

# FUSION OF GLOBAL AND LOCAL SIDE INFORMATION USING SUPPORT VECTOR MACHINE IN TRANSFORM-DOMAIN DISTRIBUTED VIDEO CODING

A. Abou-Elailah<sup>1</sup>, F. Dufaux<sup>1</sup>, J. Farah<sup>2</sup>, M. Cagnazzo<sup>1</sup>, and B. Pesquet-Popescu<sup>1</sup>

<sup>1</sup> Signal and Image processing Department, Institut TELECOM - TELECOM Paristech,  
46 rue Barrault, F - 75634 Paris Cedex 13, FRANCE

{elailah, frederic.dufaux, marco.cagnazzo, beatrice.pesquet}@telecom-paristech.fr

<sup>2</sup> Telecommunications Department, Faculty of Engineering, Holy-Spirit University of Kaslik  
P.O. Box 446, Jounieh, Lebanon  
joumanafarah@usek.edu.lb

## ABSTRACT

Side information has a strong impact on the performance of Distributed Video Coding. Commonly, side information is generated using motion compensated temporal interpolation. In this paper, we propose a new method for the fusion of local and global side information using Support Vector Machine. The global side information is generated at the decoder using global motion parameters estimated at the encoder using Scale-Invariant Feature Transform. Experimental results show that the proposed approach can achieve a PSNR improvement of up to 1.7 dB for a GOP size of 2 and up to 3.78 dB for larger GOP sizes, with respect to the reference DISCOVER codec.

**Index Terms**— Distributed Video Coding, Support Vector Machine, Classification, Side Information, Rate-Distortion Performance

## 1. INTRODUCTION

Distributed Video Coding (DVC) is a paradigm especially fitted for emerging applications such as wireless video surveillance, multimedia sensor networks, wireless PC cameras, and mobile cameras phones. In the video coding standards that are produced by MPEG and ITU-T, motion estimation and compensation are performed at the encoder. In contrast DVC, the correlation among successive frames is exploited at the decoder, while allowing a low encoding complexity. In other words, DVC allows shifting the encoder complexity to the decoder.

From information theory, the Slepian-Wolf theorem for lossless compression [1] states that it is possible to encode correlated sources (let us call them X and Y) independently and decode them jointly, while achieving the same rate bounds which can be attained in the case of joint encoding and decoding. The Wyner-Ziv (WZ) theorem [2] extends the Slepian-Wolf one to the case of lossy compression of X

when Side Information (SI) Y is available at the decoder.

Based on these theoretical results, practical implementations of DVC have been proposed [3, 4]. DISCOVER codec [5, 6] is based on transform domain WZ coding, is one of the most efficient and popular existing architectures. In this codec, the images of the sequence are split into two sets of frames, key frames (KFs) and Wyner-Ziv frames (WZFs). The Group of Pictures (GOP) of size  $n$  is defined as a set of frames consisting of one KF and  $n - 1$  WZFs. The KFs are independently encoded and decoded using Intra coding techniques such as H.264/AVC Intra mode. The WZFs are separately transformed using a  $4 \times 4$  integer Discrete Cosine Transform (DCT). The obtained coefficients are uniformly quantized. A systematic channel code such as Turbo code or Low-Density Parity Check Accumulate (LDPCA) code is applied on the resulting quantized coefficients. Only the parity bits are kept, and sent to the decoder while the systematic bits are discarded.

At the decoder, the reconstructed reference frames are used to compute the SI, which is an estimation of the WZF being decoded. This estimation can be seen as a noisy version of the original WZF. Motion-Compensated Temporal Interpolation (MCTI) [7] is used to generate the SI in the DISCOVER codec. The channel decoder corrects the DCT coefficients of the SI using the parity bits requested by the decoder through the feedback channel. Finally, reconstruction and inverse  $4 \times 4$  integer DCT are applied to obtain the decoded WZF.

In this paper, we propose a new fusion method to combine two SI using Support Vector Machine (SVM). The first SI is generated using MCTI as in DISCOVER codec and is referred to as MCTI SI. The second one is generated by applying global motion parameters on the decoded reference frames [8], and is referred to as Global Motion Compensation SI (GMC SI). In this context, the objective is to optionally fuse MCTI SI and GMC SI to reach the best Rate-Distortion (RD) performance. For this purpose, both SI are considered as two classes, and SVM classifier is applied on a block basis



decoder side. In this paper, at the border of the image, GMC SI is always taken for the estimated SI in all methods.

