

# Sliding adjustment for 3D video representation

Franck Galpin and Luce Morin

Irisa/INRIA Rennes, Université de Rennes 1,  
Campus de Beaulieu, 35042 Rennes Cedex, FRANCE

Email: Franck.Galpin@irisa.fr Luce.Morin@irisa.fr;

<http://www.irisa.fr/temics>;

Telephone: (33) 2.99.84.74.30; Fax: (33) 2.99.84.71.71

## Abstract

This paper deals with video coding of static scenes viewed by a moving camera. We propose an automatic way to encode such video sequences using several 3D models. Contrary to prior art in model-based coding where 3D models have to be known, the 3D models are automatically computed from the original video sequence. We show that several independent 3D models provide the same functionalities as one single 3D model, and avoid some draw-backs of the previous approaches. To achieve this goal we propose a novel algorithm of sliding adjustment, which ensures consistency of successive 3D models. The paper presents the method to automatically extract the set of 3D models and associate camera positions. The obtained representation can be used for reconstructing the original sequence, or virtual ones. It also enables 3D functionalities such as synthetic object insertion, lightning modification or stereoscopic visualization. Results on real video sequences are presented.

## Keywords

sliding adjustment, 3D model reconstruction, video coding, model-based coding, video manipulation

## INTRODUCTION

More and more new coding techniques include high level information in video sequence representation. This information aims to provide high level functionalities such as interactivity, video content description, video manipulation or stereo visualization.

For instance, the QuickTime-VR format provides the functionality of interactive visualization of a real static scene, by representing it as a panoramic image [1].

MPEG4 standard describes the video scene content as a set of plane objects called Video Object Plane (VOP) [2], which can be interactively moved or combined during visualization. A panoramic representation of static backgrounds is also proposed in MPEG4 with the Sprite format.

Such 2D representations do not give information on the 3D structure of the scene, and are therefore limited for video manipulation. Panoramic images provide only limited interactivity: zoom and view orientation can be changed but the view-point is fixed. With 2D representations, video manipulation such as hybrid synthetic-real video mixing, involving occlusions, shadows, lightning modification are not feasible in a realistic way. These functionalities require 3D information on the scene.

3D model-based representations for real video sequences have been studied for a long time, since they have very attractive properties. Apart from the functionalities that they

provide, they enable very low bit-rates and scalable/progressive coding [3].

3D model-based representations can be classified into explicit and implicit representations. Within the explicit representations, we can distinguish representations with known 3D models and unknown 3D models.

In 3D model-based coding with known models, a 3D model of the object in the scene is available, for instance a textured 3D triangular mesh. The video sequence is processed to compute the 3D object pose (orientation and scaling) for each frame. Sometimes, local deformations are also computed. The video sequence is represented as the 3D model and pose parameters for each frame, with optional parameters for texture and local shape deformation. This representation allows to transmit the original video at low cost and facilitates any 3D manipulation functionalities. This approach is widely used for head and shoulder video sequences coding and body animation analysis and representation [4] [5], for instance in the MPEG4-SNHC scheme [6]. Its main draw-back is that it can only be applied to video with specific contents, such as manufactured objects, head or body.

In the 3D model-based coding with unknown models, the same principle is applied but the scene contents are unknown and the 3D model shape must also be estimated from the video itself. Since shape and non-rigid motion can not be separated, this approach can only be applied to video containing one object undergoing rigid motion, or alternatively, a fixed object viewed by a moving camera. Using computer vision tools [7] [8], the 3D model shape, texture and rigid motion are estimated from the video sequence. Once the 3D model is computed, it can then be used just as a known 3D model. This approach is limited to fixed objects, but it is very attractive for applications such as realistic modeling of complex objects or interactive navigation in real environments. The most sensitive step is the estimation of camera internal parameters, or self-calibration. A theoretical solution has been established for long. Robust solutions for camera calibration and pose estimation have been proposed in recent work, and are effective for video sequence acquired with an hand-held cam-corder [9]. An other solution to deal with generic video data requires that a priori high-level information is integrated in the process through manual and sometimes expert user interaction [10].

Such an approach is thus not yet applicable to design a coding algorithm which auto-

matically produces a unique 3D model from a generic and long video sequence of a static scene without simplifying assumptions on the scene or acquisition.

Alternative approaches based on implicit 3D model-based coding have been proposed. The lightfield or the Lumigraph [11] [12] do not aim to reconstruct an explicit 3D model but provide some of the same functionalities. However, data acquisition is very constraining, as view-points must lie on a dense regular grid. Some other approaches introduce depth information in the encoded sequence, which allows stereo sequence visualization or scene manipulation [13], but which often requires stereo acquisition.

In this paper, we present an original representation of video sequences using several unknown 3D models and we propose a novel algorithm for ensuring high-level 3D functionalities using this representation. This representation can be applied in the case of a fixed scene viewed by a moving camera. We present an automatic scheme for extracting the proposed representation from the video sequence.

Some previous studies try to extract a single 3D model with a hierarchical and robust estimation of camera positions [14] [9]. A self-calibration step is also performed, allowing the reconstruction of a single 3D model. Such methods generally require a specific type of camera motion (typically a closed image sequence or an inspecting image sequence), in order to perform the self-calibration step. In this paper, we deal with video coding, assuming very long video sequences, we thus need a on-the-fly process. Moreover, we want to deal with any type of camera motion. For instance, one typical application could be navigation on a walking path where camera motion is naturally a rough frontal translation. This type of motion is closed to degenerate cases and scene points appear in a small part of the sequence. Thus we do not make the assumption that it is possible to obtain a reliable, accurate camera self-calibration for any sequences.

Unlike the classical approach of video coding with unknown models described before, we do not aim at reconstructing one single realistic 3D model of the scene. Instead, we compute a succession of 3D models, each 3D model being adequate to represent a small part of the video sequence. This approach can be viewed as an intermediate between 2D motion compensation video-coding and 3D model-based coding. Just as in the 2D approach, one 3D model can be considered as a global motion model which best fits 2D

motion in the original video sequence for a group of images (also called here a ‘‘GOP’’). Once this ‘‘motion model’’ is not valid anymore, the GOP is ended and a new 3D model is estimated for representing the next part of the video (see fig. 1). Thus, successive 3D models may contain the same parts of the 3D scene, but each model is related to a specific GOP of the video. Also GOPs size is not fixed but data driven, thus variable.

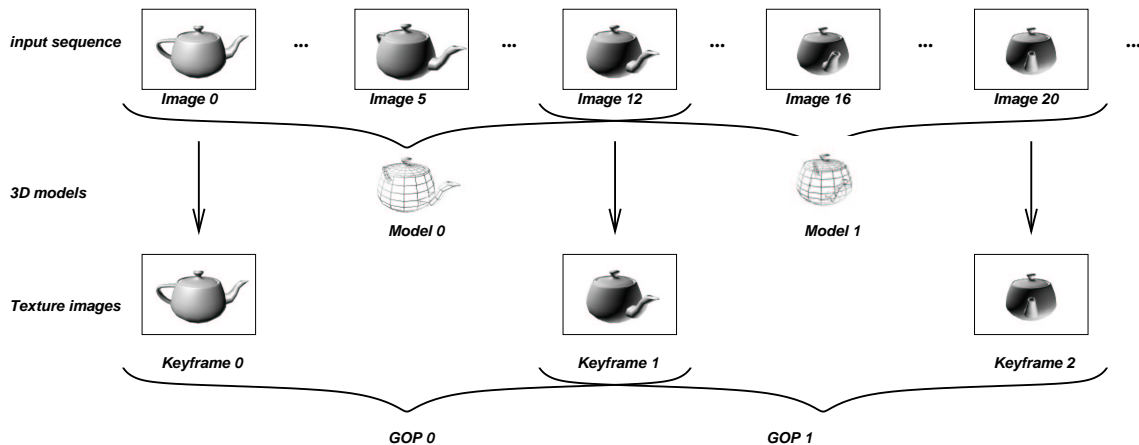


Fig. 1. Principle of video sequence representation using several overlapped 3D models.

In order to obtain the same functionalities with several 3D models as with a single 3D model, we propose a novel algorithm of sliding adjustment which ensures consistency of successive 3D models. This step enables applications such as synthetic 3D objects insertion into the video sequence, lightning modification or interactive navigation.

Using several 3D models instead of one single model has several drawbacks :

- It is of course less compact, since the models are redundant.
- It is not as simple and easy to insert synthetic data in the video sequence. With several 3D models, an adjustment step is required.
- The set of acceptable virtual views (views reconstructed from a view-point outside the acquisition set) is smaller, since the 3D models are constrained to be consistent only with a subset of original views.

In counterparts, such a representation has several advantages:

- Global consistency of estimated 3D shape and camera motion along the video sequence is no more expected. Thus, a valid representation is provided even in the case of ill-conditioned camera motion or inaccurate internal camera parameters.

