

# Vision-only mobile robot navigation with topological maps

Toon Goedemé<sup>1,2</sup> and Luc Van Gool<sup>2</sup>

<sup>1</sup>EmSD, De Nayer Technical University, WenK, Sint-Katelijne-Waver, Belgium

<sup>2</sup>ESAT - PSI - VISICS, University of Leuven, K.U.Leuven, Belgium

## Abstract

When humans navigate, they use mainly their eyes as input from the environment. Indeed, in a new place, we look around and remember the appearance of each place. We automatically organise this visual information as a kind of visual map, which can be used to recognise places later on and to find the way to a certain goal. These tasks are known respectively as localisation and path planning in mobile robotics. Experiments show that the mental map of human navigators is organised topologically, rather than metrical.

In this work we present a total navigation system for autonomous mobile robot navigation based on these insights. With only an omnidirectional camera as sensor, this system is able to build automatically and robustly accurate topologically organised environment maps of a complex, natural environment. It can localise itself using that map at each moment, both at startup phase (kidnapped robot) and if knowledge of former localisations is present. The topological nature of the map, similar to the intuitive maps humans use, is memory-efficient and enables fast and simple path planning towards a specified goal. We developed a real-time visual servoing technique to steer the system along the computed path.

The power of this approach lies in the combination of a topological world model with state-of-the-art image matching techniques based on fast wide baseline local features. Techniques as SIFT and, more recently, SURF make it possible to find image correspondences without the need for artificial markers in the scene, and also enables an efficient description of the environments. Moreover, it enables the computation of a visual similarity measure, useful to organise visual databases in a topological world model.

We applied our navigation approach on autonomous wheelchair navigation of which this text presents experimental results.

## 1 Introduction

### 1.1 Application

This chapter describes a total navigation solution for mobile robots. It enables a mobile robot to efficiently localise itself and navigate in a large man-made environment, which can



Figure 1: Left: the robotic wheelchair platform. Right: the omnidirectional camera, composed by a colour camera and an hyperbolic mirror.

be indoor, outdoor or a combination of both. For instance, the inside of a house, an entire university campus or even a small city lie in the possibilities.

Traditionally, other sensors than cameras are used for robot navigation, like GPS and laser scanners. Because GPS (and Galileo also) needs a direct line of sight to the satellites [38], it cannot be used indoors or in narrow city centre streets, i.e. the very conditions we foresee in our application. Time-of-flight laser scanners are widely applicable, but are expensive and voluminous, even when the scanning field is restricted to a horizontal plane. The latter only yields a poor world representation, with the risk of not detecting essential obstacles such as table tops.

That is why we aim at a *vision-only* solution to navigation. Vision is, in comparison with these other sensors, much more informative. Moreover, cameras are quite compact and increasingly cheap. We observe also that many biological species, in particular migratory birds, use mainly their visual sensors for navigation. We chose to use an omnidirectional camera as visual sensor, because of its wide field of view and thus rich information content of the images acquired with. For the time being, we added a range sensing device for obstacle detection, but this is to be replaced by an omnidirectional vision range estimator under development.

Our method works with *natural* environments. That means that the environment does not have to be modified for navigation in any way. Indeed, adding artificial markers to every room in a house or to an entire city doesn't seem feasible nor desirable.

In contrast to classical navigation methods, we chose a *topological* representation of the environment, rather than a metrical one, because of its resemblance to the intuitive system humans use for navigation, its flexibility, wide usability, memory-efficiency and ease for map building and path planning.

The targeted application of this research is the visual guidance of electric wheelchairs for severely disabled people. More in particular, the target group are people not able to give detailed steering commands to navigate around in their homes and local city neighbourhoods. If it is possible for them to perform complicated navigational tasks by only giving simple commands, their autonomy can be greatly enhanced. For most of them such an increase of mobility and independence from other people is very welcome.

Our test platform and camera are shown in fig. 1. Nevertheless, the price of such a robotic wheelchair is a serious issue. With our method, the only additional hardware required is a laptop (or an equivalent embedded processor), a webcam, a mirror and (for the time being) some ultrasound sensors. Because of the increased independence of the users the cost of personal helpers is reduced, making the robotic wheelchair even more economically feasible.

## 1.2 Method overview

An overview of the navigation method presented is given in fig. 2. The system can be subdivided in three parts: map building, localisation and locomotion.

The map building stage has to be gone through only once, to train the system in a new environment. The mobile system is lead through all parts of the environment, while it takes images at a constant rate (in our set-up one per second). Later, this large set of omnidirectional images is automatically analysed and converted into a topological map of the environment, which is stored in the system's memory and will be used when the system is actually in use.

The next stage is localisation. When the system is powered up somewhere in the environment, it takes a new image with its camera. This image is rapidly compared with all the images in the environment map, and an hypothesis is formed about the present location of the mobile robot. This hypothesis is refined using Bayes' rule as soon as the robot starts to move and new images come in.

When the present location of the robot is known and a goal position is communicated by the user to the robot, a path can be planned towards that goal using the map. The planned route is specified as a sequence of training images, serving as a reference for what the robot should subsequently see if on course. This path is executed by means of a visual servoing algorithm: each time a visual homing procedure is executed towards the location where the next path image is taken.

The remainder of this chapter is organised as follows. The next section gives an overview of the related work. In section 3, our core image analysis and matching technique is explained: fast wide baseline matching. The sections thereafter describe the different stages of our approach. Section 4 discusses the map building process, section 5 explains the localisation method, section 6 describes the path planning, and section 7 details the visual servoing algorithm. We conclude with an overview of experimental results (section 8) and a conclusion (section 9).

