



HAL
open science

Localization of 3D objects using model-constrained SLAM

Angélique Loesch, Steve Bourgeois, Vincent Gay-Bellile, Michel Dhome,
Olivier Gomez

► **To cite this version:**

Angélique Loesch, Steve Bourgeois, Vincent Gay-Bellile, Michel Dhome, Olivier Gomez. Localization of 3D objects using model-constrained SLAM. Machine Vision and Applications, 2018, 29 (7), pp.1041-1068. 10.1007/s00138-018-0951-x . hal-04327823

HAL Id: hal-04327823

<https://hal.science/hal-04327823v1>

Submitted on 13 May 2024

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.

Localization of 3D objects using model-constrained SLAM

Angelique Loesch · Steve Bourgeois · Vincent Gay-Bellile · Michel Dhome · Olivier Gomez

Received: date / Accepted: date

Abstract Accurate and real-time camera localization relative to an object is needed for high-quality Augmented Reality applications. However object tracking is not an easy task in an industrial context where objects may be textured or not, have sharp edges or occluding contours, be relatively small or too large to be entirely observable from one point of view.

This paper presents a localization solution built on a key-frame-based SLAM algorithm. It only uses as inputs an RGB camera and a CAD model of the object of interest. The 3D model provides an absolute constraint that reduces drastically the SLAM drift. It is based on 3D oriented contour points called edgelets, dynamically extracted from the model by Analysis-by-Synthesis technique on the graphics hardware. This model constraint is then expressed through two different formalisms in the SLAM optimization process.

The dynamic edgelet generation ensures the genericity of our tracking method, since it allow to localize polyhedral and curved objects. The proposed solution is easy to deploy, requiring no manual intervention on the model, and runs in real-time on HD video-streams. It is thus perfectly adapted for high-quality Augmented Reality experiences. Videos are available as supplementary material.

Keywords Simultaneous Localization And Mapping · Constrained bundle adjustment · Occluding contours · Memory consumption · Real-time · Augmented Reality

A. Loesch, S. Bourgeois, V. Gay-Bellile, O. Gomez
CEA LIST, Point Courier 94, Gif-sur-Yvette, F-91191 France
E-mail: {angelique.loesch, steve.bourgeois, vincent.gay-bellile, olivier.gomez2}@cea.fr

M. Dhome
Pascal Institute, Blaise Pascal University, Clermont-Ferrand France
E-mail: michel.dhome@univ-bpclermont.fr

1 Introduction

Applications such as quality control, automation of complex tasks or maintenance support with Augmented Reality (AR) could greatly benefit from visual tracking of 3D objects [5]. However, this technology is under-exploited due to the difficulty of providing deployment facility, localization quality and genericity simultaneously. Most existing solutions indeed involve a complex or an expensive deployment of motion capture sensors, or require human supervision to simplify the 3D model [31]. And finally, most tracking solutions are restricted to textured or polyhedral objects to achieved an accurate camera pose estimation [3, 36].

Tracking any object is a challenging task due to the large variety of object forms and appearances. Industrial objects may indeed have sharp edges, or occluding contours that correspond to non-static and view-point dependent edges. They may also be textured or textureless. Moreover, some applications require to take large amplitude motions as well as object occlusions into account, tasks that are not always dealt with common model-based tracking methods. These approaches indeed exploit 3D features extracted from a model, that are matched with 2D features in the image of a video-stream [13]. However the accuracy and robustness of the camera localization depend on the visibility of the object as well as on the motion of the camera.

To better constrain the localization, recent solutions rely on environment features that are reconstructed online, in addition to the model ones. These approaches combine SLAM (Simultaneous Localization And Mapping) and model-based tracking solutions by using constraints from the 3D model of the object of interest. This model can indeed be used to constrain the SLAM initialization [2] or its optimization process [33, 29, 21]. Constraining SLAM algorithms with a 3D model results in a drift free localization. However, such approaches are not generic since they are only adapted for tex-

tured or polyhedral objects.

In this paper, we propose a solution that fulfills the requirements concerning deployment facility, localization quality and genericity. This solution, based on a visual key-frame-based constrained SLAM, only exploits an RGB camera and a geometric CAD model of the object of interest. An RGB camera is indeed preferred over an RGBD sensor, since the latter imposes limits on the volume, the reflectiveness or the absorptiveness of the object, and the lighting conditions. A geometric CAD model is also preferred over a textured model since textures may hardly be considered as stable in time (deterioration, marks,...) and may vary for one manufactured object. Furthermore, textured CAD models are currently not widely spread. Contrarily to previous methods, the presented approach does not need 3D model simplification to bring out sharp edge. It deals with polyhedral and curved objects by extracting sharp, occluding contours and silhouette from a model render on GPU, and is real-time, accurate and robust to occlusion or sudden motion.

2 Related work

Several visual localization methods exist in the state of the art. This Section focuses especially on model-based tracker solutions and on constrained SLAM (C-SLAM) algorithms, adapted to object tracking.

2.1 Model-based tracking methods

Camera localization relative to an object of interest may be achieved thanks to model-based tracking methods. These solutions rely on a two-step process. First, 3D features provided by a 3D model of the object of interest are matched with their corresponding 2D features in the image. Then, the camera pose that minimizes the re-projection of these 3D features with respect to their 2D counterparts is estimated. With a textureless CAD model, these 3D features are 3D oriented surface points on the model, that result in edge points in the image. We refer to these surface points as edgelets throughout the paper. The main difficulties of model-based tracking using such features, are encountered during the 3D/2D matching step. Visual features such as 2D contour points are indeed not adapted to local discrimination (difficulty of distinguishing an edge point from another). To constrain the 3D/2D matching, the camera motion has to be small between two frames which limits the solution robustness against fast and sudden displacements. Besides, edgelets may depend of the camera point of view (silhouette, auto-occluding contours, ...) and can not be extracted from sharp edges of the model. In early works [3,37], to bypass these issues, some methods have focused on polyhedral object tracking only, or limited the camera motion to

low displacements.

However polyhedral object tracking solutions are not generic approaches since they can not deal with curved objects. When the object of interest is indeed polyhedral, surface points that generate contours in the image, mostly come from sharp edges of the model. These edgelets are thus independent of the camera point of view and can be precomputed. They are identified by a simple threshold on the dihedral angle between the model facets. However, the dihedral angle criteria is not relevant when dealing with manufactured objects with sharp edges rounded into fillets. A simplification of object models may resolve this problem by transforming round edges into sharp edges, but the simplification degree generally needs to be adjusted by an expert.

Besides, the object tracking exploitation is also reduced to expert utilization, to control the camera motion. Small amplitude displacements between images are indeed needed to facilitate 3D/2D matching between edgelets and 2D contour points, and to maintain good robustness and stability. The 3D/2D matching step is usually performed by projecting each edgelet extracted from the model and by searching locally around this projection a 2D image contour with a similar orientation. Without prediction, this projection is performed on the current image with the pose estimated on the previous one. Nevertheless, a small research area is crucial for low discriminant features such as edges to avoid matching errors and to keep real-time performances.

Several approaches work on object tracking genericity and stability to larger displacement and sudden motion.

Thus, solutions focus on tracking methods adapted to any kind of objects, both polyhedral and curved. Some model-based trackers locally approximate the 3D model surface using parametrizations based on curvature radius [16] or quadrics [26]. These parametrizations have the advantage of being differentiable and thus easily integrable in a cost function to optimize. However, for very complex objects, the number of quadrics or curvature radii may quickly increase and complicate the model parametrization. To enable computationally tractable tracking, the number of quadrics or curvature radii has to be limited which results in an approximation of object surfaces for complex objects and thus a less accurate tracking.

Other solutions exploit Analysis-by-Synthesis technique to identify the edgelets of an object model, with the use of normal and depth maps from synthetic model renders [25,27,39]. These methods are based on non photo-realistic renders to be independent of illumination conditions [8]. They have the advantage of being efficient on any object, since they consider the observation distance and exploit the model without simplification. However the resulting edgelets are only locally valid, since they are not parametrized according to the object surface. Thus, they have to be re-estimated at

each image of the sequence. This process is time consuming due to the data transfer between GPU and CPU, especially when the geometric model has many faces or when the image has a high resolution.

Other methods aim to lift the low amplitude motion constraint, in order to increase the tracking robustness to sudden motion. A first solution is to use particle filters [11] to eliminate this assumption by generating an important number of hypotheses for the current pose. Each of them is weighted based on the re-projection error between sharp edges and their closest edge in the image. This approach offers a robust tracking solution with respect to fast motions and occlusions. However, its computing cost is not compatible with a real-time execution for complex objects.

Another strategy to increase the displacement between frames, is to predict the camera pose on the current image. One possibility is to use an external sensor such as an IMU [1]. It is also possible to exploit image features that are more discriminant than edgelets, e.g. key-points, to compute a first estimation of the camera displacement between the previous image and the current one. This estimation is then refined with the model edges. In [27, 36], key-points extracted on the objects of interest are used. Since these approaches require a textured object that is always visible in the image (not occluded), an alternative solution is to exploit the object environment by extracting key-points on the whole image as proposed by constrained SLAM solutions.

2.2 Constrained SLAM solutions

2.2.1 Key-frame-based SLAM

Key-frame-based SLAM algorithms [12, 22] are iterative processes that reconstruct online 3D primitives from the environment with respect to 2D observations and estimate the camera pose for each image thanks to 2D/3D correspondences between the previously reconstructed primitives and those extracted from the images. The reconstructed map is improved with new primitives at each key-frame. To limit the error accumulation, a possible solution is to use a non-linear optimization process called bundle adjustment. It refines camera poses and 3D feature points by minimizing their re-projection errors [35]. Standard SLAM algorithms [12, 22] are often used to estimate the localization of a camera in an unknown environment. However, they are not well adapted for object tracking. Camera poses are indeed expressed in an arbitrary coordinate frame and with an arbitrary scale that is subject to drift over time. Thus, the idea to use constrained bundle adjustments that minimize simultaneously the multi-view geometry and a constraint term based on the CAD model has been developed.

2.2.2 Constrained bundle adjustment

Since standard SLAM algorithms are not appropriate for object tracking, C-SLAM methods propose to include constraints provided by the model of the object of interest. This constraint integration allows to express the SLAM reconstruction into the object-frame with the correct scale and to prevent the drift of SLAM algorithms.

Model constraints exploit an *a priori* partial knowledge of the scene geometry, usually a 3D model. They can be expressed with a metric error [29] corresponding to the 3D distance between 3D features reconstructed by SLAM and the model of the object. However combining this error with a pixel one in the constrained bundle adjustment requires an adaptive weight to deal with heterogeneous measurements. Model constraints can also be expressed through a re-projection error of 3D features reconstructed by SLAM and projected on the model while some of their DoF are fixed during the optimization [19, 34]. The constraint can furthermore be defined with a re-projection error between 3D features extracted from the model and their associated image observations. These features are considered as an absolute information, and are then fixed in the bundle adjustment process. They may come from a 3D point cloud [21] or, when objects of interest are textureless, correspond to 3D oriented points (edgelets) from edges of a CAD model [33]. This latter method is called in this article EC-SLAM (SLAM constrained to edgelets). However, although the approach of [33] is real-time and accurate, it reaches its limits for curved objects since only sharp edges are exploited.

We propose in this paper an EC-SLAM that tracks polyhedral and curved objects with robustness and accuracy. Edgelets are dynamically generated according to Analysis-by-Synthesis technique and integrated as a model constraint into a C-SLAM algorithm.

Notations Matrices are designated in this paper by sans-serif capital font such as M and vectors by bold font such as \mathbf{v} or \mathbf{V} . The projection matrix P associated to a camera is given by $P = KR^T(I_3 - \mathbf{t})$, where K is the matrix of intrinsic parameters and (R, \mathbf{t}) the extrinsic ones. \mathbf{X} corresponds to all the optimized parameters in a bundle adjustment: the extrinsic parameters of the optimized cameras $\{R_j, \mathbf{t}_j\}_{j=1}^S$ and the 3D point positions $\{\mathbf{Q}_i\}_{i=1}^{N_Q}$. $\bar{\mathbf{X}}$ corresponds to the vector concatenating all the optimized pose parameters $\bar{\mathbf{X}}_j = \{R_j, \mathbf{t}_j\}_{j=1}^S$.

3 Overview

This paper is an extension of previous works [17, 18]. It proposes a constrained SLAM solution for generic object tracking that deals with large amplitudes and occlusions.

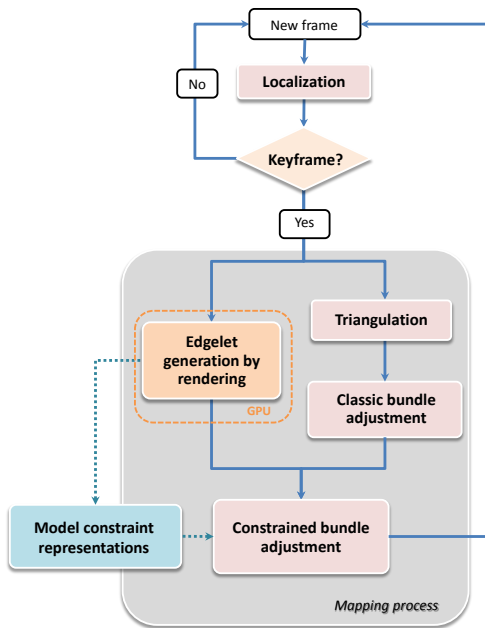


Fig. 1: Our EC-SLAM framework with a four-step mapping process: triangulation, model constraint formalism, classic and constrained bundle adjustment. Edgelet generation by rendering is performed in parallel on GPU before the model constraint representation.

Our EC-SLAM framework shown in Figure 1, is built on key-frame-based SLAM as [33] divided in a localization and a mapping threads. In the localization thread, camera pose is estimated at each frame thanks to a matching algorithm that establishes 2D/3D correspondences. The mapping process is then performed when a key-frame is detected. Generally it is decomposed into two important steps, the triangulation that extends the 3D map with new 3D features, and the refinement of camera poses and 3D feature positions through a bundle adjustment constrained to the CAD model. Our framework is however slightly different to answer the genericity and robustness issues.

3.1 Edgelet generation

In [33], the model constraint is provided by static edgelets extracted offline on the sharp edges of the CAD model. A visibility test is then performed to only get the subset of visible edgelets for each key-frame optimized in the constrained bundle adjustment.

