



HAL
open science

Iterative Multi-Planar Camera Calibration: Improving stability using Model Selection

Javier-Flavio Viguera-Gomez, Marie-Odile Berger, Gilles Simon

► **To cite this version:**

Javier-Flavio Viguera-Gomez, Marie-Odile Berger, Gilles Simon. Iterative Multi-Planar Camera Calibration: Improving stability using Model Selection. Vision, Video and Graphics - VVG'03, 2003, Bath, UK, 8 p. inria-00099483

HAL Id: inria-00099483

<https://inria.hal.science/inria-00099483v1>

Submitted on 26 Sep 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Iterative Multi-planar Camera Calibration: Improving Stability using Model Selection

J.F. Vigueras,¹ M.-O. Berger¹ and G. Simon¹

¹ LORIA/INRIA Lorraine
BP 101, 54602 Villers les Nancy, France
email: {vigueras,berger,gsimon}@loria.fr

Abstract

Tracking, or camera pose determination, is the main technical challenge in numerous applications in computer vision and especially in Augmented Reality. However, pose computation processes commonly exhibit some fluctuations and lack of precision in the estimation of the parameters. This leads to unpleasant visual impressions when augmented scenes are considered. In this paper, we propose an efficient and reliable method for real time camera tracking which avoid unpleasant statistical fluctuations. This method is based on the knowledge of a piecewise planar structure in the scene and makes use of model selection to reduce fluctuations. Videos are attached to this paper which prove the effectiveness of our approach.

1. Introduction

Augmenting real video sequences of a scene with computer generated objects is one of the main goals of many applications such as virtual museums, interactive interior design or architectural design, computer-aided repair and learning systems¹⁵. All these interactive applications require that the augmented scene is continually updated as the user moves about the real scene. Hence, one of the most basic challenge to overcome is the registration problem: the objects in the real and the virtual world must be properly aligned with respect to each other or the illusion that the two worlds coexist will be compromised.

In this paper, we address the registration problem for interactive AR applications. Such applications require sequential and real-time registration process. Though the registration problem has received a lot of attention in the computer vision community, the problem of real time registration is still far from a solved problem. Ideally, an AR system should work in all environments without the need to prepare the scene ahead of time and the user should walk anywhere he pleases. In the past, several AR systems have achieved accurate and fast tracking and registration, putting dots over objects and tracking the dots with a camera^{7,8}. However, such methods restrict the flexibility of the system. Hence, there is a need to investigate registration methods which work in un-

prepared environments and which reduce the need to know the geometry of the objects in the scene.

1.1. State of the art

Today, the approaches to sequential viewpoint computation can be divided in two main categories: model-based approach or move-matching approach. Model-based techniques rely on the identification in the images of features from the object model. Hence, a direct correspondence between the 3D object-coordinate system and each image is set up^{8,10}. This capability of treating each image independently makes such methods more appropriate for real time implementations. Another consequence of model-based tracking is the absence of drift. However, it is commonly true that few features are available for registration. Moreover, noise in the image measurements hampers their accurate detection and consequently corrupts the estimated pose. As a result, the tracking suffers from high-frequency jitter. More importantly, such methods require significant manual intervention to construct the model.

On the other hand, new move matching methods¹¹ attempt to compute the relative motion between two successive frames using planar structures. If the position of the camera in the first frame is known, the absolute position of the camera is obtained by composing each relative motion.

These systems are attractive because they do not require any knowledge on the scene. However, they can suffer from drift because errors accumulate over time.

In interactive real time applications, a good way to assess the viewpoint accuracy is to consider the visual impression of the augmented scene. Today, it appears that statistical fluctuations in the viewpoint computations lead to unpleasant jittering effects or to sliding effects in the scene. This problems are particularly conspicuous when the motion of the camera is small because noise of the extracted features lead to large fluctuations in the viewpoint computation. The problem of stabilisation was addressed in ⁴. The idea is to classify the typical movements of the camera into models (stationary, panoramic, general, zoom in) in order to fix some of the parameters assuming that their variations are due to statistical fluctuations. Of course, stability and accuracy over the remaining parameters are better because the degree of freedom of the function to be optimised is smaller. In this paper, we investigate further this idea with the following contributions: (i) we propose a method for viewpoint computation which is based on the observation of a multiplanar structure in the scene. Such structures are quite common both for indoor and outdoor scenes; (ii) following Kanatani⁴, we investigate the use of model selection to improve the stability of the computed viewpoint. Various model selection criteria are considered and tested in this study. We show that the ones which make use of the covariance on the estimated parameters give better results than the classical criteria; (iii) the effectiveness of our pose algorithm is assessed on various sequences.