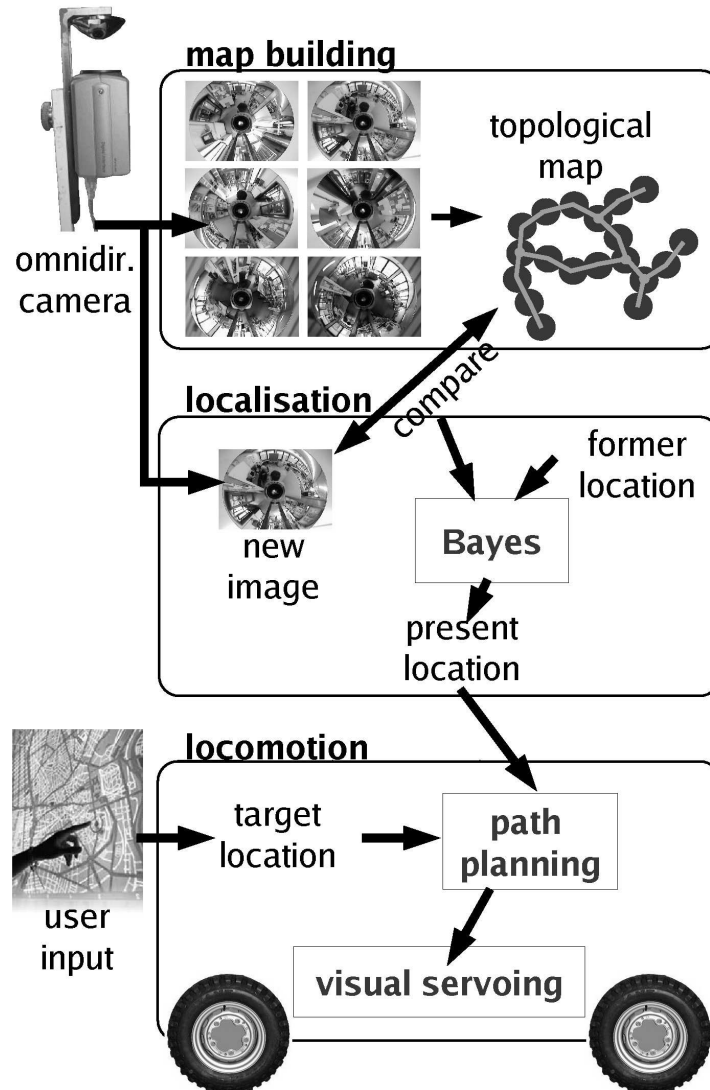


Figure 2: Overview of the navigation method

## 2 Related Work

### 2.1 Image comparison

A good image comparison method is of utmost importance in a vision-based navigation approach. Global methods compute a measure using all the pixels of the entire image. Although these methods are fast, they cannot cope with e.g. occlusions and severe viewpoint changes. On the other hand, techniques that work at a local scale, extracting and recognising *local features*, can be made robust to these effects. The traditional disadvantage of these local techniques is time complexity. In our approach, we combine novel global and local approaches resulting in fast and accurate image comparison.

#### 2.1.1 Global techniques

Many researchers use global image comparison techniques. Straightforward *global methods* like histogram-based matching, used by Ulrich and Nourbakhsh [54] don't seem distinctive enough for our application. Stricker [48] proposed a method based on the Fourier-Mellin transform to compare images. Unfortunately, the baseline can not be large which restricts that method to tracking. Another popular technique is the use of an eigenspace decomposition of the training images [21], which yields a compact database. However, these methods proved not useful in general situations because they are not robust enough against occlusions and illumination changes. That is why Jogan *et al.* [22] and Bischof *et al.* [5] developed a PCA-based image comparison that is robust against partial occlusions, respectively varying illumination.

#### 2.1.2 Local techniques

A solution to be able to cope with partial occlusions is comparing local regions in the images. The big question is how to detect these local features, also known as visual landmarks.

A simple solution to do this is by adding artificial markers to strategically chosen places in the world. To make these features easily detectable with a normal camera, they are given special (individual) photometric appearances (for instance coloured patterns [37], LEDs [1] or even 2D barcodes [41]). Using such artificial markers is perfectly possible for some applications, but often difficult. Navigation through an entire city or inside someone's house are examples of cases where pasting these markers all over the place is hardly feasible and in no case desirable.

That is why, in this project we use *natural landmarks*, extracted from the scene itself, without modifications. Moreover, the extraction of these landmarks must be automatic and robust against changes in viewpoint and illumination to ensure the detection of these landmarks under as many circumstances as possible.

Many researchers proposed algorithms for natural landmark detection. Mostly, local regions are defined around interest points in the images. The characterisation of these local regions with descriptor vectors enables the regions to be compared across images. Differ-

ences between approaches lie in the way in which interest points, local image regions, and descriptor vectors are extracted. An early example is the work of Schmid and Mohr [43], where geometric invariance was still under image rotations only. Scaling was handled by using circular regions of several sizes. Lowe *et al.* [28] extended these ideas to real scale-invariance. More general affine invariance has been achieved in the work of Baumberg [2], Tuytelaars & Van Gool [52, 53], Matas *et al.* [29], and Mikolajczyk & Schmid [31].

Although these methods are capable to find high quality correspondences, most of them are too slow to use in a real-time mobile robot algorithm. That is why we propose a much faster alternative, as explained in section 3.

## 2.2 Map structure

Many researchers proposed different ways to represent the environment perceived by vision sensors. We can order all possible map organisations by metrical detail: from dense 3D over sparse 3D to topological maps. We believe that the outer topological end of this spectrum offers the top opportunities.

### 2.2.1 Dense 3D maps

One approach is building dense 3D models out of the incoming visual data [34, 39]. Such approach has some disadvantages. It is computationally and memory demanding, and fails to model planar and less-textured parts of the environment such as walls. Nevertheless, these structures are omnipresent in our application, and collisions need to be avoided.

### 2.2.2 Sparse 3D maps

One way to reduce the computational burden is to make abstraction of the visual data. Instead of modelling a dense 3D model containing billions of voxels, a sparse 3D model is built containing only special features, called *visual landmarks*.

Examples of researchers solving the navigation problem with sparse 3D maps of natural landmarks are Se *et al.* [44] and Davison [9]. They position natural features in a metrical frame, which is as big as the entire mapped environment. Although less than the dense 3D variant, these methods are still computationally demanding for large environments since their complexity is quadratic in the number of features in the model. Also, for larger models the metric error accumulates, so that feature positions are drifting away.

### 2.2.3 Topological maps

As a matter of fact, the need for explicit 3D maps in navigation is questionable. One step further in the abstraction of environment information is the introduction of topological maps. The psychological experiments of Bühlhoff *et al.* [6] show that people rely more on a topological map than a metrical one for their navigation. In these topological maps, locally places are described as a configuration of natural landmarks. These places form the

nodes of the graph-like map, and are interconnected by traversable paths. Other researchers [24, 54, 55] also chose for topological maps, mainly because they scale better to real-world applications than metrical, deterministic representations, given the complexity of unstructured environments. Other advantages are the ease of path planning in such a map and the absence of drift.

### 2.3 Topological map building

Vale [55] developed a clustering-based method for automatic building of a topological environment map out of a set of images. Unfortunately, the latter method is only suited for image comparison techniques which are a metric function, and does not give correct results if *self-similarities* are present in the environment, i.e. places that are different but look similar.

Very popular are various *probabilistic* approaches of the topological map building problem. [40] for instance use Bayesian inference to find the topological structure that explains best a set of panoramic observations, while [46] fit hidden Markov models to the data. If the state transition model of this HMM is extended with robot action data, the latter can be modeled using a partially observable Markov decision process or POMDP, as in [23] and [51]. [56] solve the map building problem using graph cuts.

In contrast to these global topology fitting approaches, an alternative way is detecting *loop closings*. During a ride through the environment, sensor data is recorded. Because it is known that the driven path is traversable, an initial topological representation consists of one long edge between start and end node. Now, extra links are created where a certain place is revisited, i.e. an equivalent sensor reading occurs twice in the sequence. This is called a loop closing. A correct topological map results if all loop closing links are added.

Also in loop closing, probabilistic methods are introduced to cope with the uncertainty of link hypotheses and avoid links at self-similarities. [8], for instance, use Bayesian inference. [4] recently introduced Dempster-Shafer probability theory into loop closing, which has the advantage that ignorance can be modelled and no prior knowledge is needed. Their approach is promising, but limited to simple sensors and environments. In this chapter, we present a new framework for loop closing using rich visual sensors in natural complex environments, which is also based on Dempster-Shafer mathematics but uses it differently.