However, this solution tracking only polyhedral objects, we propose a different use of the 3D model by extracting dynamically edgelets with Analysis-by-Synthesis technique. In Section 4, our edgelet extraction by rendering the CAD model of the object of interest is detailed. Nevertheless, replacing

the precomputed edgelets of an EC-SLAM as [33] with a constellation generated by rendering implies modifications in the SLAM framework. In fact, since the edgelets generated from the silhouette and occluding contours of the CAD model are viewpoint dependent, they should be estimated online. Moreover, our edgelet generation process induces an additional computing cost that depends on various factors, such as the image resolution used during the rendering, the performance and current workload of GPU (it can also be used to display some content on a screen). The presented edgelet generation is then achieved in the mapping thread on GPU in parallel of other processes to prevent its computational cost from slowing down the tracking process and to maintain real-time performances. The output of this generation is finally a set of edgelets that are evenly distributed and non-ambiguous for the matching and the pose estimation steps.

3.2 Model constraint representations

The dynamic edgelets are exploited to refine the camera pose and the SLAM reconstruction in the constrained bundle adjustment as presented in Figure 1.

3.2.1 Re-projection formalism

A first possible model constraint formalism takes directly advantage of edgelet information through a re-projection error $E(\bar{\mathbf{X}})$ similar to the one proposed by [33], where our dynamic edgelets are associated to 2D contour points in the image.

In that case, the model constraint is expressed as followed:

$$E(\bar{\mathbf{X}}) = \sum_{j=1}^S H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j), \quad (1)$$

with $H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j)$ the edgelet re-projection error

$$H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j) = \sum_{i=1}^{N_M} (\mathbf{n}_{i,j} \cdot (\mathbf{m}_{i,j} - P_j \mathbf{M}_i))^2. \quad (2)$$

It measures the orthogonal distance between the projection of edgelets \mathbf{M}_i , and their corresponding edge $\mathbf{m}_{i,j}$ concatenated in the vector $\bar{\mathbf{m}}_j$, for the j^{th} key-frame. $\mathbf{n}_{i,j}$ is the normal to the projection of the edgelet direction. This re-projection formalism will be evaluated in Section 6.

However, with this constraint representation, edgelet constellations have to be stored for all the optimized key-frames. This storage increases the memory footprint of the optimization process, which grows with the duration and the resolution of the video. Besides, this formalism is time consuming since the 2D/3D correspondences have to be re-estimated at each constrained bundle adjustment on all the optimized

key-frames. Thus two other representations of the model constraint are proposed to manage this issue.

3.2.2 Output model-based poses

These other model constraint formalisms are presented in Section 5. Both correspond to hybrid model / trajectory constraints that do not exploit directly the dynamic edgelet constellations as with the re-projection formalism. They instead use the outputs of a model-based tracker (named model-based poses as opposite of the SLAM poses). These constraints are defined as pose errors more compact in memory, but expressed in pixels to be homogeneous with respect to the environment constraint (see Section 3.3) proposed by the SLAM in the constrained bundle adjustment.

The output model-based poses used in these both formalisms, are the optimal camera poses with respect to the object model that best align the dynamic edgelets and their observations in the image. Since the observations that correspond to edgelets are initially unknown, the camera pose and these observations are defined as the parameters that minimize the model-based error given by equation 2 whose the vector representation is the following:

$$H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j) = h^\top(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j)h(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j), \quad (3)$$

with $h^\top(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j) = [(\mathbf{n}_{0,j} \cdot (\mathbf{m}_{0,j} - \mathbf{P}_j \mathbf{M}_0)) \dots (\mathbf{n}_{N,j} \cdot (\mathbf{m}_{N,j} - \mathbf{P}_j \mathbf{M}_N))]$. This minimization problem is usually solved by alternating the estimation of the observations $\bar{\mathbf{m}}_j$ and the estimation of the pose parameters $\bar{\mathbf{X}}_j$. The $\bar{\mathbf{m}}_j$ are estimated by projecting the edgelets with the current pose parameters and matching them to the nearest contour with a similar orientation. The pose parameters are then refined by minimizing $H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j)$ with a Levenberg-Marquardt algorithm [14].

To ensure a good convergence, such algorithm requires an initial estimation of $\bar{\mathbf{X}}_j$ close to the solution and a configuration of edgelets/observations that constrains the 6 DoF of the camera. In our context, the first assumption is verified since the optimization process uses the SLAM pose as initial guess. On the contrary, the second assumption cannot be ensured. It is then preferable to keep the poorly constrained DoF to their initial value during the optimization process. This can be achieved by using a truncated Levenberg-Marquardt [4]. This approach relies on a different approximation of the Hessian matrix's cost function than the usual Levenberg-Marquardt. While the Hessian matrix of $H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j)$ is approximated with $W_j \approx J_j^\top J_j$ in classic Levenberg-Marquardt [28], it is approximated by $W_j = TSVD(J_j^\top J_j)$ in a truncated Levenberg-Marquardt, where $TSVD(A)$ represents the truncated SVD of the matrix A .

The resulting orientation and position parameters $\bar{\mathbf{X}}_j^c = \{R_j^c, \mathbf{t}_j^c\}$ of the optimal model-based pose are stored to be used as

a constraint for this key-frame over the rest of the tracking. The associated projection matrix is denoted \mathbf{P}_j^c . The vector $\bar{\mathbf{X}}^c$ represents the set of model-based pose parameters $\bar{\mathbf{X}}_j^c = \{R_j^c, \mathbf{t}_j^c\}$ for all the key-frames optimized in the constrained bundle adjustment. The 2D contour points associated to edgelets on key-frame j with the optimal pose $(R_j^c | \mathbf{t}_j^c)$, are concatenated in the vector $\bar{\mathbf{m}}_j^c$. Finally, we define $W_j^c = J_j^{c\top} J_j^c$ as the approximation of $H''(\bar{\mathbf{X}}_j^c, \bar{\mathbf{m}}_j^c)$.

3.3 Optimization process

Since our model constraint may present inaccuracies especially if the object is occluded, creating incorrect 2D/3D associations or erroneous model-based poses, optimizing simultaneously as [33] the model and the multi-view constraints could result in a deterioration of the multi-view geometry relationships, that may causes localization failures. Then our approach consists in first optimized the environment constraint only. We determine the optimal solution of the multi-view relationships for all the optimized key-frames. These key-frames are selected via a covisibility graph [23]. To obtain this optimal solution, a classic bundle adjustment without any model constraint is thus performed as presented in Figure 1. The cost function $G(\mathbf{X})$ to minimize is the following:

$$G(\mathbf{X}) = \sum_{i=1}^{N_Q} \sum_{j \in L_i} d^2(\mathbf{q}_{i,j}, \mathbf{P}_j \mathbf{Q}_i). \quad (4)$$

The re-projection error d of equation 4 is the euclidean distance between $\mathbf{q}_{i,j}$ the 2D observation of the i^{th} 3D point \mathbf{Q}_i in the j^{th} key-frame and the projection of this 3D point. L_i is the set of the key-frame indexes observing \mathbf{Q}_i .

Then, the model and multi-view constraints are combined in a constrained bundle adjustment. Since multi-view constraint is less robust but more accurate than the model constraint, we use a fusion strategy, which ensures that the degradation of the multi-view constraint remains small. To achieve this task, Lhuillier constrained bundle adjustment framework [15] is chosen. Thus the cost function of our constrained bundle adjustment combining multi-view and model constraints is given by:

$$F(\mathbf{X}) = \frac{\omega}{e_t - G(\mathbf{X})} + E(\bar{\mathbf{X}}), \quad (5)$$

where $E(\bar{\mathbf{X}})$ is the model constraint re-projection error or the two formalisms described in next Section 5 using the optimal model-based poses presented in Section 3.2.1. e_t is a threshold which is slightly greater than the squared re-projection error obtained after minimizing Equation 4. The fraction $\frac{\omega}{e_t - G(\mathbf{X})}$ with $\omega > 0$ corresponds to a regularization term that prevents the degradation of the multi-view relationships estimated after the classic bundle adjustment.

We can notice that in practice, with the re-projection formalism proposed by [33] and defined by equation 1, our dynamic edgelets can also be directly integrated in the constraint bundle adjustment with a classic cost function. Then this latter will be only composed of the edgelet constraint and the multi-view relationships between the images. Both terms will be expressed in pixels and optimized simultaneously. This classic cost function will be given by:

$$F(\mathbf{X}) = G(\mathbf{X}) + E(\bar{\mathbf{X}}), \quad (6)$$

with $G(\mathbf{X})$ the environment constraint presented in equation 4 and $E(\bar{\mathbf{X}})$ the re-projection formalism of the model constraint. Moreover the classic bundle adjustment optimizing the multi-view geometry relationships only, become optional.

3.4 Road map and contributions

Edgelet generation is detailed in Section 4. Additional contributions over [17, 18] include the edgelet orientation definition in Section 4.1.3, and the improvement of the CPU/GPU transfer in Section 4.1.5. Furthermore these edgelet generation steps are evaluated in Section 6. Both hybrid model / trajectory constraints are described in Section 5. Our EC-SLAM solution is evaluated on several polyhedral and curved objects, and compared with other tracking solutions on synthetic and real data to assess its genericity, robustness and accuracy. All the evaluations of Section 6 are performed with the new edgelet orientation and data transfer, since it allow more accuracy in the camera pose estimation and better performances.

Moreover, additional contributions over [17, 18] concern the evaluation of our EC-SLAM framework robustness when the model constraint is inaccurate, and the evaluation of our solution tested on public datasets. Section 7 is also a contribution of this paper where we present an Augmented Reality application involving a static curved object with movable parts.

4 Dynamic edgelet extraction

As mentioned in Section 2, tracking solutions based on Analysis-by-Synthesis technique rely on renders of 3D models to identify surface points that are likely to generate contours. For computational performance, these surface points are sampled into a set of edgelets prior to the matching process. However, in most tracking solutions, the importance of the sampling is neglected. Although it has a major impact on the matching of edgelets with 2D image contours and on the pose estimated from these correspondences, most solutions

rely on a basic constant-step sampling.

In this Section, the introduced edgelet extraction provides a constellation of edgelets adapted to both matching and pose estimation.

4.1 Edgelet generation

Similarly to [39], our virtual rendering process relies on an image space approach. This choice is motivated by the large amount of faces (hundred of thousands) of 3D models used in our applications that will result in costly rendering times for object space approach [6, 30].

Our solution aims to estimate for each pixel both its probability of being a contour and its probability of being matched with its corresponding contour in a real image.

In the following, let's consider each pixel of the image as a

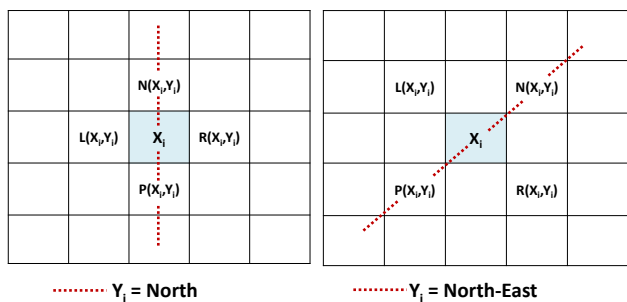


Fig. 2: Each pixel of the rendered image corresponds to the random variable X_i . $L(X_i, Y_i)$ and $R(X_i, Y_i)$ are the neighbor pixels along the normal to the direction Y_i . $N(X_i, Y_i)$ and $P(X_i, Y_i)$ are the neighbor pixels along the direction Y_i . On the left example with $Y_i = Nord$ and on the right with $Y_i = Nord - Est$.

random variable X_i , with $X_i = 1$ if the pixel is a contour, 0 otherwise. We denote Y_i the random variable that represents the 2D direction of a contour X_i , with $Y_i \in \{North, North - East, North - West, West\}$. $L(X_i, Y_i)$ (respectively $R(X_i, Y_i)$) is the function that returns the left (respectively right) neighbor of a pixel X_i along the normal to the direction Y_i , and $N(X_i, Y_i)$ (respectively $P(X_i, Y_i)$) correspond to the function that returns the next (respectively previous) neighbor of X_i along the direction Y_i . Since we are only looking to the direction that might take a contour, the neighbor pixel order (left, right and next, previous) is defined arbitrarily. Finally $Nmap(X_i)$ (respectively $Dmap(X_i)$) is the function that returns the normal (respectively the depth) of the pixel X_i .

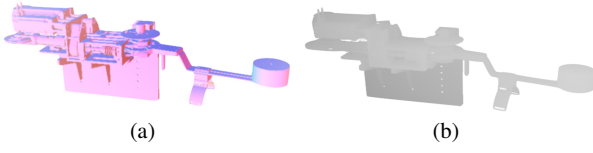


Fig. 3: Normal (left) and depth (right) maps of a rendered orthosis CAD model for a given camera point of view.

4.1.1 Edgelet probability

This first step of our edgelet generation estimates for each pixel its probability of being a contour. Usually, two types of 2D contours are distinguished: crease edges and silhouettes. On one hand, crease edges are relative to the sharp edges of the 3D model and appears as discontinuities in the normal map ($Nmap$)(see Figure 3(a)). On the other hand silhouettes correspond to the outline of the object or the auto-occluding contours and appears as discontinuities in the depth map ($Dmap$)(see Figure 3(b)). In both cases, the orientation of a 2D contour point corresponds to the 2D direction of the discontinuity in the image space.

Therefore, we define the probability of X_i to be a crease edge with direction Y_i as:

$$\mathbb{P}_{crease}(X_i|Y_i) = \max(1, crease(X_i, Y_i)), \quad (7)$$

with the angular discontinuity measure:

$$crease(X_i, Y_i) = \lambda \times (1 - Nmap(L(X_i, Y_i)) \cdot Nmap(R(X_i, Y_i))), \quad (8)$$

where $\lambda = (1 - \cos(angleMax))^{-1}$ is the normalization factor defined to reach an intensity of 1 for an angular amplitude of $angleMax$.

In a similar way, the probability of X_i to belong to the silhouette with an orientation Y_i is given by:

$$\mathbb{P}_{silhouette}(X_i|Y_i) = \max(1, silhouette(X_i, Y_i)), \quad (9)$$

where $silhouette$ is the depth discontinuity measure defined as follows:

$$silhouette(X_i, Y_i) = \frac{Laplacian(Dmap, X_i, Y_i)}{\beta \times Dmap(X_i)}. \quad (10)$$

Laplacian is the 1D laplacian oriented with respect to direction Y_i , and β is a weight factor. Since the discontinuity is normalized with respect to the observation distance, the parameter β does not depend on the dimension of the scene. Therefore, β can be easily interpreted as the minimal depth discontinuity, expressed as a ratio of observation distance, that provides a silhouette contour with a probability of 1.