In Multi-view DVC, two SI are usually generated. The first SI ( $SI_t$ ) is generated from previously decoded frames in the same view, while the second one ( $SI_v$ ) is estimated using previously decoded frames in adjacent views. The authors in [13] proposed new techniques for fusion  $SI_t$  and  $SI_v$ . Inspired from [13], a linear fusion of GMC SI and MCTI SI is proposed as follows:

$$SI(b) = \frac{f_{MCTI} \cdot (GMC SI) + f_{GMC} \cdot (MCTI SI)}{(f_{GMC} + f_{MCTI})} \quad (3)$$

This method is referred to as ‘FusLin’. Dufaux [14] proposed a solution that consists in combining  $SI_t$  and  $SI_v$  using SVM.

### 3. PROPOSED METHOD

The block diagram of our proposed codec architecture is depicted in Figure 1. It is based on the DISCOVER codec [5][6]. The shaded (green) blocks correspond to the three new modules introduced in this paper: Model, Classification, and generating SVM SI.

The block for the SI can be predicted from GMC SI or MCTI SI using SVM classifier. In this paper, we use SVM<sup>Light</sup> software implementation [15]. The training stage to generate the model is described with the classification procedure in Subsection 3.1. Finally, the proposed methods for the combination of GMC SI and MCTI SI based on the predicted value by the SVM classifier is described in Subsection 3.2.

#### 3.1. Model and Classification

First, we select the most discriminative features to be used in SVM. For this reason, three features are estimated in the proposed method as follows:

$$\begin{cases} f_1 = f_{GMC} \\ f_2 = f_{MCTI} \\ f_3 = f_{GMC} - f_{MCTI} \end{cases} \quad (4)$$

In the training stage, the first WZF is encoded using H.264/AVC Intra mode as the KFs. This frame is used to build the model for SVM. More precisely, let  $D_{GMC}$  and  $D_{MCTI}$  be the difference between the first WZF and the GMC SI and MCTI SI for the current block respectively. A block belongs to GMC SI if  $D_{GMC}$  is smaller than  $D_{MCTI}$ , and belongs to MCTI SI otherwise. Only the blocks ( $N$  blocks) which give the largest difference  $D$  ( $D = |D_{GMC} - D_{MCTI}|$ ) are taken in the training stage. This step consists in increasing accuracy and precision of the training.

Next, classification procedure is carried out on the first WZF using this model. The blocks which are well-classified are taken into account for a second learning stage, in order

to produce the final model (find the hyperplane that optimally separates the blocks of GMC SI and MCTI SI). This model will then be used in the classification procedure for all WZFs in the sequence.

In the classification, three features  $f_1$ ,  $f_2$ , and  $f_3$  are computed for each WZF using GMC SI and MCTI SI. The SVM classifier computes a predicted value for each block based on the features and the obtained model.

#### 3.2. Proposed fusion

The SVM classifier gives a decision value  $d$  for each block.  $d$  represents the distance between this block and the separating hyperplane. Based on this value, we define two fusion algorithms. The first algorithm consists of binary combination of GMC SI and MCTI SI. The second algorithm linearly combines the two SI.

**SVM binary fusion** - In this method, the value  $d$  is directly used to combine the two SI as follows:

$$SI(b) = \begin{cases} GMC SI & \text{if } d > 0 \\ MCTI SI & \text{otherwise} \end{cases} \quad (5)$$

where  $d$  represents the classification label at block  $b$ . This method is referred to as ‘SVM’.

**SVM linear fusion** - This method aims at combining linearly GMC SI and MCTI SI. The linear combination is defined as follows:

$$SI(b) = \begin{cases} GMC SI & \text{if } d > T \\ MCTI SI & \text{if } d < (-T) \\ \frac{(T+d) \cdot GMC SI + (T-d) \cdot MCTI SI}{2 \cdot T} & \text{if } |d| < T \end{cases} \quad (6)$$

where  $T$  represents a threshold. This method is referred to as ‘SVMLin’.

**Oracle fusion** - This method is impractical, but it aims at estimating the upper bound limit that can be achieved by combining GMC SI and MCTI SI, using the original WZF. This fusion is defined as follows:

$$SI(b) = \begin{cases} GMC SI & \text{if } D_{GMC} < D_{MCTI} \\ MCTI SI & \text{otherwise} \end{cases} \quad (7)$$

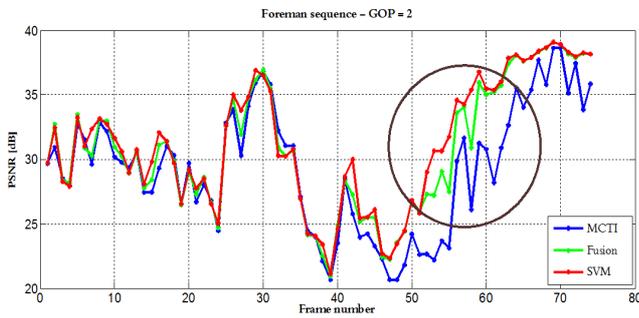
This method is referred to as ‘Oracle’.