### 2.4 Visual Servoing

As explained in section 6, the execution of a path using such a topological environment map boils down to a series of visual servoing operations between places defined by images.

Cartwright and Collett [7] proposed the so-called bearing-only 'snapshot' model, inspired by the visual homing behaviour of insects such as bees and ants. Their proposed algorithm consists of the construction of a home vector, computed as the average of landmark displacement vectors. Franz *et al.* [14] analysed the computational foundations of this method and derived its error and convergence properties. They conclude that every visual homing method based solely on bearing angles of landmarks like this one, inevitably

depends on basic assumptions such as equal landmark distances, isotropic landmark distribution or the availability of an external compass reference. Unfortunately, because none of these assumptions generally hold in our targeted application we propose an alternative approach.

If both image dimensions are taken into account, not limiting the available information to the bearing angle, the most obvious choice is working via epipolar geometry estimation (e.g. [3, 52]). Unfortunately, in many cases this problem is ill conditioned. A workaround for planar scenes is presented by Sagüés [42], who opted for the estimation of homographies. Svoboda [49] proved that motion estimation with omnidirectional images is much better conditioned compared to perspective cameras. That is why we chose a method based on omnidirectional epipolar geometry. Other work in this field is the research of Mariottini *et al.* [30], who split the homing procedure in a rotation phase and a translation phase, which can not be used in our application because of the non-smooth robot motion.

### 3 Fast wide baseline matching

The novel technique we use for image comparison is *fast wide baseline matching*. This key technique enables extraction of natural landmarks and image comparison for our map building, localisation and visual servoing algorithms.

We use a combination of two different kinds of wide baseline features, namely a rotation reduced and colour enhanced form of Lowe's *SIFT* features [28], and the *invariant column segments* we developed [16]. These techniques extract local regions in each image, and describe these regions with a vector of measures which are invariant to image deformations and illumination changes. Across different images, similar regions can be found by comparing these descriptors. This makes it possible to find correspondences between images taken from very different positions, or under different lighting conditions. The crux of the matter is that the extraction of these regions can be done beforehand on each image separately, rather than during the matching. Database images can be processed off-line, so that the images themselves do not have to be available at the time of matching with another image.

#### 3.1 Camera motion constraint

The camera we use is a catadioptric system, consisting of an upward looking camera with a hyperboloidal mirror mounted above it. The result is a field of view of  $360^\circ$  in horizontal direction and more than  $180^\circ$  in vertical direction. The disadvantage is that these images contain severe distortions, as seen for instance in fig. 5.

We presume the robot to move on one horizontal plane. The optical axis of the camera is oriented vertically. In other words, allowed movements consist of translations in the plane and rotation around a vertical axis. Figure 3 shows an illustration on this.

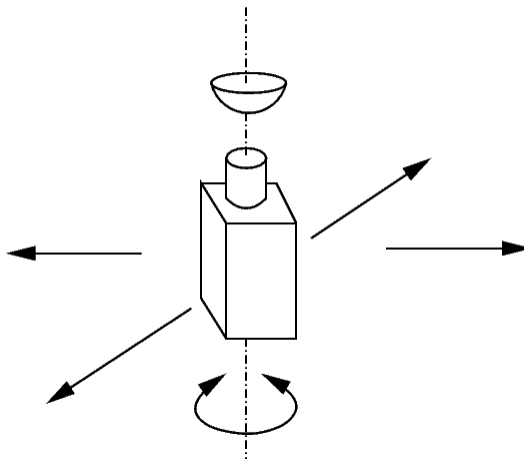


Figure 3: Illustration of the allowed movements of the camera.

### 3.2 Rotation reduced and colour enhanced SIFT

David Lowe presented the *Scale Invariant Feature Transform* [28], which finds interest points around local extrema in a scale-space of difference-of-Gaussian (DoG) images. The latter tend to correspond to blobs which contrast with their background. A dominant gradient orientation and scale factor define an image patch around each interest point so that a local image descriptor can be found as a histogram of normalised gradient orientations. SIFT features are invariant to rotation and scaling, and robust to other transformations.

A reduced form of these SIFT features for use on mobile robots is proposed by Ledwich and Williams [26]. They used the fact that rotational invariance is not needed for a camera with a motion constraint as in fig. 3. Elimination of the rotational normalisation and rotational part of the descriptor yields a somewhat less complex feature extraction and more robust feature matching performance.

Because the original SIFT algorithm works on greyscale images, some mismatches occur at similar objects in different colours. That is why we propose an outlier filtering stage using a colour descriptor of the feature patch based on global colour moments, introduced by Mindru *et al.* [32]. We chose three colour descriptors:  $C_{RB}$ ,  $C_{RG}$  and  $C_{GB}$ , with

$$C_{PQ} = \frac{\int PQ \, d\Omega \int d\Omega}{\int P \, d\Omega \int Q \, d\Omega}, \quad (1)$$

where  $P, Q \in \{R, G, B\}$ , i.e. the red, green, and blue colour bands, centralised around their means. After matching, the correspondences with Euclidean distance between the colour description vectors above a fixed threshold are discarded.

We tested these algorithms on the image pair in fig. 5. With the original SIFT algorithm, the first 13 matches are correct. Using our rotation reduced and colour enhanced algorithm, we see that up to 25 correct matches are found without including erroneous ones.

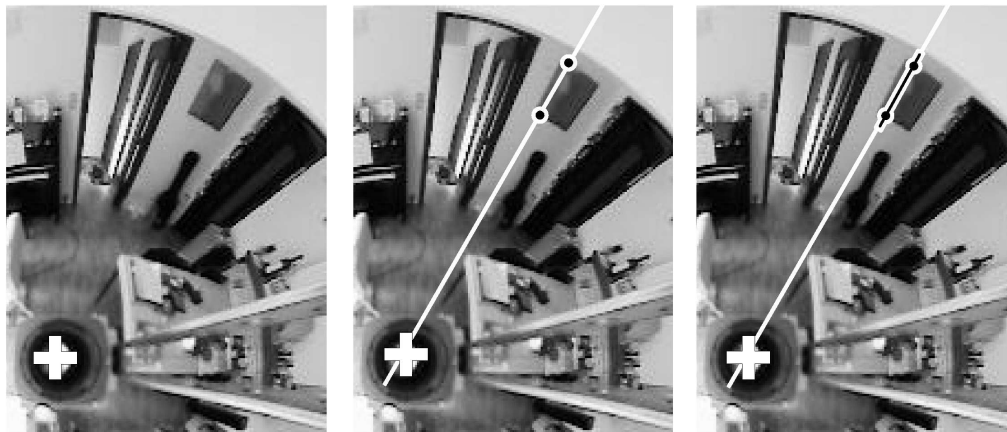


Figure 4: Illustration of the invariant column segment extraction algorithm: (left) part of the original image, the white cross identifies the projection centre, (centre) local maxima of the gradient for one radius, (right) one pair of maxima defines a column segment.

### 3.3 Invariant column segments

We developed wide baseline features which are specially suited for mobile robot navigation. There we exploited the special camera motion (section 3.1) and the fact that man-made environments contain many vertical structures. Examples are walls, doors, and furniture. These don't have to be planar, so cylindrical elements like pillars do comply too. Vertical lines in the world always project to radial lines in the omnidirectional image for the constrained camera motions.