Consequently, the probability of being a contour is defined as:

$$\mathbb{P}_{contour}(X_i) = \max_{Y_j}(\max(\mathbb{P}_{crease}(X_i|Y_j), \mathbb{P}_{silhouette}(X_i|Y_j))) \quad (11)$$

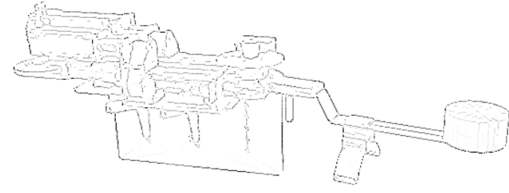


Fig. 4: Matching map of the rendered orthosis CAD model for a given camera point of view. More the edgelet pixel is dark, more the edgelet is likely to match to a 2D contour point.

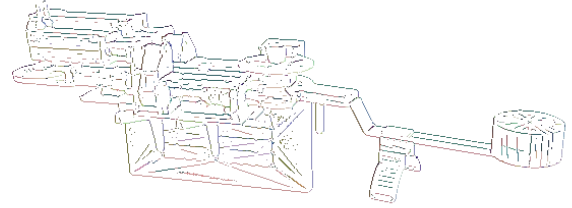


Fig. 5: Orientation map of the rendered orthosis CAD model for a given camera point of view. The 3D orientation is expressed through the R, G and B channels of the rendered image.

and the 2D direction associated to X_i is defined as:

$$Dir_{2D}(X_i) = \arg \max_{Y_j}(\max(\mathbb{P}_{crease}(X_i|Y_j), \mathbb{P}_{silhouette}(X_i|Y_j))). \quad (12)$$

4.1.2 Matching probability

This second step estimates the probability that each virtual contour is matched to its corresponding contour in the real image. Usually, during the matching process, the edgelets are associated to the nearest contour located along the normal of the contour. Consequently, the probability of correctly matching an edgelet can be assessed from the number of contours encountered within this 1D neighborhood.

Therefore, we propose to value this probability from the probability map estimated at the previous step. In the following, let's consider X_n to be the pixels of the 1D neighborhood of a pixel X_i still along the normal of the contour. And let N be the random variable representing the number of contours within the neighborhood. We define the probability of a pixel to be both a contour and correctly matched in a real image according to the Expected Value of N as:

$$\mathbb{P}_{match}(X_i) = \frac{\mathbb{P}_{contour}(X_i)}{1 + \mathbb{E}(N)}. \quad (13)$$

This matching probability is represented in Figure 4 for each edgelet detected on the rendered orthosis model.

4.1.3 3D edgelet orientation

The 3D edgelet orientation is in theory the same one that the edge they belong. This orientation can be easily determined from the 3D positions of two adjacent contour points in the edgelet probability map, method exploited in our previous works. However, the rasterization step of the rendering process creates aliasing and damage the orientation estimation. The erroneous edgelet orientation is then able to disturb the 2D/3D matching process and decrease the tracking accuracy. In order to deal with this issue, we propose a new solution in the image space based on the hypothesis that our 3D model mesh is an orientable 2-manifold surface. Our approach locally approximates the model surface by planes whom the intersections correspond to crease edges, or silhouettes (including auto-occluding contours). Thus the idea is to identify the two planes whose the intersection forms an edge that approximates locally the 3D model contour. The intersection direction will then corresponds to the 3D edgelet orientation.

Each pixel X_i that might be an edgelet, belongs to the intersection between two planes according to our definition. That's why an analysis of pixels in the adjacent neighborhood of X_i in the contour normal direction is performed, to establish which of them can be associated to the model planes creating the contour. The two pixel selection differs if X_i is a crease edge or a silhouette. On one hand, if the contour point X_i is part of a crease edge, the plane pixels are chosen by taking the adjacent pixels from both sides of X_i in the direction $Dir_{2D}(X_i)$. Then the 3D direction of the plane intersection and by extension the 3D direction of the contour point X_i is given by:

$$Dir_{3D}(X_i) = Nmap(L(X_i, Dir_{2D}(X_i))) \wedge Nmap(R(X_i, Dir_{2D}(X_i))). \quad (14)$$

On the other hand, if the contour point X_i belongs to the silhouette, the two pixel selection has to take into account that the object has auto-occluding contours. Indeed one of the planes creating the edge where X_i stands, is not visible on the front-face rendering. To visualize this occulted plane, a back-face rendering is achieved. The first neighbor pixel belonging to the hidden plane corresponds to the pixel X_i in terms of 2D position but in the back-face rendering. The second plane pixel is choosing in the front-face rendering as for the crease edges. Its selection depends of its distance to the camera position $dist_{cam}$. Thus, 3D direction of the plane intersection and the contour point X_i is defined as:

$$Dir_{3D}(X_i) = Nmap_{back}(X_i) \wedge Nmap(C(X_i, Dir_{2D}(X_i))), \quad (15)$$

with $Nmap_{back}$ the function that returns the back-face rendering normal of the pixel X_i . $C(X_i, Dir_{2D}(X_i))$ is the neigh-

bor pixel in the front-face rendering the closest of the camera. An orientation map of the rendered orthosis model is shown in Figure 5.

4.1.4 Edgelet sampling

An edgelet map \mathbb{P}_{match} (see Figure 4), that provides for each pixel its probability of being an edgelet correctly matched to its corresponding contour, has been estimated previously. This edgelet map is then exploited by the sampling step to provide a set of edgelets relevant to the matching and pose estimation processes.

On one side, to be pertinent to the matching process, this set of edgelets must maximize the expectation of the matching success. On the other side, to be suitable to the pose estimation process, the matches estimated from this set of edgelets must constrain the 6 DoF of the camera pose. Particularly, a set of matches that is unevenly distributed in the image, both in term of position and contour orientation, is not a relevant configuration for pose estimation.

Our sampling process (see Figure 6) relies on a division of both the 2D position and 2D orientation spaces of the edgelets projections. The 2D position space is divided in a regular grid of $N \times N$ buckets, and each bucket is divided in 4 angular sectors. Then, a set of edgelets is sampled for each angular sector of each spatial bucket, in order to obtain an evenly distribution in 2D space. The sampling itself is achieved with respect to a sampling probability affected to each edgelet. To define this sampling probability, let $\{Z_i\}$ be the set of edgelets of an angular sector of a spatial bucket. The probability $\mathbb{P}_{sampling}$ of an edgelet $Z \in \{Z_i\}$ is defined as follows:

$$\mathbb{P}_{sampling}(Z) = \frac{\mathbb{P}_{match}(Z)}{\sum_i \mathbb{P}_{match}(Z_i)}. \quad (16)$$

With this sampling strategy, inside an angular sector of a spatial bucket, edgelets with a high matching probability \mathbb{P}_{match} are more likely to be sampled. Moreover, the random nature of the sampling prevents the local agglomeration of edgelets.

4.1.5 Implementation details

While the edgelet sampling step is achieved on CPU, the edgelet generation step is performed on GPU. Our implementation relies on multiple rendering passes with GLSL shaders. The first pass computes the depth map $Dmap$ and the normal map $Nmap$. The second and third one compute the back-face normal map $Nmap_{back}$ and the 3D orientation Dir_{3D} . The fourth one computes the edgelet map $\mathbb{P}_{contour}$ and the 2D direction Dir_{2D} . The last one computes the matching probability map \mathbb{P}_{match} . Since the sampling is achieved

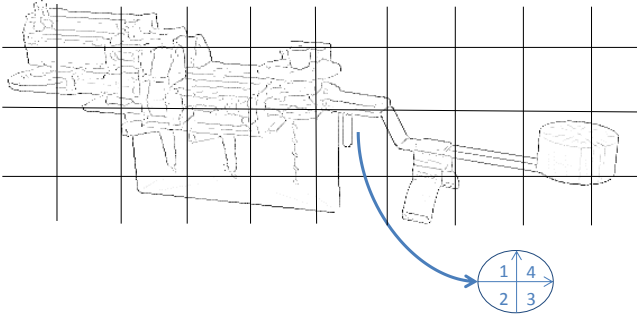


Fig. 6: Edgelet sampling method. The rendered image is divided in buckets, also divided in 4 angular sectors to obtain an edgelet set evenly distributed in term of 2D position and orientation.

on CPU, the textures containing all the edgelet information have to be transferred to CPU memory. However in each map, most of the pixels are not relevant and correspond to the background or the model surfaces instead of edgelet data. That is why, instead of transferring the entire textures asynchronously by GPU as in earlier works, a filter is proposed as a pre-processing in order to only transfer the edgelet pixels to CPU. This new filter is a rendering pass based on a geometry shader and a transform feedback process.

A set of vertices is generated on GPU through a Vertex Buffer Object (VBO) in order to optimize computation times. This set contains as many vertices as pixels in one texture. The geometry shader reduces the set size by storing only the vertices representing edgelet pixels with their 2D coordinates as attributes. Moreover the geometry shader adds more attributes to each edgelet vertex, as 2D and 3D positions, normal and depth values, 3D and 2D orientations and matching probability values of each edgelet. Then the vertices are saved in an other VBO thanks to the Transform Feedback process with an interleaved data mode, in order to simplify the edgelet reading on CPU. The new set of vertices is organized on GPU similarly as the sequential reading on CPU and the reading performance is optimal. Since all the edgelet information are stored contiguously, edgelets with their attributes can be indeed read in a sequential way. In this filter step, the rasterization process coming usually after the geometry shader in a rendering pass, is deleted since it is time consuming and not relevant in our case. The vertices containing the edgelet data are then directly transferred to CPU in an asynchronous way after the geometry shader execution, with the help of Pixel Buffer Object (PBO).

Moreover, to reduce the amount of data to transfer, the probability is encoded on an unsigned byte (the value 255 corresponding to 1).

5 Hybrid model/trajectory constraint

Edgelets extracted from the CAD model on GPU as described in the previous Section, are exploited to expressed the model constraint of our EC-SLAM algorithm. They can be directly used in the model constraint representation as proposed in the re-projection formalism presented in Section 3.2.1, or dedicated to the optimization of the model-based poses exploited into both formalisms of our model constraint in the present Section.

We describe two approaches slightly different that uses the 3D model to estimate model-based poses as presented in Section 3.2.2. These model-based poses are exploited as hybrid model/trajectory constraints in the constrained bundle adjustment in order to decrease the memory footprint of the optimization process. However, using the outputs of a model-based tracker as constraints implies to deal with heterogeneous error terms between the environment constraint and the model one. We then propose two hybrid model / trajectory constraints to combine the benefits of both kinds of constraints and keep error terms homogeneous.

Different pose distances are introduced in this Section, that measure the deviation of pose parameters from $\bar{\mathbf{X}}^c$ by their impact on the model constraint.

5.1 Main idea

We propose to use a difference of pixel errors as a pose distance. More precisely, it consists in measuring the difference between the re-projection error (see equation 2) corresponding to a key-frame pose $(R_j | \mathbf{t}_j)$ and the one corresponding to the model-based pose with the parameters $\bar{\mathbf{X}}_j^c$ associated to the same j^{th} key-frame. The corresponding distance is then defined by the following equation:

$$E(\bar{\mathbf{X}}) = \sum_{j=1}^S H(\bar{\mathbf{X}}_j, \bar{\mathbf{m}}_j^c) - \sum_{j=1}^S H(\bar{\mathbf{X}}_j^c, \bar{\mathbf{m}}_j^c), \quad (17)$$

Notice that to ensure the pixel error to be minimal at the pose $(R_j^c | \mathbf{t}_j^c)$, $E(\bar{\mathbf{X}})$ is evaluated by using the 2D contour points $\bar{\mathbf{m}}_j^c$ associated to edgelets for the key-frame j . We define the vector representation of our hybrid error $E(\bar{\mathbf{X}})$ as follow:

$$E(\bar{\mathbf{X}}) = H(\bar{\mathbf{X}}, \bar{\mathbf{m}}^c) - H(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c), \quad (18)$$

with $\bar{\mathbf{m}}^c$ the vector concatenating the 2D contour points associated with the parameters $\bar{\mathbf{X}}^c$ for all the edgelets and for all the optimized key-frames. An illustration of this hybrid error for a given key-frame and a given edgelet is represented in Figure 7. With this definition we ensure that $E(\bar{\mathbf{X}})$ is minimal when the model constraint $H(\bar{\mathbf{X}}, \bar{\mathbf{m}}^c)$ is minimal (*ie.* at the pose parameters $\bar{\mathbf{X}}^c$). Besides, the pose defined by the parameters $\bar{\mathbf{X}}^c$, the 2D contour points $\bar{\mathbf{m}}^c$, and

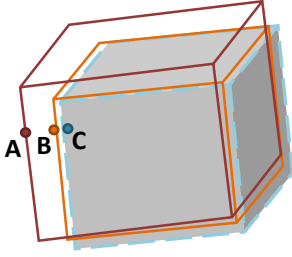


Fig. 7: Point A (respectively B) is an edgelet extracted from a cube model and projected according to a SLAM key-frame pose in red (respectively according to the model-based pose in orange). Point C is the 2D contour point associated to the edgelet B . The error $E(\bar{\mathbf{X}})$ of equation 18 can be interpreted as the difference $d_{ortho}(A,C) - d_{ortho}(B,C)$, with d_{ortho} the orthogonal distance.

consequently the pixel error associated to the parameters $\bar{\mathbf{X}}^c$ are constant. Since the 2D contour points $\bar{\mathbf{m}}^c$ are no longer re-estimated, it is not required to store the contour images. Furthermore, contrary to the re-projection formalism described in Section 3.2.1, the edgelet/contour matching step is not necessary anymore, which allows to reduce computation time. Only the vector of the 2D contour points $\bar{\mathbf{m}}^c$ and, in case of dynamic edgelets, the edgelet constellation are memorized for each key-frame. Compared to re-projection formalism, it implies a large reduction of memory footprint and computation cost. However, this memory consumption can remain important if the number of edgelets and key-frames is high. In the following Subsections, we introduce two approximations of $E(\bar{\mathbf{X}})$ that reduce the memory footprint.

5.2 First approximation

This first approximation relies on the hypothesis that the model-based tracker has converged to a local minimum. Thus it is reasonable to consider that the distance between the projection of the edgelets with respect to the parameters $\bar{\mathbf{X}}^c$ and their corresponding 2D contour points is negligible. In Figure 7, this assumption corresponds to approximate $d_{ortho}(A,C) - d_{ortho}(B,C)$ by $d_{ortho}(A,B)$.

