

# Optimal Mode Selection for a Pulsed-Quality Dual Frame Video Coder

Athanasios Leontaris, Vijay Chellappa, and Pamela C. Cosman\*

**Abstract**—A dual frame video coder employs two past reference frames for motion compensated prediction. Compared to conventional single frame prediction, the dual frame encoder can have advantages both in distortion-rate performance and in error resilience. In previous work, it was shown that optimal mode selection can enhance the performance of a dual frame encoder. In another strand of previous work, it was shown that uneven assignment of quality to frames, to create high-quality long-term reference frames, can enhance the performance of a dual frame encoder. In this letter, we combine these two strands and demonstrate the performance advantages of optimal mode selection among high-quality frames for video transmission over noisy channels.

**Index Terms**—Video compression, mode switching, H.264, dual frame buffer, high quality updating, per-pixel estimation, multiple frame prediction.

## I. INTRODUCTION

Traditionally, hybrid video codecs employ motion-compensated prediction to compress an input raw video stream. A block of pixels in the current frame is predicted from a displaced block in a previous frame. A motion vector points to the coordinates of the displaced block. The difference (error signal) between the original block and its prediction is compressed and transmitted along with the corresponding displacement (motion) vectors. This approach has formed the cornerstone of modern video coding algorithms such as MPEG-4 [1] and H.263+ [2].

Performance was improved when the search for the best prediction included additional past frames apart from the previous one. Examples of this multiple reference frame approach are [3], [4], [5], and it has recently been standardized as part of the H.264/AVC video coding standard [6]. To counter the increased memory and computational cost, the number of reference frames can be constrained to be small. Fukuhara *et al.* used only two reference frames, with encouraging results [7]. The first reference buffer contained the previous frame, and the second one contained a reference frame from the distant past that was periodically updated. We refer to this as *dual frame* coding. In [8], the authors use a linearly weighted combination of two frames, primarily to enhance the error robustness of the codec. Budagavi *et al.* [9] used Markov

chain analysis to prove that multiple frames increase error robustness.

A novel algorithm for estimating distortion due to packet losses was introduced in [10], for conventional single frame encoding. This distortion estimation was extended in [11], to allow dual frame coding with rate-distortion optimal mode switching. For packet loss scenarios, the performance increased noticeably. In other recent work [12], periodic high-quality frames were generated, and retained as the long-term frames in a dual frame encoding approach.

In this letter, we combine the high-quality construction and buffering of long-term frames from [12] with the optimal mode selection approach from [11] and show that the combination provides a significant advantage for lossy packet networks. For reasons of simplicity and reduced computational complexity, we made use of half-pel motion compensation to allow easier calculation of the distortion estimate, and the loop filter was disabled as well. Both greater motion vector accuracy and loop filtering could be enabled in our scheme, with some minor modifications and approximations (for the second moments) in the distortion estimation. The letter is organized as follows: Section 2 describes how distortion was estimated and optimal mode selection was performed for a dual frame encoder. Section 3 discusses the High Quality (HQ) updating approach. In Section 4, we provide experimental results for the combination of these two approaches. The letter is concluded with Section 5.

## II. OPTIMAL MODE SELECTION FOR A DUAL FRAME CODER

Dual frame motion compensation is depicted in Fig. 1, and works as follows. While encoding frame  $n$ , the encoder and decoder both maintain two reference frames in memory. The short-term reference is frame  $n - 1$ . The long-term reference can be selected in a number of ways; we used *jump updating* in which the long-term reference frame varies from as recent as frame  $n - 2$  to as old as frame  $n - N - 1$ . When encoding frame  $n$ , if the long-term reference frame is  $n - N - 1$ , then, when the encoder moves on to encoding frame  $n + 1$ , the short-term reference frame slides forward by one to frame  $n$ , and the long-term reference frame jumps forward by  $N$  to frame  $n - 1$ . The long-term reference frame then remains static for  $N$  frames, and then jumps forward again. We refer to  $N$  as the updating parameter. This approach was first adopted in [7].

