

# SPATIAL SCALABILITY AND COMPRESSION EFFICIENCY WITHIN A FLEXIBLE MOTION COMPENSATED 3D-DWT

*Nagita Mehrseresht and David Taubman*

The University of New South Wales, Sydney, Australia

## ABSTRACT

We investigate the implications of the conventional “t+2D” MC 3D-DWT structure for spatial scalability, and propose a more flexible “2D+t+2D” structure. An initial  $P$  levels of spatial wavelet decomposition are followed by  $T$  levels of motion compensated temporal decomposition, applied separately to each spatial resolution level. A further  $S - P$  levels of spatial decomposition are applied to the resulting subbands. By adjusting  $P$ , the structure allows us to trade energy compaction with the potential for artifacts at reduced spatial resolutions. This allows us to study the interaction between scalability and compression efficiency. We show that the “t+2D” structure ( $P = 0$ ) necessarily maximizes compression efficiency, while allowing misaligned spatial aliasing artifacts to arise at reduced resolutions. These artifacts can be removed by increasing the value of  $P$ , at an inevitable cost in compression efficiency. We show how this cost can be minimized.

## 1. INTRODUCTION

Recently a large amount of effort has been invested in wavelet-based fully scalable video coding. Interactive multimedia, video surveillance and video delivery over heterogeneous networks and/or in error prone environments can be named as some of the applications which stand to benefit from this research. The predictive feedback paradigm inherent in traditional video compression algorithms is incompatible with the requirement of highly scalable compression. Instead, the preferable paradigm is that of feed-forward compression in which a spatio-temporal (3D) discrete wavelet transform (DWT) is followed by embedded quantization and coding. To exploit interframe redundancy, the transform must compensate for the motion between frames.

Without motion compensation, the 3D-DWT can be formed by separable extension of the 1D-DWT to the spatial and temporal dimensions. In this case, spatial and temporal DWT stages can be performed in any desired order, without altering the final subbands. With the introduction of motion compensation, however, this commutative property is lost, principally because the spatial wavelet transform is not shift invariant. By changing the order of the spatial and motion compensated temporal DWT (MC TDWT) operators, a family of different MC 3D-DWTs may be formed.

Early work on scalable MC 3D-DWT followed the so-called “t+2D” paradigm, in which MC TDWT is applied to the full resolution video frames, followed by a spatial DWT (SDWT) of the temporal subbands [1][2]. New schemes proposed in [3][4] also fall into this category. In [5][6] an alternate, “2D+t” approach may be found, in which MC TDWT is applied to the spatial subbands generated by spatial decomposition of the original video frames.

Until this time, very little effort has been invested in directly comparing different structures for MC 3D-DWT. In this paper

we propose a flexible “2D+t+2D” structure for MC 3D-DWT, in which any desired number of spatial decomposition levels are performed on the original full resolution video, followed by MC TDWT of the spatial subbands. Further levels of spatial decomposition may be performed on the resulting subbands to provide additional levels of spatial scalability and energy compaction. Using this structure, we can flexibly choose the particular MC 3D-DWT which is most suited to the desired application. Perhaps more significant than the structure itself is the insight which it provides into the interaction between spatial aliasing, scalability and energy compaction. We first suggested the “2D+t+2D” structure in [7], along with some preliminary results. The present paper is concerned with a more careful investigation into the interaction between spatial scalability and compression efficiency.

## 2. MOTION COMPENSATED TEMPORAL LIFTING

We build upon the LIMAT [4] framework for MC TDWT, using the deformable mesh motion model in connection with the bi-orthogonal 5/3 kernel, due to its superior performance for temporal filtering. However, our approach is suitable for other motion models and wavelet kernels. Using this framework, MC TDWT is accomplished through a sequence of temporal lifting steps. To exploit inter-frame dependence, we compensate for motion *inside* each lifting step [4]. Let  $\mathcal{W}_{k_1, k_2}(f_{k_1})$  denote a motion compensated mapping of frame  $f_{k_1}$  onto the coordinate system of frame  $f_{k_2}$ . Using this notation, we can implement the MC lifting steps for the 5/3 analysis filters by

$$h_k = f_{2k+1} - \frac{1}{2}[\mathcal{W}_{2k, 2k+1}(f_{2k}) + \mathcal{W}_{2k+2, 2k+1}(f_{2k+2})] \quad (1)$$

$$l_k = f_{2k} + \frac{1}{4}[\mathcal{W}_{2k-1, 2k}(h_{k-1}) + \mathcal{W}_{2k+1, 2k}(h_k)]. \quad (2)$$