Under such hypothesis, the 2D contour points $\bar{\mathbf{m}}_j^c$ of equation 17 are replaced by the projection of their corresponding 3D edgelets with respect to the parameters $\bar{\mathbf{X}}_j^c$. The resulting approximation of $E(\bar{\mathbf{X}})$ is given by:

$$E(\bar{\mathbf{X}}) = \sum_{j=1}^S H(\bar{\mathbf{X}}_j, \bar{\mathbf{M}}_j^c), \quad (19)$$

where $\bar{\mathbf{M}}_j^c = \{P_j^c \mathbf{M}_i\}_{i=1}^{N_M}$ is the vector concatenating edgelets \mathbf{M}_i projected according to P_j^c .

With this approach, it is not necessary to keep in memory the 2D associated contours, since only the model-based pose parameters are exploited. In the case of polyhedral object, the memory footprint of the resulting constrained bundle adjustment is very small, since only a set of edgelets shared by all the key-frames has to be stored as proposed by [33]. However, the memory consumption is higher with dynamic edgelets extracted from occluding contours as presented in Section 4, since they depend on the point of view. They are generated online and stored for each key-frame. That is why a second approximation of equation 18 is proposed to reduce memory consumption.

5.3 Second approximation

This second approximation relies on the hypothesis that the poses of the EC-SLAM key-frames are located in the neighborhood of their corresponding model-based poses. This approximation is valid in practice since the model-based poses are obtained by refining the SLAM poses with equation 2. Under this hypothesis, equation 18 can be approximated with a second order Taylor expansion around the model-based poses $\bar{\mathbf{X}}^c$:

$$E(\bar{\mathbf{X}}) \approx E(\bar{\mathbf{X}}^c) + E'(\bar{\mathbf{X}}^c)\delta + \frac{1}{2}\delta^\top E''(\bar{\mathbf{X}}^c)\delta, \quad (20)$$

where $\delta = \bar{\mathbf{X}} - \bar{\mathbf{X}}^c$ is the difference between pose parameters (positions and orientations) provided by the SLAM and model-based tracking algorithms. $\bar{\mathbf{X}}^c$ is obtained after the model-based refinement and corresponds to the minima of the model-based cost functions (equation 2) for all the optimized key-frames. Thus $H'(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c)$ (equation 3) is null. Consequently the first derivative $E'(\bar{\mathbf{X}}^c) = H'(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c) = 0$. Besides the second derivative $E''(\bar{\mathbf{X}}^c) = H''(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c)$ can be approximated by the $S \times S$ diagonal matrix W^c with S the number of optimized key-frames. The diagonal term of this matrix is $W_j^c = J_j^{c\top} J_j^c$, $j \in [0..S]$, which corresponds to the approximation of $H''(\bar{\mathbf{X}}_j^c, \bar{\mathbf{m}}_j^c)$ introduced in Section 3.2.2. This matrix is estimated during the model-based refinement step for the key-frame j . W_j^c is determined only once per key-frame and stored in addition to the optimal pose parameters $\bar{\mathbf{X}}_j^c$. The hybrid error defined in equation 20 becomes:

$$E(\bar{\mathbf{X}}) = E(\bar{\mathbf{X}}^c) + \frac{1}{2}\delta^\top W^c \delta \quad (21)$$

Since $E(\bar{\mathbf{X}}^c) = H(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c) - H(\bar{\mathbf{X}}^c, \bar{\mathbf{m}}^c) = 0$ according to equation 18, the hybrid error definition becomes:

$$E(\bar{\mathbf{X}}) = \frac{1}{2}\delta^\top W^c \delta \quad (22)$$

This definition of our hybrid constraint presents several advantages. W^c converts the difference between the pose parameters to an error in pixels and makes this error homogeneous with the multi-view term of the constrained bundle

adjustment (equation 5). It also allows not to store contour images for each key-frame contrary to equation 2. Moreover, neither the edgelets, nor their associated observations intervene, allowing not to store any of them. Only the model-based pose $\{R_j^c, \mathbf{t}_j^c\}$ and the 6×6 matrix W_j^c are exploited in the optimization process for each key-frame.

This hybrid constraint results in a bundle adjustment with a memory footprint invariant to the video resolution and the model complexity.

6 Experimental results

In this entire Section, the EC-SLAM of [33] exploiting static edgelets extracted from sharp edges of CAD models, and using a re-projection error as presented in Section 3.2.1 is called rSEC-SLAM. On the contrary, our SLAM solutions exploiting dynamic edgelets extracted from sharp and occluding contours (see Section 4) are called DEC-SLAM. The approach using the re-projection error is named rDEC-SLAM. If the first hybrid formalism is used (see Section 5.2), our solution is called rhDEC-SLAM. Finally, our DEC-SLAM exploiting the second hybrid representation (see Section 5.3) is named phDEC-SLAM.

We first evaluate in Section 6.2 our EC-SLAM framework exploiting dynamic edgelets and using the second hybrid constraint presented in Section 5.3. Results are only presented with the phDEC-SLAM, since it has similar results in term of accuracy as the other ones. Secondly, we evaluate in Section 6.3 the different model constraint formalisms described in Sections 3.2.1, 5.2 and 5.3 in term of computation time, memory consumption and accuracy. Our DEC-SLAM is also compared on these criteria with the solution rSEC-SLAM of [33] exploiting static edgelets. Finally, in Section 6.4, our DEC-SLAM using the three model constraint formalisms is evaluated on two public benchmarks: CoRBS [38] and ICL-NUIM [7].

Evaluations are performed on synthetic and real data with objects that have different natures to assess the genericity of the proposed solution. In the experiments, we use a laptop with an Intel (R) Core (TM) i7-4800MQ CPU @ 2.70GHz processor and a NVIDIA GeForce GT 730M graphics hardware.

6.1 Synthetic data presentation

The synthetic sequences used for the experimental results, are generated with a resolution of 640×480 . Only the tracked object is modified. The camera trajectory and the object environment are the same for all the sequences. The camera / object distances do not exceed 8 meters and the object of interest is not always entirely in the camera field of view.

The object environment is composed of 4 walls covered with brick texture and an untextured ground. Several objects, with a volume of about 4m^3 have been tested:

- a torus (Figure 9(a)), a simple curved object with a 13K face model.
- a dragon (Figure 9(b)) that presents both sharp edges and occluding contours. The associated model has 100K faces.
- a part of a robotic exoskeleton arm, referred to as an orthosis, (Figure 9(c)) with many sharp edges. The CAD model of this object has 264K faces.
- a dwarf (Figure 12(a)), a curved object with a 8K face model.
- a bypass (Figure 17(a)), a curved object from the chemical industry, and mainly composed by pipes. Its CAD model used during the tracking has 152K faces. The bypass environment is slightly different from the others. For some robustness evaluation (see Section 6.2.4), the bypass can be occluded by an orthosis as shown in Figure 11(b).

The EC-SLAM initialization for the synthetic sequences is obtained by the ground truth pose.

Quantitative results are obtained by measuring the difference between the ground truth and the estimated camera positions. These localization errors are expressed as a percentage of the camera / object distance. An other possible measure used in these experimental results is an orientation error expressed in radians between the the camera orientation and the ground truth one. A mean 2D error expressed in pixels is also computed and corresponds to the mean 2D distance between edgelets projected according to the estimated camera poses and edgelets projected according to the ground truth camera poses.

6.2 DEC-SLAM framework evaluation

Evaluations of our DEC-SLAM framework through phDEC-SLAM solution, concerns the edgelet orientation definition (presented in Section 4.1.3), the sampling strategy (Section 4.1.4), the transfer time between GPU and CPU of edgelet information (Section 4.1.5), and the exploitation of the Lhuillier cost function in our constrained bundle adjustment to deal with erroneous constraint (Section 3). A first comparison with a model-based tracker algorithm is also performed to demonstrate the robustness and the stability of our solution when sudden motions occur. A second comparison with the rSEC-SLAM of [33] is achieved to evaluate the use of dynamic edgelets against static ones extracted offline on sharp edges only. Finally, our DEC-SLAM algorithm is tested on several real sequences tracking any kind of objects, to assess its genericity.

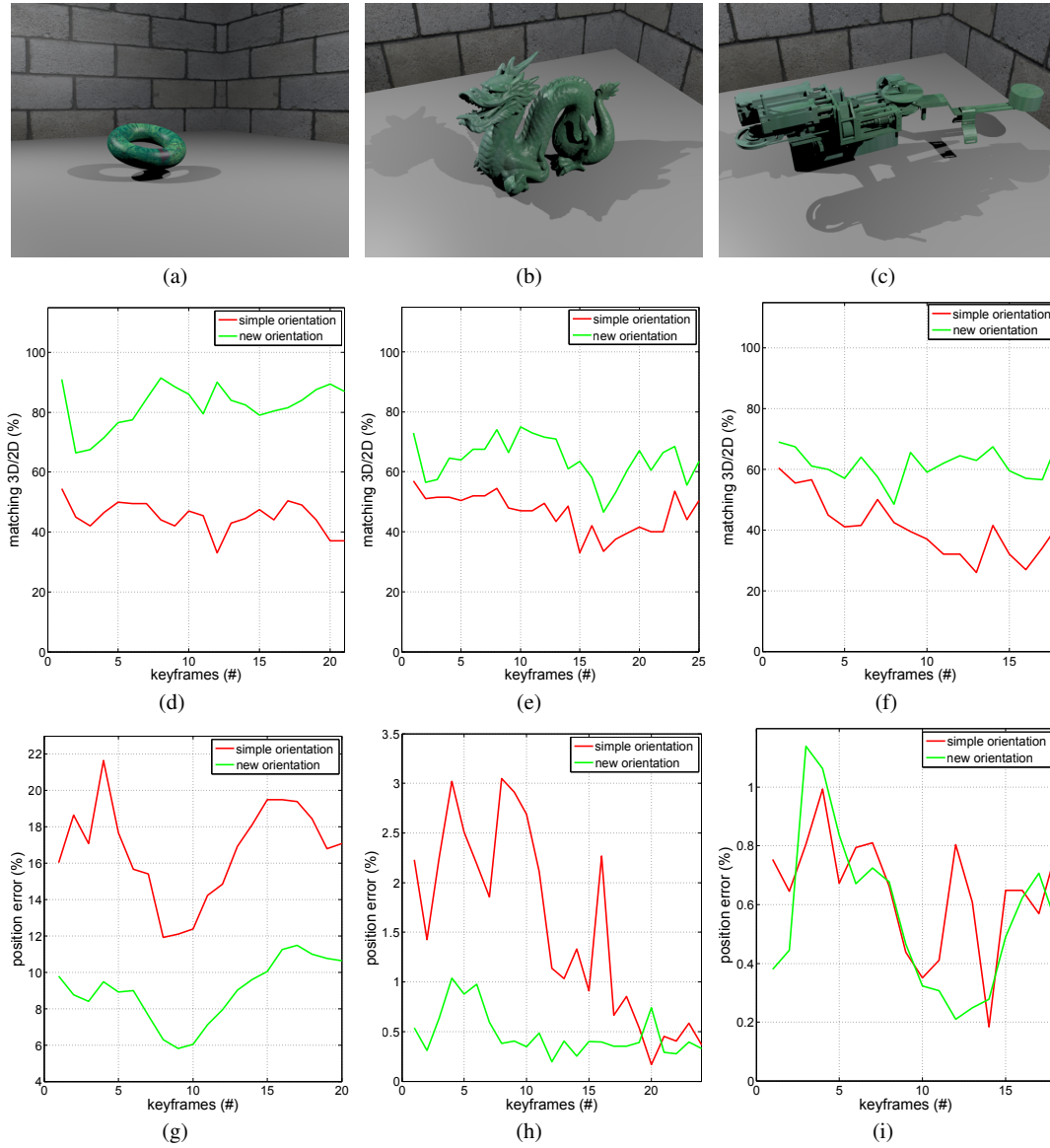


Fig. 9: Comparison of DEC-SLAM with simple 3D orientation definition and the solution we proposed that uses dynamic generated edgelets with new orientations. In green our solution exploiting the presented edgelet generation and in red the other method.

6.2.1 Edgelet orientation

The 3D edgelet orientation described in Section 4.1.3 is compared to a simple orientation definition based on the 3D positions of two adjacent contour points of the rendered model in the image space.

The edgelet orientation is used in two processes of DEC-SLAM algorithm: the 2D/3D matching step when edgelets are associated to image contours, and then the camera pose estimation based on the minimization of the orthogonal distance between projected 3D edgelets and their associated 2D

contours (for the re-projection formalism or the optimization of the output model-based poses for the hybrid constraint representations). Thus the 3D orientations are first evaluated in a qualitative way, then their impact on the rate of 2D/3D matches and on the accuracy of the camera position are evaluated.

Figure 8 shows the 3D orientations obtained for a simple object as the torus. Figure 8(a) corresponds to the simple 3D orientation, whereas Figure 8(b) represents our 3D edgelet orientation for a given point of view. With our method, the 3D directions seem to correspond to the expected ones. On the contrary, with the simple definition, the 3D orientation

are not well estimated because of the aliasing effect from the rasterization step of the rendering process.

Figure 9 presents the impact of our new 3D orientation definition. The evaluation is performed on the torus (Figure 9(a)), the dragon (9(b)) and the orthosis (9(c)) sequences in order to confirm the genericity of our method. Our DEC-SLAM is compared to one with simple 3D edgelet orientation. The impact on the 2D/3D matching step is shown on Figures 9(d), 9(e) and 9(f) representing the rate of edgelets associated to image contours. Indeed, the well estimation of the 3D edgelet direction increases the probability that a projected edgelet matches with an image contour with a similar orientation. The 2D/3D matching rate is presented in percent according to the number of extracted edgelets for each key-frame. For each sequence, the 2D/3D matching mean rate is higher for our solution with 82.2% against 45.0% for the torus sequence, 64.2% against 46.3% for the dragon sequence and 63.1% against 41.4% for the orthosis sequence. The gap between the two methods is more important when the object is curved, the aliasing effect being more significant.