## 4. EXPERIMENTAL RESULTS

In order to evaluate the performance of the proposed methods, we performed extensive simulations, adopting the same test conditions as described in DISCOVER [5, 6], *i.e.* test video sequences are at QCIF spatial resolution and sampled at 15 frames/sec. The obtained results of the proposed methods SVM (Eq. 5) and SVMLin (Eq. 6) are compared to the DISCOVER codec, to the binary fusion (Eq. 2), to the linear fusion (Eq. 3), and to ‘Oracle’ fusion (Eq. 7).

SI Average PSNR [dB]							
Method	MCTI	GMC	Fusion	FusLin	SVM	SVMLin	Oracle
<b>GOP = 2</b>							
Stefan	22.57	25.88	26.27	26.19	26.45	<b>26.54</b>	27.21
Foreman	29.31	30.70	30.77	30.97	31.21	<b>31.30</b>	31.90
Bus	24.72	22.99	26.96	26.83	26.92	<b>27.18</b>	27.94
Coastguard	31.43	29.28	32.02	31.95	32.11	<b>32.23</b>	32.62
<b>GOP = 4</b>							
Stefan	21.28	25.27	25.33	25.23	25.59	<b>25.66</b>	26.47
Foreman	27.58	29.62	29.24	29.47	29.77	<b>29.87</b>	30.72
Bus	23.48	22.41	25.93	25.88	25.91	<b>26.14</b>	26.91
Coastguard	29.85	28.19	30.78	30.76	30.90	<b>31.03</b>	31.46
<b>GOP = 8</b>							
Stefan	20.64	24.85	24.79	24.71	25.06	<b>25.15</b>	25.99
Foreman	26.24	28.62	28.08	28.30	28.68	<b>28.79</b>	29.69
Bus	22.53	21.84	24.95	24.95	24.95	<b>25.17</b>	25.90
Coastguard	28.75	27.50	29.85	29.87	29.97	<b>30.10</b>	30.60

**Table 1.** SI average PSNR for a GOP size equal to 2, 4, and 8 (QI = 8).



**Fig. 2.** PSNR of MCTI SI, Fusion, and the proposed method SVM for Foreman sequence for a GOP size of 2.

#### 4.1. SI performance assessment

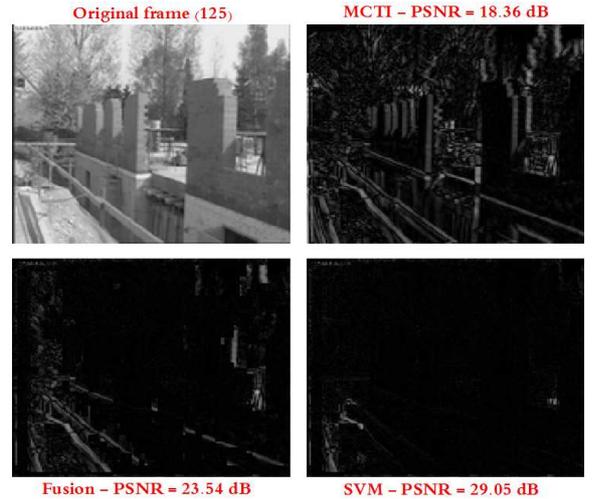
Figure 2 shows the SI PSNR for Foreman sequence, for a GOP size of 2. The proposed method (SVM) allows a consistent improvement, and achieves a gain up to 4.4 dB for some frames

Figure 3 shows the visual difference of the SI for Foreman (frame number 125), for a GOP size of 8. the SI obtained by MCTI technique is not good as shown in this figure (top-right - 18.36 dB). On the contrary, the SI obtained by the proposed method (SVM) is significantly better than the SI estimated by both MCTI and Fusion. The gain up to 10.7 dB compared to MCTI, and up to 5.5 dB compared to the previous fusion (Fusion method) for this frame.

Table 1 shows the average PSNR of the SI for the different methods, different sequences, and different GOP sizes. The proposed technique (SVMLin) leads to best SI quality for all test sequences.