- Such a representation is very well suited for large environments, where the observed scene part progressively changes along the sequence.
- The representation is robust to small violations of the rigidity constraint, for instance presence of small moving objects or specular surfaces. Such data is temporarily taken into account by the 3D models, as a change in the model geometry.
- 3D models are constructed sequentially. Streaming of the representation can be easily achieved for communication and compression purposes. This point is particularly important for very long sequences (several thousands of images) for which the automatic reconstruction of a single 3D model is very complex and computationally intensive.

This representation is thus very well suited for representing large natural scenes such as the ones acquired by outdoors walk through in cities, parks, with uncontrolled camera motion. Applications concern virtual tours in realistic environments, with possibility of scene manipulation and interactive navigation.

The paper is organized as follows. We first present the principle of video representation using several 3D models, the coding scheme and visualization procedure. We explain how this representation can be used to regenerate the initial video sequence, as well as virtual ones. The motion estimation step is briefly described because it uses classical techniques. We then present in more details the automatic selection of variable sized GOPs, and the sliding adjustment algorithm. Finally, we validate the method on real applications such as interactive navigation, virtual lightning, synthetic object insertion or stereoscopic visualization. Examples of the obtained results on real video sequences are presented and discussed.

## I. VIDEO CODING USING SEVERAL UNKNOWN 3D MODELS

Our approach is quite similar to 3D model based coding using a single unknown 3D model: shape and texture are estimated from the video sequence itself, using shape from motion techniques. Camera motion for each frame is also estimated from the video sequence. The following assumptions are made: we use perspective projection model and we assume that the observed scene is fixed, contains mostly lambertian objects and is not entirely planar. The same scheme can be applied to a fixed background if moving objects have been segmented out from the video, as an alternative to the MPEG4-Sprite mode for

instance.

In our approach, camera internal parameters are not necessarily known. If not provided they are affected arbitrary values. Camera motion is not constrained. We yet assume that camera motion is not a pure rotation around optical center.

Instead of computing one single model for the whole sequence, several 3D models are computed for the same 3D scene. Each 3D model is relative to a GOP (see figure 1). More precisely, for a given GOP, the 3D model and associated camera positions are estimated from the images in the GOP. At the decoder, the estimated model is projected onto the estimated camera positions to reconstruct these images. The coding scheme is thus sequential, as in classical motion-compensation video-coders. The 3D model is expected to best fit the 2D information in the GOP, by minimizing a cost function based on MSE. Thus 3D shape may not be realistic as long as it allows reconstruction of the original sequence with minimum distortion. Subsequently, 3D models for successive GOPs may represent the same 3D object but they may have different shape, different texture and even different scale.

Only a subset of images are used for 3D reconstruction of the 3D models. These images are called *keyframes*. Two successive keyframes are used to compute one 3D model. They are the first image and last image of the GOP associated with this 3D model. One keyframe is used as texture image for defining the 3D model texture. All keyframes are thus part of the representation, as texture images of the 3D models. Successive GOPs overlap by one keyframe. Keyframes can thus be reconstructed at the decoder either using one 3D model or the other. This is important for smooth transition between GOPs during visualization (note that what we call a keyframe is different from keyframes delimiting shots in video structure analysis).

3D model reconstruction from two views extracted from a video is known to be very sensitive to the choice of these two views: several criteria must be verified in order for the estimation to be geometrically and numerically stable. This is the reason why GOP size can not be fixed. GOP size varies depending on data driven keyframe choice. We propose a robust method to select keyframes, which is described in section IV.

For simple visualization of the original sequence, 3D models can be completely indepen-

dent: they can have different scales and they can be described in different reference frames. This is still true for visualization along a virtual path, as long as this path contains all the camera positions associated with the keyframes. At these specific viewpoints, transition between 3D models is smooth, because both models, though different in 3D space, project onto the same 2D image: the common keyframe.

However, independent 3D models are not suited for other 3D augmented reality functionalities, like inserting objects or lights. The 3D object position should be specified in a reference frame valid along the whole sequence. With independent 3D models for each GOP, this is not feasible. We thus propose an accurate method to set a 3D model compatible with the previous and next models. This method is derived from the classical bundle adjustment method, but it is specifically adapted to sequential processing and minimum distortion coding purposes. This algorithm is called *sliding adjustment* and it will be described in section V.

The general coding scheme is shown in figure 2.

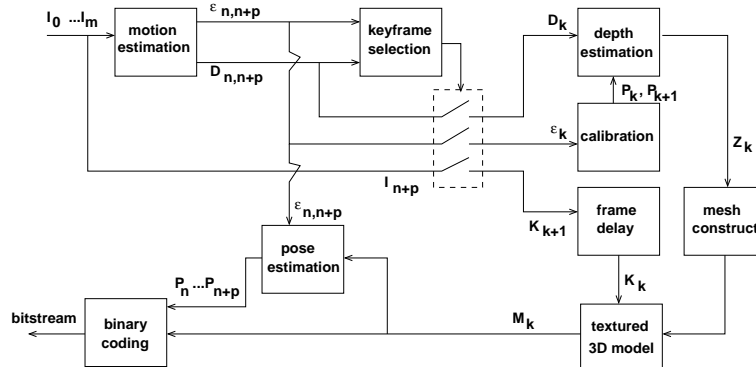


Fig. 2. Bloc-diagram of the 3D encoder

## II. BACKGROUND

In this section, we briefly remind the principle of 3D model reconstruction from 2 views and we set the notations used in the paper.

Let consider two images  $I$  and  $I'$  viewed by cameras  $\mathcal{C}$  and  $\mathcal{C}'$ . We denote  $(R, t)$  the relative rigid transformation between cameras, and  $A$  the matrix containing the internal parameters of the cameras.

Let  $\mathcal{E}$  be a set of matched points between images  $I$  and  $I'$ . For  $m_i$  a point of  $\mathcal{E}$  in  $I$  and

$m'_i$  a point of  $\mathcal{E}$  in  $I'$ , we have the following relation in homogeneous coordinates:

$$\begin{aligned}\tilde{m}_i &= P.\tilde{M}_i \\ \tilde{m}'_i &= P'.\tilde{M}_i\end{aligned}\tag{1}$$

where  $P = A.(I_3|0)$  and  $P' = A(R|t)$ .

If  $m_i$  and  $m'_i$  are matched, then the epipolar constraint is expressed in homogeneous coordinates as:

$$\tilde{m}'_i{}^T.F.\tilde{m}_i = 0\tag{2}$$

where  $F$  denotes the well-known fundamental matrix.  $F = A^{-T}.[t]_{\times}.R.A^{-1}$  and  $[.]_{\times}$  is the matrix associated with the cross-product. We denote  $E = [t]_{\times}.R$  the essential matrix associated with  $F$ .

We define the *epipolar residual* associated with  $F$  as the sum:

$$\frac{1}{2} \sum_i d(m_i, F.m'_i) + d(m'_i, F^T.m_i)\tag{3}$$

computed on all matched points  $m_i, m'_i$  in  $\mathcal{E}$ . Nullity of the epipolar residual means that the symmetrical epipolar constraint defined by  $F$  is verified for all points in  $\mathcal{E}$ . If this residual is small (i.e. with sub-pixel value), we then say that  $F$  is consistent with  $\mathcal{E}$ .

In the following, we denote  $K_k$  the keyframe images. For a given keyframe image  $K_k$ ,  $R_k, t_k$  denote the camera motion parameters and  $O_k$  denotes the center of projection for the corresponding camera, i.e.  $O_k = -R_k^{-1}.t_k$ .

Keyframes  $K_k$  and  $K_{k+1}$  delimit a GOP and  $\mathcal{M}_k$  denotes the 3D model associated to this GOP.  $\mathcal{E}_k$  denotes a set of points matched between  $K_k$  and  $K_{k+1}$ .

We also denote  $\vec{u}(m)$  the unitary tangent vector of the view-line passing through pixel  $m$  in image  $K_k$  (i.e. line  $(O_k, m)$ ).

### III. 2D AND 3D MOTION ESTIMATION

#### A. 2D motion estimation

##### A.1 Dense motion estimation

Motion estimation is performed between two current images  $I = I_n$  and  $I' = I_{n+p}$ . Usually  $I = K_k$  is the last selected keyframe and  $I'$  is a subsequent image in the video, which is evaluated as a potential next keyframe  $K_{k+1}$ .

Motion estimation is provided by previously developed algorithms. We use a mesh-based motion estimator based on a multi-resolution scheme over hierarchical meshes. It allows dense estimation between current images  $I$  and  $I'$ , by successive estimation/relaxation steps between successive images [15]. This motion estimator provides a dense motion field between  $I$  and  $I'$  (see fig. 8(a)). For each pixel  $m_i$  in image  $I$ , its 2D displacement is a 2D vector denoted  $D(m_i)$  and its corresponding position in image  $I'$  is thus  $m'_i = m_i + D(m_i)$ .

## A.2 Sparse point matching

We also compute a set of matched points  $\mathcal{E}$  between current images  $I$  and  $I'$  using the motion field. These points are chosen among the vertices of the mesh used in motion estimation. We select 200 vertices uniformly scattered in the image and which get highest scores for Harris and Stefen detector [16] in image  $I$ . More details about motion estimation can be found in [17] [18].

