

MODEL-AIDED CODING: USING 3-D SCENE MODELS IN MOTION-COMPENSATED VIDEO CODING

Peter Eisert, Thomas Wiegand

Telecommunications Laboratory
University of Erlangen-Nuremberg
{eisert,wiegand}@LNT.de

Bernd Girod

Information Systems Laboratory
Stanford University
girod@ee.stanford.edu

ABSTRACT

We show that traditional waveform-coding and 3-D model-based coding are not competing alternatives but should be combined to support and complement each other. Both approaches are combined such that the generality of waveform coding and the efficiency of 3-D model-based coding are available where needed. The combination is achieved by providing the block-based video coder with a second reference frame for prediction which is synthesized by the model-based coder. Since the coding gain of this approach is directly related to the quality of the synthetic frame, we have extended the model-aided coder [1] to cope with illumination changes and multiple objects. Remaining model failures and objects that are not known at the decoder are handled by standard block-based motion-compensated prediction. Experimental results show that bit-rate savings of up to 45 % are achieved at equal average PSNR when comparing the model-aided codec to TMN-10, the test model of the H.263 standard.

1. INTRODUCTION

In recent years, several video coding standards such as H.261, H.263 [2], MPEG-1, and MPEG-2 have been introduced, which mainly address the compression of generic video data for digital storage and communication services. These schemes are designed on the basis of the statistics of the video signal without knowledge of the semantic content and can therefore robustly be used for arbitrary scenes. The design of model-based codecs [3] is based on the semantics of the scene. Hence, if the semantic information of the scene can be exploited, higher coding efficiency may be achieved by model-based video codecs. Such a 3-D model-based codec is restricted to scenes that can be composed of objects that are known by the decoder. One typical class of scenes are head-and-shoulder sequences which can be frequently found in applications such as video-telephone or video-conferencing systems. For head-and-shoulder scenes, bit-rates of about 1 kbit/s with acceptable quality can be achieved [4]. This has also motivated the recently determined *Synthetic and Natural Hybrid Coding* (SNHC) part of the MPEG-4 standard [5].

The combination of traditional hybrid video coding methods with model-based coding for higher coding efficiency has been proposed by several researchers [6, 7]. In these approaches, the mode decision is done for an entire frame and therefore the information from the 3-D model cannot be exploited if parts of the frame cannot be described by the model-based coder.

In [1] we have presented an extension of an H.263 video coder [2] that utilizes information from a model-based coder. Instead

of exclusively predicting the current frame of the video sequence from the previous decoded frame, prediction from the synthetic frame of the model-based coder is additionally allowed. The *model-aided coder* decides which prediction is efficient in terms of rate-distortion performance. Hence, the coding efficiency does not decrease below H.263 in the case the model-based coder cannot describe the current scene.

In this paper, we extend the model-aided coder [1] to exploit the efficiency of the model-based coder also for more sophisticated video sequences with changing lighting conditions or multiple objects. Parameters describing the illumination in the scene are estimated together with motion and deformation of the objects resulting in more accurate model frames. Experimental results demonstrate that the improved rate-distortion performance of the model-aided codec can also be measured for head-and-shoulder sequences with multiple objects and varying illumination.

2. VIDEO CODING ARCHITECTURE

Figure 1 shows the architecture of the proposed model-aided video coder (MAC). This figure depicts the well-known hybrid video

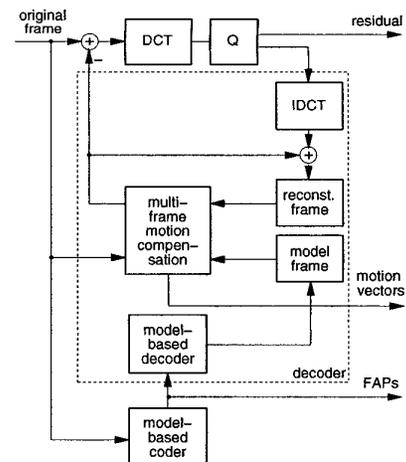


Fig. 1. Structure of the model-aided video coder. Traditional block-based MCP from the previous decoded frame is extended by prediction from the current model frame.