In dual frame motion compensation, each macroblock can be encoded in one of three coding modes: intra coding, inter coding using the short-term buffer (inter-ST-coding), and inter

A. Leontaris, V. Chellappa, and P. C. Cosman are with the Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA 92093-0407 USA. (E-mail: {aleontar,vchellap,pcosman}@code.ucsd.edu Tel: 858-822-0157 FAX: 858-822-3426)

EDICS: 2.IMMD Image and Multidimensional Signal Processing, 2.MMSP Multimedia Signal Processing.

coding using the long-term buffer (inter-LT-coding). In [11], the choice among these three was made using an extended version of the ROPE algorithm, as described briefly below.

Using the notation from [10], we use  $f_n$ ,  $\hat{f}_n$ , and  $\tilde{f}_n$  to denote the original frame  $n$ , the encoder reconstruction of the compressed frame, and the decoder version (possibly error concealed) of the frame, respectively. We assume that the long-term frame buffer was updated  $m$  frames ago. Thus, it contains  $\hat{f}_{n-m}$  at the transmitter and  $\tilde{f}_{n-m}$  at the receiver. The expected distortion for pixel  $i$  in frame  $n$  is:

$$d_n^i = E\{(f_n^i - \tilde{f}_n^i)^2\} = (f_n^i)^2 - 2f_n^i E\{\tilde{f}_n^i\} + E\{(\tilde{f}_n^i)^2\} \quad (1)$$

Calculation of  $d_n^i$  requires the first and second moments of the random variable of the estimated image sequence  $\tilde{f}_n^i$ . We have two separate inter modes, the *inter-ST* and *inter-LT*. Let  $i$  denote the pixel in the current frame,  $k$  denote the pixel in the previous frame that is associated with pixel  $i$  in the current frame using error concealment, and  $j$  denote the pixel in the reference frame (either ST or LT) that is the prediction of pixel  $i$  in the current frame derived using the motion vector. Let  $p$  denote the packet erasure rate (which equals the pixel loss probability for our variable-length packets which contain a single horizontal group of blocks), and  $q = 1 - p$ .

The two required moments for a pixel in an *intra*-coded macroblock (MB) are [10]:

$$E\{\tilde{f}_n^i\} = q(\hat{f}_n^i) + pqE\{\tilde{f}_{n-1}^k\} + p^2E\{\tilde{f}_{n-1}^j\} \quad (2)$$

$$E\{(\tilde{f}_n^i)^2\} = q(\hat{f}_n^i)^2 + pqE\{(\tilde{f}_{n-1}^k)^2\} + p^2E\{(\tilde{f}_{n-1}^j)^2\} \quad (3)$$

The first and second moments of  $\tilde{f}_n^i$  for a pixel in an *inter*-coded MB are:

$$E\{\tilde{f}_n^i\} = q(\hat{e}_n^i + E\{\tilde{f}_{n-g}^j\}) + pqE\{\tilde{f}_{n-1}^k\} + p^2E\{\tilde{f}_{n-1}^j\} \quad (4)$$

$$E\{(\tilde{f}_n^i)^2\} = q((\hat{e}_n^i)^2 + 2\hat{e}_n^i E\{\tilde{f}_{n-g}^j\} + E\{(\tilde{f}_{n-g}^j)^2\}) + pqE\{(\tilde{f}_{n-1}^k)^2\} + p^2E\{(\tilde{f}_{n-1}^j)^2\} \quad (5)$$

where these equations are valid for the *inter-LT* mode if  $g = m$  and they are valid for the *inter-ST* mode if  $g = 1$ . Using these equations, the encoder can estimate recursively the per-pixel distortion of the reconstructed video at the decoder. For more details we refer the reader to [10] and [11].