## B. Camera parameters estimation

The parameters needed for reconstruction of the 3D model are internal and external camera parameters for each keyframe.

### B.1 Internal parameters

Internal camera parameters need not to be accurate more than by an order of magnitude (this is one of the advantages of the representation with several models rather than one model). We thus arbitrarily choose parameters which are equal for all keyframes: we fix optical center  $(u_0, v_0)$  at image center and we assume square pixels. An order of magnitude for the focal length may be provided by previous calibration or constructor data. If not, focal length is arbitrarily fixed to 500 (this is what is done for all the presented results). An other solution would be to estimate internal parameters through self-calibration techniques directly from the video sequence. However previous studies have shown that it is a highly unstable procedure with general acquisition conditions [19]. It is thus not adapted to automatic coding schemes.

## B.2 External parameters

External camera parameters  $(R_k, t_k)$  are estimated from the set of matched points  $\mathcal{E}_k$ . The goal is to obtain a set of internal and external parameters which is consistent with  $\mathcal{E}_k$ , i.e. the epipolar residual associated with fundamental matrix  $F_k = A^{-T} \cdot [t_k]_{\times} \cdot R_k \cdot A^{-1}$  must be of sub-pixel value. For the sake of simplicity, index  $k$  is omitted in the remaining of the section; all parameters implicitly refer to the current GOP between  $K_k$  and  $K_{k+1}$ .

We first estimate fundamental matrix  $F$  by minimizing the epipolar residual for all points in  $\mathcal{E}$ , using a classical median least squares algorithm [20]. The obtained fundamental matrix is denoted  $F_m$ . The essential matrix  $E_m$  is obtained from  $F_m$  as  $E_m = A^T \cdot F_m \cdot A$ . A first set of camera parameters, denoted  $R_c$  and  $t_c$ , are then computed from  $E_m$  using a state-of-the-art decomposition method [21].  $R_c$  and  $t_c$  minimize:

$$\|E_m - E_c\|_f \quad (4)$$

where  $\|\cdot\|_f$  denotes the Frobenius distance and  $E_c = R_c \cdot [t_c]_{\times}$ . This estimation thus minimizes a matrix distance between estimated essential matrix  $E_c$  and essential matrix related to matched points  $E_m$ . However, it does not ensure that the camera parameters are compatible with the set of matched points  $\mathcal{E}$ . Indeed the corresponding matrix  $F_c = A^{-T} \cdot [t_c]_{\times} \cdot R_c \cdot A^{-1}$  is not similar to  $F_m$  and the epipolar residual (3) associated with  $F_c$  is large. In other words,  $R_c$ ,  $t_c$  and  $F_c$  are not consistent with  $\mathcal{E}$ .

As explained before, the epipolar residual associated with  $F_c$  also assesses the projection error for points in  $\mathcal{E}$ , when 3D reconstruction is performed using  $A$ ,  $R_c$  and  $t_c$ .

This is a very important criterion to take into account as we are looking for a 3D model and camera position which enable to re-create the original images between 2 keyframes. Moreover, the next step of sliding adjustment is sensitive to initialization, so a refinement step is performed on the  $(R_c, t_c)$  parameters as follows.

We first compute the 3D position of each point in  $\mathcal{E}$  and we then estimate the pose of this 3D points set with Dementhon algorithm [22]. This technique computes the camera position  $(R_d, t_d)$  which minimizes projection error for a given set of 3D points and their corresponding image points. We obtain a new pair  $(R_d, t_d)$  which is compatible both with  $A$  and  $\mathcal{E}$ . This is verified by computing the epipolar residual for  $F_d = A^{-T} \cdot [t_d]_{\times} \cdot R_d \cdot A^{-1}$ .

Figure 3 shows a comparative plot of the epipolar residuals associated with  $F_c$  and  $F_d$ , as a function of frame number, for the *stairway* sequence. We can see that the proposed refinement greatly improves consistency of camera parameters with image data. Moreover it shows to be a more robust technique, as it often provides reasonable solution (epipolar residual smaller than 1 pixels) when the decomposition method provides an invalid solution (epipolar residual greater than 4 pixels).

At this point the estimated parameters  $R_d, t_d$  are compatible with  $\mathcal{E}$ .

Once camera parameters are estimated, results are used to decide whether frame  $I'$  should be used as keyframe image  $K_{k+1}$ . This is done by an automatic keyframe selection algorithm described in section IV. If  $I'$  is detected as an invalid  $K_{k+1}$ , a further frame is chosen as a candidate and motion estimation steps are started again for the same GOP. If  $I'$  is a valid keyframe  $K_{k+1}$ , 3D model  $\mathcal{M}_k$  associated with the current GOP is computed from dense motion field  $D_k$  and camera parameters  $R_k$  and  $t_k$ . To achieve global consistency of the 3D models, a sliding adjustment procedure is performed. It provides the 3D model scale and camera parameters  $R_k, t_k$  which are consistent with the previous GOP. Sliding adjustment will be described in details in section V. The whole procedure is then started again for next GOP, with  $I = K_{k+1}$  as first image in the GOP.

#### IV. KEYFRAME SELECTION

We propose a simple method to select keyframes on the fly in order to insure a valid 3D model reconstruction. The first selected keyframe  $K_0$  is the first image of the video sequence  $I_0$ . The other keyframes are selected while processing and coding the sequence. Suppose keyframe  $K_k$  has been selected and we are looking for the next keyframe  $K_{k+1}$ . The selection is made by the verification of 3 criteria estimated from the dense motion estimation and matched points.

##### A. Selection criteria

The following criteria are computed for each image  $I'$  following  $I = K_k$  in order to decide if  $I'$  will be the next keyframe.

- $C_1 \quad \bar{D} > S_m$ : where  $\bar{D}$  is the average apparent motion and  $S_m$  a static threshold fixed to 10 pixels,

- $C_2$ : the percentage of outgoing points is less than a threshold  $S_o$ , i.e. no more than  $S_o$  percent of points from  $I$  in  $\mathcal{E}$  are not present in  $I'$ , with  $S_o$  of typical value 30%,
- $C_3$   $\sum_i d^2(m_i, F_d.m'_i) + d^2(m'_i, F_d^T.m_i) < S_f$ : where  $F_d$  is the fundamental matrix computed from the estimated camera motion (see below),  $(m_i, m'_i) \in \mathcal{E}$  the set of matched points between  $I$  and  $I'$ , and  $S_f$  a static threshold typically fixed to 0.5 pixels.

The first criterion tends to favor a good precision for the depth field: the best precision on depth is clearly achieved with perpendicular lines of views.  $C_1$  is a necessary criterion for a significant change in camera viewpoints, but is not sufficient, since  $C_1$  can be defeated with a large rotation component.

The second criterion  $C_2$  insures that the two keyframes share a large part of the scene. This is necessary because the 3D model contains only points viewed in both images.

The third criterion  $C_3$  ensures a valid 3D model reconstruction, by testing the epipolar residual. This criterion has two means. First it allows to detect ill-conditioned configurations. In such cases, due to numerical errors, the estimated 3D model, camera motion and fundamental matrix are not consistent with the image data and motion field. The epipolar residual is then very large. The second mean of  $C_3$  is to ensure that the 3D model projects onto image  $I'$  with sub-pixel error. The epipolar residual is the average 2D projection error for points in  $\mathcal{E}$ . Thus it evaluates the ability of the 3D model to accurately represent image  $I'$ .

These three criteria are used as follows: for a fixed image  $I = I_n$ , successive images  $I_{n+1}, I_{n+2}, \dots$  are examined until  $C_1$  and  $C_2$  are verified. Following images are considered as candidates for  $K_{k+1}$ . For each candidate image  $I'$ , camera motion and fundamental matrix are estimated in order to evaluate  $C_3$ . We then select as  $K_{k+1}$  the last candidate  $I'$  before  $C_2$  is false or before  $C_3$  is not verified for more than 2 successive frames. One or two successive frames  $I'$  with true  $C_1$  and  $C_2$  and false  $C_3$  is considered due to instable numeric estimation of  $F$ . This is the reason why the GOP is not ended before 3 or more successive images do not verify  $C_3$ .

## B. Validation

This approach has been validated on several video sequences with various camera motions; we show the results on two typical cases. Figure 4 presents the evolution of the

3 criteria on the two test sequences:  $C_1$  on top,  $C_2$  in the middle,  $C_3$  at the bottom. Horizontal lines show the threshold values for each criterion, and vertical dotted lines indicate keyframes. Left column refers to the *street* sequence, and right column refers to the *stairway* sequence (see section VI-C for images from the original video sequences).

Thresholds  $S_o$  and  $S_f$  are manually fixed close to typical values in order to obtain large GOPs. The following parameters were used to encode the *street* sequence:  $S_m = 10$ ,  $S_o = 40\%$  and  $S_f = 0.35$ . Epipolar residuals are computed only when average motion is greater than 10 pixels. We can notice that the 3 criteria allow to select keyframes mostly on outgoing points percentage, because this sequence have a quite stable motion.