Here, these new wide baseline features are described for the use on omnidirectional images. More details and how they can be used on perspective camera images are described in [16].

The extraction process of the wide baseline features starts as illustrated in figure 4. We stress that every step is invariant to changes in viewpoint and illumination. Along every line through the centre of the original image (left), we look for points having a local maximum gradient value (centre). Every consecutive pair of gradient maxima along the line defines the begin and end of a new invariant column segment (right).

We characterise the extracted column segments with a descriptor that holds information about colour and intensity properties of the segment. This 10-element vector includes:

- Three *colour invariants*. To include colour information in the descriptor vector, we compute the colour invariants, based on generalised colour moments (equation 1), over the column segment. To include information about the close neighbourhood of the segment, the line segment is expanded on both sides with a constant fraction of the segment length (in our experiments 0.2). Figure 4 (right) shows this.

- *Seven intensity invariants.* To characterise the intensity profile along the column segment, the best features to use are those obtained through the Karhunen-Löve transform (PCA). But because all the data is not known beforehand this is not practical. As is well known [20], the Fourier coefficients can sometimes offer a close approximation of the KL coefficients. In our method, because it is computationally less intensive and gives real output values, we choose to use the seven first coefficients of the *discrete cosine transform* (DCT), instead of Fourier. We compute seven DCT coefficients:

$$F(u) = \sqrt{\frac{2}{N}} \sum_{i=0}^{N-1} \Lambda(i) \cos \left[ \frac{\pi u}{2N} (2i + 1) \right] f(i) \quad (2)$$

Where  $u = 1 \dots 7$ ,  $f(i)$  is the normalised intensity profile along the line segment, and  $\Lambda(i) = \frac{1}{\sqrt{2}}$  if  $i = 0$ , else  $\Lambda(i) = 1$ . Because of the normalisation, the first coefficient ( $u = 0$ ) is zero, so we don't use it. The DCT computations in our algorithm are executed fast using the 8-point 1D method described in [27].

In many cases there are horizontally constant elements in the scene. This leads to many very resembling column segments next to each other. To avoid matching over and over again very similar line segments, we first do a clustering of the line segments in each image. This clustering can be implemented very efficiently, because the line segments are already sorted in order of extraction. As a clustering measure we use the Mahalanobis distance of the descriptor vectors, extended with the horizontal distance between the line segments. In each cluster a prototype segment is chosen for use in the matching procedure.

### 3.4 Matching

These two kinds of local wide baseline features are very suited to quickly find correspondences between two widely separated images. A correspondence pair is established if for both features the other is the closest to it in the feature space, for the entire data set. Also, this match must be at least a fixed ratio better than the second best match. To be able to cope with different ranges of the elements of the descriptor vectors, distances are computed using the Mahalanobis measure (where we assume the cross-correlations to be zero):

$$d_{ij} = \sqrt{\sum_k \frac{(x_{ik} - x_{jk})^2}{\sigma_k^2}} \quad (3)$$

To speed up the matching, a Kd-tree of the reference image data is built. We used the on-line available package ANN (Approximate Nearest Neighbour) by Mount *et al.* [33], for this.

Fig. 5 shows the matching results on a pair of omnidirectional images. As seen in these examples, the SIFT features and the column segments are complementary, which pleads for the combined use of the two. The computing time required to extract features in two  $320 \times 240$  images and find correspondences between them is about 800 ms for the enhanced

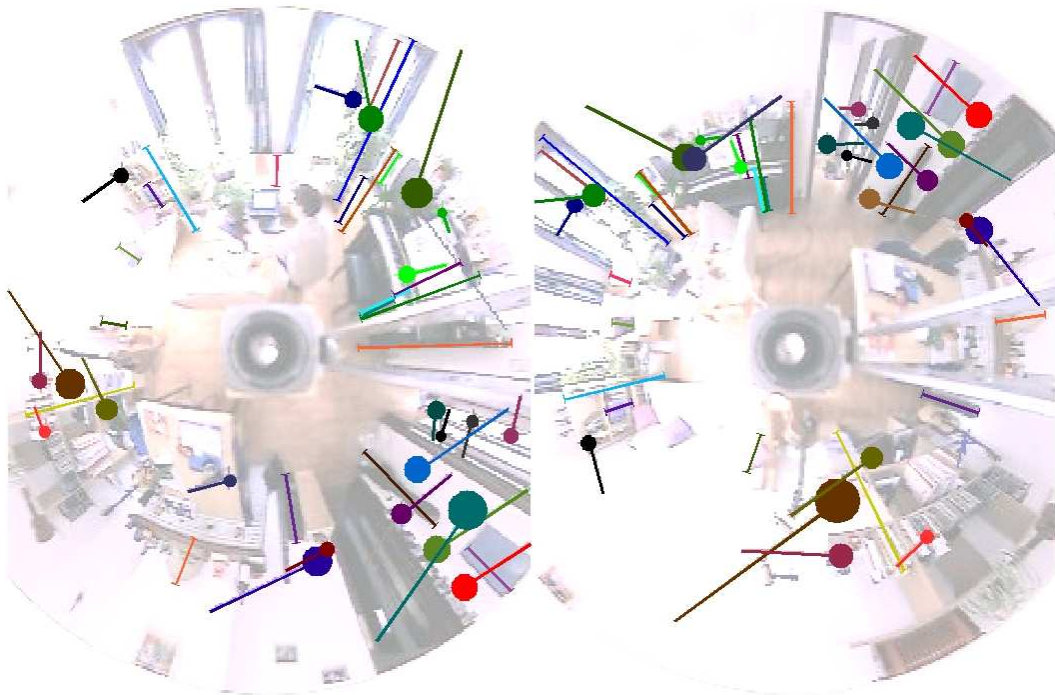


Figure 5: A pair of omnidirectional images with colour-coded corresponding column segments (radial lines) and SIFT features (circles with tail).

SIFT features and only 300 ms for the vertical column segments (on a 800 MHz laptop). Typically 30 to 50 correspondences are found.

For the description of a feature only the descriptors are used in the end, and not the underlying pixel data. As a result, the memory requirements for storing the reference images of entire environments can be kept limited.

## 4 Map Building

The navigation approach proposed is able to automatically construct a topological world representation out of a sequence of training images. During a training tour through the entire environment, omnidirectional images are taken at regular time intervals. The order of the training images is known. Section 4.1 describes the map structure targeted. In section 4.2, the image comparison technique based on fast wide baseline features is described which is used by the actual map building algorithm, presented in section 4.4. The mathematical formalism used in this algorithm, Dempster-Shafer probability theory, is briefly introduced in section 4.3.

## 4.1 Topological Maps

To be of use in the following parts of the navigation method, the topological map must describe all ‘places’ in the environment and the possible connections between these places. The topology of the world, being the maze of streets in a city or the structure of a house, must be reflected in the world model. The question remains what exactly is meant with such a *place* and how to delimit it. In a building, a place can be defined as a separate room. But that is not such a good definition either. What to do with long corridors, and outdoors with city streets?