Finally the accuracy of our solution exploiting the new 3D edgelet orientation definition is compared to a DEC-SLAM with a simple edgelet orientation. It is evaluated through a position error. The quantitative results are presented in Figure 9(g), 9(h) and 9(i). On the torus sequence, for both methods, the position error is important. The optimization process is indeed less constrained due to the object symmetry. However, our 3D edgelet orientation allows a camera position more accurate since the mean position error is around 9% whereas the solution with the simple orientation definition has a mean position error of more than 16%. For the dragon sequence and the orthosis one, the tracking is also more accurate with the new 3D orientation. The mean position error is around 0.5% for the dragon sequence using the proposed 3D orientation definition, against 1.5% with the simple method. For the orthosis sequence, the object of in-

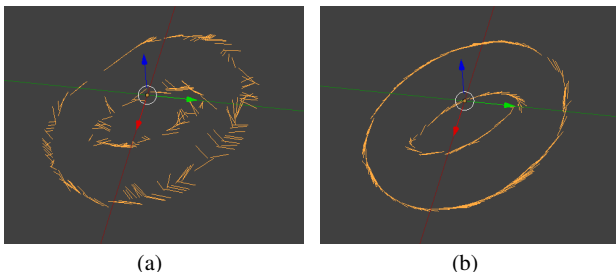


Fig. 8: Representation of 3D edgelet orientations of the torus rendered model for a given camera point of view with simple orientation on the left and with our proposed orientation on the right.

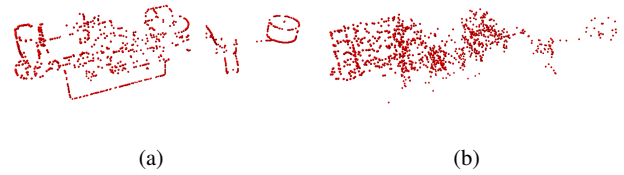


Fig. 10: Set of edgelets extracted with our sampling method on the left and with a random sampling solution on the right.

	min (%)	max (%)	mean (%)	std. (%)
Random sampling	0.9228	19.0420	8.4383	8.4943
Proposed sampling	0.5788	10.3170	3.5466	3.4458

Table 1: Localization errors (%) with different sampling strategies.

terest is polyhedral and the aliasing effect is less perceptible. Thus both methods are accurate with a mean position error less than 1%.

The DEC-SLAM with the proposed 3D orientation allows a higher 2D/3D matching rate and a higher or similar accuracy than the one with simple edgelet orientation.

6.2.2 Sampling strategy

The proposed sampling strategy described in Section 4.1.4 is compared with a random sampling. This evaluation is performed on the orthosis sequence where the proposed solution is ran a hundred times with the two sampling strategies (400 edgelets are used). The mean errors over the 100 trials and over all the images of the orthosis sequence are estimated, as well as the min/max errors and the standard deviation (std.). Results are reported in Table 1. Our proposed sampling better constrains the matching and the pose estimation, halving the localization errors. Edgelets resulting of the different sampling methods are also shown in Figure 10. Our sampling allows a more evenly distribution.

6.2.3 Edgelet data transfer time

Our method of edgelet transfer from GPU to CPU detailed in Section 4.1.5 is compared to a simple asynchronous texture transfer using Pixel Buffer Object.

In the presented approach, the transfer begins with the synchronization between GPU and CPU for the data reading access. Then the filtering pass follows, storing only the edgelet data. The transfer ends with the reading of the VBO containing edgelet information. With the simple texture trans-

	SD sequence (ms)		HD sequence (ms)	
	mean	std.	mean	std.
Texture transfer	3.55	0.17	17.07	3.75
Proposed solution	0.56	0.10	1.03	0.23

Table 2: Computation time for the edgelet data transfer between GPU and CPU. The comparison is made between a texture transfer and the proposed method.

fer method, the transfer also begins with the synchronization between GPU and CPU. However, there is no filter performed and the transfer ends when all the textures containing edgelet information are read.

Table 2 presents the mean and standard deviation (std.) of transfer and reading times on the whole key-frames of the dragon sequence. Computation time are expressed in milliseconds. To evaluate the impact of our optimization, transfer computation time is evaluating on a Standard Definition (SD) sequence with a 640×480 resolution and a High Definition (HD) sequence with a 1280×960 resolution.

Table 2 shows an important time saving since transfer time between the simple texture transfer and our proposed solution decrease of over 84% for the SD video and of 94% for the HD one. Moreover, we can see that the simple texture transfer approach is almost five times more time consuming between the SD and the HD sequence, whereas our proposed transfer method is only almost twice more time consuming. Edgelet transfer time with the simple texture transfer method depends directly of the sequence definition since all the pixels of all the textures are transferred to CPU. On the contrary, our method proposes to only transfer the edgelet pixels through an optimized array. Thus the presented approach runs faster and is less affected by the sequence resolution.

6.2.4 Lhuillier cost function

In this experiment, the robustness of our solution exploiting inaccurate constraints is evaluated. Particularly, we seek to compare our constrained bundle adjustment using Lhuillier cost function and another one using a classic cost function optimizing simultaneously the environment and the model constraints as in [33]. The hybrid constraint may indeed be erroneous, especially if the output of the model-based tracker is inaccurate when the object is occluded or small in the image.

The experiment consists in occluding an object of interest to disturb the model-based tracking and obtain erroneous model-based poses for the constrained bundle adjustment. Figure 11 shows the position errors and 2D errors with our

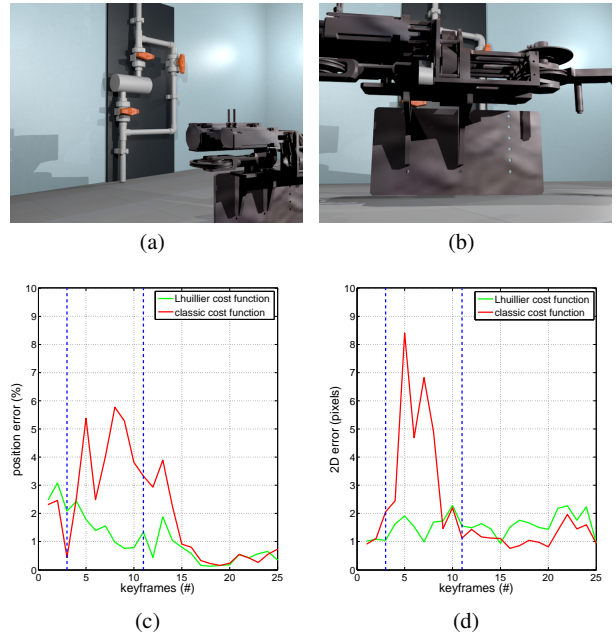


Fig. 11: Comparison between Lhuillier cost function and classic cost function on a synthetic sequence. The bypass is the object of interest and is occluded by an orthosis in a part of the video. (a) and (b) present some frames of this sequence. (c) shows the position errors of our phDEC-SLAM with a classic cost function (green) or Lhuillier cost function (red) in the constrained bundle process. The blue dots correspond to the interval where the bypass is occluded and the model-based poses are erroneous. (d) present the 2D errors of both methods.

phDEC-SLAM approach and with a phDEC-SLAM constrained with a classic cost function. The two object trackers are experimented on a synthetic sequence where a bypass is momentarily occluded by an orthosis as shown in Figures 11(a) and 11(b). When the bypass is occluded, the model-based pose can reach a position error of almost 22%. However this erroneous constraint does not degrade the camera pose since after the optimization process exploiting Lhuillier cost function, the mean position error is less than 1% (see Figure 11(c)) and the mean 2D error is 1.56 pixels. However with the use of a classic cost function, the localization is less accurate when the model-based pose is erroneous. The mean localization error is more than 2% and the mean 2D error is around 2 pixels.

6.2.5 Comparison to model-based tracking

We compare our DEC-SLAM with model-based tracking in term of robustness against large displacements. The latter uses edgelets extracted with the proposed solution described in Section 4. For the matching step, the pose estimated on

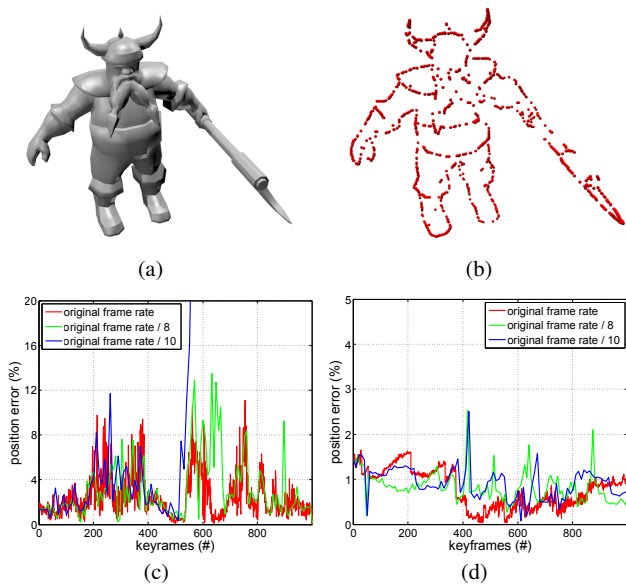


Fig. 12: Comparison of model-based tracking and our DEC-SLAM solution on the dwarf sequence. (a) Illustration of the dwarf. (b) Edgelets extracted with the solution described in Section 4 for a given camera pose. Localization errors resulting from model-based tracking (c) and from the proposed DEC-SLAM (d) on the dwarf sequence with the original frame rate (red), a frame rate divided by 8 (green) and by 10 (blue).

the previous image is used to project each edgelet and the nearest 2D image contour with a similar orientation is selected as its correspondent. The comparison is made on the dwarf synthetic sequence described in Section 6.1. The two algorithms are evaluated by varying the frame rate of the camera on the dwarf sequence. The travel speed of the camera is unchanged.

Three scenarios are tested: the original frame rate and frame rates divided by 8 and by 10. Figure 12(c) and 12(d) represent the position errors for the model-based tracking and the proposed solution respectively. The DEC-SLAM succeeds on the three scenarios. The mean error stays in the interval $[0.78\%, 0.92\%]$ and is almost constant for all frame rates. On the other hand, model-based tracking fails on the sequence with the lowest frame rate due to too large amplitude motions. The mean position error increases with the displacement amplitude, its values are 2.41% and 6.65% for the sequences with the original frame rate and with a frame rate divided by 8, respectively. Moreover, the standard deviation is important (2.76% vs 0.33% for our solution) due to tracking instabilities (jitter). The proposed solution is thus more accurate and more robust to large displacements than model-based tracking.

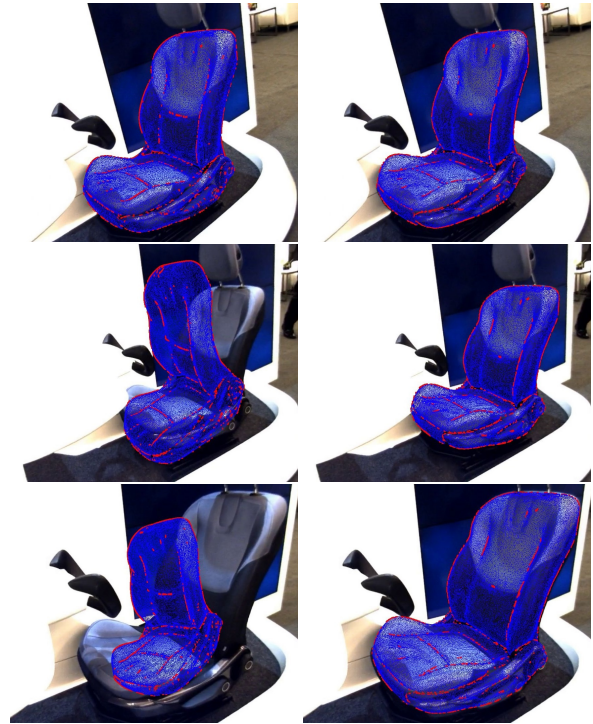


Fig. 13: Comparison of the proposed solution (right) with model-based tracking (left) on the car seat sequence where sudden motions occur.

This robustness comparison is also made on a real sequence to obtain qualitative results, as shown in Figure 13. The tracked object corresponds to a car seat that is fully curved and whose CAD model has 79K faces.

Thanks to the SLAM prediction based on key-points, the generated edgelets at each key-frame are extracted from an accurate pose contrary to the model-based approach, that uses the pose estimated on the previous frame. Moreover, the pose predicted by the SLAM also facilitates the 3D/2D matching step, especially when large amplitude motion occurs. Thus our solution results in a stable (jitter-free) localization while model-based tracking presents some localization instabilities and fails to track when the motion amplitude is large as illustrated in Figure 13. Additional data are given in the first part of Online Resource 1, with a comparison video on the car seat sequence.

6.2.6 Comparison to rSEC-SLAM

We compare our proposed solution with the rSEC-SLAM algorithm described in [33] in terms of localization accuracy and tracking genericity. The dragon and the orthosis sequences are used for the evaluation. Both objects have sharp edges, but they are not predominant for the dragon. Fig. 14(a) and 14(b) represent position errors on the two sequences. For the dragon sequence, the proposed solution is

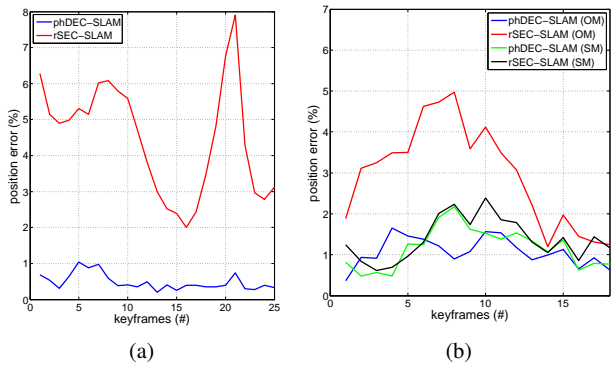


Fig. 14: Comparison of the rSEC-SLAM of [33] and our phDEC-SLAM, on the dragon (a) and the orthosis (b) sequences. Localization errors for the orthosis sequence are measured with the use of a simplified model (SM) and the original one (OM).

5 times more accurate than the one proposed by [33]. In fact, the edgelets obtained with the dihedral criteria are not evenly distributed on the object whereas the dynamic edgelets generated with the solution we proposed (see Section 4) better represents the silhouette and object contours. For the orthosis sequence, the original model and a simplified version are used. With the proposed solution, the localization errors are almost the same whatever the model complexity. On the other hand, the localization accuracy of the rSEC-SLAM is affected by the model complexity. Mean errors values are of 1.18% and 2.36% with the simplified and the original model respectively. In fact, without model simplification, edgelet detection with the dihedral criteria results in a noisy constellation. The proposed solution does not require any model simplification to generate an evenly distributed edgelet constellation and thus to achieve accurate tracking.

Figures 15(a) and 15(b) illustrate on a real bypass sequence the same issue with a set of 2000 edgelets dynamically generated with the solution described in Section 4 and extracted with a dihedral criteria, respectively. A constellation of edgelets obtained via the dihedral criteria is not evenly distributed, whereas the edgelets obtained with the proposed solution are distributed on the whole object surface, including the pipes.