The following parameters were used to encode the *stairway* sequence:  $S_m = 10$ ,  $S_o = 30\%$  and  $S_f = 0.4$ . This sequence is more instable than the previous one, due to unstabilized camera motion during acquisition. However, we can notice that the three criteria allow to select keyframes despite the instability of epipolar geometry.

The GOP size varies according to video contents. For the *street* sequence, where camera motion is homogeneous, GOP size is quite stable, with a value around 40 frames in a GOP. For the *stairway* sequence, GOP size values between 5 to 30 frames with a typical value of 25 frames. GOPs are adapted to scene content and camera motion so that a single 3D model can accurately represent the frames in the GOP.

## V. SLIDING ADJUSTMENT

At this point, we have a set of camera parameters which are independent, i.e. a computed 3D model  $\mathcal{M}_k$  whose geometry is optimal for the GOP  $k$  between  $K_k$  and  $K_{k+1}$ . Since we want local consistency between each successive 3D models, we use a sliding window to compute the camera positions in order to increase the consistency of a pair (camera, 3D model) with its neighbors.

### A. Initialization

As the proposed sliding adjustment is solved using a non-linear optimization procedure, initial values for the estimated parameters must be provided. These value should be close enough to the solution to allow the sliding adjustment to converge toward an acceptable solution. The camera position is first initialized with the previously computed  $(R_d, t_d)$

parameters. Each 3D model has its own scale factor because camera translation and 3D model are defined up to a scale factor  $\alpha$ , as shown by the following equation:

$$\begin{aligned} \forall M = (x, y, z) \in \mathcal{M}_k, \\ \tilde{m} = A.(R|t).(x, y, z, 1) \Leftrightarrow \tilde{m} = A.(R|\alpha.t).(\alpha.x, \alpha.y, \alpha.z, 1) \end{aligned} \quad (5)$$

Consistent scales for successive models are also estimated by the sliding adjustment procedure. Initial values thus have to be provided as well. We describe here how the initial scale of each 3D model is set to be similar to its previous 3D model.

Scale of first 3D model  $\mathcal{M}_0$  is not modified and is taken as a reference for the whole sequence. We compute its gravity center  $G_0$  from the set of matched points  $\mathcal{E}_0$  (points lying at infinity are not taken into account). Assuming that gravity center  $G_k$  and scale factor  $\alpha_k$  have been computed for model  $\mathcal{M}_k$ , the following steps are then iteratively performed:

- track points in  $\mathcal{E}_k$  from keyframe  $K_{k+1}$  to keyframe  $K_{k+2}$ ,
- compute  $G'_{k+1}$ , the gravity center of the 3D points reconstructed using these points,
- compute scale factor for  $\mathcal{M}_{k+1}$ :  $\alpha_{k+1} = |G_k.R_k + t_k - O_{k+1}| / |G'_{k+1} - O_{k+1}|$ ,
- rescale 3D model  $\mathcal{M}_{k+1}$  and associated matrices  $P_{k+1}, P_{k+2}$  using  $\alpha_{k+1}$ ,
- compute new gravity center  $G_{k+1}$  from the set of points  $\mathcal{E}_{k+1}$

At the end of this process, we obtain a set of 3D models  $\mathcal{M}_k$  which are defined in the same basis and which have a consistent scale from one 3D model to the next 3D model in the stream.

### B. Algorithm

Our algorithm is based on bundle adjustment [23] but it is adapted to our application, namely:

- each set  $\mathcal{E}_k$  generates a 3D model  $\mathcal{M}_k$  whereas in bundle adjustment all 3D points are merged into one single model,
- local consistency is performed on a sliding temporal window: for a given keyframe  $K_k$ , only neighboring images and camera positions are taken into account,
- some 3D points are constrained to stay on their view-line for coding purpose.

We now define some useful notations:

- $P_k = A.(R_k|t_k)$  is the projection matrix which perfectly projects the 3D model  $\mathcal{M}_k$  on image  $K_k$ . This matrix is also the projection matrix for the last image of previous GOP:  $P_k$  projects the 3D model  $\mathcal{M}_{k-1}$  on  $K_k$  with an error due to the imperfections of 3D model  $\mathcal{M}_{k-1}$ .
- $\{M_i^k\}$  is a set of 3D points computed with projection matrices  $P_k$  and  $P_{k+1}$ , from the set of robust points  $\mathcal{E}_k$  extracted in image  $K_k$  and matched in image  $K_{k+1}$ .  $\{M_i^k\}$  can be seen as a subset of 3D model  $\mathcal{M}_k$ .
- $m_i^k$ : this is a point in  $\mathcal{E}_k$  extracted in keyframe  $K_k$ .
- $m_i^{k,l}$ : this is a point extracted in keyframe  $K_k$  and tracked till keyframe  $K_l$  by summation of estimated motion fields.
- $\hat{m}_i^{k,l} = P_l.M_i^k$ : this is a point of  $\{M_i^k\}$  projected in image  $K_l$  with projection matrix  $P_l$ .

For each keyframe  $K_k$ , a cost function  $f = f_1 + f_2 + f_3$  is minimized. The considered costs are 2D residual errors when projecting the 3D models onto keyframes inside the sliding window. When processing current keyframe  $K_k$ , the following parameters have been estimated at previous iteration on  $K_{k-1}$ : point sets  $\{M_i^{k-1}\}$ ,  $\{M_i^k\}$ , projection matrices  $P_{k-1}$ ,  $P_k$  and  $P_{k+1}$ . Among them  $\{M_i^{k-1}\}$ ,  $P_{k-1}$  and  $P_k$  are final values, whereas  $\{M_i^k\}$  and  $P_{k+1}$  values are estimated again in the current  $K_k$  process.

We now describe each term in the cost function and we give its geometrical interpretation.

- The first term  $f_1$  ensures that model  $\mathcal{M}_k$  projects correctly onto keyframe  $K_{k+1}$ , by finding the best projection matrix  $P_{k+1}$ .  $P_k$  is expected to perfectly project  $\mathcal{M}_k$  onto  $K_k$ , thus point  $M_i^k$  must not be modified on image  $K_k$ . A 3D point  $M_i^k$  thus has only one degree of freedom: moving along its view-line. This constraint writes:

$$M_i^k = O_k + \lambda_i^k \cdot \vec{u}(m_i^k) \quad (6)$$

Under this constraint, the cost function  $f_1$  is defined as:

$$f_1(P_{k+1}, \{M_i^k\}) = \sum_i \|m_i^{k,k+1} - \hat{m}_i^{k,k+1}\|^2 = \sum_i \|m_i^{k,k+1} - P_{k+1}.M_i^k\|^2 \quad (7)$$

$f_1$  is a function of both  $P_{k+1}$  matrix, and point set  $\{M_i^k\}$ . The unknown parameters in  $f_1$  are  $\{\lambda_i^k\}$  which define  $\{M_i^k\}$  and  $(R_{k+1}, t_{k+1})$  which define  $P_{k+1}$ .

- $f_2$  ensures consistency of keyframe  $K_{k+1}$  with 3D model  $\mathcal{M}_{k-1}$ . The set of points  $\{M_i^{k-1}\}$  has been computed on a previous step as well as projection matrices  $P_{k-1}$  and  $P_k$ , and they are thus fixed parameters. We search the best projection matrix  $P_{k+1}$  for both  $\mathcal{M}_{k-1}$  and  $\mathcal{M}_k$ , i.e. we add a new cost function:

$$f_2(P_{k+1}) = \sum_i \|m_i^{k-1,k+1} - \hat{m}_i^{k-1,k+1}\|^2 = \sum_i \|m_i^{k-1,k+1} - P_{k+1} \cdot M_i^{k-1}\|^2 \quad (8)$$

Minimizing the cost function  $f_1 + f_2$  insures that 3D model  $\mathcal{M}_k$  is consistent with the previous model  $\mathcal{M}_{k-1}$ .

