all fractional-pel displacements, distortion in the transform domain is estimated by using the Walsh-Hadamard transform, while the rate of the motion vector candidates is estimated by using a fixed, pre-calculated table. This search process takes place for each of the 4 macroblock partitions and each reference frame, and the cost of the overall best motion vector candidate(s) of all 4 macroblock modes is finally compared against the cost of the intra mode decision to choose the macroblock mode with minimum cost.

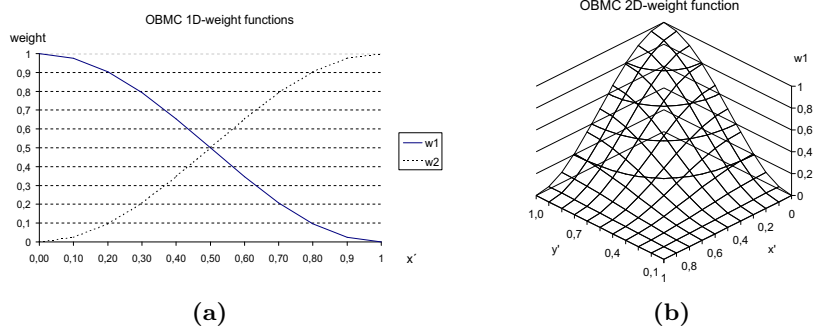


Fig. 2. (a) 1-D profile of 2-D weighting functions along the horizontal or vertical axes of two neighboring overlapping blocks. (b) 2-D weighting function

The prediction error luminance (chrominance) signal is formed by the weighted sum of the differences between all 16×16 (8×8) overlapping blocks from the current frame and their related overlapping blocks with displaced locations in the reference frame, which have been estimated in the BME stage for the corresponding core blocks. In the case of an intra macroblock, we compute the weighted sum of the differences between the overlapping blocks of the current intra blocks and its related DC-values. As a weighting function w , we used the 'raised cosine', as shown in Fig. 2. For a support of $N \times N$ pels, it is given by

$$w(n, m) = w_n \cdot w_m, \quad w_n = \frac{1}{2} \left[1 - \cos \frac{2\pi n}{N} \right] \quad \text{for } n = 0, \dots, N. \quad (1)$$

In our presented approach, we choose $N = 16$ ($N = 8$) for the luminance (chrominance, resp.) in Eq. (1), which results in a 16×16 (8×8) pixel support centered over a "core" block of size 8×8 (4×4) pels for the luminance (chrominance, resp.). For the texture interpolation of sub-pel positions, the same filters as specified in TML8 [2] have been used.

4 Wavelet Transform

In wavelet-based image compression, the so-called 9/7-wavelet with compact support [3] is the most popular choice. Our proposed coding scheme, however, utilizes a class of biorthogonal wavelet bases associated with infinite impulse response (IIR) filters, which was recently constructed by Petukhov [9]. His approach relies on the construction of a dual pair of rational solutions of the matrix equation

$$M(z)\tilde{M}^T(z^{-1}) = 2I, \quad (2)$$

where I is the identity matrix, and

$$M(z) = \begin{pmatrix} h(z) & h(-z) \\ g(z) & g(-z) \end{pmatrix}, \quad \tilde{M}(z) = \begin{pmatrix} \tilde{h}(z) & \tilde{h}(-z) \\ \tilde{g}(z) & \tilde{g}(-z) \end{pmatrix}$$

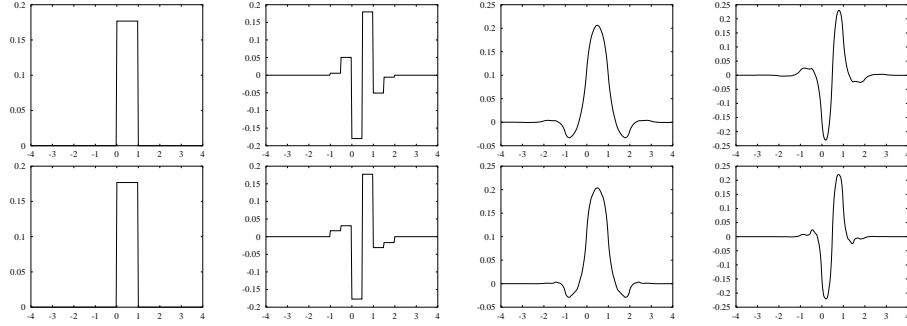


Fig. 3. From left to right: scaling function of analysis, analyzing wavelet, scaling function of synthesis, and synthesizing wavelet used for I-frame coding (*top row*) and P-frame coding (*bottom row*)

are so-called ‘modulation matrices’.