The accuracy of EC-SLAM is directly affected by the quality of the edgelet constellation as seen in Figure 15(c) and 15(d) on the real bypass sequence. The DEC-SLAM is more accurate than the rSEC-SLAM proposed in [33] since their edgelets constellation poorly constraints camera poses. Additional data are given in the second part of Online Resource 1, with a comparison video on the bypass sequence.

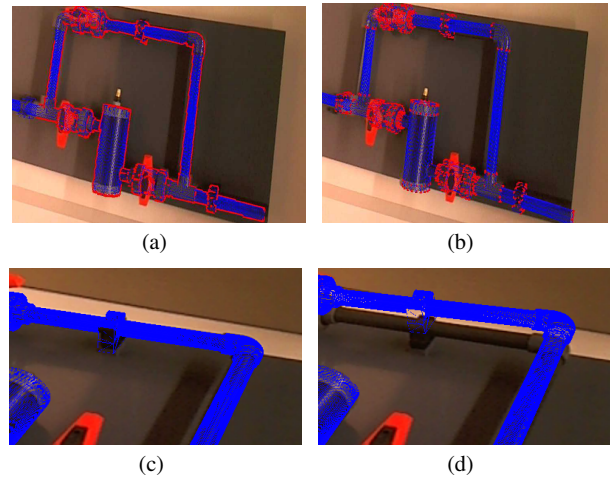


Fig. 15: Comparison between DEC-SLAM (left) and the rSEC-SLAM of [33] (right) on the bypass sequence. (a) Dynamic generated edgelets for that point of view (b) Static edgelets generated with a dihedral criteria. (c) and (d) Model re-projections on a bypass portion.

6.2.7 Tracking genericity evaluation

To finally assess the genericity of our DEC-SLAM method, experiments are carried on several real sequences with different kind of objects. Our solution presents similar accurate and robust tracking results on these objects of interest whatever the model constraint formalism. Then Figure 24 shows 3D model projections on the real scene according to the estimated camera poses with phDEC-SLAM solution only.

Our proposed solution has been successfully tested with objects for which a CAD model is available, for example with a sport car (Figure 24(g)), a bypass from the chemical industry (Figure 15) and a seat from a car OEM (Figure 13). The genericity of our solution is also demonstrated by tracking objects like a metal statue of dragon (Figure 24(d)), and a Raving Rabbid (Figure 24(a)) using 3D models reconstructed by photogrammetry. Industrial objects like a real-sized car (Figure 24(j)), a car cylinder head (Figure 24(m)), or an other kind of orthosis (Figure 24(p)) with a known CAD model have been tracked successfully.

These objects might be polyhedral (the sport car, the cylinder head or the orthosis) or curved (the Raving Rabbid, the dragon or the car seat), but they also can present sharp and occluding contours at the same time (the bypass or the real-sized car). Some are textured (the dragon), while others are absolutely textureless (the Raving Rabbid).

Our solution is robust to occlusions, thanks to the proposed DEC-SLAM framework, which guarantees the multi-view relationships to be well estimated. Additional data are given in the third part of Online Resource 1, with a compilation

	Polyhedral object (sport car)		
	rSEC-SLAM [33]	rhDEC-SLAM	phDEC-SLAM
3 optimized key-frames	47.9	27.9	23.6
10 optimized key-frames	130.5	77.9	76.2
30 optimized key-frames	216.6	203.9	188.7

Table 3: Computation time in milliseconds of the mapping process for the rhDEC-SLAM, phDEC-SLAM and rSEC-SLAM on the sport car sequence. It is given for 3,10 and 30 optimized key-frames.

	Curved object (dragon statue)		
	rDEC-SLAM	rhDEC-SLAM	phDEC-SLAM
3 optimized key-frames	97.8	94.8	94.1
10 optimized key-frames	172.1	170.6	130.3
30 optimized key-frames	385.3	330.0	261.3

Table 4: Computation time in milliseconds of the mapping process for the rDEC-SLAM, rhDEC-SLAM and phDEC-SLAM on the dragon sequence. It is given for 3,10 and 30 optimized key-frames.

video that shows our tracking solution on several polyhedral and curved objects.

6.3 EC-SLAM comparison

The different formalisms of the presented solution are compared in terms of computational cost, memory consumption and accuracy in this Section.

6.3.1 Computation time evaluation

In this Subsection, the different model constraint formalisms integrated in our DEC-SLAM framework are evaluated in term of computational costs. Since these solutions differ only by their bundle adjustment, this experiment consists in comparing the processing time required by the mapping process described in Section 3. The mapping process is achieved once per new key-frame, then we compare the median processing time of these executions. Also, since the computation time depends on the number of optimized key-frames, the different solutions are compared for a set of 3, 10 and 30 optimized key-frames.

This experiment is achieved on two real sequences recorded with a HD resolution (1280×1024). The two objects used are a metal statue of dragon (see Figure 24(d)) with a 3D model reconstructed by photogrammetry and a sport car (see Figure 24(g)), whose the CAD model is available. The SLAM

initialization on these real videos is performed manually and followed by a contour refinement process.

The sport car illustrates the performances for a polyhedral object, while the dragon sequence illustrates the performances for a curved object. For the car sequence, the performances of our rhDEC-SLAM and phDEC-SLAM are compared to the performances of rSEC-SLAM of [33]. For the dragon sequence, the rDEC-SLAM, rhDEC-SLAM and phDEC-SLAM are compared in term of computational cost.

The computation time of the mapping processes is presented in Tables 3 and 4. The classic bundle adjustment optimizing the multi-view constraint and the dynamic edgelet generation (and optimal output model-based poses for both hybrid constraint representations) are running in parallel, to optimize the computation time. Moreover, even if rDEC-SLAM and rSEC-SLAM have not the model-based pose refinement step in their framework (contrary to the rhDEC-SLAM and phDEC-SLAM), the 2D/3D associations between edgelets and their 2D counterparts have to be re-estimated at each constrained bundle adjustment on all the optimized key-frames. This process participates in slowing down its execution. Computation time is globally more important with the tracking of curved objects. It is due to the edgelet generation process, which has to be performed online at the last key-frame since the occluding contours depend on the camera point of view (see Section 3.1). On the contrary, the edgelet extraction is performed offline only once for polyhedral objects. Tables 3 and 4 show that the exploitation of any hybrid constraint representation allows to reduce computation time. Compared to the use of re-projection formalism, the computation time decreases of around 44% for the

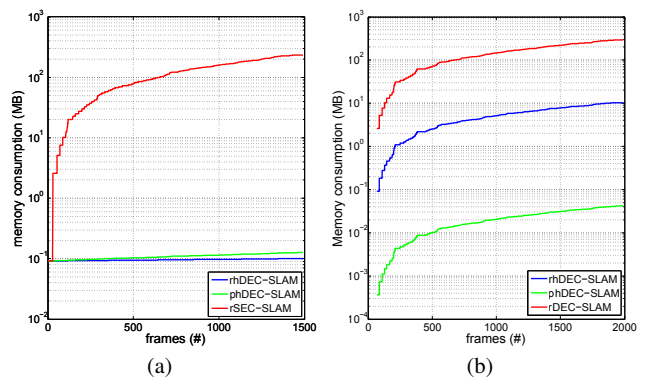


Fig. 16: The memory footprint associated to the car sequence (a) and the dragon one (b), is drastically decreased with the use of the rhDEC-SLAM and phDEC-SLAM (see Section 5) compared to rSEC-SLAM of [33] or our rDEC-SLAM. The memory consumption is expressed in MB with a logarithmic scale on the y axis.

sport car sequence and around 8% on the dragon one even if a model-based pose refinement is performed in addition to the constrained bundle adjustment with the use of hybrid constraint (with 3 and 10 key-frames in the bundle adjustment). The computation time decrease of around 10% and 23% with the use of 30 key-frames.

6.3.2 Memory consumption comparison

Optimizing a SLAM reconstruction with a bundle adjustment constrained to a 3D model implies the storage of additional data for each key-frame. In this Section, the memory footprint of these data is evaluated, depending on the use of the chosen model constraint formalism and the edgelet generation process (static or dynamic). These evaluations are conducted on the two real sequences introduced in Section 6.3.1. The edgelet extraction is performed offline for the car sequence and online for the dragon statue. Our rhDEC-SLAM and phDEC-SLAM are compared to [33] for the sport car and with our rDEC-SLAM for the dragon. The results are summarized in Figures 16(a) and 16(b).

As explained in Section 3.2, rDEC-SLAM or rSEC-SLAM require to store edgelets and contour images for each key-frame optimized in the constrained bundle adjustment. Consequently the memory footprint increases every time a new key-frame is created. The impact is even greater as the video resolution (for curved objects) and the number of edgelets are high. For the real sequences, the number of edgelets projected and associated to 2D contour points is set to 2000, and 94 (respectively 113) key-frames are created over the car sequence (respectively dragon sequence).

Thus, as shown in Figures 16(a) and 16(b), the memory footprint for rSEC-SLAM is near 235 MB at the end of the sport car sequence and more than 292 MB for the rDEC-SLAM at the end of the dragon one. The latter sequence requires more memory since edgelets extracted from occluding contours are stored for each key-frame, which is not required with polyhedral objects. The memory footprint drastically decreases with the use of a hybrid constraint formalisms, since contour images are not required anymore in the con-

strained bundle adjustment. Pose parameters of the model-based tracker outputs are henceforth stored instead. In addition to the pose parameters storage, our rhDEC-SLAM requires to save sets of edgelets. Since they are extracted only once with polyhedral object, the memory consumption is limited to 0.1 MB at the end of the car sequence. However, generating and accumulating sets of edgelets at each key-frame for curved objects, increase the memory footprint to 10 MB at the end of the dragon sequence. Finally the phDEC-SLAM allows a memory footprint close to 0.12 MB (respectively 0.04 MB) at the end of the sport car sequence (respectively dragon sequence), only pose parameters and matrices W^c of the model-based tracker outputs (see Section 5.3) being stored. We can notice that with polyhedral objects, the memory consumption is higher for the phDEC-SLAM (0.12 MB) than for the rhDEC-SLAM (0.10 MB). Indeed, in addition to the pose parameters stored for both hybrid constraint representations, for the rhDEC-SLAM, a unique set of edgelets is stored for the whole sequence, whereas for the phDEC-SLAM, matrices W_j^c are stored at each key-frame creation.

This latter allows our memory consumption to be invariant to the number of edgelets and their generation, invariant to the resolution of the tracking video. Thus, if polyhedral (respectively curved) objects are tracked, the rhDEC-SLAM (respectively phDEC-SLAM) is more suitable. Additional data are given in the fourth part of Online Resource 1, with a comparison video on the sport car and dragon sequences.

6.3.3 Accuracy evaluation

Our proposed DEC-SLAM algorithm is evaluated in terms of accuracy and robustness on polyhedral and curved objects and also compare to rSEC-SLAM described in [33].

The accuracy of our DEC-SLAM solution is evaluated on two synthetic sequences. The first one corresponds to the bypass sequence (see Figure 17(a)) presented in Section 6.1. The second sequence has a sedan car as object of interest. Its CAD model has 50K faces. The sedan car environment is composed of city buildings as seen in Figure 17(b).

Figures 18(a), 18(b) and 18(c) present the localization error in position and orientation, and the 2D error on the bypass sequence. Figures 18(d), 18(e) and 18(f) present the same errors on the sedan car sequence.

The rDEC-SLAM and the phDEC-SLAM provide similar accuracies on both polyhedral and curved objects, contrary to the rhDEC-SLAM that presents some lacks of accuracy. On the sedan sequence, the median pixel error is 1.38 for the rhDEC-SLAM, 1.12 for the phDEC-SLAM, 1.28 for the rDEC-SLAM and 1.63 for the rSEC-SLAM of [33]. On the bypass sequence, the median pixel error is 1.39 for the rhDEC-SLAM, 0.99 for the phDEC-SLAM, 0.91 for the rDEC-SLAM and 0.66 for the rSEC-SLAM of [33]. This

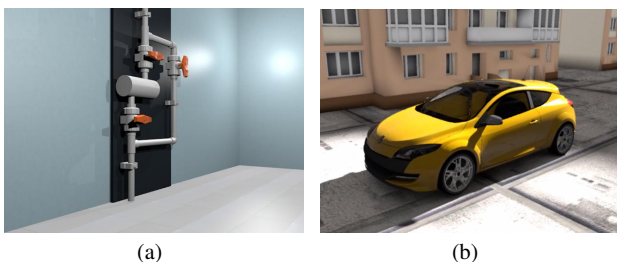


Fig. 17: Images from the bypass (a) and sedan car (b) synthetic sequences.

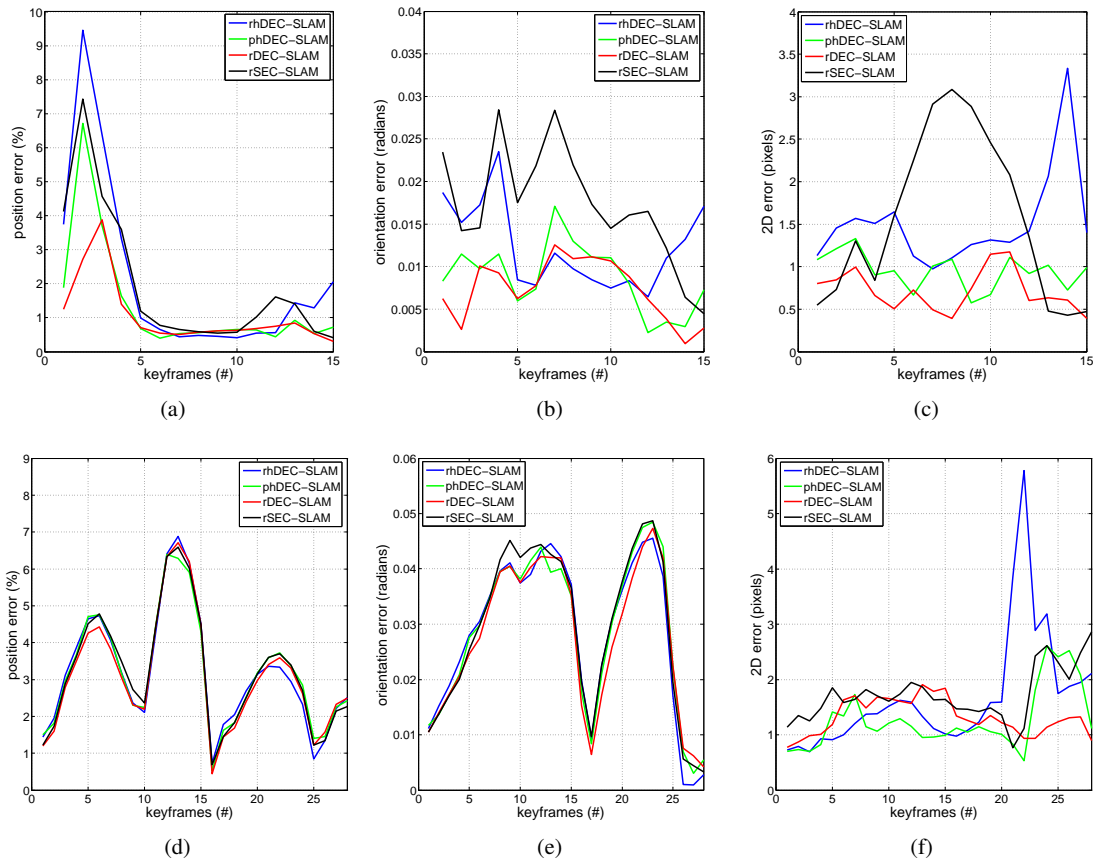


Fig. 18: Accuracy comparison between static edgelet cSLAM [33] (black), dynamic edgelet cSLAM (red) and the first hybrid (blue) and the second (green) cSLAM on the bypass (top) and sedan car (bottom) sequences.

latter is less accurate when the object of interest is curved. Edgelets obtained with the dihedral criteria are not evenly distributed on the object whereas the dynamic edgelets generated with the proposed solution, better represents the silhouette and object contours.