- Finally we want consistency of 3D model  $\mathcal{M}_k$  with next 3D model  $\mathcal{M}_{k+1}$ . This is done by estimating  $P_{k+2}$  which best projects points  $\{M_i^k\}$  and  $\{M_i^{k+1}\}$  on image  $K_{k+2}$ , under the constraint that  $\{M_i^{k+1}\}$  project perfectly on  $K_{k+1}$ . This is ensured by minimizing cost function  $f_3$ :

$$\begin{aligned} f_3(P_{k+2}, \{M_i^k\}, \{M_i^{k+1}\}) &= \sum_i \|m_i^{k,k+2} - \hat{m}_i^{k,k+2}\|^2 + \sum_i \|m_i^{k+1,k+2} - \hat{m}_i^{k+1,k+2}\|^2 \\ &= \sum_i \|m_i^{k,k+2} - P_{k+2} \cdot M_i^k\|^2 + \\ &\quad \sum_i \|m_i^{k+1,k+2} - P_{k+2} \cdot M_i^{k+1}\|^2 \end{aligned} \quad (9)$$

under the constraint:

$$M_i^{k+1} = O_{k+1} + \lambda_i^{k+1} \cdot \vec{u}(m_i^{k+1}) \quad (10)$$

The final cost function becomes:

$$f(P_{k+1}, \{M_i^k\}, P_{k+2}, \{M_i^{k+1}\}) = f_1(P_{k+1}, \{M_i^k\}) + f_2(P_{k+1}) + f_3(P_{k+2}, \{M_i^k\}, \{M_i^{k+1}\}) \quad (11)$$

Equation (11) is a large non-linear system with the following characteristics:

- 6 unknown parameters for the projection matrix  $P_{k+1}$ : 3 for translation  $t_{k+1}$ , and 3 for rotation  $R_{k+1}$ .
- 6 unknown parameters for projection matrix  $P_{k+2}$
- $\text{Card}(\{M_i^k\})$  unknown parameters for the set of points  $\{M_i^k\}$ ,
- $\text{Card}(\{M_i^{k+1}\})$  unknown parameters for the set of points  $\{M_i^{k+1}\}$ ,
- $2 \cdot \text{Card}(\{M_i^k\})$  equations for the constraint  $f_1$  (one equation on  $x$  axis and another on  $y$  axis),

- $2 \cdot \text{Card}(\{M_i^{k-1}\})$  equations for the constraint  $f_2$ ,
- $2 \cdot \text{Card}(\{M_i^k\}) + 2 \cdot \text{Card}(\{M_i^{k+1}\})$  equations for the constraint  $f_3$ ,

One must notice that this large system is a very sparse system: a given 3D point  $M_i$  interferes in 4 or less equations. This system is then solved using a classical non linear estimation algorithm which deals with large sparse systems. We have used the MinPack [24] package implementation. At the end of the minimization step  $k$ , the final 3D model  $\mathcal{M}_k$  is computed using the 2 projection matrices  $P_k$  and  $P_{k+1}$ , and projection matrix  $P_{k+2}$  is taken as an initial value for the next step  $k+1$ .

Figure (5) presents PSNR value obtained with the *stairway* video sequence reconstructed using previous 3D model  $\mathcal{M}_{k-1}$  instead of  $\mathcal{M}_k$  for each GOP  $k$ . Figure (6) shows the PSNR obtained when using next 3D model. They show that the sliding adjustment increases the ability of next and previous 3D models to represent the current GOP. The expected consistency of each 3D model with previous and next 3D models is thus achieved.

### C. Extended sliding adjustment

The presented method can be extended to take into account more than 3 successive models. This is done by adding cost functions similar to  $f_3$ . Consider that we want to increase consistency with the  $p$  next keyframes (in the previous section  $p = 1$ ). We generalize function  $f_3$  into  $g_n$ , which gives the contribution of keyframe  $K_{k+n}$  into the total cost function:

$$g_n(P_{k+n+1}, \{M_i^k\}, \dots, \{M_i^{k+n}\}) = \sum_{q=0}^n \sum_i \|m_i^{k+q, k+n+1} - P_{k+n+1} \cdot M_i^{k+q}\|^2 \quad (12)$$

under the constraints:

$$M_i^{k+q} = O_{k+q} + \lambda_i^{k+q} \cdot \vec{u}(m_i^{k+q}) \quad q = 0, n \quad (13)$$

For  $n = 1$ ,  $g_1$  is the contribution of  $K_{k+2}$  and is equal to  $f_3$ .

The final cost function which takes into account  $p$  keyframes and 3D models is:

$$g(P_{k+1}, \{M_i^k\}, P_{k+2}, \{M_i^{k+1}\}, \dots, P_{k+p+1}, \{M_i^{k+p}\}) = f_1(P_{k+1}, \{M_i^k\}) + f_2(P_{k+1}) + \sum_{n=1}^p g_n(P_{k+n+1}, \{M_i^k\}, \dots, \{M_i^{k+n}\}) \quad (14)$$

The extension to  $p$  3D models ensures consistency of the 3D models for a longer time in the video. However, without a precise set of camera parameters, a larger window may degrade the estimation results.

## VI. RESULTS

### A. Textured 3D model generation

We just present the final step for 3D models generation. The visualization step is also adapted to our representation with local 3D models.

For each pair of successive keyframe images  $K_k$  and  $K_{k+1}$  a 3D model  $\mathcal{M}_k$  is computed. Since we have a dense motion field  $D_k$ , camera internal parameters  $A$  and camera motion parameters  $R_k, t_k$  and projection matrices  $P_k, P_{k+1}$  from sliding adjustment, this is fairly simple. For any pixel  $m$  in  $K_k$ , its corresponding position is given by  $m' = m + D_k(m)$ , and the 3D point is recovered by solving projection equations (1). A dense depth map is then constructed (see fig. 8(b)). In order to have a 3D model which can be easily visualized, only vertices of a regular 2D triangular mesh are reconstructed. The reconstructed points define a continuous 3D triangular mesh. This mesh is textured using image  $K_k$ .

This 3D model is simply described in a format Rec3D quite similar to VRML format. A Rec3D file is then generated, which can be interactively visualized in real time with classical 3D rendering libraries like OpenGL [25].

### B. Adapted visualization through 3D model fading

The proposed representation with the set of 3D models  $\{\mathcal{M}_k\}$  can reconstruct the input video sequence. However, some artifacts appears in the reconstructed video sequence, in particular when we switch from one 3D model to another. They are mostly due to illumination changes between 2 keyframes, occluded areas and to the accuracy of the 3D model to reconstruct the GOP. In order to take into account such problems, we propose an original 3D model fading technique. This technique is an approximation of real 3D model morphing. It is a simple way to take into account not only texture changes but also the geometric changes from one 3D model to the next in the stream. A two passes rendering is performed, a first pass using the current 3D model, and a second using the next 3D model. The resulting reconstructed sequences are then blended with respect to

the  $\alpha$  factor defined by:

$$\alpha(C) = |O_k - C| / |O_k - O_{k+1}| \quad (15)$$

where  $O_k$  and  $O_{k+1}$  are camera centers for first and last image of the GOP and  $C$  is the current camera position. Thus the  $\alpha$  factor balances the reconstructed sequences contribution from first to last image of the GOP, in proportion with the distance of the current camera position from the keyframes camera positions.

### C. Results on real video sequences

The proposed algorithm has been implemented and tested on several real video sequences. Reconstruction of the original video sequence with the 3D models stream is easily done, but is not discussed here since it is similar to classical image interpolation. So, we test the accuracy of the 3D models with classical applications of 3D model manipulation: free navigation, illumination changes, augmented reality, stereo visualization. We perform visual quality estimation. The corresponding video sequences can be found at <http://www.irisa.fr/temics/Demos/3D4> showing some applications of the representation (compression aspects are also shown but not discussed here).

The first is a *street* video sequence which is a walk in a city with a global translation along z-axis (fig. 7). Sequence has been acquired with a mechanical stabilizer, internal parameters are unknown and fixed (see III-B), the focal distance is approximated to 500. The dense motion field between image 0 (keyframe 0) and image 40 (keyframe 1) is shown on figure 8(a). The vector scale is 0.25 and the average motion is 41 pixels. The corresponding depth map on image 0 is shown on figure 8(b) (the furthest areas from the camera are in dark). We see that the depth map is quite regularized and reconstructed the global shape of the street with details on the relief (pot plant, street lamp, etc.). The corridor shape of the street is more visible on the figure 9, where we clearly see the planarity of the floor and the right angle with the 2 walls on the sides with the floor.

The second test sequence is the *stairway* sequence which is a walk with a global translation along x-axis (fig. 14). The video sequence has been acquired with a simple hand-held camcorder, internal parameters are unknown and fixed, the focal distance is again approximated to 500. This video sequence is quite harder than the previous one due to uniform

textures and water in the fountain which confuses the motion estimation. The figure 15(a) and 15(b) show motion field and depth map for *stairway* video sequence between image 71 and image 83: vector scale is 0.25 and the average motion is 40.2 pixels. We see that we obtain a valid scene geometry: we find the trees and the shape of the stairway. The figure 16 shows in details the geometry of the scene as depth maps.

### C.1 Free navigation

Free navigation in the representation is performed by just specifying new camera positions. The figure 10 shows a texture image (used for the texture mapping) and a virtual view generated with this image and the corresponding 3D model: we simulate a virtual walker which performs few steps on the left and turn the head to the right. We see that the window and the gate in the background are occulted by the left wall, according to the scene geometry. The figure 11 shows 2 other virtual views taken far away from original viewpoints: (a) shows the image produces with a large rotation to the left and (b) a view in details of the scooter. Texture stretching is visible on surfaces which are not visible with a frontal view or which are in occulted areas.