Equations (1) and (2) are commonly known as “*prediction*” and “*update*” steps, respectively.

Regardless of what motion model and interpolation operators are used for  $\mathcal{W}$ , the temporal transform can be trivially inverted by reversing the order of the lifting steps and replacing addition with subtraction, as indicated by equations (3) and (4)

$$f_{2k} = l_k - \frac{1}{4}[\mathcal{W}_{2k-1, 2k}(h_{k-1}) + \mathcal{W}_{2k+1, 2k}(h_k)] \quad (3)$$

$$f_{2k+1} = h_k + \frac{1}{2}[\mathcal{W}_{2k, 2k+1}(f_{2k}) + \mathcal{W}_{2k+2, 2k+1}(f_{2k+2})]. \quad (4)$$

## 3. FLEXIBLE “2D+T+2D” STRUCTURE

The purpose of MC 3D-DWT analysis is to generate a spatio-temporal multi-resolution representation for the original video sequence. In the sequel, we consistently use the term “*resolution*



first instance by discarding up to the first  $P$  spatial resolution levels. Lower spatial resolutions may be obtained (when  $S > P$ ), in a similar manner to the way in which the spatial scalability is achieved with the “t+2D” structure. Thus, for resolution reduction factors up to  $2^P$ , the low resolution video sequences are identical to those obtained by extracting reduced spatial resolutions from the original video frames. For larger reduction factors, however, *non-aligned spatial aliasing artifacts* can appear within the spatially scaled video, in regions where the motion model fails.

To understand the origin and nature of these non-aligned spatial aliasing artifacts, consider the case  $P = 0$  and suppose we reconstruct the video at half its original spatial resolution, writing  $f'_k$  for the resulting video frames. Ignoring the typically small differences between  $f'_{2k}$  and  $\mathcal{A}_L(f_{2k})$  (there is no difference if update steps are skipped), we obtain

$$\begin{aligned} f'_{2k+1} &\approx \mathcal{A}_L(h_k) + \frac{1}{2} \mathcal{W}_{2k,2k+1}^1 \circ \mathcal{A}_L(f_{2k}) + \dots \\ &\approx \mathcal{A}_L(h_k) + \frac{1}{2} \mathcal{A}_L \circ \mathcal{W}_{2k,2k+1} \circ \mathcal{S}(\mathcal{A}_L(f_{2k}), 0) + \dots \\ &= \mathcal{A}_L(f_{2k+1}) - \frac{1}{2} \mathcal{A}_L \circ \mathcal{W}_{2k,2k+1} [f_{2k} - \mathcal{S}(\mathcal{A}_L(f_{2k}), 0)] - \dots \\ &= \mathcal{A}_L(f_{2k+1}) - \frac{1}{2} \mathcal{A}_L \circ \mathcal{W}_{2k,2k+1} \circ \mathcal{S}(0, \mathcal{A}_H(f_{2k})) - \dots \end{aligned}$$

In the above, we use  $\mathcal{A}_L$  to denote one level of spatial wavelet analysis, retaining only the LL subband, and we write “...” to represent contributions from  $f_{2k+2}$ ; these are analogous to those from  $f_{2k}$ . Evidently, the difference between  $f'_{2k+1}$  and  $\mathcal{A}_L(f_{2k+1})$  arises from the spatial aliasing components which remain after synthesizing the high-pass subbands of  $f_{2k}$  (resp.  $f_{2k+2}$ ), motion compensating them, and taking their low-pass subbands.