The method for multiplanar viewpoint computation is given in section 2. Section 3 exhibit results and strategies for model selection. Finally, various snapshots of augmented scenes are provided.

2. Multiplanar viewpoint computation

2.1. Overview

This section gives an overview of our registration method. The equations of the planes used by the registration process are given by the user. In our approach, the intrinsic parameters are supposed constant and are computed beforehand. The first camera pose is also estimated. Often, this estimation is obtained by using a poster in the scene.

Once the preprocessing stage has been achieved, the registration follows a four step loop: key-points are extracted and matched from frame to frame. Then, for each model, the projection matrix is computed using constraints induced by the homographies. Finally the right motion model is selected and the viewpoint is computed accordingly. In the following, the main steps of this algorithm are described with further details (Fig. 1).

Initialisation stage:

1. Give the equation of the observed planes used for registration,

2. Compute the projective matrix for the first frame P^0 ,

Computation of the projective matrix P^i for $i > 0$:

1. Compute the set of matched key-points between images $i - 1$ and i for each observed plane.

2. For each motion model, compute P_i from P_{i-1} using the constraints induced by the homographies

3. Select the best model according to the selection criterion which is a tradeoff between accuracy and simplicity of the model

4. Compute the motion using the selected model

Figure 1: Overview of the multi-planar tracking method

2.2. Planar viewpoint computation

We assume that the position, orientation, and the internal parameters of the camera are known for the first image. Then all the images of the sequence may be related to the preceding one by setting up correspondences between points.

We know that given two projection matrices $P_1 = [I|0]$ and $P_2 = [A|a]$ and a plane defined by the vector v such that $v^T X + 1 = 0$, the corresponding homography matrix ⁵ can be expressed as:

$$H = K_2(A - av^T)K_1^{-1} \quad (1)$$

where $K_i = \begin{pmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{pmatrix}$ is the matrix of intrinsic parameters for the image i .

Given the intrinsic parameters and a set of matched points (x_j, x'_j) on the considered plane between two images, a classical procedure to get the viewpoint parameters is to minimise the mean error of the matched points with respect to the transformation ⁴, thus:

$$A, a = \arg \text{Min} \left(J(\mathbf{A}, \mathbf{a}) = \frac{1}{N} \sum_{j=1}^N \|x'_j - Z(Hx_j)\|^2 \right)$$

where $Z(\cdot)$ denotes normalisation to make the third component 1.

2.3. Multi-planar calibration

Our experiments proved that the accuracy of the single plane registration method is not sufficient to obtain a good visual impression of the augmented scene. Indeed the accuracy depends on the relative position of the camera and of the observed plane and also on the number of matched points. Moreover, as sequential viewpoint computation is considered, errors are accumulated over time and the viewpoint pa-

parameters tend to diverge from the real ones especially when large sequences are considered.

That is the reason why we suggest to use several planes because it brings more information about the tridimensional space and reduces considerably the variability of the estimated calibration parameters. It will also help us to handle large environments for AR applications.

When several planes are considered, the function to be optimised is:

$$J(\mathbf{A}, \mathbf{a}) = \frac{1}{N_1 + \dots + N_n} \sum_{k=1}^n \sum_{j=1}^{N_k} \|x'_{kj} - Z(H_k x_{kj})\|^2$$

where n is the number of planes, N_k the number of points belonging to the plane k , H_k the corresponding homography.

A typical method used to minimise this non-linear function is Newton iterations but it is very sensitive to the initial estimation. Thus, we use the Levenberg-Marquardt method being more stable and almost as fast as the Newton method.

2.4. Results

To prove the effectiveness of the approach, we considered a **synthetic image sequence** using the model of our three-plane calibration target. Different motions were considered: x and y translations, panoramic motion, and stationary motion. Gaussian noise with covariance matrix $\sigma^2 \mathbf{I}$ ($\sigma = 0.5$ pixels) was added to the image points.