That is why we define a place with regard to the needs of the localisation and locomotion algorithms. To be able to get a sufficiently detailed localisation output, the sampling of places must be dense enough. For the locomotion algorithm, which performs visual homing between two places at a time, the distance between these places must be not too big to ensure errorless motion. On the other hand, a compact topological map with fewer places requires less memory and enables faster localisation and path planning.

We discuss the image comparison method used in the map building algorithm before deciding on this place definition, as this comparison will lie at its basis.

## 4.2 Image Comparison Measure

The main goal of this section is to determine for each arbitrary pair of images a certain similarity measure, which tells how visually similar the two images are. Our image comparison approach consists of two levels, a global and a local comparison of the images. We first compare two images with a coarse but fast global technique. After that, a relatively slower comparison with more precision based on local features only has to be carried out on the survivors of the first stage.

### 4.2.1 Global colour similarity measure

To achieve a fast global image similarity measure between two images, we compute the same moments we used for the local features (equation 1) over the entire image. These moments are invariant to illumination changes, i.e. offset and scaling in each colour band. The Euclidean distance between two sets of these image colour descriptors gives a visual dissimilarity measure between two images.

With this dissimilarity measure, we can clearly see for instance the difference of images taken in different rooms. Because images taken in the same room but at different positions have approximately the same colour scheme, a second dissimilarity measure based on local features is needed to distinguish them.

### 4.2.2 Local measure based on matches

First, we search for feature matches between the two images, using the techniques described in section 3. The dissimilarity measure is taken to be inversely proportional to the *number of matches*, relative to the average number of features found in the images. Also the difference

in relative configuration of the matches is taken into account. Therefore, we first compute a global angular alignment of the images by computing the average angle difference of the matches. The dissimilarity measure  $D_m$  is now also made proportional to the average angle difference of the features after this global alignment:

$$D_m = \frac{1}{N} \cdot \frac{n_1 + n_2}{2} \cdot \frac{\sum |\theta_i|}{N} \quad (4)$$

$$= \frac{(n_1 + n_2) \sum |\theta_i|}{2N^2}, \quad (5)$$

where  $N$  corresponds to the number of matches found,  $n_i$  the number of extracted features in image  $i$ ,  $\theta_i$  the angle difference for one match after global alignment.

### 4.2.3 Combined dissimilarity measure

We combine these two measures: only those pairs of images who have a colour dissimilarity under a predefined threshold are candidates for computing a matching dissimilarity.

This combined *visual* distance between two images is related to the physical distance between the corresponding viewpoints, but is certainly not a linear measure for it. As a matter of fact, the disparity and appearance difference of features is also related to the distance of the corresponding natural landmark to the cameras. Therefore, in large spaces (halls, market squares), a certain visual distance will be corresponding to a much larger physical distance, compared to the same visual distance between two images in a small space.

With this visual distance, the place definition problem in section 4.1 can be addressed on the basis of a constant visual distance between places instead of a constant physical distance.

## 4.3 Introduction to Dempster-Shafer theory

The proposed topological map building algorithm relies on Dempster-Shafer theory [10,45] to collect evidence for each loop closing hypothesis. Therefore, a brief overview of the central concepts of Dempster-Shafer theory is presented in this section.

Dempster-Shafer theory offers an alternative to traditional probabilistic theory for the mathematical representation of uncertainty. The significant innovation of this framework is that it makes a distinction between multiple types of uncertainty. Unlike traditional probability theory, Dempster-Shafer defines two types of uncertainty:

- *Aleatory Uncertainty* – the type of uncertainty which results from the fact that a system can behave in random ways (a.k.a. stochastic or objective uncertainty)
- *Epistemic Uncertainty* – the type of uncertainty which results from the lack of knowledge about a system (a.k.a. subjective uncertainty or ignorance)

This makes it a powerful technique to combine several sources of evidence to try to prove a certain hypothesis, where each of these sources can have a different amount of knowledge (*ignorance*) about the hypothesis. That is why Dempster-Shafer is typically used for sensor fusion.

For a certain problem, the set of mutually exclusive possibilities, called the *frame of discernment*, is denoted by  $\Theta$ . For instance, for a single hypothesis  $H$  about an event this becomes  $\Theta = \{H, \neg H\}$ . For this set, traditional probability theory will define two probabilities  $P(H)$  and  $P(\neg H)$ , with  $P(H) + P(\neg H) = 1$ . Dempster-Shafer's analogous quantities are called *basic probability assignments* or *masses*, which are defined on the *power set* of  $\Theta$ :  $2^\Theta = \{A | A \subseteq \Theta\}$ . The *mass*  $m : 2^\Theta \rightarrow [0, 1]$  is a function meeting the following conditions:

$$m(\emptyset) = 0 \quad \sum_{A \in 2^\Theta} m(A) = 1. \quad (6)$$

For a single hypothesis  $H$ , the power set becomes  $2^\Theta = \{\emptyset, \{H\}, \{\neg H\}, \{H, \neg H\}\}$ . A certain sensor or other information source can assign masses to each of the elements of  $2^\Theta$ . Because some sensors do not have knowledge about the event (e.g. it is out of the sensor's field-of-view), they can assign a certain fraction of their total mass to  $m(\{H, \neg H\})$ . This mass, called the *ignorance*, can be interpreted as the probability mass assigned to the outcome ' $H$  OR  $\neg H$ ', i.e. when the sensor does not know about the event, or is—to a certain degree—uncertain about the outcome. This means also that using this technique no prior probability function is needed, *no knowledge* can be expressed as total ignorance.

Sets of masses about the same power set, coming from different information sources can be combined together using *Dempster's rule of combination*:

$$m_1 \oplus m_2(C) = \frac{\sum_{A \cap B = C} m_1(A)m_2(B)}{1 - \sum_{A \cap B = \emptyset} m_1(A)m_2(B)} \quad (7)$$

This combination rule is useful to combine evidence coming from different sources into one set of masses. Because these masses can not be interpreted as classical probabilities, no conclusions about the hypothesis can be drawn from them directly. That is why two additional notions are defined, *support* and *plausibility*. They are computed as:

$$Spt(A) = \sum_{B \subseteq A} m(B) \quad Pls(A) = \sum_{A \cap B \neq \emptyset} m(B) \quad (8)$$

These values define a confidence interval for the real probability of an outcome:  $P(A) \in [Spt(A), Pls(A)]$ . Indeed, due to the vagueness implied in having non-zero ignorance, the exact probability cannot be computed. But, decisions can be made based on the lower and upper bounds of this confidence interval.