coding loop that is extended by a model-based coder. The model-based coder is running simultaneously to the hybrid coder, generating a synthetic model frame. This model frame is employed as a

second reference for block-based motion-compensated prediction (MCP) in addition to the previous reconstructed reference frame. For each macroblock, the video coder decides which of the two frames to use for MCP. The bit-rate reduction for the proposed scheme arises from those parts in the image that are well approximated by the model frame. For these blocks, the bit-rate required for transmission of the motion vector and DCT-coefficients for the residual coding is often highly reduced. For more details about the rate-distortion optimized mode decision and the changes made to the H.263+ syntax, see [1].

3. MODEL-BASED CODEC

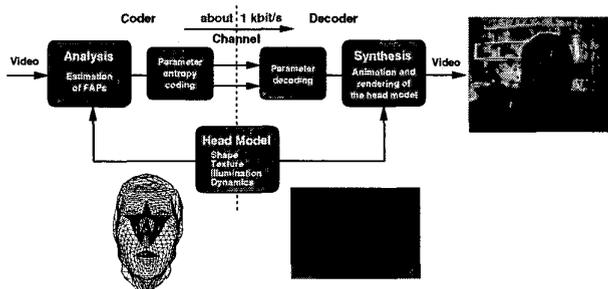


Fig. 2. Basic structure of the model-based codec.

The structure of the model-based codec is depicted in Fig. 2. The encoder analyzes the incoming frames and estimates the parameters of 3-D motion and deformation for all objects in the scene. The deformations for the head model are represented by a set of facial animation parameters (FAPs) according to the MPEG-4 standard [5]. Motion and deformation of other objects are parameterized similarly. All parameters are entropy-encoded and transmitted through the channel. The information from the 3-D models and the facial expression synthesis are incorporated into the parameter estimation. The 3-D models describe the shape, texture, and the motion constraints of the objects. For synthesis of facial expressions, the transmitted FAPs are used to deform the 3-D head model. The other objects are similarly moved and deformed in the virtual scene. Finally, individual video frames are approximated by simply rendering the 3-D scene.

In our model-based coder all parameters are estimated simultaneously using a hierarchical optical flow based method [4]. In the optimization, an analysis-synthesis loop is employed. The mean squared error between the rendered scene and the current video frame is minimized by estimating changes of the FAPs and the parameters for the other objects. To simplify the optimization in the high-dimensional parameter space, a linearized solution is directly computed using information from the optical flow and motion constraints from the models. For more details about the parameter estimation and the generation of model frames, please refer to [4, 8].

3.1. Illumination Compensation

The model-aided coder presented in [1] is not capable of representing lighting changes correctly since the texture is not updated periodically. Therefore, the coding gain is much smaller for video sequences with varying illumination. In order to exploit the information from the model frame also for this class of sequences, we add

an illumination component to the scene model that describes the photometric properties of object surfaces and light sources. This way, the lighting in the model frame can be compensated towards the original frame by changing the parameters of the photometric model.

The incident light in the original scene is assumed to consist of ambient light and a directional light source with illumination direction \mathbf{l} . The object surface is modeled by Lambertian reflection, and thus the relation between the video frame intensity I and the corresponding value I_{model} from the head model is

$$\begin{aligned} I^R &= I_{model}^R (c_{amb}^R + c_{dir}^R \cdot \max\{-\mathbf{n} \cdot \mathbf{l}, 0\}) \\ I^G &= I_{model}^G (c_{amb}^G + c_{dir}^G \cdot \max\{-\mathbf{n} \cdot \mathbf{l}, 0\}) \\ I^B &= I_{model}^B (c_{amb}^B + c_{dir}^B \cdot \max\{-\mathbf{n} \cdot \mathbf{l}, 0\}). \end{aligned} \quad (1)$$

c_{amb} and c_{dir} control the intensity of ambient and directional light, respectively [9] and the surface normal \mathbf{n} is derived from the 3-D head model. The Lambertian model is applied to all object pixels in the image. Each pixel contributes 3 equations for the 3 RGB color components with a common direction of the incident light. 8 parameters (ambient light: 3, directional light: 3, illumination direction: 2) characterize the current illumination situation for the entire object. By estimating these parameters with a linear least-squares estimator as shown in [9], we are able to compensate the dominant brightness differences of corresponding points in the synthesized model frame and the camera frame. This improves the quality of the model frame used for prediction in the model-aided coder significantly if the illumination in the scene changes.