In Figure 2, we compare the actual translation coordinates T_z and the computed coordinates when a single plane, 2 planes, and 3 planes, are used. The viewpoint is found computing the 6 extrinsic parameters and fixing all the intrinsics to a pre-calibrated value. These results show that using a single plane, the estimated viewpoint is very unstable and the estimated coordinates are lacking of precision. By adding a second and a third plane both precision and regularity are improved considerably.

2.5. Improving the robustness of the viewpoint computation

It is well known that false matchings $x_j \leftrightarrow x'_j$ can severely disturb the viewpoint computation process, and RANSAC algorithm is classically used for every visible plane in order to discard false matchings on each one. However such an approach considers the planes independently and does not take into account the multi-planar model of the scene. We propose here two methods to cope with this problem:

Method 1 We first use an iterative method to refine the set of inliers. As the multi-planar model of the scene is available, the homography induced by each plane can be obtained from the projection matrix (eq. 1) computed with the multi-planar algorithm. This allows us to get the set of inliers compatible with the estimated homographies. The projection matrix is then computed from this new set of inliers, which is

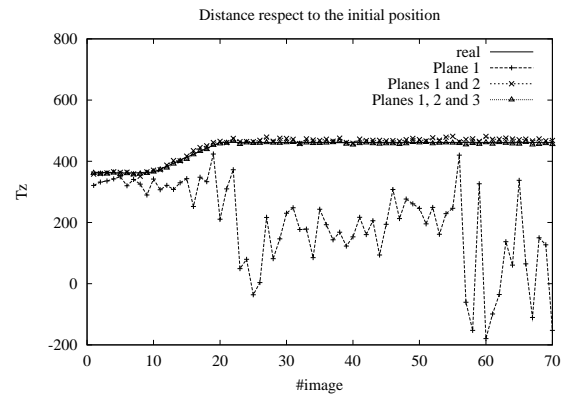


Figure 2: The computed Z-translation through the sequence using one, two or three planes.

in better agreement with the scene geometry, and the process is iterated until convergence. This way, false matchings may be removed and new matchings can be added.

However, this approach may fail if a small number of points is available on a given plane. In this case, the RANSAC algorithm is not always able to select the right points, then the first estimate of the viewpoint may be erroneous and also the obtained inliers set is. To cope with this problem, we suggest an approach which fully integrates the multi-planar model in the robust estimation process:

Method 2

1. Four point correspondences are randomly chosen in the full set of matched points (the union of the matchings for all visible planes).
2. The viewpoint is estimated from these four correspondences using the multi-planar method.
3. The homography induced by each plane is computed from the projection matrix and from the known plane equations using eq. 1.
4. A new set of inliers is obtained for each plane. This is the union of the correspondences which are in agreement with the computed homography in each plane.
5. Repeat 1 to 4, L times (the number of samples is chosen according to the law $L = \log(1-p)/\log(1-(1-\epsilon)^4)$, where $p = 0.99$ and the proportion of outliers is at most $\epsilon = 0.3$ in our experiments.)
6. Select the parameters corresponding to the biggest set of inliers.

A **turntable sequence** was considered to prove the efficiency of these two iterative methods (and classical RANSAC over planes), and to assess the accuracy of the viewpoint algorithm. The table we have used does micro-metrical precision 1-D translations and rotations around one axis; hence, it is possible to know the real displacements of the camera being fixed to it. In this experiment, we consider a

closed sequence which is described in Fig. 3.a and b. Fig. 3.c exhibits the computed translation along the Z axis when the three described methods are used. As the sequence is closed, a good way to assess the accuracy is to check if the final position is the same as the initial one. Fig. 3.c clearly shows that using method 1, the final position is very close to the initial one. For method 2, the difference between the initial and the final position is smaller than using the classical approach. To consider the influence of the method on the visual impression of the augmented scene, a cube is added to the scene and is shown in the first image of the sequence in Fig. 4.a. The three other images show the augmented scene in the final position, which is the same as the initial one, when the three matching methods are used: classical, method 1 and method 2 (Fig. 4.b,c and d). These snapshots proved that sliding effects occurred when the classical matching method is used. The use of methods 1 and 2 clearly improved the viewpoint accuracy with noticeable better results for method 2. However, as the computational cost of method 2 is very high, we prefer to use the method 1 which has a good compromise between computational time and accuracy.