The only major difference of the mode selection in this work, compared to [11], is that we use H.264 here, and therefore need to change the approximation used for half-pixel motion vectors, to account for the 6-tap interpolation filters used in H.264.

Given the distortion estimate, the encoder switches between intra, inter-ST or inter-LT coding on a macroblock basis, in an optimal fashion for a given bit rate and packet loss rate. The goal is to minimize the total distortion subject to a bit rate constraint. Individual macroblock contributions to this cost are additive, thus it can be minimized on a macroblock basis. Therefore, the encoding mode for each MB is chosen by minimizing  $\min_{(mode)} J_{MB} = \min_{(mode)} (D_{MB} + \lambda R_{MB})$  where

$D_{MB} = \sum_{i \in MB} d_n^i$  and  $R_{MB}$  denote per MB distortion and rate, respectively, and  $\lambda$  is the Lagrange multiplier. The coding mode (*intra*, *inter-ST* and *inter-LT*) is chosen to minimize the Lagrangian cost. Because of the uneven quality levels (discussed in the next section) being assigned in the current paper, contrary to what was done in [11], we do not optimize over the Quantization Parameter (QP). We instead use the one chosen for that particular frame (or MB), manually or through a rate allocator.

### III. HIGH QUALITY UPDATING

A question that arises when designing such a system is the choice of the optimal update parameter  $N$  for a given image sequence, frame rate and bit rate. It depends heavily on the sequence's characteristics, such as occlusion effects and scene changes. Thus, the problem is one of computer vision and image sequence characterization. An optimal solution will require significant computational resources. In [12] it was proposed not to attempt to *select* an optimal frame to be buffered as long-term, but rather to reformulate the problem as one of *constructing* a good frame explicitly. In this approach, every  $N$  frames, one frame is coded with additional bit rate at the expense of other regular frames. This frame is then buffered and used as the long-term reference frame for the subsequent  $N$  frames.

A second issue is how to allocate bit rate to the long-term frame. In this work, as in [12], the rate allocation was heuristic. We used a quantization parameter for the long-term frame that was lower by 7 compared to the quantization parameter for the regular frames. One plausible *optimal* solution to the problem would be to consider a block of  $N$  frames, and optimize the QP or rate of the long-term frame and the QP or rate of the rest over *all*  $N$  frames by applying rate-distortion optimization for all possible combinations of rates/QPs. However, the computational complexity would be immense, and there would be a delay of  $N$  frames.

The actual transmission of extra bits for the periodic high-quality frames can be accomplished in two ways. One can incur extra delay for the high-quality frames, by using extra transmission time for them, sending more information at the same average bit rate. Alternatively, the sender could use extra channel bandwidth for a short period of time, to send the extra bits for the high-quality frames.

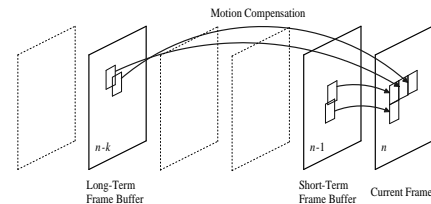


Fig. 1. Dual frame buffer scheme.

It was found in [12] that using a dual frame encoder with periodic high-quality frames (pulsed quality) that could be used as long term frames provided about a 0.6 dB advantage over a regular dual frame encoder where the long-term reference frames have the same quality as other frames.

#### IV. EXPERIMENTAL RESULTS

The previous work with pulsed-quality dual frames considered only transmission over noiseless channels [12]. In this letter, our goal is to use the pulsed-quality creation of long-term reference frames, but for packet erasure channels, and to use the optimal mode selection approach of [11] to optimally select among *intra* coding, *inter-low-quality* coding, and *inter-high-quality* coding based on estimates of the distortion that arises both from packet erasures and from the difference in low and high quality frames.