### C.2 Lightning modification

Lightning modification is performed with classical illumination algorithm (see [25] for details). Contrary to global illumination changes of a video sequence, the illumination takes into account 3D information such as distance from the light to the surfaces, giving a better realism. Lightning modification on the street video sequence are presented on figure 12: a headlight have been added to change illumination in the scene. The figure shows 2 images of the reconstructed video sequence. We notice the illumination is consistent with the geometry, i.e. decreasing with the distance from the light to the surface. We also notice that illumination is invariant from one 3D model (a) to another (b) and some artifacts on the top street lamp due to 3D mesh continuity.

### C.3 Object insertion

Insertion of virtual objects in the video sequence taking into account depth and occlusions is easy with the representation, contrary to classical 2D object insertion. The figures

13 and 17 show a virtual sphere added in the sequence *street* and *stairway*. The sphere is placed to intersect the scene, showing the occlusions accuracy. We can also notice the good stability of the sphere's position along the video sequence. Moving objects can also be inserted, taking into account depth information given by the representation.

#### C.4 Stereo sequence generation

The generation of stereoscopic video sequences just requires to reconstructed the scene twice (for the left and the right eyes) with a small shift. The figure 18 shows such a pair of images extracted from the *stairway* video sequence. The video sequence is visualized on specific device or with eyes' defocus technique for still images. Stereo visualization has been successfully tested on stereo display.

## VII. CONCLUSION

We have proposed an automatic scheme to extract a stream of 3D models from a video sequence of a fixed scene. The presented scheme offers a good compromise for 3D reconstruction from any video sequence of static scenes when the reconstruction of a unique 3D model is not possible or desirable: this is the case for very long video sequences (computation complexity) which require on-the-fly analysis, or when camera motion is not appropriate for a unique 3D reconstruction (forward translation).

We have proposed a simple technique for automatic video sequence clustering which allows to reconstruct several 3D models. These 3D models are then computed using a sliding adjustment which allows to keep most of the functionalities of an unique 3D model. We have validated our approach on several video sequences.

Such approaches could be extended to any video sequences (not only static scene) where objects segmentation are known. Moreover, as in MPEG4 Sprite coding, very low bitrate coding might be achieved with such a representation. We thus plan to study the coding performance of our approach compared to standard video coders.

## VIII. ACKNOWLEDGMENT

We thank Stphane Pateux for the fruitful discussions and his cooperation regarding the use of his algorithms. This work was supported by the RNRT project V2NET.

## REFERENCES

- [1] Shenchang Eric Chen, “QuickTime VR — an image-based approach to virtual environment navigation,” *Computer Graphics*, vol. 29, no. Annual Conference Series, pp. 29–38, 1995.
- [2] ISO/IEC JTC1/SC29/WG11, “Generic coding of audio-visual objects : Visual,” Information Technology N2802, ISO/IEC, July 99, Final Proposed Draft Amendment 1.
- [3] Huang T. S. and Tang L., “Very low bit-rate video compression using 3d models,” in *Proc. IWSNHC3DI*, 1995.
- [4] B. Girod and al., “3d image models and compression - synthetic hybrid or natural fit?,” in *Proc. ICIP*, Oct. 1999.
- [5] F. Prêteux and M. Malciu, “Model-based head tracking and 3d pose estimation,” in *Visual Conference on Image Processing*, San Jose, californie, 1998, pp. 94–110.
- [6] ISO/IEC JTC1/SC29/WG11, “Mpeg-4 animation framework extension (afx) vm 3.0,” Tech. Rep. N4020, ISO/IEC, Mar. 2001.
- [7] H.-Y. Shum and R. Szeliski, “Stereo reconstruction from multiperspective panoramas,” in *Seventh International Conference on Computer Vision*, Sept. 1999.
- [8] M. Pollefeys, R. Koch, M. Vergauwen, B. Deknuydt, and L. Van Gool, “three-dimensional scene reconstruction from images,” in *proceedings SPIE Electronic Imaging, Three-Dimensional Image Capture and Applications III*, 2000.
- [9] D. Nistèr, “Reconstruction from uncalibrated sequences with a hierarchy of trifocal tensors,” in *Proc. ECCV*, Dublin, Ireland, 2000, pp. 649–663.
- [10] Realviz, “Realviz,” URL <http://www.realviz.com>, 2001.
- [11] S.J. Gortler, R. Grzeszczuk, R. Szeliski, and M.F. Cohen, “The lumigraph,” in *Proc. SIGGRAPH*, Aug. 1996.
- [12] R. Kochand M. Pollefeys and L. Van Gool, “Realistic 3d scene modeling from uncalibrated image sequences,” in *Proc. ICIP*, Kobe Japan, Oct. 1999, Invited contribution to special session on Image Analysis and Synthesis.
- [13] J.-R. Ohm and K. Müller, “Incomplete 3d - multiview representation of video objects,” *IEEE Transaction on Circuits and Systems for Video Technology*, vol. 9, no. 2, pp. 389–400, mar 1999.
- [14] Andrew W. Fitzgibbon and Andrew Zisserman, “Automatic camera recovery for closed or open image sequences,” in *ECCV (1)*, 1998, pp. 311–326.
- [15] G. Marquant and S. Pateux, “Mesh and ”crack lines”: Application to object-based motion estimation and higher scalability,” in *Proc. ICIP*, Sept. 2000.
- [16] C. Harris and M. Stephens, “A combined corner and edge detector,” in *Proc. 4th Alvey Vision Conf.*, 1988.
- [17] Franck Galpin and Luce Morin, “Computed 3d models for very low bitrate video coding,” in *Proceedings of the IEEE conference on Visual Communications and Image Processing, VCIP'2001*, Jan. 2001, vol. 4310.
- [18] Franck Galpin and Luce Morin, “Video coding using streamed 3d representation,” *Calculateurs parallèles, Réseaux et systèmes répartis, numéro spécial : Image et vidéo, Editeur Hermes*, vol. 12, no. 3-4, pp. 431–442, 2000.
- [19] Sylvain Bougnoux, “From projective to euclidean space under any practical situation, a criticism of self-calibration,” in *IEEE International Conference on Computer Vision*, 1998, pp. 790–796.
- [20] P.H.S. Torr and D.W.Murray., “The development and comparison of robust methods for estimating the fundamental matrix,” *IJCV*, vol. 24, no. 3, Sept. 1997.
- [21] O. D. Faugeras, Q.T. Luong, and S.J. Maybank, “Camera self calibration: Theory and experiments,” in *Proc. ECCV*, Santa Margherita Ligure, Italy, June 1992.

- [22] D.F. DeMenthon and L.S. Davis, “Model-based object pose in 25 lines of code,” *International Journal of Computer Vision*, vol. 15, pp. 123–141, June 1995.
- [23] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon, “Bundle adjustment – a modern synthesis,” in *Proc. Vision Algorithms*, Oct. 1999.
- [24] J. Moré, B. Garbow, K. Hillstom, and M. Report, “User guide for minpack-i,” Tech. Rep. ANL-80-74, Argonne National Laboratory, 1980.
- [25] Jackie Neider, Tom Davis, and Mason Woo, *OpenGL Programming Guide*, Addison-Wesley, Reading MA, 1993.

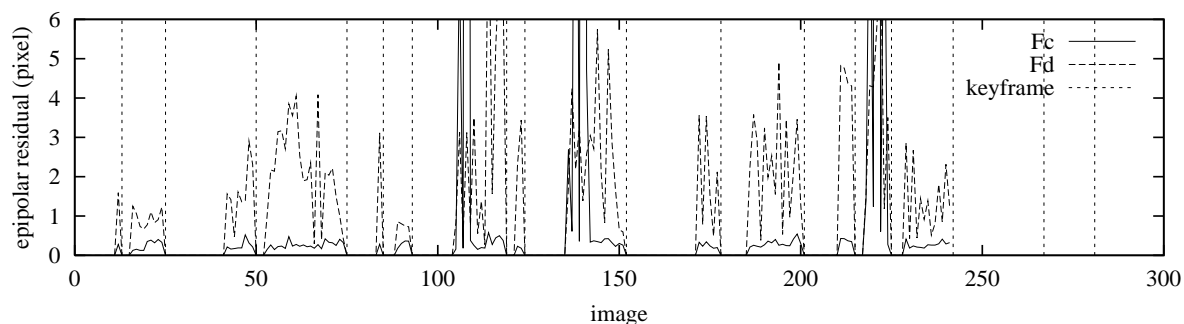


Fig. 3. Epipolar residuals along the *stairway* video sequence after the calibration step ( $F_c$ ) and after the localization step ( $F_d$ ).