3. Use of stabilisation methods

3.1. Aims

Even when the precision of the viewpoints is improved by considering several planes, fluctuations in the parameters are often observed and may lead to unpleasant visual impressions such as jittering or sliding when augmented scenes are considered. These fluctuations are especially conspicuous when the camera motion is small because of noise and imprecision in computing the points coordinates. In the past, several papers used Kalman filtering for prediction and stabilisation task. However, the use of a Kalman filter is not always advantageous for AR. This is because a low order dynamical model of human motion may not be always appropriate except under very constraint scenarios.

Following Matsunaga⁴ and Torr¹³ we investigate the use of motion model selection to reduce fluctuations of the camera parameters and to improve the visual impression of the augmented scene. The underlying idea in model selection is as follows: a higher order motion model fits any data set more accurately than a lower order model. However, high order models fit part of the random noise they are supposed to remove. Thus, a high order model, although accurate, is less stable to random perturbations in the data. A good motion model must strike the right balance between accuracy and stability. The model selection principle demands that the model should explain the data very well and at the same time have a simple structure.

3.2. State of the art

Many model selection criteria for balancing the residual and the degree of freedom of the model have been proposed in

the literature³. All of them are the sum of an accuracy criterion and of a term which is a measure of the complexity of the model. Most of them are based on statistical and information-theoretic criteria. Among them, the most widely used criterion are the geometric Akaike's criterion and the minimum description length (MDL) criterion. The AIC criterion can be viewed as an approximation of an entropy criterion (the Kullbak-Leibler distance) whereas the MDL criterion try to choose the model that minimises the number of bits required to express the model:

$$G_{AIC} = \hat{J} + 2k\epsilon^2$$

$$G_{MDL} = \hat{J} - k\epsilon^2 \log \epsilon^2$$

where k is the degree of freedom of the motion. The square noise level ϵ^2 can be estimated from the residuals \hat{J} (the one corresponding to the highest order model^{4 2}). Whatever the considered criterion, the use of a too complex model is penalised with respect to simpler model.

Kanatani⁴ previously applied this technique to the calibration problem by using a single plane, specifically, a calibration pattern. In this seminal work, Kanatani classifies the movements in six types, specifically those with fixed focal length are four:

| Movement | Known parameters | Variables |
|----------------|--|------------------------------|
| stationary | $\mathbf{A}_i = \mathbf{A}_{i-1}, \mathbf{a}_i = \mathbf{a}_{i-1}$ | — |
| panoramic | $\mathbf{a}_i = \mathbf{a}_{i-1}$ | \mathbf{A}_i |
| t —predicted | $t_i = 2t_{i-1} - t_{i-2}, \mathbf{a}_i = \mathbf{a}_{i-1}$ | \mathbf{A}_i |
| general model | — | $\mathbf{A}_i, \mathbf{a}_i$ |

In the t —predicted model, the camera position is linearly extrapolated and the optimisation is only performed on the rotation.

In Kanatani's approach, only two criteria are considered: G_{AIC} and G_{MDL} . However, there are many other criteria, especially those which make use of the covariance matrix or the information matrix on the estimated parameters.

3.3. Our approach to model selection

We suggest to use together the model selection strategy and the multi-planar calibration in order to improve the stability and the accuracy of the estimated parameters. There are different branches using model selection, but there is no such successful criterion in general, as can be seen in some reviews comparing some of them for different problems: finding the polynom's degree², surface merging³, type of motion¹², detection of geometric primitives⁶.

For this reason, we compare different model selection criteria (Table 1). These criteria use the same accuracy measure: the residual error evaluated at the maximum likelihood parameters. But, the complexity term is different for each model and depends on different assumptions over the parameters and their distribution. In this work, we especially

| Frame | Movement |
|-----------|---------------------------|
| 0 – 5 | Stationary |
| 5 – 25 | Rotation +10° |
| 25 – 65 | Rotation –20° |
| 65 – 75 | Translation 10cm. (left) |
| 75 – 85 | Stationary |
| 85 – 115 | Rotation +15° |
| 115 – 125 | Translation 10cm. (right) |
| 125 – 135 | Rotation –5° |
| 135 – 140 | Stationary |