In [9], a one-parametric family of filters h_a , g_a , \tilde{h}_a and \tilde{g}_a satisfying Eq. (2) was constructed:¹

$$h_a(z) = \frac{1}{\sqrt{2}}(1 + z), \tag{3}$$

$$\tilde{h}_a(z) = \frac{(2 + a)(z^{-1} + 3 + 3z + z^2)(z^{-1} + b + z)}{4\sqrt{2}(2 + b)(z^{-2} + a + z^2)}, \tag{4}$$

$$g_a(z) = \frac{(2 + a)(z^{-1} - 3 + 3z - z^2)(-z^{-1} + b - z)}{4\sqrt{2}(2 + b)}, \tag{5}$$

$$\tilde{g}_a(z) = \frac{1}{\sqrt{2}} \frac{1 - z^{-1}}{z^{-2} + a + z^2}, \tag{6}$$

where $b = \frac{4a-8}{6-a}$, $|a| > 2$, $a \neq 6$.

To adapt the choice of the wavelet basis to the nature and statistics of the different frame types of intra and inter mode, we performed a numerical simulation on this one-parametric family of IIR filter banks yielding the optimal value of $a = 8$ for intra frame mode and $a = 25$ for inter frame mode in Eqs. (3)–(6). Graphs of these optimal basis functions are presented in Fig. 3. Note that the corresponding wavelet transforms are efficiently realized with a composition of recursive filters [9].

5 Pre-Coding of Wavelet Coefficients

For encoding the quantized wavelet coefficients, we follow the conceptual ideas initially presented in [6] and later refined in [7]. Next, we give a brief review of the involved techniques. For more details, the readers are referred to [6, 7].

¹ h_a and g_a denote low-pass and high-pass filters of the decomposition algorithm, respectively, while \tilde{h}_a and \tilde{g}_a denote the corresponding filters for reconstruction.

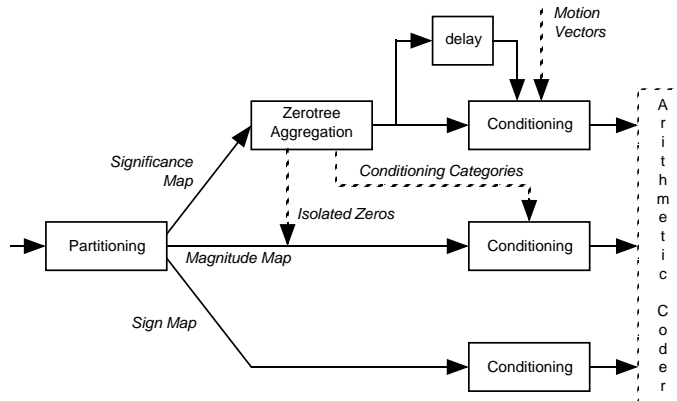


Fig. 4. Schematic representation of the pre-coder used for encoding the quantized wavelet coefficients

5.1 Partitioning

As shown in the block diagram of Fig. 4, an initial ‘partitioning’ stage divides each frame of quantized coefficients into three sub-sources: a significance map, indicating the position of significant coefficients, a magnitude map holding the absolute values of significant coefficients, and a sign map with the phase information of the wavelet coefficients. Note that all three sub-sources inherit the subband structure from the quantized wavelet decomposition, so that there is another partition of each sub-source according to the given subband structure.

5.2 Zerotree Aggregation

In a second stage, the pre-coder performs an ‘aggregation’ of insignificant coefficients using a quad-tree related data structure. These so-called *zerotrees* [4, 6] connect insignificant coefficients, which share the same spatial location along the multiresolution pyramid. However, we do not consider zero-tree roots in bands below the maximum decomposition level. In inter-frame mode, coding efficiency is further improved by connecting the zerotree root symbols of all three lowest high-frequency bands to a so-called ‘integrated’ zerotree root which resides in the LL-band.