3.2. Multiple Object Motion

So far, the model-based coder lacks the generality to cope with multiple object motion and deformation. For example, the sequence *Clapper Board* (Fig. 4) shows a clap moving in front of a person occluding most parts of the face. In order to exploit the model frame also for multiple object sequences, some modification to the parameter estimator are necessary. Two different cases are distinguished: first, only the head and shoulder part is modeled in the synthetic scene and, second, all objects are described by a 3-D model.

If no 3-D model exists for the additional objects, the model frame does not show them and cannot be expected to improve the prediction in the corresponding area. The rate-distortion decision of the multi-frame predictor, however, ensures that the coding efficiency does not decrease below H.263 even for this case. On the other hand, the model frame can still contribute to the prediction of those parts in the image that are not occluded or uncovered. This requires the motion estimator to determine the parameters also from partly occluded objects. The occluded parts are detected using image gradients and intensity differences between model and camera frame. They are classified as outliers in the over-determined system of equations and not used for parameter estimation. Additionally, only those FAPs are estimated that are influenced by a sufficient number of equations. Otherwise they remain constant until they are uncovered again.

Higher coding gains can be obtained if the additional objects are also modeled in the synthetic scene. The parameter estimation is performed in the same way for all objects and only the description for shape and motion/deformation constraints is adapted to the individual object. For the clap in Fig. 4, e.g., a planar triangular mesh is extracted from the first frame showing the entire object.

Five animation parameters are estimated: translation in three directions, rotation in the image plane and opening of the clap. All pixels showing a particular object contribute to the corresponding system of linear equations for the estimation of animation parameters. The classification of the pixels to individual objects is determined exploiting knowledge from the synthetic 3-D scene.

4. EXPERIMENTAL RESULTS

Experiments are conducted with the two self-recorded natural CIF sequences *Clapper Board* and *Illumination*. Rate-distortion curves are measured by varying the DCT quantizer parameter over values 10, 15, 20, 25, and 31. Bit-streams are generated that are decodable producing the same PSNR values as at the encoder. The data for the first intra-coded frame and the initial 3-D model are excluded from the results thus simulating steady-state behavior, i.e., we compare the inter-frame coding performance of both codecs excluding the transition phase at the beginning of the sequence. To specify the coding performance of the proposed model-aided codec (MAC), we compare it to the H.263 test model, TMN-10 (Annexes D, F, I, J, and T enabled).

For the special case of head-and-shoulder sequences, bit-rate savings of 35 % at the low bit-rate end corresponding to a coding gain of 2-3 dB PSNR are reported [1, 8]. If the lighting in the scene changes this coding gain is reduced, since the model frames no longer represent the original video frames correctly. The additional estimation of the lighting situation, however, allows to adapt the illumination condition in the synthetic scene to the real world. The effectiveness of the illumination estimation is illustrated in

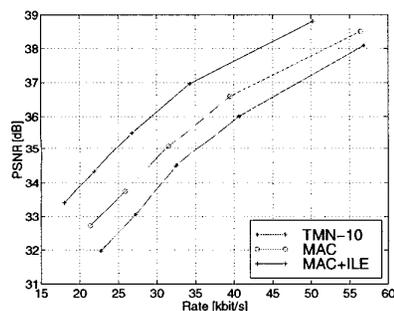


Fig. 3. Rate-distortion plot for the sequence *Illumination* illustrating the achieved improvement when using an illumination estimator (ILE).

Fig. 3 for the sequence *Illumination*. During the acquisition of this sequence, one natural light source was moved to alter the illumination conditions. Two experiments are performed. For the first one, only the FAPs are estimated to create a model frame. For the second experiment, we additionally estimate the illumination parameters and generate motion- and illumination-compensated model frames. As shown in Fig. 3, the gain in PSNR for the model-aided coder compared to the TMN-10 is about 1 dB if no illumination compensation is performed. An additional gain of about 1.5 dB is achieved when exploiting illumination information.