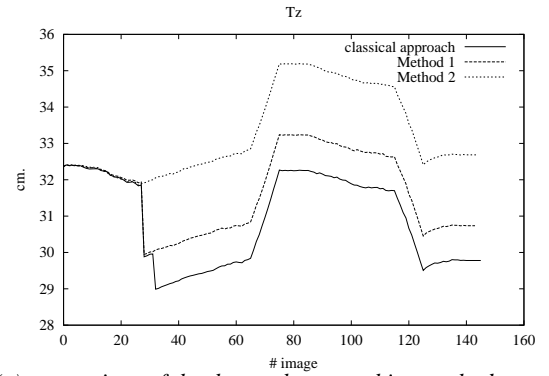
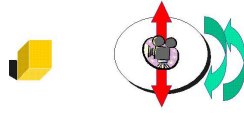


Figure 3: Turntable sequence: (a) and (b): actual motion of the camera, (c) comparison of the three robust matching methods on the turntable sequence.

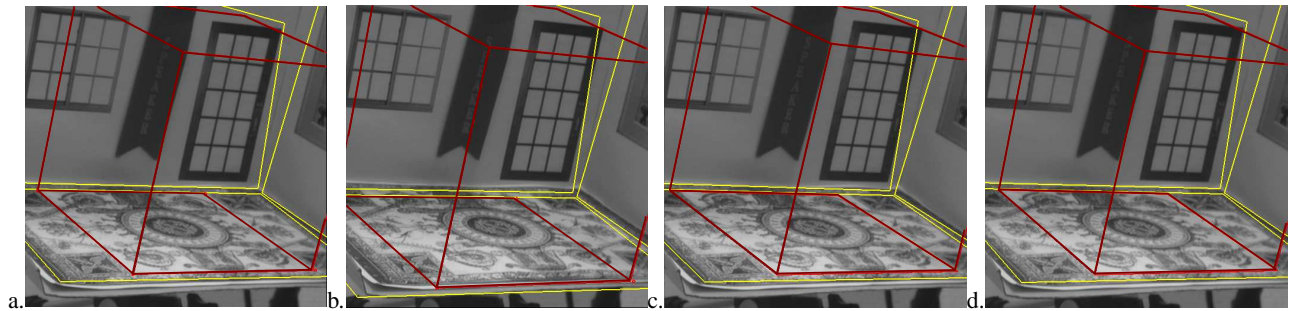


Figure 4: Comparison of the final position when the three matching methods are used: initial position (a); final position when the classical method (b), method 1 (c) and method 2 (d) are used. Compare the position of the cube corners over the carpet [red] and the edges of the three planar regions [yellow].

| Criterion | Complexity term |
|-------------------------------|---|
| Akaike's AIC ¹ | $2k$ |
| Bozdogan's CAIC ² | $k(\log n + 1)$ |
| Bozdogan's CAICF ² | $k(\log n + 2) + \log \mathbf{I}(\theta_k) $ |
| Schwarz's BIC ⁹ | $2k \log n$ |
| MDL | $1/2 k \log(n)$ |
| Kanatani's gMDL ⁴ | $-k \log \epsilon^2$ |
| BMSC-RISS ³ | $k/2 \log_2 * (\theta_k^T \mathbf{I}(\theta_k) \theta_k) + \log_2(V_k)$ |

Table 1: The criteria considered in our study (k is the size of the model and n is the number of data).

considered criteria which involve the covariance matrix on the estimated parameters ($V(\theta_k)$) and the Fisher information matrix ($\mathbf{I}(\theta_k) = V(\theta_k)^{-1}$). Indeed, often, criteria such as AIC are only asymptotic approximations of another ones which includes the covariance or the information matrix. So, we hope that such criteria will improve the model selection.

In general, we show and compare diverse approaches for

every stage of the algorithm, pointing out that a good automatic recognition of the kind of movement (by model selection) and an efficient tool to find the real correspondences improve the stability and precision of the viewpoint parameters. We also describe the behaviour for different combinations of these approaches which are suitable to real-time AR applications. We assess these approaches by using both synthetic and real data measuring the precision of the viewpoint estimations, and show examples to compare the visual effect in augmented sequences.