#### 4.4 Map Building Algorithm

We apply this mathematical theory on the topological map building problem posed. Out of a series of omnidirectional images, acquired at constant rate during a tour through the environment. Firstly, these images are clustered into *places*. Then, loop closing hypotheses are

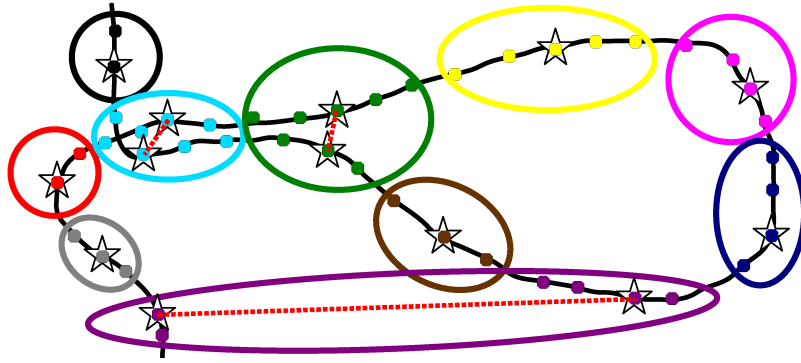


Figure 6: Example for the image clustering and hypothesis formulation algorithms. Dots are image positions, black is exploration path, clusters are visualised with ellipses, prototypes of (sub)clusters with a star. Hypotheses are denoted by a dotted red line.

formulated between visually similar places of which evidence is collected using Dempster-Shafer theory. Once decisions are made about these hypotheses, the correct topology of the world is known.

This technique makes it possible to cope with *self-similar* environments. Places that look alike but are different will more likely get their link hypothesis rejected.

#### 4.4.1 Image clustering

The dots in the sketch figure of 6 denote places where images are taken. Because they were taken at constant time intervals and the robot did not drive at a constant speed, they are not evenly spread. We perform agglomerative clustering with complete linkage based on the combined visual distance (see section 4.2.3) on all the images, yielding the ellipse shaped clusters in fig. 6. The black line shows the exploration path as driven by the robot.

#### 4.4.2 Hypothesis formulation

As can be seen in the lower part of fig. 6, not all image groups nicely cover one distinct place. This is due to *self-similarities*, or distinct places in the environment that are different but look alike and thus yield a small visual distance between them.

For each of the clusters, we can define one or more subclusters. Images within one cluster which are linked by exploration path connections are grouped together. For each of these *subclusters* a prototype image is chosen as the *medoid*<sup>1</sup> based on the visual distance, denoted as a star in the figure.

For each pair of these subclusters within the same cluster, we define a *loop closing hypothesis*  $H$ , which states that if  $H = true$ , the two subclusters describe the same physical

<sup>1</sup>The medoid of a cluster is computed analogous to the centroid, but using the median instead of the average.

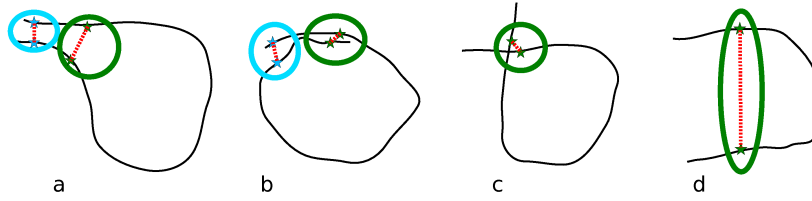


Figure 7: Four topological possibilities for one hypothesis.

place and must be merged together. We will use Dempster-Shafer theory to collect evidence about each of these hypotheses.

#### 4.4.3 Dempster-Shafer evidence collection

For each of the hypotheses defined in the previous step, a decision must be made if it was correct or wrong. Figure 7 illustrates four possibilities for one hypothesis. We observe that a hypothesis has more chance to be true if there are more hypotheses in the neighbourhood, like in case *a* and *b*. If no neighbouring hypotheses are present (*c,d*), no more evidence can be found and no decision can be made based on this data.

We conclude that for a certain hypothesis, a neighbouring hypothesis adds evidence to it. It is clear that, the further away this neighbour is from the hypothesis, the less certain the given evidence is. We chose to model this subjective uncertainty by means of the ignorance notion in Dempster-Shafer theory. That is why we define an *ignorance function*  $\xi$  containing the distance between two hypotheses  $H_a$  and  $H_b$ :

$$\xi(H_a, H_b) = \begin{cases} 1 - \sin\left(\frac{d_H(H_a, H_b)\pi}{2d_{th}}\right) & (d_H \leq d_{th}) \\ 0 & (d_H > d_{th}) \end{cases} \quad (9)$$

where  $d_{th}$  is a distance threshold and  $d_H(H_a, H_b)$  is the sum of the distances between the two pairs of prototypes of both hypotheses, measured in number of exploration images.

To gather *aleatory* evidence, we look at the visual similarity of both subcluster prototypes, normalised by the standard deviation of the intra-subcluster visual similarities:

$$s_V(H_a) = \frac{s_V(\text{prot}_{a1}, \text{prot}_{a2})}{\sigma_{\text{subclus}}(s_V)}, \quad (10)$$

Where the visual similarity  $s_V$  is the inverse of the visual distance, defined in equation 4.

Each neighbouring hypothesis  $H_b$  yields the following set of Dempster-Shafer masses, to be combined with the masses of the hypothesis  $H_a$  itself:

$$\begin{aligned} m(\{\emptyset\}) &= 0 \\ m(\{H_a\}) &= s_V(H_b)\xi(H_a, H_b) \\ m(\{\neg H_a\}) &= (1 - s_V(H_b))\xi(H_a, H_b) \\ m(\{H_a, \neg H_a\}) &= 1 - \xi(H_a, H_b) \end{aligned} \quad (11)$$

Hypothesis masses are initialised with the visual similarity of its subcluster prototypes and an initial ignorance value (0.25 in our experiments), which models its influenceability by neighbours.

#### 4.5 Hypothesis decision

After combination of each hypothesis's mass set with the evidence given by neighbouring hypotheses (up to a maximum distance  $d_{th}$ ), a decision must be made if this hypothesis was correct and thus if the subclusters must be united into one place or not.

Unfortunately, as stated above, only positive evidence can be collected, because we can not gather more information about totally isolated hypotheses (like  $c$  and  $d$  in fig. 7). This is not too bad, because of different reasons. Firstly, the chance for correct, but isolated hypothesis (case  $c$ ) is low in typical cases. Also, adding erroneous loop closings ( $c$  and  $d$ ) will yield an incorrect topological map, whereas leaving them out will keep the map useful for navigation, but a bit less complete. Of course, new data about these places can be acquired later, during navigation.

It is important to remind oneself that the computed Dempster-Shafer masses can not directly be interpreted as probabilities. That is why we use equation 8 to compute the support and plausibility of each hypothesis after evidence collection. Because these values define a confidence interval for the real probability, a hypothesis can be accepted if the lower bound (the support) is greater than a threshold.

After this decision, a final topological map can be built. Subclusters connected with accepted hypotheses are merged into one place, and a new medoid is computed as prototype of it. For hypotheses that are not accepted, two distinct places should be constructed.

## 5 Localisation

When the system has learnt a topological map of an environment, this map can be used for a variety of navigational tasks, firstly *localisation*. For each arbitrary new position in the known environment, the system can find out *where* it is. The output of this localisation algorithm is a *location*, which is—opposed to other methods like GPS—not expressed as a metric coordinate, but as one of the topological places defined earlier in the formerly explained map building stage.

The training set doesn't need to cover every imaginable position in the environment, in contrast to [25]. A relatively sparse coverage is sufficient to localise every possible position. That is because the image comparison method we developed is based on wide baseline techniques and hence can recognise scene items from substantially different viewpoints.