The dual frame buffer with high quality updating was implemented with the H.264/AVC reference software version JM 7.4. The encoder was modified, using functionality built into the standard. The resulting video codec produces a fully standard compliant H.264 bitstream. Since we deal with standard QCIF imagery at  $176 \times 144$  pixels, we took each slice (packet) to contain 11 MBs, so each variable-length packet is equivalent to a horizontal slice of MBs.

A multiple frame buffer of size two was employed. In the case of regular updating, the first frame is the previous and the second a long-term one. When high quality updating is used, the sole difference is that the frames selected to be buffered as long-term ones have been explicitly coded with a lower QP compared to all others. In our simulations we code the long-term frames with a QP that is lower by 7 compared to the general QP for the entire sequence. The long-term frame was updated every 10 frames (updating parameter  $N = 10$ ). This fixed value was selected experimentally as in [12] and is a compromise between the optimal values for many sequences.

Packet losses corrupt the bit stream. The loss of a packet translates to losing a horizontal slice of MBs. Packets can be decoded independently of one another. Error concealment is applied by using the median of the motion vectors of the three upper MBs to conceal from the previous frame. If the upper slice has been lost as well, then we just copy the co-located MB from the previous frame. The error concealment is modeled within the distortion estimation equations.

A hundred different random patterns were used to obtain the displayed results. Since the regular dual frame coder takes no account of the possibility of packet losses in choosing between the coding modes, it underutilizes the intra mode, and performs very poorly in a lossy environment. Thus, we also investigated the performance of random intra refresh algorithms.

The QP value was selected empirically in [12] after extensive experimentation for standard video sequences, and, although it cannot be said to be optimal, it was found to be an acceptable compromise for a range of image sequences. Constraining the QP selection and comparing against intra-coded long-term frames, we found that the present scheme outperformed intra-LTs for most of the cases by up to 1dB.

Fig. 3(a) illustrates the system's performance for the "carphone" image sequence, which consists of a man talking on his videophone in a moving car. For a packet loss rate of 10% the "zero intra" coding proves extremely error-prone. It has very few intra-coded MBs, and a  $p = 10\%$  packet loss ratio drops the PSNR down to 21dB from more than 32dB. Increasing the allocated rate has no effect since the absence of intra-coded MBs totally compromises the bistream's resilience.

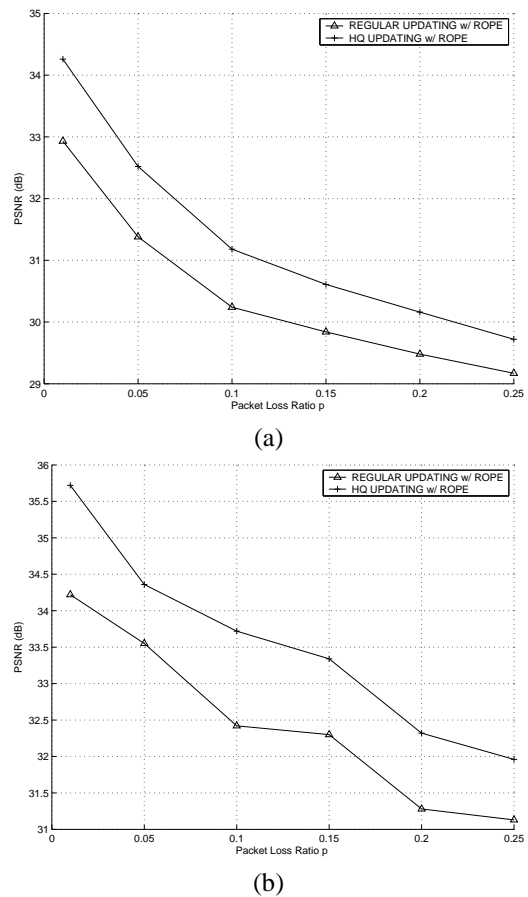


Fig. 2. PSNR performance vs. packet loss rate. (a) Image Sequence "carphone" QCIF at 10fps,  $N = 10$ , 122.5kbps. (b) Image Sequence "mother-daughter" QCIF at 10fps,  $N = 10$ , 34.4kbps.