3.4. Experimental results

3.4.1. Predicted model

The possibility of using linearly predicted models was described in ⁴. The main problem of sequential calibration is that some variations are very small, and some subsequences may seem piecewise linear even if they are not. The consequence is that for some noise level, the predicted model is often preferred to the general one because its complexity is simpler than the real model. This may lead to divergence after some images as exhibited in the experiment conducted

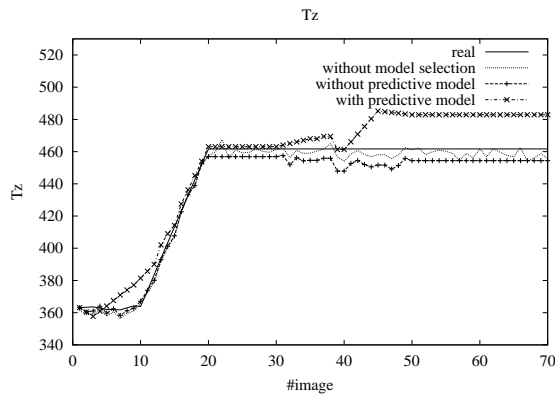


Figure 5: The computed Z-translation using model selection with and without t -predicted model

on the **synthetic sequence** (Fig. 5). In this case, we use the CAICF criterion, but the behaviour is similar for the other ones. In our approach, it is important both to select always the right model and to compute the parameters accurately. This is the reason why we decide to use just three models for the fixed focal length calibration, and in the experiments that follow.

3.4.2. Comparing selection criteria

In order to compare the selection criteria, we use the **synthetic sequence** corrupted with various noise levels. For each image i we consider the model selected between frame $i - 1$ and i . As an iterative procedure is used, we use as initialisation the actual viewpoint for frame $i - 1$ in order to avoid drift problem. As the actual model is known, we show in Table 2 and 3 the percentage of correct model choice obtained for each criterion and for each model.

These tables proved that for a moderate noise level ($\sigma = 0.3$ pixels), most of the criteria perform well. We see that some criteria prefer the more general model (AIC, gMDL) while some others always chose the simpler one (BIC). It can also be noticed that some criteria are more sensitive to noise and select a wrong model often than others. In general, the criteria based on Akaike's Information Criterion (AIC, CAIC and CAICF) performs well. However, AIC tend to admit more overfitting than the CAIC or CAICF and the stabilisation performance is then reduced. The experiments we performed on Bayesian criteria such as BMSC-RISS are not convincing. First, this criterion tend to admit underfitting for the general motion when the noise level is relatively high ($\sigma = 1.0$). Second, the results seem to depend tightly on the a priori probability on the various motion models.

That is the reason why the experiments in the following are done using the CAICF as selection criterion, because it performs well and it considers also that the nature of the

| Motion | Criterion | $\sigma = 0.3$ | | |
|---------|-----------|----------------|---------|---------|
| | | Underfit | correct | Overfit |
| Static | AIC | - | 83.1% | 16.9% |
| | CAIC | - | 98.7% | 1.3% |
| | CAICF | - | 100.0% | 0.0% |
| | BIC | - | 100.0% | 0.0% |
| | gMDL | - | 77.5% | 22.5% |
| | MDL | - | 83.2% | 16.8% |
| | BMSC-RISS | - | 86.3% | 13.7% |
| Pan | AIC | 0.0% | 85.3% | 14.7% |
| | CAIC | 0.0% | 99.3% | 0.7% |
| | CAICF | 0.0% | 98.7% | 1.3% |
| | BIC | 0.0% | 100.0% | 0.0% |
| | gMDL | 0.0% | 84.7% | 15.3% |
| | MDL | 0.0% | 85.4% | 14.6% |
| | BMSC-RISS | 0.0% | 100.0% | 0.0% |
| General | AIC | 0.0% | 100.0% | - |
| | CAIC | 1.5% | 98.5% | - |
| | CAICF | 1.3% | 98.7% | - |
| | BIC | 5.4% | 94.6% | - |
| | gMDL | 0.0% | 100.0% | - |
| | MDL | 0.0% | 100.0% | - |
| | BMSC-RISS | 5.8% | 94.2% | - |

Table 2: Percentage of good model selection for various criteria, noise level = 0.3 .

parameters is different by including the Fisher's information matrix in the complexity term.

If the noise level increases, some criteria tend to select a simpler model, specially the panoramic model even when the real one varies in translation and in rotation. Often, small translations are mistaken by panoramic motions because the direction of motion of the points is the same and the residual error is also similar, but the complexity of the model to detect a translation is bigger (because it is just included in the general model) than to detect a panoramic movement, the first model has 6 degrees of freedom, rather than 3 for the second one.