Actually, two localisation modes exist. When starting up the system, there is no *a priori* information on the location. Every location is equally probable. This is called *global localisation*, alias the *kidnapped robot problem*. Traditionally, this is known to be a hard problem in robot localisation. In contrast, if there is knowledge about a former localisation not too

long ago, the locations in the proximity of that former location have a higher probability than others further away. This is called *location updating*.

We propose a system that is able to cope with both localisation modes. A probabilistic approach is taken. Instead of making a hard decision about the location, a probability value is given to each location at each time instant. The Bayesian approach we follow is explained in the next subsection.

## 5.1 Bayesian Filtering

Define  $x \in X$  a place of the topological map.  $Z$  is the collection of all omnidirectional images  $z$ , so that  $z(x)$  corresponds to the training observation at place  $x$ . At a certain time instant  $t$ , the system acquires a new image  $z_t$ . The goal of the localisation algorithm is to reveal the place  $x_t$  where this image was taken.

We define the *Belief function*  $Bel(x, t)$  as the probability of being at place  $x$  at time  $t$ , given all previous observations. So,

$$Bel(x, t) = P(x_t | z_t, z_{t-1}, \dots, z_0), \quad (12)$$

for all  $x \in X$ . In the *kidnapped robot* case, there is no knowledge about previous observations hence  $Bel(x, t_0)$  is initialised equal for all  $x$ .

Using Bayes' rule, we find:

$$Bel(x, t) = \frac{P(z_t | x_t, z_{t-1}, \dots, z_0) P(x_t | z_{t-1}, \dots, z_0)}{P(z_t | z_{t-1}, \dots, z_0)}. \quad (13)$$

Because the denominator of this fraction is not dependent on  $x$ , we replace it by the normalising constant  $\eta$ . If we know the current location of the system, we assume that future locations do not depend on past locations. This property is called the *Markov Assumption*:  $P(z_t | x_t, z_{t-1}, \dots, z_0) = P(z_t | x_t)$ . Using it, together with the probabilistic sum rule, equation 13 yields:

$$Bel(x, t) = \eta P(z_t | x_t) \sum_{x_{t-1} \in X} [P(x_t | x_{t-1}) Bel(x, t-1)] \quad (14)$$

This allows us to calculate the belief recursively based on two variables: the *next state density* or *motion model*  $P(x_t | x_{t-1})$  and the *sensor model*  $P(z_t | x_t)$ .

## 5.2 Motion Model

The motion model  $P(x_t | x_{t-1})$  explicits the probability of a transition from one place  $x_{t-1}$  to another  $x_t$ . It seems logical to assume that a transition in the time instant  $(t-1) \rightarrow t$  between places that are far from each other is less probable than between places close to each other. We model this effect with a Gaussian:

$$P(x_t | x_{t-1}) = \frac{1}{\beta_x} e^{-dist(x_{t-1}, x_t) / \sigma_x^2}. \quad (15)$$

In this equation, the function  $dist(x_1, x_2)$  corresponds to a measurement of the distance between the two places. We approximate it as the minimum number of place transitions needed to go from  $x_1$  to  $x_2$  on the topological map, computed with the Dijkstra algorithm [12]. In equation 15,  $\beta_x$  is a normalisation constant, and  $\sigma_x^2$  is the variance of the distances, measured on the map data. Once the topological map is known, the complete motion model can be computed off-line for usage during localisation.

### 5.3 Sensor Model

The entity  $P(z_t|x_t)$ , called the sensor model, is the probability of acquiring a certain observation  $z_t$  if the location  $x_t$  is known. This is related to the visual dissimilarity of that observation  $z_t$  and the training observation  $z(x_t)$  at location  $x_t$ . The probability of acquiring an image at a certain place that differs much from the training image taken at that place has a low probability. We model this sensor model also by a Gaussian:

$$P(z_t|x_t) = \frac{1}{\beta_z} e^{-diss(z_t, z(x_t))/\sigma_z}. \quad (16)$$

This time, the function  $diss(z_1, z_2)$  refers to the visual dissimilarity explained in section 4.2.3.  $\sigma_z$  is the standard deviation of it, measured on the data (or for large maps a local subset of the data) and  $\beta_z$  is a normalisation constant.

Unlike the motion model, the sensor model cannot be computed beforehand. It depends on the newly incoming query image data. Every location update step the visual dissimilarities of the query image with many database images must be computed. This validates our efforts to make the computation of the visual dissimilarity measure as fast as possible.

## 6 Path planning

With the method of the previous section, at each time instant the most probable location of the robot can be found, from which a path to a goal can be determined. How the user of the system, for instance the wheelchair patient, gives the instruction to go towards a certain goal is highly dependent on the situation. For every disabled person, for instance, an individual interface must be designed adapted to his/her possibilities.

We assume a certain goal is expressed as a certain place of the topological map, e.g. as a voice command  $\ll Kitchen! \gg$ . From the present pose, computed by the localisation algorithm, a path can be easily found towards it using Dijkstra's algorithm [12]. This path is expressed as a series of topological places which are traversed, see fig. 8 (d).

## 7 Visual servoing

The algorithm described in this section makes the robot move along a path, computed by the previous section. Such a path is given as a sparse set of prototype images of places. The physical distance between two consecutive path images is variable (1 to 5 metres in our

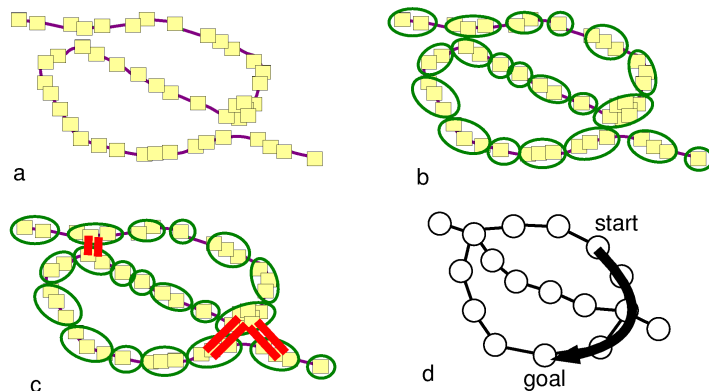


Figure 8: Topological map building. (a) Squares denote places where training images were taken. (b) Place grouping of images. (c) Cross-place matching. (d) In the resulting topological map, a path is found between a certain start and goal place.

tests), but the visual distance is constant, such that there are enough local feature matches as needed by this algorithm.

It is easy to see that following such a sparse visual path boils down to a succession of *visual homing* operations. First, the robot is driven towards the place where the first image on the path is taken. When arrived, it is driven towards the next path image, and so on. Because a smooth path is desired for the application, the motion must be continuous without stops at path image positions.

We tackle this problem by estimating locally the spatial structure of the wide baseline features using epipolar geometry. Hence, at this point we bring in some 3D information. This may seem at odds with our topological approach, but the depth maps are very sparse and only calculated locally so that errors are kept local, don't suffer from error build-up, and are efficient to compute.

Fig. 9 offers an overview of the proposed method. Each of the *visual homing* operations is performed in two phases, an initialisation phase (section 7.1) and an iterated motion (section 7.2) phase.