5.3 Conditional Coding

The final ‘conditioning’ part of the pre-coding stage supplies the elements of each source with a ‘context’, *i.e.*, an appropriate model for the actual coding process in the arithmetic coder. Fig. 5 (a) shows the prototype template used for conditioning of elements of the significance map. In the first part, it consists of a causal neighborhood of the actual coding event C , which depends on the scale and orientation of a given band. Except for the lowest frequency bands, the

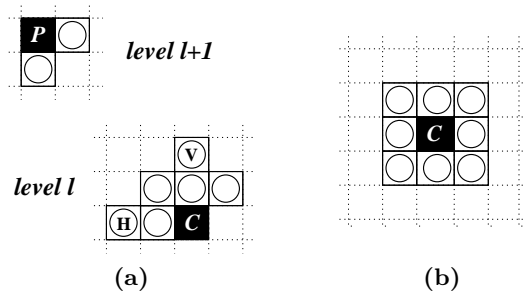


Fig. 5. (a) Two-scale template (*white circles*) with an orientation dependent design for conditional coding of an event C of the significance map; V , H : additional element used for vertical and horizontal oriented bands, respectively. (b) 8-neighborhood of significance used for conditioning of a given magnitude C

template uses additional information of the next upper level (lower resolution) represented by the neighbors of the parent P of C , thus allowing a 'prediction' of the non-causal neighborhood of C . The processing of the lowest frequency band depends on the intra/inter decision. In intra mode, mostly non-zero coefficients are expected in the LL-band, so there is no need for coding a significance map. For P-frames, however, we indicate the significance of a coefficient in the LL-band by using the four-element kernel of our prototype template (Fig. 5 (a)), which is extended by the related entry of the significance map belonging to the previous P-frame.

The processing of subbands is performed band-wise in the order from lowest to highest frequency bands and the partitioned data of each band is processed such that the significance information is coded (and decoded) first. This allows the construction of special conditioning categories for the coding of magnitudes using the local significance information. Thus, the actual conditioning of magnitudes is performed by classifying magnitudes of significant coefficients according to the local variance estimated by the significance of their 8-neighborhood (cf. Fig. 5 (b)). For the conditional coding of sign maps, we are using a context built of two preceding signs with respect to the orientation of a given band [7].

For coding of the LL-band of I-frames, the proposed scheme uses a DPCM-like procedure with a spatially adaptive predictor and a backward driven classification of the prediction error using a six-state model.

6 Binarization and Adaptive Binary Arithmetic Coding

All symbols generated by the pre-coder are encoded using an adaptive binary arithmetic coder, where non-binary symbols like magnitudes of coefficients or motion vector components are first mapped to a sequence of binary symbols by means of the unary code tree. Each element of the resulting "intermediate" code-word given by this so-called *binarization* will then be encoded in the subsequent process of binary arithmetic coding.

At the beginning of the overall encoding process, the probability models associated with all different contexts are initialized with a pre-computed start distribution. For each symbol to encode the frequency count of the related binary decision is updated, thus providing a new probability estimate for the next coding decision. However, when the total number of occurrences of symbols related to a given model exceeds a pre-defined threshold, the frequency counts will be scaled down. This periodical rescaling exponentially weighs down past observations and helps to adapt to non-stationarities of a source.

For intra and inter frame coding we use separate models. Consecutive P-frames as well as consecutive motion vector fields are encoded using the updated related models of the previous P-frame and motion vector field, respectively. The binary arithmetic coding engine used in our presented approach is a straightforward implementation similar to that given in [13].

7 Experimental Results

7.1 Test Conditions

To illustrate the effectiveness of our proposed coding scheme, we used an improved MPEG-4 coder [10] as a reference system. This coder follows a rate-distortion (R-D) optimized encoding strategy by using a Lagrangian cost function, and it generates bit-streams compliant with MPEG-4, Version 2 [1]. Most remarkable is the fact that this encoder provides PSNR gains in the range from 1.0-3.0 dB, when compared to the MoMuSys reference encoder (VM17) [10]. For our experiments, we used the following encoding options of the improved MPEG-4 reference coder:

- $\frac{1}{4}$ -pel motion vector accuracy enabled
- Global motion compensation enabled
- Search range of ± 32 pels
- 2 B-frames inserted (*IBBPBBP...*)
- MPEG-2 quantization matrix used

For our proposed scheme, we have chosen the following settings:

- $\frac{1}{4}$ -pel motion vector accuracy
- Search range of ± 32 pels around the motion vector predictor
- No B-frames used (*IPPP...*)
- Five reference pictures were used for all sequences except for the 'News'-sequence (see discussion of results below)

Coding experiments were performed by using the test sequences 'Foreman' and 'News' both in QCIF resolution and with 100 frames at a frame rate of 10 Hz. Only the first frame was encoded as an I-frame; all subsequent frames were encoded as P-frames or B-frames. For each run of a sequence, a set of quantizer parameters according to the different frame types (I,P,B) was fixed.