6.4 Evaluation on public datasets

To assess the robustness and accuracy of our DEC-SLAM approach, we run our method on two public datasets. Even if these benchmarks are more adapted for RGBD solutions, EC-SLAM are able to localize accurately the camera without depth information. The SLAM initialization on these sequences is obtained by the first ground truth pose.

6.4.1 CoRBS dataset

Our DEC-SLAM solution is evaluated and compared on the Comprehensive RGBD Benchmark for SLAM (CoRBS) dataset

that proposes real sequences with different objects of interest. It provides the real depth and RGB frames thanks to the use of a Kinect v2, with a ground truth trajectory of the camera and a ground truth 3D model of the scene. The CoRBS dataset is composed of 5 sequences named kt0, kt1, kt2, kt3 and kt4 for each object of interest (a desk and a wooden human manikin). An overview of these two objects with their reconstructed 3D models is presented in Figure 19. This benchmark is interesting to evaluate our DEC-SLAM with the different model constraint formalisms, since the tracked objects are polyhedral and complex like the desk, or curved as the wooden human manikin. The sequences have a 640×480 definition and a frequency of 30Hz.

We evaluate our rDEC-SLAM, rhDEC-SLAM and phDEC-SLAM solutions and compare it to RGBD C-SLAM [20] (only the kt0 sequence is available for this approach) and rSEC-SLAM [33]. Table 5 describes the Absolute Trajectory Error (ATE) inspired by [32] and expressed through the Root-Mean-Square-Error metric (RMSE). This error quantifies the accuracy of the entire trajectory for all the five desk sequences. The rDEC-SLAM solution presents the best re-

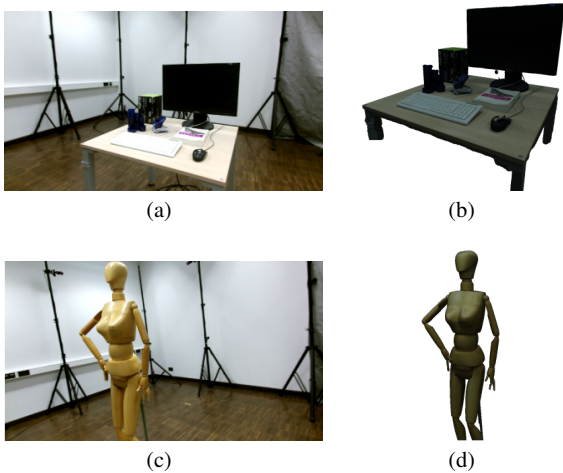


Fig. 19: CoRBS sequences [38]. Frames from the different sequences where the desk and human manikin are tracked (left). Their associated reconstructed 3D model (right).

sults in term of accuracy for the first three sequences and is very close of the most accurate approach on the last two sequences. For the kt0 one, RGBD C-SLAM has the same RMSE of 0.013 meters as our rDEC-SLAM. However the other model constraints formalisms integrated into our DEC-SLAM solution provide also accurate localization with a trajectory error around 0.020 meters for this sequence. On this polyhedral object, rSEC-SLAM of [33] reaches good performances, especially for the kt1 and kt3 sequences where it gets the smallest trajectory error. Our rhDEC-SLAM and phDEC-SLAM also track with accuracy the desk with similar localization error. The first one has the best result for the kt1 and kt4 sequences, the second one following just near by. This latter has the best RMSE for the kt2 sequences.

Table 6 describes the localization error for all the five human manikin sequences. The object of interest is a curved one. Then rSEC-SLAM [33] and RGBD C-SLAM are less accurate, few sharp edges existing on the 3D model. Edgelets extracted with the dihedral criteria are not well distributed on the model and do not well constrain the camera pose estimation. The rDEC-SLAM approach obtains the best trajectory errors except for the kt1 sequence where the phDEC-SLAM precedes it with a RMSE of 0.146 meters.

The proposed DEC-SLAM method obtain accurate results on CoRBS sequences with all the model constraint formalisms, the objects of interest being polyhedral or curved.

6.4.2 ICL-NUIM dataset

The Imperial College London and National University of Ireland Maynooth (ICL-NUIM) dataset [7] is a benchmark

		Desk sequence					
		Error (m)	kt0	kt1	kt2	kt3	kt4
RGBD C-SLAM [20]	RMSE	0.013	-	-	-	-	-
	Std.	0.005	-	-	-	-	-
	Min	-	-	-	-	-	-
	Max	0.073	-	-	-	-	-
rSEC-SLAM [33]	RMSE	0.020	0.010	0.601	0.287	0.354	
	Std.	0.006	0.003	0.238	0.126	0.256	
	Min	0.006	0.001	0.047	0.053	0.017	
	Max	0.048	0.026	1.011	0.841	1.244	
rDEC-SLAM	RMSE	0.013	0.010	0.600	0.294	0.360	
	Std.	0.005	0.004	0.238	0.128	0.257	
	Min	0.003	0.001	0.042	0.061	0.017	
	Max	0.041	0.025	1.007	0.858	1.246	
rhDEC-SLAM	RMSE	0.020	0.010	0.603	0.288	0.351	
	Std.	0.006	0.004	0.238	0.126	0.257	
	Min	0.006	0.001	0.053	0.055	0.016	
	Max	0.049	0.025	1.016	0.854	1.260	
phDEC-SLAM	RMSE	0.021	0.011	0.600	0.293	0.358	
	Std.	0.012	0.004	0.237	0.129	0.254	
	Min	0.001	0.001	0.042	0.066	0.026	
	Max	0.063	0.030	1.003	0.854	1.222	

Table 5: Trajectory error for the CoRBS Desk sequences [38].

		Human sequence					
		Error (m)	kt0	kt1	kt2	kt3	kt4
RGBD C-SLAM [20]	RMSE	0.036	-	-	-	-	-
	Std.	0.017	-	-	-	-	-
	Min	-	-	-	-	-	-
	Max	0.208	-	-	-	-	-
rSEC-SLAM[33]	RMSE	0.032	0.223	0.077	0.740	0.147	
	Std.	0.014	0.120	0.047	0.349	0.074	
	Min	0.002	0.015	0.006	0.052	0.017	
	Max	0.067	0.619	0.218	1.572	0.472	
rDEC-SLAM	RMSE	0.019	0.158	0.027	0.510	0.086	
	Std.	0.008	0.087	0.011	0.252	0.048	
	Min	0.001	0.016	0.003	0.024	0.005	
	Max	0.055	0.440	0.059	1.282	0.248	
rhDEC-SLAM	RMSE	0.046	0.150	0.062	0.760	0.272	
	Std.	0.021	0.072	0.032	0.308	0.076	
	Min	0.003	0.026	0.032	0.141	0.130	
	Max	0.095	0.317	0.174	1.524	0.475	
phDEC-SLAM	RMSE	0.030	0.146	0.029	0.516	0.127	
	Std.	0.013	0.073	0.012	0.252	0.074	
	Min	0.003	0.011	0.002	0.043	0.012	
	Max	0.068	0.290	0.070	1.310	0.431	

Table 6: Trajectory error for the CoRBS Human manikin sequences [38].

created for a different context than object tracking evaluation. The scene indeed corresponds to a living room as seen in Figure 20(a), a much larger object of interest than the ones we exploit previously. However, even if this dataset is out of our context and our DEC-SLAM is not necessarily well adapted for this kind of scene when the object of interest is too large for the camera field of view, it is able to run on these sequences. The 3D model of the scene is indeed available (Figure 20(b)). The ICL-NUIM sequences have a 640×480 resolution and are recorded at 30 Hz. Images are synthetic but with real world lightning conditions. The sequences are also noisy according to an RGB noise model to simulate the one providing of real camera.

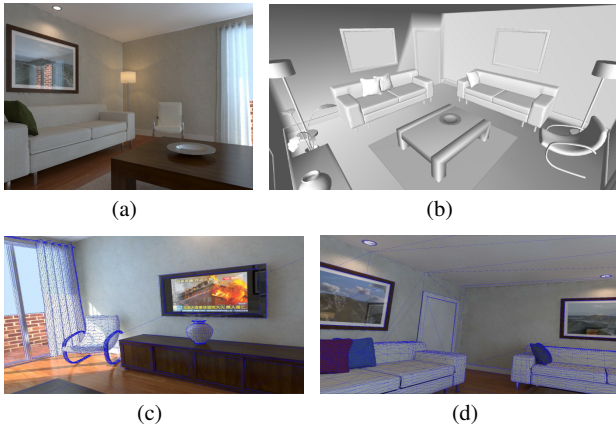


Fig. 20: ICL-NUIM sequences [7]. (a) A frame from the kt1 sequence. (b) The 3D model of the living room used as a constraint in our solution. (c) and (d) 3D model projected on the kt1 and kt2 sequences with the rDEC-SLAM.

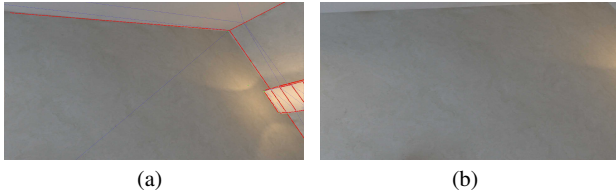


Fig. 21: Our DEC-SLAM is not able to estimate camera poses when any model information or environment key-points are available to constrain our method. Without depth information, our solution fails when the camera is looking to planar textureless regions. (a) and (b) are frames from the ICL-NUIM kt0 sequence [7]. The tracking is lost when only walls are visible.

Table 7 describes the trajectory accuracy expressed with the RMSE metric of [32] as for the CoRBS dataset, for the living room sequences kt1 and kt2. We compare our rDEC-SLAM, rhDEC-SLAM and phDEC-SLAM to the rSEC-SLAM of [33] and RGBD SLAM approaches such as DVO [10], FOVIS [9], ICP of KinectFusion [24] and finally, RGBD C-SLAM [20].

Trajectory errors for the kt0 and kt3 living room sequences are not presented in this paper. Our DEC-SLAM solution does not exploit RGBD frames and are not able to estimate the camera pose when this latter is looking at a planar textureless region of the scene (walls) as it happens with these both sequences (see Figure 21).

However accurate results are obtained on the kt1 and kt2 sequences with our DEC-SLAM and the different model constraint formalisms. On kt1, the rDEC-SLAM and the phDEC-SLAM have the second best trajectory error with a RMSE of 0.018 meters behind ICP solution [24] (RMSE equal to

	Error (m)	kt1	kt2
DVO [10]	RMSE	0.125	0.473
	Std.	0.037	0.175
	Min	0.051	0.137
	Max	0.200	0.834
FOVIS [9]	RMSE	1.868	1.495
	Std.	0.871	0.504
	Min	0.333	0.270
	Max	3.039	2.773
ICP [24]	RMSE	0.005	0.010
	Std.	0.002	0.004
	Min	0.001	0.004
	Max	0.011	0.015
RGBD C-SLAM [20]	RMSE	0.025	0.023
	Std.	0.015	0.011
	Min	-	-
	Max	0.087	0.093
rSEC-SLAM [33]	RMSE	0.024	0.010
	Std.	0.011	0.004
	Min	0.003	0.001
	Max	0.089	0.025
rDEC-SLAM	RMSE	0.018	0.010
	Std.	0.013	0.004
	Min	0.000	0.001
	Max	0.094	0.025
rhDEC-SLAM	RMSE	0.068	0.048
	Std.	0.038	0.017
	Min	0.001	0.005
	Max	0.157	0.077
phDEC-SLAM	RMSE	0.018	0.011
	Std.	0.012	0.005
	Min	0.001	0.001
	Max	0.091	0.034

Table 7: Trajectory error for the kt1 and kt2 ICL-NUIM sequences [7].

0.005 meters). On kt2, the rDEC-SLAM has the best trajectory error with the ICP approach and the rSEC-SLAM [33] that corresponds to a RMSE equal to 0.010 meters. Our phDEC-SLAM as the second best trajectory accuracy with a RMSE of 0.011 meters.

Even if we do not exploit depth information, our DEC-SLAM approach presents similar accurate results compared to [24] and superior results to [10,9,20]. Moreover, the mean 2D error of our DEC-SLAM on both sequences is around 4 pixels, that is small enough for RA application. The 3D model of the scene is globally well re-projected on the images as seen in Figures 20(c) and 20(d).

7 Application

DEC-SLAM is useful to Augmented Reality (AR) applications for maintenance support, automation of complex tasks or other quality controls. Our method is able to localize complex static objects, but also objects with simple dynamic parts thanks to dynamic edgelet generation (see Section 4). Thus more complex scenarios can be proposed to guide in-

dustrial people with AR applications. In the industry domain, objects of interest may be composed of several parts, curved or not. For example, in the petrochemical industry and in a maintenance support scenario, a valve may be open or close during the sequence, or in a part assembly scenario, elements may be piled up successively. In these cases, the objects of interest have different states. To give suitable information, the AR application has then to localize these object parts according to the camera pose and their current state.

7.1 Automatic object part state estimation

In order to estimate the state of a movable object part, the idea is to exploit 3D contours extracted from the CAD model as presented in Section 4.