Attempting to protect the bistream by employing some intra-coded MBs, we used a random intra refresh update. We experiment by forcing 20, 33 and 45 random intra-coded MBs *per-frame*. Performance increases with the transmission rate, due to the added protection of the intra MBs, particularly as the number of intra-coded MBs increases.

However, we observe that providing high-quality to the long-term frames does less well than regular quality for these heuristic intra refresh approaches. This is likely due to the fact that the high-quality frames are depriving other frames of their share of rate, and the random intra refresh macroblocks also deprive other MBs of their share of rate; the competing effects of these heavy rate users hurt the final performance. Similar conclusions can be drawn in Fig. 3(b).

However, with the use of rate-distortion optimal mode selection among *intra* coding, *inter-low-quality* coding, and *inter-high-quality* coding, the pulsed-quality dual frame approach outperforms regular dual frame updating by as much as 1-1.5dB throughout Figs. 3(a)-(b). For Fig. 3(b) in particular, the performance gap increases with the bit rate.

We then investigated the performance for varying packet loss ratios. Figs. 2(a)-(b) demonstrate that HQ updating outperforms regular updating for packet loss ratios ranging from 1% to 25%. The performance gain varies from 0.6-1.5dB depending on the image sequence. It is in general higher

for low-motion sequences such as “mother-daughter” but can reach 1 dB for active sequences such as “carphone”. QPs were chosen to achieve the same ( $\pm 5\%$ ) total bit rate for the graph points.

We experimented with using more than 2 reference frames, and found, as in [13], that expanding the reference buffer size beyond 2 frames produces sharply diminishing returns. In particular, we tried (1) two long-term (high-quality) frames plus one short-term frame, (2) two short-term frames and one long-term (high quality) frame, and (3) one long-term and four short-term frames. In all cases, the gains over the dual frame high-quality case were quite small (on the order of 0.1-0.2 dB). It appears that, for the sequences we tried, the immediate past frame captures most of the benefit that short term references can provide (high correlation with current frame), and a single high-quality frame captures most of the benefit that the high-quality long-term past can provide.

Subjective quality evaluation showed that the periodic coding of frames at high quality is only rarely noticeable for the error-free case, and for the error-prone case it is completely masked by the error concealment and propagation.

## V. CONCLUSION

In conclusion, the dual-frame coder with pulses of high quality provided to the long-term frame, when used in conjunction with random intra refresh, performs *less well* than the regular dual frame coder where long-term frames are chosen from among regular quality frames. The optimal mode selection does significantly better than random intra refresh for both regular quality and pulsed-quality dual frame coders, and the pulsed-quality approach performs *much better* than the regular dual frame approach when used with the optimal mode selection. This result says that the method of creating or choosing a long-term reference frame (pulsed quality or regular quality) and the method of choosing, for each macroblock, whether or not to use that long-term reference frame (or use the short-term or intra mode) can work together synergistically or can oppose each other. The gains in performance ranged from 0.6 to 1.6dB.

The results point to the superiority of high quality (pulsed quality) over regular updating for lossy packet network video transmission with a dual frame coder. This gain comes at trivial extra computational and implementation cost and can be easily deployed in a standard compliant H.264 codec.

## VI. ACKNOWLEDGEMENTS

This work was supported in part by the National Science Foundation, the Office of Naval Research, and the CoRe program of the State of California.

## REFERENCES

- [1] T. Sikora, “The MPEG-4 video standard verification model,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 7, no. 1, pp. 19–31, Feb. 1997.
- [2] G. Côté, B. Erol, M. Gallant, and F. Kossentini, “H.263+: Video coding at low bit rates,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 8, no. 7, pp. 849–865, Nov. 1998.
- [3] M. Gothe and J. Vaisey, “Improving motion compensation using multiple temporal frames,” in *Proc. IEEE Pacific Rim Conference on Communications, Computers and Signal Processing*, vol. 1, May 1993, pp. 157–160.