In a second experiment, the influence of unknown objects in the scene is investigated. Fig. 4 shows the first frames of the head-and-shoulder sequence *Clapper Board*. During the first 50 frames, the face is occluded by an object that cannot be represented by the 3-D models available at the decoder. As a result, the correspond-



Fig. 4. Frames 0, 11, 22, 33, 44, and 55 of the sequence *Clapper Board*.

ing model frames do not contain this additional object as shown in Fig. 7 a). Since prediction from the previous decoded frame and residual coding provides us with robustness against model failures, the model-aided coder represents the entire frame correctly (Fig. 7 b)). The coding efficiency of the model-aided coder, however, drops down during the first frames as shown in the temporal evolution of the PSNR in Fig. 5. If the face is visible again the model frame can be exploited and the PSNR recovers showing high coding gains. The overall rate-distortion performance for the

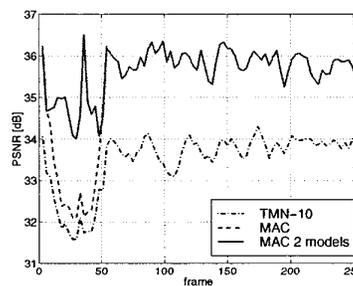


Fig. 5. Temporal evolution of PSNR for the sequence *Clapper Board*. Both coders use a DCT quantizer parameter of 31.

entire sequence is depicted in Fig. 6. Bit-rate savings of 33 % at the low bit-rate end are achieved. The quality of the reconstructed

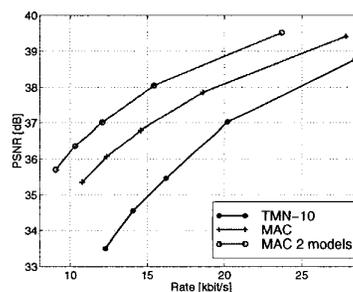


Fig. 6. Rate-distortion plot for sequence *Clapper Board*.

frames is illustrated in Fig. 7. Image b) shows frame 54 encoded with the model-aided coder, while c) corresponds to the TMN-10 coder at the same bit-rate. Figure 7 d) shows a frame from the TMN-10 coder that has the same PSNR as the model-aided frame. Even though the PSNR is the same, the subjective quality of the