The appearance of spatial aliasing terms is not surprising. Any form of spatial resolution reduction, including that associated with  $\mathcal{A}_L(f_{2k+1})$ , incurs some aliasing. In fact, subject to accurate motion modeling, it can be shown [9] that  $f'_{2k+1}$  typically has even less total aliasing power than  $\mathcal{A}_L(f_{2k+1})$ ! Wherever the motion model fails, however, the aliasing terms produced by  $\mathcal{A}_L \circ \mathcal{W}_{2k,2k+1} \circ \mathcal{S}(0, \mathcal{A}_H(f_{2k}))$  do not line up with the scene features (typically edges) which generated them. The resulting, non-aligned aliasing artifacts, can be visually disturbing.

## 5. LEAKAGE COMPENSATION AND EFFICIENCY

In the special case of ideal bandlimited spatial subband filters, and pure translational motion, the synthesized baseband signal  $\mathcal{S}(0, f_{u,p})$  in equation (5) is perfectly bandlimited. This, coupled with the fact that  $\mathcal{W}_{u,v}^{p-1}$  is an LTI operator, means that no information is lost when  $\mathcal{A}_H$  is applied to  $\mathcal{W}_{u,v}^{p-1} \circ \mathcal{S}(0, f_{u,p})$ . In this special case, it is not hard to show [9] that the temporal subbands produced using equations (6) and (7) are identical to those produced by the “t+2D” structure.

In practice, the subband filters are not perfectly bandlimited and real motion can be locally expansive or contractive. These factors generally result in the loss (or leakage) of information during the implementation of  $\widehat{\mathcal{W}}_{u,v}^p$ . Of course, this does not affect the invertibility of the complete transform, which is trivially obtained by replacing  $\mathcal{W}_{u,v}$  with  $\widehat{\mathcal{W}}_{u,v}^p$  in equations (3) and (4). It does, however, reduce the energy compaction and hence compression efficiency associated with the motion compensated temporal transform.

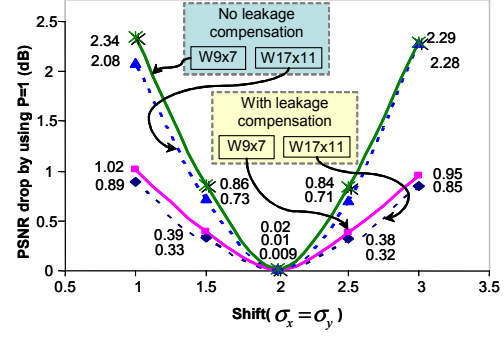


Fig. 2. Loss in compression performance when  $P = 1$  relative to  $P = 0$ , using a synthetic video sequence with ideal motion.

To see this, consider the prediction step in equation (6), with  $t = 1$  so that  $\mathbf{1}_p^{t-1} = \mathbf{f}_p$ . Assuming a perfect motion model,  $f_{2k+1,p}$  would be perfectly compensated if we could motion compensate the original high resolution frames  $f_{2k}$  and  $f_{2k+2}$ , finding spatial resolution  $p$  of the result. Instead of motion compensating the original frames, however, we are only using the information contained in  $f_{2k,p}$  and  $f_{2k+2,p}$ . In so doing, we are unable to exploit the motion induced leakage of information from lower frequency subbands of  $\mathbf{f}$  to  $\mathbf{f}_p$  (call this “type I leakage”), and we are unable to exploit the motion induced leakage of information from higher frequency subbands of  $\mathbf{f}$  to  $\mathbf{f}_p$  (call this “type II leakage”). As a result, the energy of the high-pass temporal subbands  $\mathbf{h}_p^t$  is generally higher than in the case of the “t+2D” structure, leading to reduced compression efficiency.