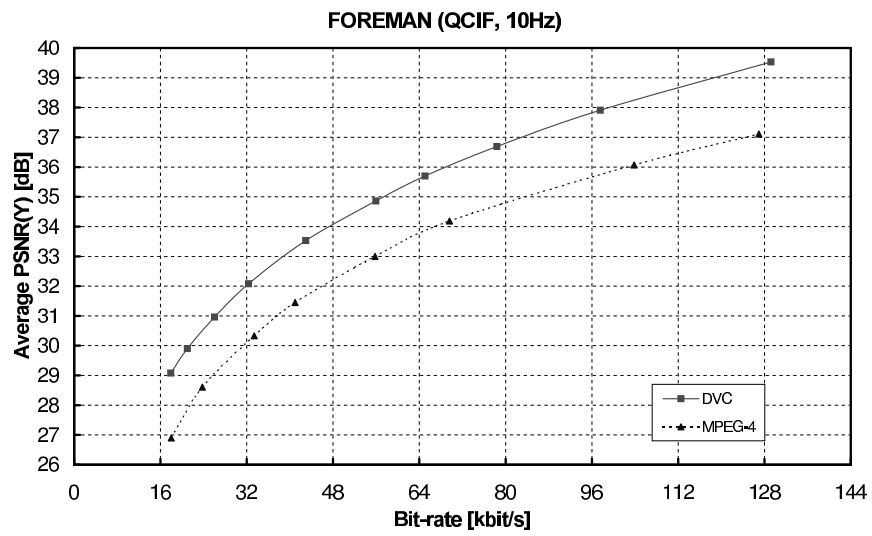


Fig. 6. Average Y-PSNR against bit-rate using the QCIF test sequence 'Foreman' at a frame rate of 10 Hz

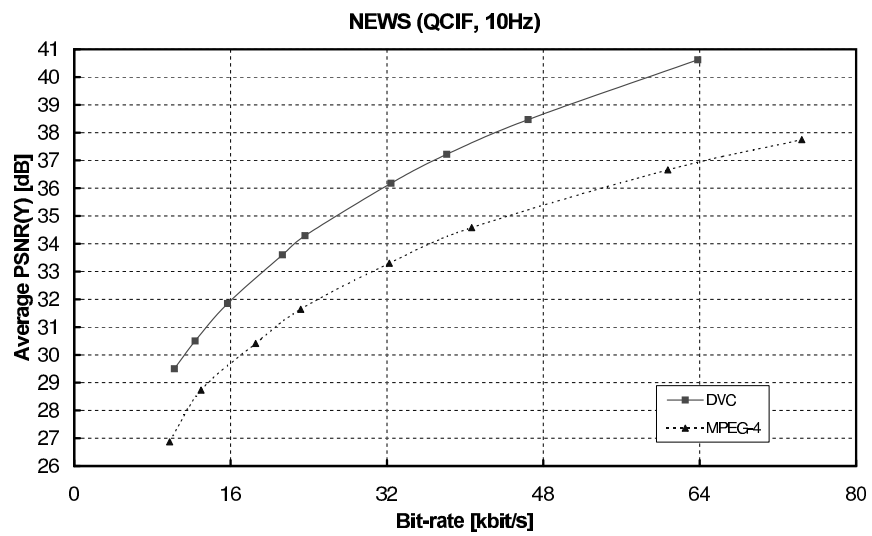


Fig. 7. Average Y-PSNR against bit-rate using the QCIF test sequence 'News' at a frame rate of 10 Hz

7.2 Test Results

Figs. 6–7 show the average PSNR gains obtained by our proposed DVC scheme relative to the MPEG-4 coder for the test sequences 'Foreman' and 'News', respectively. For the 'Foreman'-sequence, significant PSNR gains of 2.0–2.5 dB on the luminance component have been achieved (cf. Fig. 6). Figure 8 shows a comparison of the visual quality for a sample reconstruction at 32 kbit/s. The results we obtained for the "News"-sequence show dramatic PSNR improvements of about 2.5–3.5 dB. To demonstrate the ability of using some *a priori* knowledge about the scene content, we checked for this particular sequence in addition to the five most recent reference frames one additional reference frame 50 frames back in the past according to the repetition of parts of the scene content. By using the additional reference frame memory for this particular test case, we achieved an additional gain of about 1.5 dB PSNR on the average compared to the case where the reference frame buffer was restricted to the five most recent reference frames only.