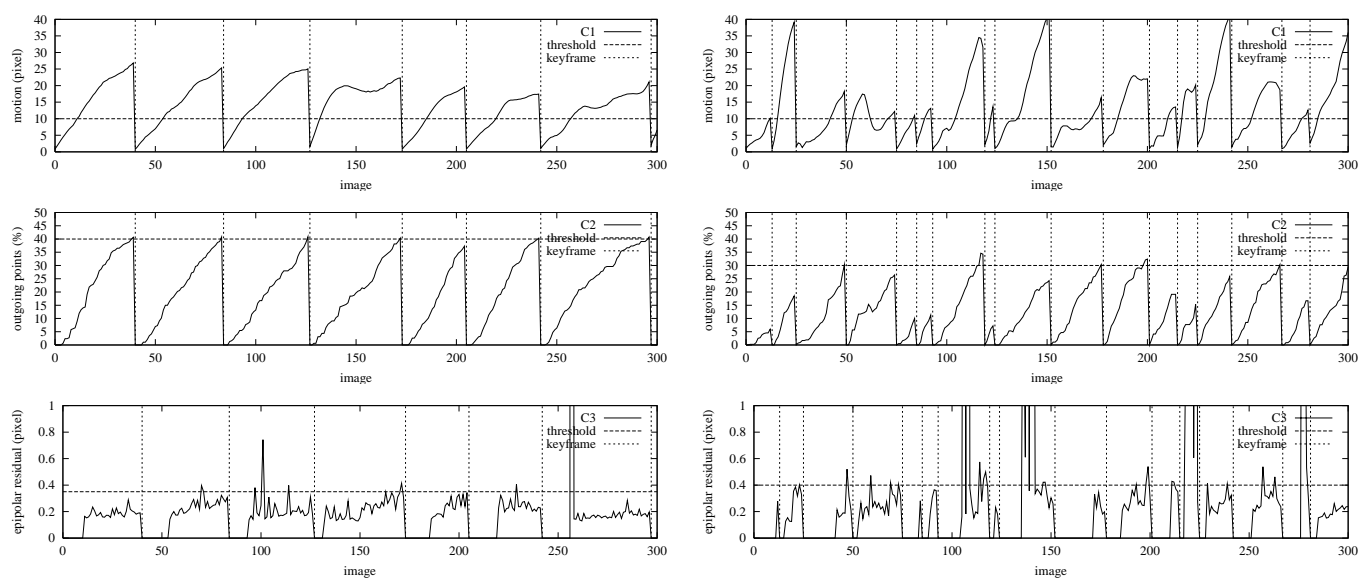


Fig. 4. Figures shows the evolution of the 3 criteria on the *street* video sequence (left) and *stairway* video sequence (right).

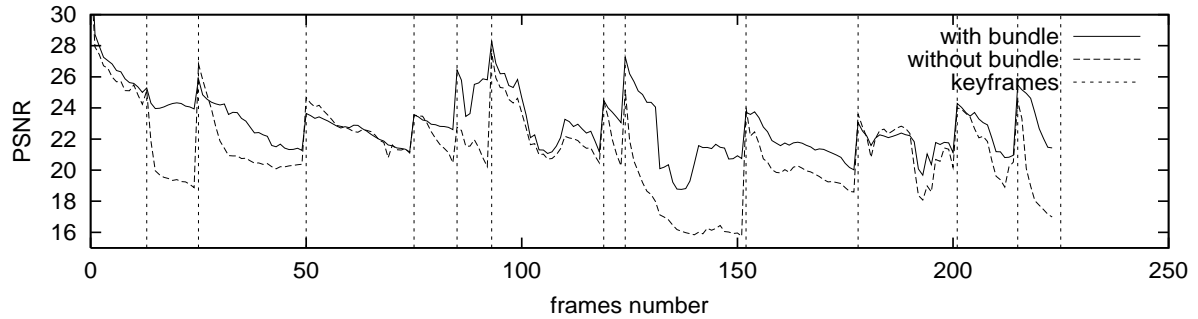


Fig. 5. PSNR value for reconstruction of *stairway* video sequence. Each GOP  $k$  is reconstructed with the previous corresponding 3D model  $\mathcal{M}_{k-1}$ . Method with sliding bundle adjustment have a better quality, showing a better consistency of previous 3D model with the current GOP.

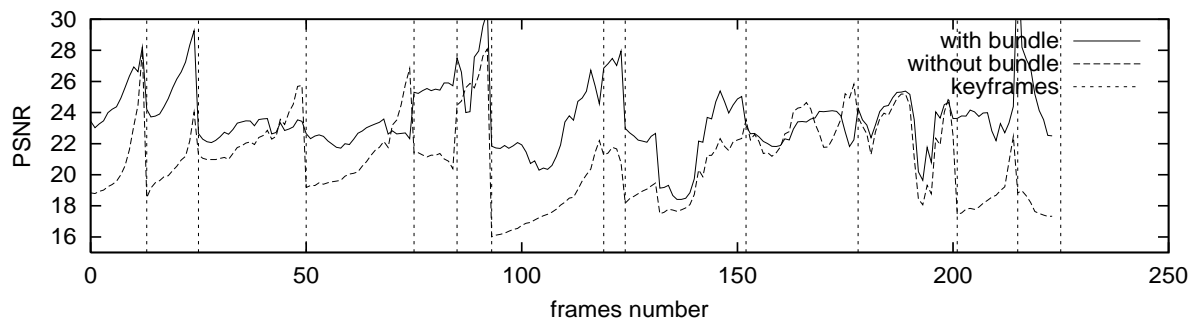


Fig. 6. PSNR value for reconstruction of *stairway* video sequence. Each GOP  $k$  is reconstructed with the next corresponding 3D model  $\mathcal{M}_{k+1}$ . Method with sliding bundle adjustment have a better quality, showing a better consistency of next 3D model with the current GOP.



Fig. 7. Original *street* sequence. From top-left to bottom-right, image 0, 40, 80, 120, 160, 200

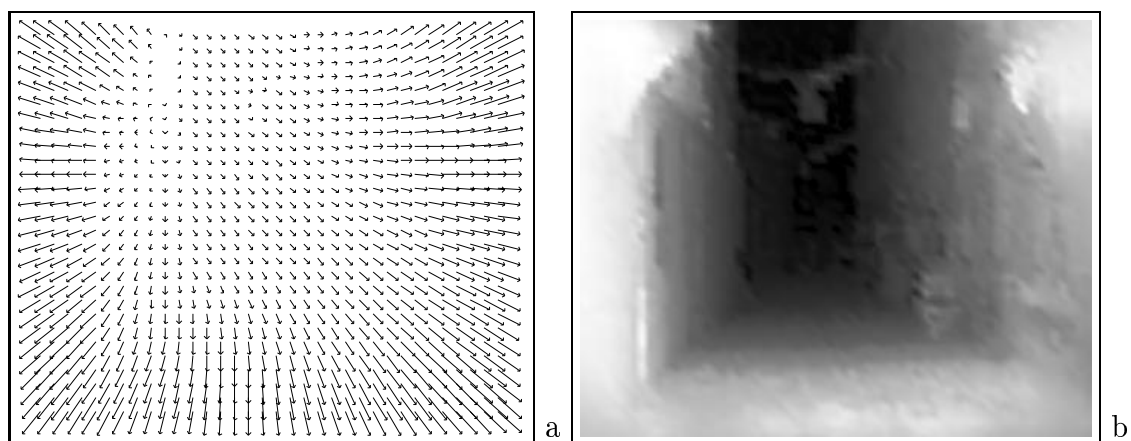


Fig. 8. (a) Dense motion field between image 0 and 40 in *street* sequence. (b) Depth field extract from dense motion field.

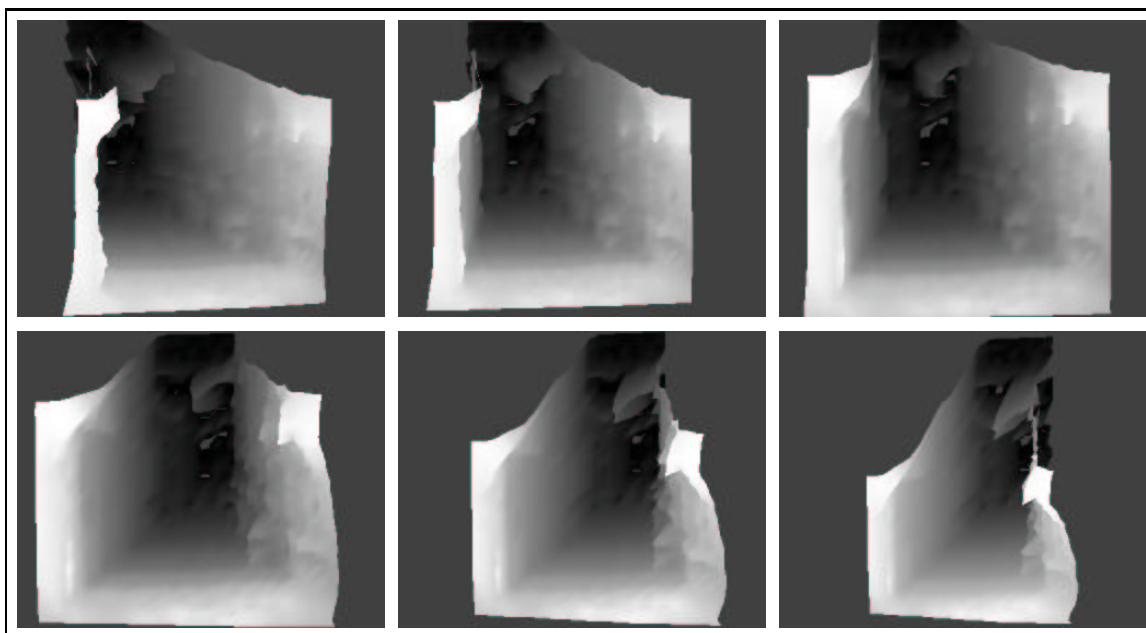


Fig. 9. Rotation around a 3D model automatically extracted from *street* sequence. The images are the depth map of each views.