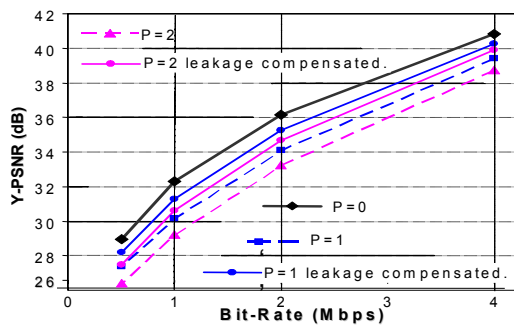
Fig. 2 provides an experimental demonstration of these effects, in the context of a synthetic video sequence created by cropping a sequence of shifted CIF resolution windows from a much larger image. The horizontal axis in the figure represents the amount of shift between successive frames. The vertical axis measures the drop in compression efficiency (measured in PSNR) of a complete compression system (see Section 6) when going from  $P = 0$  to  $P = 1$  level of pre-temporal spatial decomposition. The PSNR values correspond to 1Mbps compressed bit-rate. As expected, there are no leakage effects when the shift is an even number of pixels, since the spatial subband transform is periodically shift invariant with period 2. Leakage effects are greatest when the shift is an odd number of pixels. Also as expected, the use of longer subband filters (17x11 vs. 9x7), with sharper frequency responses, can reduce the aliasing leakage.

The lower curves in Fig. 2 are obtained by compensating for type I leakage effects. We cannot compensate for type II leakage effects without using information from higher frequency spatial subbands in the MC TDWT of  $\mathbf{f}_p$ . Such information is not available when inverting the temporal transform at a reduced spatial resolution. We can, however, borrow information from lower spatial resolution levels to compensate for type I leakage. This does not sacrifice the desirable property that spatial scalability, with reduction factors up to  $2^P$ , produces video sequences identical to those obtained by reducing the spatial resolution of the original video frames.

To compensate for type I leakage, we modify equation (5) to

$$\widehat{\mathcal{W}}_{u,v}^p(l_{u,p}^t) = \mathcal{A}_H \circ \mathcal{W}_{u,v}^{p-1} \circ \mathcal{S}(l_u^t \widehat{LL}_p, l_{u,p}^t) \quad (8)$$

Here,  $l_u^t \widehat{LL}_p$  may be understood as an estimate of the missing low-



**Fig. 3.** Luminance PSNR of the “Stefan” sequence at full resolution using  $T = 3$ ,  $S = 5$  and  $P = \{0, 1, 2\}$ .

pass subband. The idea is to emulate the behaviour of applying the baseband MC TDWT directly to the LL video sequence produced by a  $(p - 1)$ -level SDWT, followed by application of the  $\mathcal{A}_H$  operator. It is possible to do this perfectly for  $t = 1$ . More generally, however, we must restrict our estimate of  $\widehat{l_u^t LL_p}$  to use information only from the corresponding temporal resolution, at lower spatial resolutions. We do this by synthesizing the lower resolution subbands available from corresponding lifting steps at the same temporal resolution, i.e.,

$$\widehat{l_u^t LL_p} = \begin{cases} l_{u,P+1}^t & p = P \\ \mathcal{S}(l_{u,P}^t, l_{u,P+1}^t) & p = P - 1 \geq 1 \\ \mathcal{S}(l_{u,P-1}^t, \mathcal{S}(l_{u,P}^t, l_{u,P+1}^t)) & p = P - 2 \geq 1 \\ \dots & \dots \end{cases}$$

## 6. EXPERIMENTS WITH REAL VIDEO SEQUENCES

In this section, we compare the compression efficiency of the MC 3D-DWT when different values of  $P$  are used. For an SNR (quality/bit-rate) scalable video compression scheme, the spatio-temporal subbands must further be subjected to embedded quantization and coding. This allows the quality of the overall video sequence to be adjusted at each spatio-temporal resolution by discarding appropriate subsets from the embedded bitstream. We use the EBCOT algorithm [10] embodied within the JPEG2000 image compression standard to create an efficient embedded codestream for spatio-temporal subbands.