However it is not enough to only compare the rate of 2D/3D matching between projected edgelets and image contours for each potential state to estimate the accurate current one. Projected edgelets and 2D image contours are not perfectly fused when the model is rendered in the current state (see Figure 22(a)). The camera pose may be inaccurate, the 3D model may not be identical to the object of interest particularly when it is created by photogrammetry. The edgelet extraction and the 2D contour detection may also be inaccurate. The object state estimation is then not obvious. These contour alignment issues have to be taking into account to estimate correctly which edgelets from the different rendered models correspond to the current state, and thus to propose a robust state estimation.

Our state estimation approach is based on a score for each potential state (open, close or half open for example in the case of a valve), computed as a weighted sum of 2D/3D matches between edgelets and 2D contours. The highest score



Fig. 22: Super-imposition between projected edgelets from the current state (red) and 2D image contours (blue) on the left, and between projected edgelets from a false potential state (green) and 2D contours on the right. To estimate the accurate state of the valve, false positive association as between edgelets from the potential state and the static part of the object have to be considered.

		Other potential states			
		2D/3D matching	Same orientation	No match	Wrong orientation
One potential state	Same orientation		0	+	+
	No match		-	0	+
	Wrong orientation		-	-	0

Table 8: Weights for the state estimation.

defines the potential state as the accurate current state. Since some associations must be irrelevant and disturb the state estimation, our method does not look for matching an edgelet to a contour pixel. It aims to associate each 2D contour pixel to edgelets extracted from the different rendered models based on the potential states. In a Region Of Interest (ROI), 2D contours of the static part as well as the movable one are analyzed in order to help the state estimation. The 2D/3D matching step is performed for the pixel with edgelets from all the potential model renders and the weight according to the association status is presented in the Table 8 that follows. Three different status may happen when we want to associate a pixel with an edgelet. For a given 2D contour pixel, a match is found with a projected edgelet having the same orientation, a match is found with a projected edgelet having a wrong orientation, or no match is found. The weight according to a 2D/3D association between the given pixel and an edgelet from a potential state depends of the 2D/3D association of this same pixel with edgelets from the other potential states. The score of each potential state corresponds to the sum of all the 2D/3D correspondences weighted according to their status for every 2D contour pixels. Thus a pixel with a 2D/3D association having the same status for each potential state, is not relevant to estimate if a state is actually the accurate current state. Its weight is then null. A 2D/3D correspondence with the same orientation between a pixel and a projected edgelet will have a positive weight if it exists only for one potential state. It will have a negative weight for the other potential states where it does not exist. The table also shows that if a noisy contour pixel has no match for a potential state but has a 2D/3D association with a wrong orientation on the other potential states, this correspondence is relevant and a positive weight will be given for this no match.

7.2 Integration into a AR application

In an AR application, the state estimation is performed during a live sequence in addition to the other tasks as the camera localization and the AR effect displays. It is also im-

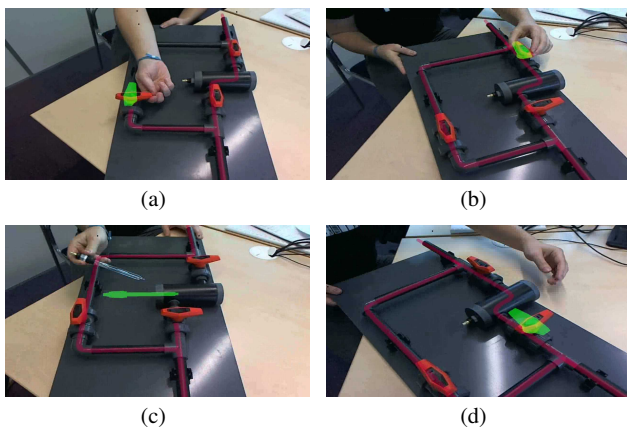


Fig. 23: Maintenance support on a bypass system with AR effect displays.

portant to have a real-time running application for the user comfort.

To display virtual information in the AR application, the camera pose has to be known. This camera localization is also in our case needed to estimate the state of the movable parts. Then the AR application is running on a multi-thread architecture with a thread for the localization process corresponding to our DEC-SLAM, an other one for the display and a last one for the state estimation.

The presented application aims to guide a user in order to replace a pH sensor. The sensor is integrated into a bypass system. This one allows to deviate the liquid contained into the pipes in order to change the pH sensor without leak. The maintenance support scenario is the following:

- Open the bypass in order to let the liquid circulate.
- Close the entrance of the main pipe to deviate the fluid.
- Close the pipe exit.
- Change the pH sensor located in the main pipe.
- Open the main pipe exit.
- Open the main pipe entrance.
- Close the bypass.

The application has to detect the bypass state (open/close), validate automatically the different steps of the scenario and show the next move to the user. The application user can see through a tablet screen different AR effects like the fluid circulation according to the valve state in red, or actions to do on the valves and the pH sensor in green (see Figure 23). Additional data are given in the fifth part of Online Resource 1.

8 Conclusion and Perspectives

In this article, a real-time solution for camera localization relative to any industrial object is proposed. A key-frame-

based SLAM algorithm is presented with a model constraint improving the tracking accuracy.

This constraint is represented through dynamic edgelets extracted by rendering on the graphic hardware. They are sampled to guarantee homogeneous spatial and angular distributions and to prevent the inclusion of ambiguous edgelets in matching steps of the DEC-SLAM algorithm. Two different formalisms in addition to the re-projection error representation, are proposed to integrate the model constraint in the optimization process. They correspond to hybrid model/trajectory constraint expressions that exploit the output of a model-based tracker, optimized thanks to the dynamic edgelets. They reduce the memory footprint of the DEC-SLAM while achieving similar accuracy. Our solution has been tested on several polyhedral and curved objects to attest its genericity. Experimental results demonstrate that our solution is as robust and accurate as rSEC-SLAM [33] exploiting static edgelets extracted offline from sharp model edges and using a re-projection error formalism. In addition, the hybrid constraint representations are able to reduce drastically the memory consumption and the computation time. Our results illustrate the high accuracy reached that allows convincing Augmented Reality applications.

As further work, we aim to focus on larger object to track as power plants or pipelines, or more rooms with available 3D models as with the ICL-NUIM dataset. Dealing with big objects may be indeed a challenging task. Whereas our DEC-SLAM solution proposes an accurate localization with local optimization when tracked objects are relatively small, it is not necessarily the case when objects are too large for the camera field of view. Potentially significant drifts over time may occur because of not well constrained DoF. Localization error may then accumulate itself, needing local and global corrections. However with the use of our hybrid constraint formalisms light in memory consumption and with the use of a model-based tracking of the large object of interest, a global optimization of the trajectory may be possible when localization error occurs.

Acknowledgements This work was partly funded by the french research program FUI through the project NASIMA. The authors would also like to thank their project partners Diotasoft and Faurecia for providing the car seat sequence.

References

1. Bleser, G., Stricker, D.: Advanced tracking through efficient image processing and visualinertial sensor fusion. *Computers and Graphics* **3**(1), 59–72 (2009)
2. Bleser, G., Wuest, H., Stricker, D.: Online camera pose estimation in partially known and dynamic scenes. In: *International Symposium on Mixed and Augmented Reality* (2006)
3. Drummond, T., Cipolla, R.: Real-time visual tracking of complex structures. *PAMI* **24**(7), 932–946 (2002)



Fig. 24: Accurate localization results on a Raving Rabbid, a dragon, a sport car, a real-sized car, a car cylinder head, and an orthosis with our phDEC-SLAM. Results are similar with rDEC-SLAM and rhDEC-SLAM. The accuracy can be appreciated by the projection of the 3D blue models on the object of interest in the images.

4. Finsterle, S., Kowalsky, M.: A truncated levenberg marquardt algorithm for the calibration of highly parameterized nonlinear models. *Computers & geosciences* **37**(6), 731–738 (2011)
5. Gay-Bellile, G., Bourgeois, S., Tamaazousti, M., Naudet-Collette, S., Knodel, S.: A mobile markerless augmented reality system for the automotive field. In: *International Symposium on Mixed and Augmented Reality Workshop* (2012)
6. Hajagos, B., Szcsi, L., Csbfalvi, B.: Fast silhouette and crease edge synthesis with geometry shaders. In: *Spring Conference on Computer Graphics* (2013)
7. Handa, A., Whelan, T., McDonald, J., Davison, A.: A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In: *International Conference on Robotics and Automation*. Hong Kong, China (2014)
8. Hertzmann, A.: Introduction to 3d non-photorealistic rendering: Silhouettes and outlines. In: *Special Interest Group on Computer GRAPHics and Interactive Techniques* (1999)
9. Huang, A.S., Bachrach, A., Henry, P., Krainin, M., Maturana, D., Fox, D., Roy, N.: Visual odometry and mapping for autonomous flight using an rgb-d camera. In: *International Symposium on Robotics Research* (2011)
10. Kerl, C., Sturm, J., Cremers, D.: Robust odometry estimation for rgb-d cameras. In: *International Conference on Robotics and Automation* (2013)
11. Klein, G., Murray, D.: Full-3d edge tracking with a particle filter. In: *British Machine Vision Conference* (2006)
12. Klein, G., Murray, D.: Parallel tracking and mapping for small AR workspaces. In: *International Symposium on Mixed and Augmented Reality* (2007)
13. Lepetit, V., Fua, P.: Monocular model-based 3d tracking of rigid objects: A survey. *Foundations and Trends in Computer Graphics and Vision* **1**(1), 1–89 (2006)
14. Levenberg, K.: A method for the solution of certain non-linear problems in least squares. *Quarterly of Applied Mathematics* **2**(2), 164–168 (1944)
15. Lhuillier, M.: Incremental fusion of structure-from-motion and gps using constrained bundle adjustments. *Pattern Analysis and Machine Intelligence* **34**(12), 2489–2495 (2012)
16. Li, G., Tsin, Y., Genc, Y.: Exploiting occluding contours for real-time 3d tracking: A unified approach. In: *International Conference on Computer Vision* (2007)
17. Loesch, A., Bourgeois, S., Gay-Bellile, V., Dhome, M.: Generic edgelet-based tracking of 3d objects in real-time. In: *Intelligent RObots and Systems* (2015)
18. Loesch, A., Bourgeois, S., Gay-Bellile, V., Dhome, M.: A hybrid structure / trajectory constraint for visual slam. In: *3D Vision* (2016)
19. Lothe, P., Bourgeois, S., Dekeyser, F., Royer, E., Dhome, M.: Towards geographical referencing of monocular slam reconstruction using 3d city models: Application to real-time accurate vision-based localization. In: *Computer Vision and Pattern Recognition* (2009)
20. Melbouci, K., Collette, S.N., Gay-Bellile, V., Ait-aider, O., Dhome, M.: Model based rgbd slam. In: *International Conference on Image Processing* (2016)
21. Middelberg, S., Sattler, T., Untzelmann, O., Kobbelt, L.: Scalable 6-dof localization on mobile devices. In: *European Conference on Computer Vision* (2014)
22. Mouragnon, E., Lhuillier, M., Dhome, M., Dekeyser, F., Sayd, P.: Real time localization and 3d reconstruction. In: *Computer Vision and Pattern Recognition* (2006)
23. Mur-Artal, R., Montiel, J.M.M., Tardos, J.D.: Orb-slam: a versatile and accurate monocular slam system. *Transactions on Robotics* **31**(5), 1147–1163 (2015)
24. Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Fitzgibbon, A.: Kinectfusion: Real-time dense surface mapping and tracking. In: *International Symposium on Mixed and Augmented Reality* (2011)
25. Nienhaus, M., Doellner, J.: Edge-enhancement - An algorithm for real-time non-photorealistic rendering. *Journal of WSCG* **11**(2) (2003)
26. Oikawa, M.A., Taketomi, T., Yamamoto, G., Fujisawa, M., Amano, T., Miyazaki, J., Kato, H.: Local quadrics surface approximation for real-time tracking of textureless 3d rigid curved objects. In: *Symposium on Virtual and Augmented Reality* (2012)
27. Petit, A., Marchand, E., Kanani, K.: Combining complementary edge, point and color cues in model-based tracking for highly dynamic scenes. In: *International Conference on Robotics and Automation* (2014)
28. Press, W.H.: *Numerical recipes 3rd edition: The art of scientific computing*. pp. 282–283. Cambridge university press (2007)
29. Ramadasan, D., Chevaldonne, M., Chateau, T.: Dcslam: A dynamically constrained real-time slam. In: *International Conference on Image Processing* (2015)
30. Raskar, R.: Hardware support for non-photorealistic rendering. In: *Special Interest Group on computer GRAPHics and Interactive Techniques Workshop on Graphics hardware* (2001)
31. Stanimirovic, D., Damasky, N., Webel, S., Koriath, D., Spillner, A., Kurz, D.: [poster] a mobile augmented reality system to assist auto mechanics. In: *International Symposium on Mixed and Augmented Reality* (2014)
32. Sturm, J., Engelhard, N., Endres, F., Burgard, W., Cremers, D.: A benchmark for the evaluation of rgb-d slam systems. In: *Intelligent Robots and Systems* (2012)
33. Tamaazousti, M., Gay-Bellile, V., Collette, S.N., Bourgeois, S., Dhome, M.: Real-time accurate localization in a partially known environment: Application to augmented reality on textureless 3d objects. In: *International Symposium on Mixed and Augmented Reality Workshop* (2011)
34. Tamaazousti, M., Gay-Bellile, V., Naudet-Collette, S., Bourgeois, S., Dhome, M.: Nonlinear refinement of structure from motion reconstruction by taking advantage of a partial knowledge of the environment. In: *Conference on Computer Vision and Pattern Recognition* (2011)
35. Triggs, B., McLauchlan, P., Hartley, R., Fitzgibbon, A.: Bundle adjustment - a modern synthesis. In: *International Conference on Computer Vision* (2000)
36. Vacchetti, L., Lepetit, V., Fua, P.: Combining edge and texture information for real-time accurate 3d camera tracking. In: *International Symposium on Mixed and Augmented Reality* (2004)
37. Vacchetti, L., Lepetit, V., Fua, P.: Stable real-time 3d tracking using online and offline information. *Pattern Analysis and Machine Intelligence* **26**(10), 1385–1391 (2004)
38. Wasenmuller, O., Meyer, M., Stricker, D.: CoRBS: Comprehensive rgb-d benchmark for slam using kinect v2. In: *Winter Conference on Applications of Computer Vision* (2016). URL <http://corbs.dfki.uni-kl.de/>
39. Wuest, H., Wientapper, F., Stricker, D.: Adaptable model-based tracking using analysis-by-synthesis techniques. In: *Computer Analysis of Images and Patterns* (2007)