#### 4.2. Rate Distortion Performance

The RD performance is shown for the Stefan, Foreman, Bus, and Coastguard sequences in Table 2, in comparison to the DISCOVER codec, using the Bjontegaard metric [16], for a



**Fig. 3.** Visual difference of the SI estimated by MCTI, Fusion, and the proposed method SVM for frame number 125 of Foreman sequence, for a GOP size of 8 (QI = 8).

GOP size equal to 2, 4 and 8.

The proposed method SVMLin always achieves a gain compared to the other methods for Foreman, Bus and Coastguard sequences, for all GOP sizes. For Stefan sequence, the proposed method SVM achieves the best gain for all GOP sizes.

It is clear that the performance of the proposed fusion becomes close to that of ‘Oracle’ fusion, for all test sequences. The difference between them is small than 0.5 dB for all GOP sizes.

The gains become even more significant for a GOP size equal to 8. In fact, for SVM, we obtain a bit reduction up to  $-52.46\%$ , which corresponds to an improvement of 3.78 dB on the decoded frames w.r.t. DISCOVER codec for Stefan sequence. For Foreman sequence, the proposed method SVM-Lin allows a gain of up to 2.01 dB, with a rate reduction of 34.20%, compared to the DISCOVER codec, while the previous method ‘Fusion’ allows a gain up to 1.26 dB, with a rate reduction of 22.77%, compared to the DISCOVER codec, for this sequence.

## 5. CONCLUSION

A new technique based on the SVM for the fusion of global and local SI is proposed in this paper. Experimental results show that our proposed method can achieve a gain in RD performance up to 1.7 dB for a GOP size of 2 and 3.78 dB for longer GOP sizes, compared to DISCOVER codec, especially when the video sequence contains high motion.

Method	GMC	Fusion	FusLin	SVM	SVMLin	Oracle
<b>GOP = 2</b>						
<b>Stefan</b>						
$\Delta_R$ [%]	-25.59	-24.49	-21.38	<b>-25.70</b>	-25.45	-27.43
$\Delta_{PSNR}$ [dB]	1.70	1.61	1.37	<b>1.70</b>	1.68	1.84
<b>Foreman</b>						
$\Delta_R$ [%]	-8.90	-7.90	-9.46	-11.31	<b>-12.02</b>	-14.30
$\Delta_{PSNR}$ [dB]	0.53	0.46	0.55	0.68	<b>0.72</b>	0.86
<b>Bus</b>						
$\Delta_R$ [%]	5.02	-13.42	-10.05	-13.05	<b>-14.09</b>	-17.09
$\Delta_{PSNR}$ [dB]	-0.25	0.80	0.59	0.79	<b>0.84</b>	1.03
<b>Coastguard</b>						
$\Delta_R$ [%]	9.97	-4.94	-3.71	-5.70	<b>-6.32</b>	-8.20
$\Delta_{PSNR}$ [dB]	-0.46	0.25	0.18	0.28	<b>0.31</b>	0.42
<b>GOP = 4</b>						
<b>Stefan</b>						
$\Delta_R$ [%]	-45.52	-43.12	-37.55	<b>-45.09</b>	-44.51	-47.81
$\Delta_{PSNR}$ [dB]	3.16	2.94	2.46	<b>3.13</b>	3.07	3.38
<b>Foreman</b>						
$\Delta_R$ [%]	-22.77	-16.03	-18.58	-23.58	<b>-24.61</b>	-29.85
$\Delta_{PSNR}$ [dB]	1.33	0.90	1.05	1.38	<b>1.43</b>	1.78
<b>Bus</b>						
$\Delta_R$ [%]	-2.74	-25.80	-21.74	-26.08	<b>-26.99</b>	-31.37
$\Delta_{PSNR}$ [dB]	0.16	1.52	1.26	1.54	<b>1.60</b>	1.90
<b>Coastguard</b>						
$\Delta_R$ [%]	6.64	-16.34	-14.43	-18.45	<b>-19.28</b>	-24.01
$\Delta_{PSNR}$ [dB]	-0.29	0.67	0.58	0.77	<b>0.81</b>	1.04
<b>GOP = 8</b>						
<b>Stefan</b>						
$\Delta_R$ [%]	-53.02	-50.35	-44.18	<b>-52.46</b>	-51.99	-55.90
$\Delta_{PSNR}$ [dB]	3.83	3.55	2.98	<b>3.78</b>	3.73	4.11
<b>Foreman</b>						
$\Delta_R$ [%]	-32.68	-22.77	-26.16	-32.82	<b>-34.20</b>	-39.86
$\Delta_{PSNR}$ [dB]	1.93	1.26	1.45	1.93	<b>2.01</b>	2.42
<b>Bus</b>						
$\Delta_R$ [%]	-11.49	-32.33	-28.55	-32.14	<b>-33.24</b>	-38.56
$\Delta_{PSNR}$ [dB]	0.58	1.88	1.62	1.89	<b>1.96</b>	2.34
<b>Coastguard</b>						
$\Delta_R$ [%]	-7.95	-28.14	-26.50	-31.64	<b>-32.45</b>	-39.02
$\Delta_{PSNR}$ [dB]	0.27	1.20	1.09	1.37	<b>1.41</b>	1.76