Fig. 10. (a) Texture image used with 3D model 4 extracted from the *street* video sequence , (b) a virtual view generated with this texture image: the virtual walker performs few steps on the left and turn the head to the right.



Fig. 11. Two virtual views of *street* video sequence generated far away from original viewpoints.



Fig. 12. Image 0 and 40 from the *street* sequence: a headlight is added to reconstruct the original video sequence.



Fig. 13. A virtual sphere is added: (a) view with extracted 3D model 1, (b) view with extracted 3D model 2. The occlusions are taken into account, and the sphere position is quite stable from one 3D model to the next one.



Fig. 14. Original *stairway* sequence. From top-left to bottom-right, image 50, 75, 100, 125, 150, 175

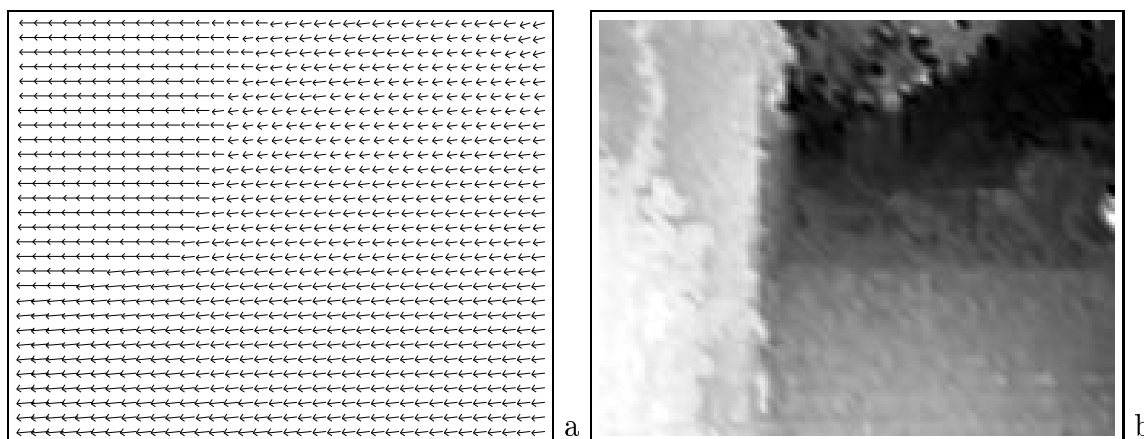


Fig. 15. (a) Dense motion field between image 71 and 83 in *street* sequence. (b) Depth field extract from dense motion field.

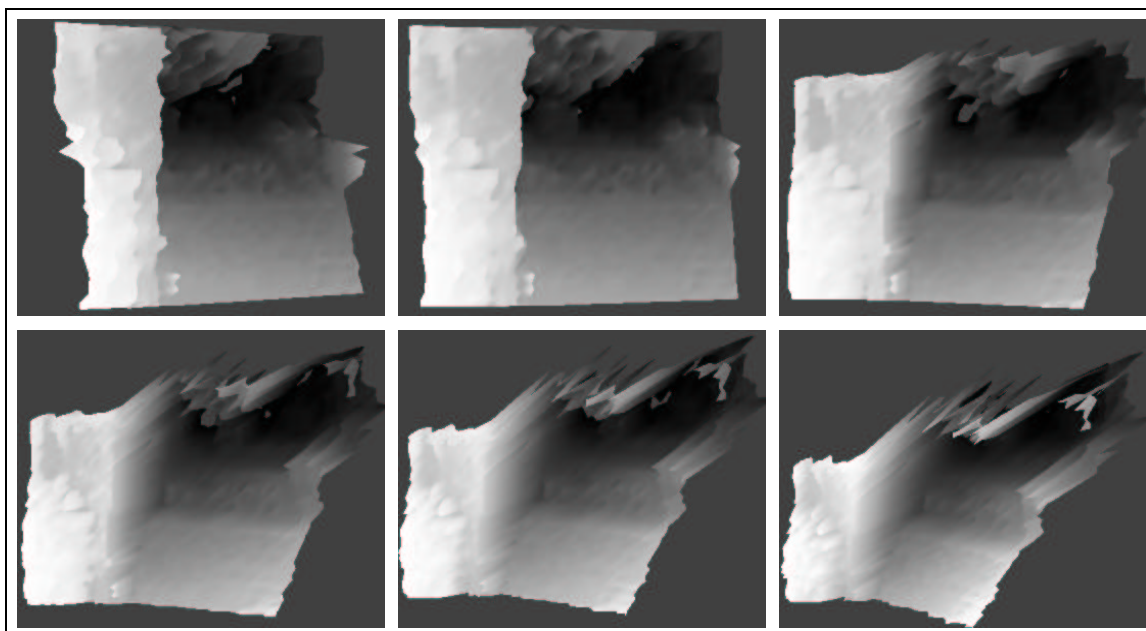


Fig. 16. Rotation around a 3D model automatically extracted from *stairway* sequence. The images are the depth map of each views.

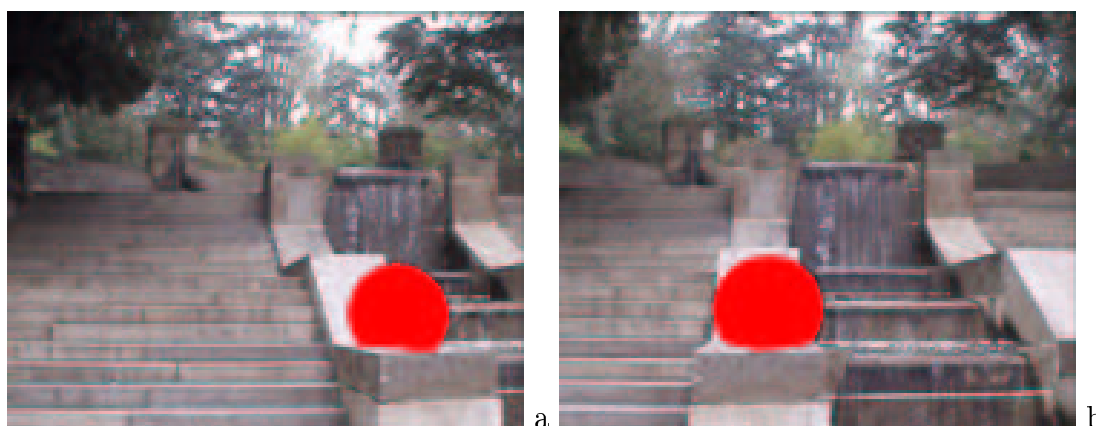


Fig. 17. A virtual sphere is added in the scene. (a) view with extracted 3D model 10, (b) view with extracted 3D model 15. The occlusions are taken into account, and the sphere position is stable from one 3D model to the next one.

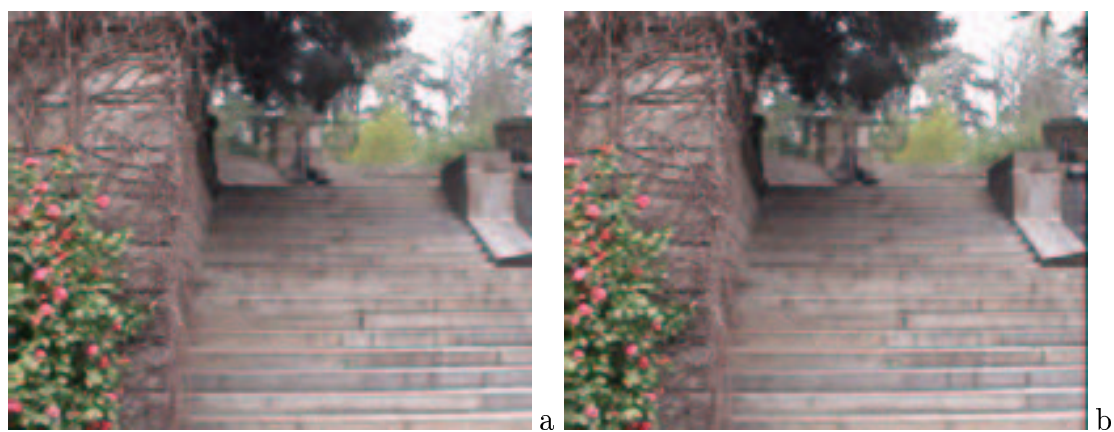


Fig. 18. Left (a) and right (a) views of an image extracted from the reconstructed stereoscopic video sequence.