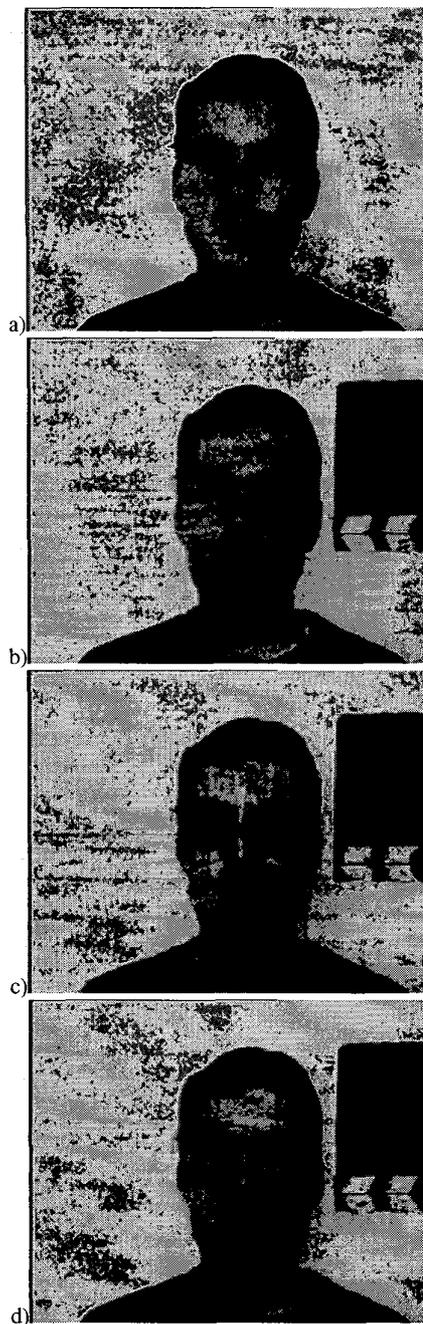


Fig. 7. Frame 54 of the sequence *Clapper Board*. a) Model frame without clapper board; b) MAC, 36.1 dB, 2900 bits; c) TMN-10 with same average bit-rate as MAC, 32.7 dB, 3200 bits; d) TMN-10 with same average PSNR as MAC, 36.0 dB, 5800 bits.

reconstructed frame from the model-aided coder is higher since facial features are reproduced more accurately and with less artifacts. The difference is even more striking when viewing motion sequences¹.

¹<http://www.LNT.de/eisert/mac.html>

In a third experiment, the clap in the sequence *Clapper Board* is described by an additional 3-D model placed in the synthetic scene. This model is manually acquired using the texture from one frame that shows the entire clap. Motion and deformation parameters are estimated for both objects using the approach in Section 3. With these parameters, model frames are generated that represent all objects in the scene. Running the model aided coder with these model frames results in a much higher PSNR for the first frames compared to the case when using only the head model. This is illustrated in the upper curve of Fig. 5. At the low bit-rate end, an average gain of 4.6 dB is achieved for the first 60 frames when using two models while the use of a single head model results in a gain of 1.2 dB PSNR. The overall rate-distortion performance for the entire sequence is depicted in Fig. 6. Bit-rate savings of 45 % corresponding to a coding gain of about 3.5 dB are achieved.

5. CONCLUSIONS

The proposed model-aided codec which combines model-based video coding with block-based motion-compensated prediction yields a superior video coding scheme for head-and-shoulder sequences. The coding efficiency of this codec can be further increased by improving the 3-D models describing the scene. We have extended the model-aided codec to exploit the efficiency of the model-based codec also for more sophisticated video sequences with changing lighting conditions or multiple objects. Parameters describing the illumination in the scene are estimated together with motion and deformation of the objects resulting in more accurate model frames. Experiments have shown that bit-rate savings of up to 45 % can be achieved at equal average PSNR.

6. REFERENCES

- [1] P. Eisert, T. Wiegand, and B. Girod, "Rate-distortion-efficient video compression using a 3-D head model," in *Proc. International Conference on Image Processing (ICIP)*, Kobe, Japan, Oct. 1999, vol. 4, pp. 217–221.
- [2] ITU-T Recommendation H.263 Version 2 (H.263+), "Video Coding for Low Bitrate Communication," Jan. 1998.
- [3] D. E. Pearson, "Developments in model-based video coding," *Proceedings of the IEEE*, vol. 83, no. 6, pp. 892–906, June 1995.
- [4] P. Eisert and B. Girod, "Analyzing facial expressions for virtual conferencing," *IEEE Computer Graphics and Applications*, vol. 18, no. 5, pp. 70–78, Sep. 1998.
- [5] *ISO/IEC FDIS 14496-2, Generic Coding of audio-visual objects: (MPEG-4 video), Final Draft International Standard, Document N2502*, 1999.
- [6] M. F. Chowdhury, A. F. Clark, A. C. Downton, E. Morimatsu, and D. E. Pearson, "A switched model-based coder for video signals," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 4, no. 3, pp. 216–227, June 1994.
- [7] H. G. Musmann, "A layered coding system for very low bit rate video coding," *Signal Processing: Image Communication*, vol. 7, no. 4-6, pp. 267–278, Nov. 1995.
- [8] P. Eisert, T. Wiegand, and B. Girod, "Model-aided coding: A new approach to incorporate facial animation into motion-compensated video coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 10, no. 3, pp. 344–358, Apr. 2000.
- [9] P. Eisert and B. Girod, "Model-based coding of facial image sequences at varying illumination conditions," in *Proc. 10th Image and Multidimensional Digital Signal Processing Workshop IMDSP '98*, Alpbach, Austria, Jul. 1998, pp. 119–122.