3.4.3. The Turntable Sequence

We demonstrate the effectiveness of the approach on the turntable sequence, which was described in Section 2.5. Fig. 6 shows the distance from the current camera position to the initial camera position computed with and without model selection. We can notice that when model selection is used, the trajectory is more stable. As the total translation of the turntable is perfectly known (10cm), we can also estimate the accuracy of the process by comparing the estimated translation with and without model selection. When model selection is used the estimated total translation is 9.82 cm,

| Motion | Criterion | $\sigma = 1.0$ | | |
|---------|-----------|----------------|---------|---------|
| | | Underfit | correct | Overfit |
| Static | AIC | - | 83.7% | 16.3% |
| | CAIC | - | 100.0% | 0.0% |
| | CAICF | - | 100.0% | 0.0% |
| | BIC | - | 100.0% | 0.0% |
| | gMDL | - | 0.0% | 100.0% |
| | MDL | - | 83.8% | 16.2% |
| | BMSC-RISS | - | 100.0% | 0.0% |
| Pan | AIC | 0.0% | 86.7% | 13.3% |
| | CAIC | 0.0% | 100.0% | 0.0% |
| | CAICF | 0.0% | 97.3% | 2.7% |
| | BIC | 0.0% | 100.0% | 0.0% |
| | gMDL | 0.0% | 0.0% | 100.0% |
| | MDL | 0.0% | 88.0% | 12.0% |
| | BMSC-RISS | 0.0% | 100.0% | 0.0% |
| General | AIC | 11.5% | 88.5% | - |
| | CAIC | 24.1% | 75.9% | - |
| | CAICF | 20.3% | 79.7% | - |
| | BIC | 33.6% | 66.4% | - |
| | gMDL | 0.0% | 100.0% | - |
| | MDL | 11.5% | 88.5% | - |
| | BMSC-RISS | 36.6% | 63.4% | - |

Table 3: Percentage of good model selection for various criteria, noise level= 1.

whereas it is around 6.56cm without model selection. In addition, as the sequence is closed, the drift can be used to assess the two methods. The total drift of the camera position when model selection is used is 0.14 cm. Without model selection, the drift is 3.26 cm. During the stationary and the rotating motion, the distance between the current and the initial position is constant. When model selection is used, this distance is really constant, whereas it is not without model selection. Two videos using the same room with more abrupt motions are attached to this paper: the sequence *roomWithoutMS.mpg* exhibits the augmented scene and the selected model when no motion selection is used; the video *roomWithMS.mpg* exhibits results when motion selection is used. In the last video, the symbol in the upper-left corner of the images indicates the selected model. The red cross indicates the stationary model, the green circle corresponds to the panoramic rotation, and the blue square to the general model.

These results clearly demonstrate that model selection produces a smoother trajectory and a better visual impression. They also prove that the use of model selection improve the accuracy of the viewpoints and reduces noticeably the drift problems that are common at long sequences.

In order to quantify the time needed for viewpoint com-

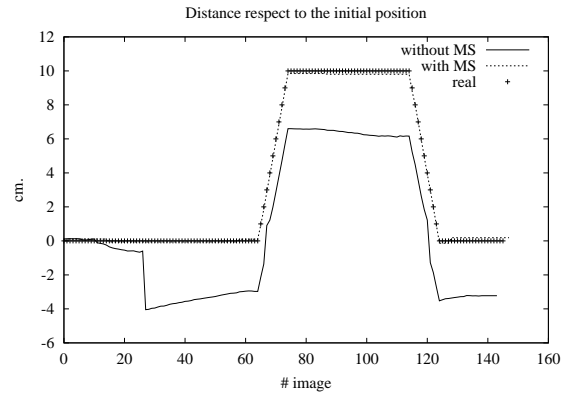


Figure 6: The distance between the current viewpoint and the initial one for the turntable sequence.

putation, table 4 gives the times needed for the different steps of viewpoint computation for one image of the calibration target: extraction and matching steps (we use the MIC algorithm¹⁴ to extract the key points), robust matching using the RANSAC algorithm with method 1 and 2, and viewpoint computation using model selection. About 500 key-points were extracted from each image. After retaining only the points which belongs to three planes, only 100 points are considered in the viewpoint computation process. The full process is about 64 ms when method 1 is used and 124 ms for method 2. This means that we can handle 16 images per second with method 1 and 8 images with method 2.