**Table 2.** Rate-distortion performance gain for *Stefan*, *Foreman*, *Bus*, and *Coastguard* sequences towards DISCOVER codec, using Bjontegaard metric, for a GOP size of 2, 4, and 8.

## 6. REFERENCES

- [1] J.D. Slepian and J.K. Wolf, "Noiseless coding of correlated information sources," *IEEE Transactions on Information Theory*, vol. IT-19, pp. 471–480, July 1973.
- [2] A. Wyner and J. Ziv, "The rate-distortion function for source coding with side information at the decoder," *IEEE Transactions on Information Theory*, vol. 22, pp. 1–10, July 1976.
- [3] R. Puri and K. Ramchandran, "PRISM: A video coding architecture based on distributed compression principles," *EECS Department, University of California, Berkeley, Tech. Rep. UCB/ERL M03/6*, 2003.
- [4] B. Girod, A. Aaron, S. Rane, and D. Rebello-Monedero, "Distributed video coding," *Proceedings of the IEEE*, vol. 93, pp. 71–83, Jan. 2005.
- [5] X. Artigas, J. Ascenso, M. Dalai, S. Klomp, D. Kubasov, and M.Ouaret, "The DISCOVER codec: Architecture, techniques and evaluation," in *Proc. of Picture Coding Symposium*, Lisboa, Portugal, Oct. 2007.
- [6] "Discover project," <http://www.discoverdvc.org/>.
- [7] C. Brites, J. Ascenso, and F. Pereira, "Improving transform domain Wyner-Ziv video coding performance," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2006, vol. 2, pp. 525–528.
- [8] A. Abou-Elailah, F. Dufaux, M. Cagnazzo, B. Pesquet-Popescu, and J. Farah, "Fusion of global and local motion estimation for distributed video coding," *IEEE Transactions on Circuits and Systems for Video Technology (Submitted)*.
- [9] J. Ascenso, C. Brites, F. Dufaux, A. Fernando, T. Ebrahimi, F. Pereira, and S. Tubaro, "The VISNET II DVC Codec: Architecture, Tools and Performance," in *Proc. of the 18th European Signal Processing Conference (EUSIPCO)*, 2010.
- [10] A. Abou-Elailah, J. Farah, M. Cagnazzo, B. Pesquet-Popescu, and F. Dufaux, "Improved side information for distributed video coding," in *3rd European Workshop on Visual Information Processing (EUVIP)*, Paris, France, July 2011, pp. 42 – 49.
- [11] T. Maugey, C. Yaacoub, J. Farah, M. Cagnazzo, and B. Pesquet-Popescu, "Side information enhancement using an adaptive hash-based genetic algorithm in a Wyner-Ziv context," in *IEEE International Workshop on Multimedia Signal Processing*, Saint-Malo, France, Oct. 2010, pp. 298 –302.
- [12] F. Dufaux and T. Ebrahimi, "Encoder and decoder side global and local motion estimation for distributed video coding," in *IEEE International Workshop on Multimedia Signal Processing (MMSP)*, 2010, pp. 339 – 344.
- [13] T. Maugey, W. Miled, M. Cagnazzo, and B. Pesquet-Popescu, "Fusion schemes for multiview distributed video coding," in *17th European Signal Processing Conference (EUSIPCO)*, Scotland, Aug. 2009.
- [14] F. Dufaux, "Support vector machine based fusion for multi-view distributed video coding," in *17th International Conference on Digital Signal Processing (DSP)*, Corfu, Aug. 2011, pp. 1 –7.
- [15] "SVM implementation," [http://www.cs.cornell.edu/People/tj/svm\\_light/](http://www.cs.cornell.edu/People/tj/svm_light/).
- [16] G. Bjontegaard, "Calculation of average PSNR differences between RD-curves," in *VCEG Meeting*, Austin, USA, Apr. 2001.