The standard CIF resolution test sequence “Stefan” is compressed using MC 3D-DWT with  $T = 3$  levels of MC TDWT and  $S = 5$  levels of 9/7 SDWT. Fig. 3 illustrates the luminance (Y) PSNR of the reconstructed full resolution sequence at different bit-rates, when different numbers of pre-temporal spatial decomposition levels  $P$  are used. This figure confirms the results obtained in Fig. 2, although that figure was obtained with a simulated source and ideal motion, having abnormally high energy compaction which tends to exaggerate relative losses due to leakage. As before, we see that  $P = 0$  (“t+2D” structure) yields the highest compression efficiency. Moreover, the efficiency continues to drop as  $P$  is increased. Fig. 3 also confirms the benefits of leakage compensation.

Our choice of the “Stefan” sequence here is motivated by the fact that it exhibits complex motion, with some motion model failure. In this case,  $P > 0$  is required if non-aligned aliasing artifacts are to be avoided at reduced spatial resolutions.

## 7. CONCLUSIONS

In this paper we proposed a flexible structure for MC 3D-DWT, which enables us to trade energy compaction with the potential for artifacts at reduced spatial resolutions by choosing the number of pre-temporal spatial decomposition levels  $P$ . We found that compression performance of the MC 3D-DWT is maximized by using the “t+2D” structure ( $P = 0$ ). On the other hand, we also showed that such a structure can lead to visually annoying artifacts, in places where the motion model fails. Using the flexible structure, we can avoid these artifacts by selecting an appropriate value for  $P$ . We found that the performance loss associated with  $P > 0$  is mainly caused by spatial frequency leakage during motion compensation. We showed that this phenomenon can be reduced by leakage compensation and, to a certain extent, by the selection of subband filters with better frequency roll-off performance. Since space is limited, in this paper we report results only for the “Stefan” sequence; however, we have tested the proposed structure with a wide range of synthetic and standard CIF and SIF resolution video sequences.

## 8. REFERENCES

- [1] K. Ohm, “Three-dimensional subband coding with motion compensation,” *IEEE Trans. Image Proc.*, vol. 3, no. 5, pp. 559–571, Sept. 1994.
- [2] S. Choi and J. Woods, “Motion compensated 3d subband coding of video,” *IEEE Trans. Image Proc.*, vol. 8, pp. 155–167, Feb 1999.
- [3] V. Bottreau, M. Benetiere, B. Felts, and B. Pesquet-Popescu, “A fully SCalable 3d subband video codec,” *IEEE Int. Conf. Image Proc.*, pp. 1017–1020, 2001.
- [4] A. Secker and D. Taubman, “Lifting based invertible motion adaptive transform, LIMAT, framework for highly scalable video compression,” *IEEE Trans. Image Proc.*, vol. 12, pp. 1530–1542, Dec. 2003.
- [5] J. Woods and G. Linlienfield, “A resolution and frame-rate scalable Subband/Wavelet video coder,” *IEEE Tran. Circ. and Syst. for Video Tech.*, vol. 11, no. 9, pp. 1035–1044, Sept 2001.
- [6] Y. Andreopoulos, M. Schaar, A. Munteanu, J. Barbarien, P. Schelkens, and J. Cornelis, “Complete-to-overcomplete discrete wavelet transforms for scalable video coding with MCTF,” *SPIE Visual Comm. Image Proc.*, pp. 719–731, 2003.
- [7] D. Taubman, “Successive refinement of video: Fundamental issues, past efforts and new directions,” *SPIE Visual Comm. and Image Prog.*, pp. 649–663, 2003.
- [8] N. Mehrseresht and D. Taubman, “Adaptively weighted update steps in motion compensated lifting based scalable video compression,” *IEEE Int. Conf. Image Proc.*, pp. 771–774, 2003.
- [9] —, “A flexible structure for fully scalable motion compensated 3d-DWT with emphasis on the impact of spatial scalability,” *submitted to IEEE Tran. Image Proc.*, 2004.
- [10] D. Taubman, “High performance scalable image compression with EBCOT,” *IEEE Trans. Image Proc.*, vol. 9, no. 7, pp. 3445–3462, July 2000.