| | |
|---|--------|
| MIC | 15 ms |
| Matching | 9 ms |
| RANSAC Method 1 | 25 ms |
| RANSAC Method 2 | 85 ms |
| Viewpoint estimation with model Selection | 15 ms |
| Total Method 1 | 64 ms |
| Total Method 2 | 124 ms |

Table 4: Computation rates obtained on a Pentium IV 2. Ghz

3.4.4. The snooker sequence

This large sequence was taken using a hand held camera in the hall of our laboratory. The user was free to move anywhere he wanted. Due to the brightness of the floor, some sheets of paper were put down on it to make easier the tracking process. During the sequence, two panoramic motions were realized (see *hallCamera.mpg*), one with a tripod and the other without a tripod. Both are correctly labelled by the model selection process as can be seen on the video (*hallTrack.mpg*). The set of inliers is also visible on this video. Finally, some snapshots of the scene augmented



Figure 7: Snapshots of the scene augmented with the snooker.

with a **snooker** are shown in Fig. 7 and in video *hallAugmented.mpg* proving the effectiveness of our method.

4. Conclusion

We proposed in this paper several improvements to view-point computation for multi-planar environments. The use of model selection with various criterion proved that criterion involving information on the covariance of the estimated parameters are well suited to stabilisation. Videos attached to this paper proved that this method significantly improved the visual impression of the augmented scene. We now investigate if these criteria are well suited when varying focal lens are considered. Our first experiments proved that the use of non nested models is more difficult to handle. We also plan to investigate the influence of the accuracy on the first view-point on the whole process. Indeed, it appears that good registration results can only be obtained if good intrinsic camera parameters are available.

References

1. H. Akaike. A new look at the statistical model identification. *IEEE Trans Aut Ctrl*, 19(6):716–723, 1974.
2. H. Bozdogan. Model Selection and Akaike’s Information Criterion (AIC); The General Theory and its Analytical Extensions. *Psychometrika*, 52(3):345–370, 1987.
3. K. Bubna et C.V. Stewart. Model selection and surface merging in reconstruction algorithms. In *Proceedings of 6th International Conference on Computer Vision, Bombay (India)*, pages 895–902, 1998.
4. K. Kanatani C. Matsunaga. Calibration of a Moving Camera Using a Planar Pattern: Optimal Computation, Reliability Evaluation and Stabilization by Model Selection. In *Proceedings of 6th European Conference on Computer Vision, Trinity College Dublin (Ireland)*, pages 595–609, 2000.
5. R. I. Hartley et A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521623049, 2000.
6. K. Kanatani. Model Selection for Geometric Inference. In *Proceedings of 5th Asian Conference on Computer Vision, Melbourne, Australia*, pages 23–25, 2002.
7. J. P. Mellor. Realtime Camera Calibration for Enhanced Reality Visualization. In *Proceedings of Computer Vision, Virtual Reality, and Robotics in Medicine’95 (CVRMed’95)*, pages 471–475, April 1995.
8. U. Neumann et Y. Cho. A selftracking augmented reality system. In *Proceedings of the ACM Symposium on Virtual Reality Software and Technology*, pages 109–115, 1996.
9. G. Schwarz. Estimating the Dimension of a Model. *The Annals of Statistics*, 6(2):461–464, 1978.
10. G. Simon et M.-O. Berger. A Two-stage Robust Statistical Method for Temporal Registration from Features of Various Type. In *ICCV’98, Bombay (India)*, pages 261–266, January 1998.
11. G. Simon, A. Fitzgibbon, and A. Zisserman. Markerless tracking using planar structures in the scene. In *Proc. International Symposium on Augmented Reality*, pages 137–146, October 2000.
12. P.H.S. Torr. An Assessment of Information Criteria for Motion Model Selection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition, Puerto Rico, PR (USA)*, pages 47–52, June 1997.
13. P.H.S. Torr, A.W. Fitzgibbon, et A. Zisserman. Maintaining multiple motion model hypotheses over many views to recover matching and structure. In *ICCV’98, Bombay (India)*, pages 485–491, 1998.
14. M. Trajkovic et Mark Hedley. Fast corner detection. *Image and Vision Computing*, (16):75–87, 1998.
15. J. Vallino. *Interactive Augmented Reality*. Thèse de doctorat, University of Rochester, December 1998.