- [4] N. Vasconcelos and A. Lippman, “Library-based image coding,” in *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. v, Apr. 1994, pp. 489–492.
- [5] T. Wiegand, X. Zhang, and B. Girod, “Long-term memory motion-compensated prediction,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 9, no. 1, pp. 70–84, Feb. 1999.
- [6] T. Wiegand, “Final draft international standard for joint video specification H.264,” JVT of ISO/IEC MPEG and ITU-T VCEG, JVT-G050, Mar. 2003.
- [7] T. Fukuhara, K. Asai, and T. Murakami, “Very low bit-rate video coding with block partitioning and adaptive selection of two time-differential frame memories,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 7, no. 1, pp. 212–220, Feb. 1997.
- [8] C.-S. Kim, R.-C. Kim, and S.-U. Lee, “Robust transmission of video sequence using double-vector motion compensation,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 11, no. 9, pp. 1011–1021, Sept. 2001.
- [9] M. Budagavi and J. D. Gibson, “Multiframe video coding for improved performance over wireless channels,” *IEEE Trans. Image Processing*, vol. 10, no. 2, pp. 252–265, Feb. 2001.
- [10] R. Zhang, S. L. Regunathan, and K. Rose, “Video coding with optimal inter/intra-mode switching for packet loss resilience,” *IEEE J. Select. Areas Commun.*, vol. 18, no. 6, pp. 966–976, June 2000.
- [11] A. Leontaris and P. C. Cosman, “Video compression with intra/inter mode switching and a dual frame buffer,” in *Proc. IEEE Data Compression Conference*, Mar. 2003, pp. 63–72.
- [12] V. Chellappa, P. C. Cosman, and G. M. Voelker, “Dual frame motion compensation with uneven quality assignment,” in *Proc. IEEE Data Compression Conference*, Mar. 2004.
- [13] A. Leontaris and P. C. Cosman, “Video compression for lossy packet networks with mode switching and a dual-frame buffer,” *IEEE Trans. Image Processing*, 2004.

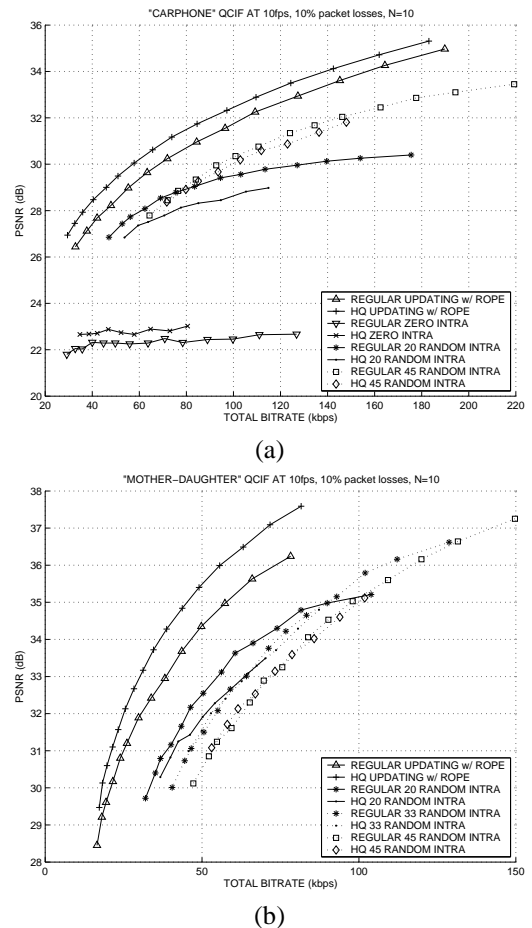


Fig. 3. Packet loss ratio  $p = 10\%$ . (a) Image Sequence “carphone”. (b) Image Sequence “mother-daughter”.