8 Conclusions and Future Research

The coding strategy of DVC has proven to be very efficient. PSNR gains of up to 3.5 dB relative to an highly optimized MPEG-4 coder have been achieved. However, it should be noted that in contrast to the MPEG-4 coding system, no B-frames were used in the DVC scheme, although it can be expected that DVC will benefit from the usage of B-frames in the same manner as the MPEG-4 coder, *i.e.*, depending on the test material, additional PSNR improvements of up to 2 dB might be achievable. Another important point to note, when comparing the coding results of our proposed scheme to that of the highly R-D optimized MPEG-4 encoder used for our experiments, is the fact that up to now, we did not incorporate any kind of high-complexity R-D optimization method. We even did not optimize the motion estimation process with respect to the overlapped motion compensation, although it is well known, that conventional block motion estimation is far from being optimal in an OBMC framework [8]. Furthermore, we believe that the performance of our zerotree-based wavelet coder can be further improved by using a R-D cost function for a joint optimization of the quantizer and the zerotree-based encoder. Thus, we expect another significant gain by exploiting the full potential of encoder optimizations inherently present in our DVC design. This topic will be a subject of our future research.

References

1. ISO/IEC JTC1SC29 14496-2 MPEG-4 Visual, Version 2.
2. Bjontegaard, G. (ed.): H.26L Test Model Long Term Number 8 (TML8), ITU-T SG 16 Doc. VCEG-N10 (2001)
3. Cohen, A., Daubechies, I., Feauveau, J.-C.: Biorthogonal Bases of Compactly Supported Wavelets, *Comm. on Pure and Appl. Math.*, Vol. 45 (1992) 485–560



Fig. 8. Comparison of subjective reconstruction quality: Frame no. 22 of 'Foreman' at 32 kbit/s. (a) DVC reconstruction (b) MPEG-4 reconstruction. Note that the MPEG-4 reconstruction has been obtained by using a de-blocking filter

4. Lewis, A., Knowles, G.: Image Compression Using the 2D Wavelet Transform, *IEEE Trans. on Image Processing*, Vol. 1, No. 2 (1992) 244–250
5. Marpe, D., Blättermann, G., Wiegand, T.: Adaptive Codes for H.26L, ITU-T SG 16 Doc. VCEG-L13 (2001)
6. Marpe, D., Cycon, H. L.: Efficient Pre-Coding Techniques for Wavelet-Based Image Compression, *Proceedings Picture Coding Symposium 1997*, 45–50
7. Marpe, D., Cycon, H. L.: Very Low Bit-Rate Video Coding Using Wavelet-Based Techniques, *IEEE Trans. on Circuits and Systems for Video Technology*, Vol. 9, No. 1 (1999) 85–94
8. Orchard, M. T., Sullivan, G. J.: Overlapped Block Motion Compensation: An Estimation-Theoretic Approach, *IEEE Trans. on Image Processing*, Vol. 3, No. 5 (1994) 693–699
9. Petukhov, A. P.: Recursive Wavelets and Image Compression, *Proceedings International Congress of Mathematicians 1998*
10. Schwarz, H., Wiegand, T.: An Improved MPEG-4 Coder Using Lagrangian Coder Control, ITU-T SG 16 Doc. VCEG-M49 (2001)
11. Wiegand, T., Zhang, X., Girod, B.: Long-Term Memory Motion-Compensated Prediction, *IEEE Trans. on Circuits and Systems for Video Technology*, Vol. 9, No. 1 (1999) 70–84
12. Wiegand, T., Lightstone, M., Mukherjee, D., Campbell, T. G., Mitra, S. K.: Rate-Distortion Optimized Mode Selection for Very Low Bit Rate Video Coding and the Emerging H.263 Standard, *IEEE Trans. on Circuits and Systems for Video Technology*, Vol. 6, No. 2 (1996) 182–190
13. Witten, I., Neal, R., Cleary, J.: Arithmetic Coding for Data Compression, *Communications of the ACM*, Vol. 30 (1987) 520–540