## 7.1 Initialisation phase

From each position within the reach of the next path image (the *target* image), a visual homing procedure can be started. Our approach first establishes wide baseline local feature correspondences between the present and the target image, as described in section 3. That information is used to compute the epipolar geometry, which enables us to construct a local map containing the feature world positions, and to compute the initial homing vector.

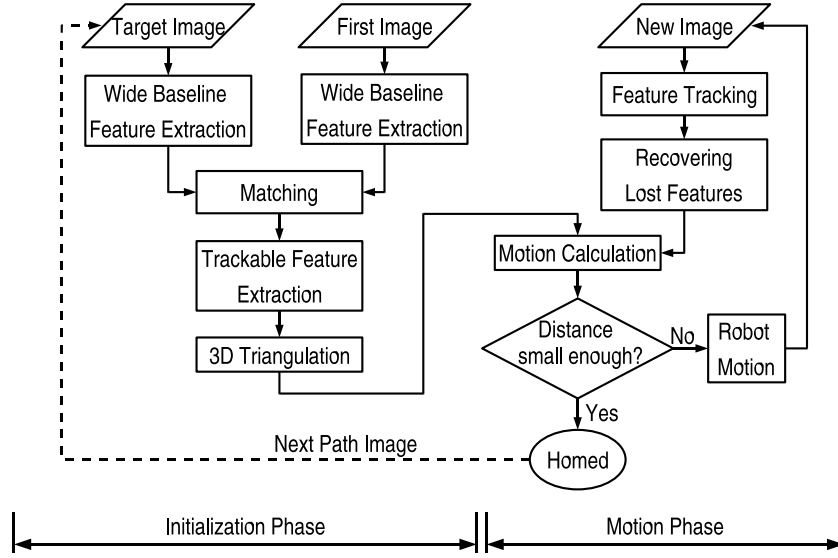


Figure 9: Flowchart of the proposed algorithm for visual servoing along a path.

### 7.1.1 Epipolar geometry estimation

Our calibrated single-viewpoint omnidirectional camera is composed of a hyperbolic mirror and a perspective camera. As imaging model, we use the model proposed by Svoboda and Pajdla [49] (which is less general, but less complicated than the one by Geyer and Daniilidis [15]). This enables the computation of the epipolar geometry based on 8 point correspondences. In [50], Svoboda describes a way to robustly estimate the *essential matrix*  $E$ , when there are outliers in the correspondence set. The essential matrix is the equivalent of the fundamental matrix in the case of known internal camera calibration,  $F = K^{-T}EK^{-1}$ .

Svoboda’s so-called *generate-and-select* algorithm to estimate  $E$  is based on repeatedly solving an overdetermined system built from the correspondences that have a low ‘*outlierness*’ and evaluating the *quality measure* of the resulting essential matrix. Because our tests with this method did not yield satisfactory results, we implemented an alternative method based on the well-known Random Sample Consensus (RANSAC [13]) paradigm.

The set-up is sketched in fig. 10. One visual feature with world coordinates  $\mathbf{X}$  is projected via point  $\mathbf{u}$  on the first mirror to point  $\mathbf{p}$  in the image plane of the first camera. In the second camera, the mirror point is called  $\mathbf{v}$  and the image plane point  $\mathbf{q}$ . For each of the correspondences, the mirror points  $\mathbf{u}$  and  $\mathbf{v}$  can be computed as

$$\mathbf{u} = \mathcal{F}(K^{-1}\mathbf{p})K^{-1}\mathbf{p} + \mathbf{t}_C, \quad (17)$$

with  $\mathbf{t}_C = [0, 0, -2e]^T$  and

$$\mathcal{F}(\mathbf{x}) = \frac{b^2(e x_1 + a \|\mathbf{x}\|)}{b^2 x_1^2 - a^2 x_2^2 - a^2 x_3^2}. \quad (18)$$

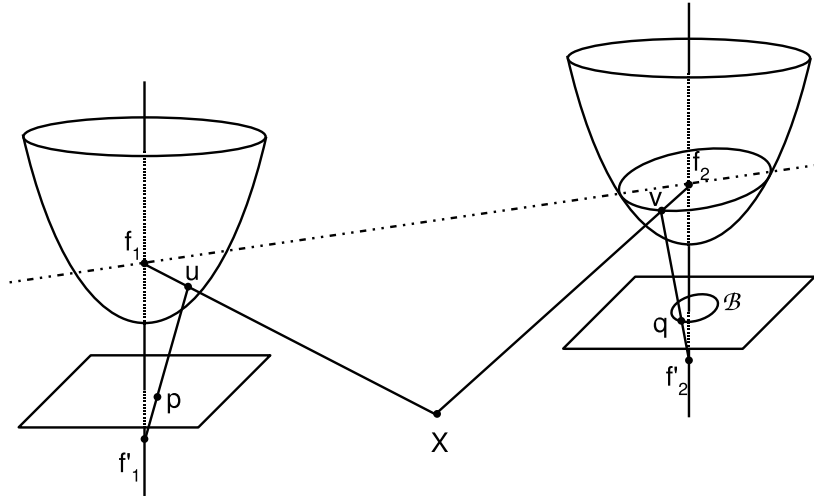


Figure 10: Projection model for a pair of omnidirectional images.

In these equations  $K$  is the internal calibration matrix of the camera, and  $a$ ,  $b$  and  $e$  are the parameters of the hyperbolic mirror, with  $e = \sqrt{a^2 + b^2}$ .

If  $E$  is the *essential matrix*, for all correspondences  $\mathbf{v}^T E \mathbf{u} = 0$ . This yields for each correspondence pair one linear equation in the coefficients of  $E = [e_{ij}]$ . The essential matrix can be computed as the solution of the homogeneous system

$$A\mathbf{e} = \mathbf{0}, \quad (19)$$

with  $\mathbf{e} = [e_{11}, e_{12}, e_{13}, e_{21}, \dots, e_{33}]$  and the rows of  $A$  equal to

$$\mathbf{a}_i = [v_{i1}u_{i1}, v_{i1}u_{i2}, v_{i1}u_{i3}, \dots, v_{i3}u_{i3}]. \quad (20)$$

For each random sample of 8 correspondences, an  $E$  matrix can be calculated. This is repeatedly done and for each  $E$  matrix candidate the inliers are counted. A correspondence is regarded an inlier if the second image point  $\mathbf{q}$  lies within a predefined distance from the epipolar *ellipse*, defined by the first image point  $\mathbf{q}$ . This epipolar ellipse  $\mathcal{B}$  with equation  $\mathbf{x}^T B \mathbf{x} = 0$  is computed with  $B =$

$$\begin{bmatrix} -4t^2a^2e^2 + r^2b^4 & rsb^4 & rtb^2(-2e^2 + b^2) \\ rsb^4 & -4t^2a^2e^2 + s^2b^4 & stb^2(-2e^2 + b^2) \\ rtb^2(-2e^2 + b^2) & stb^2(-2e^2 + b^2) & t^2b^4 \end{bmatrix} \quad (21)$$

and  $[r, s, t] = E\mathbf{u} = E(\mathcal{F}(K^{-1}\mathbf{p})K^{-1}\mathbf{p} + \mathbf{t}_C)$ . Fortunately, this ellipse becomes a circle when the motion is in one plane, so that the distance from a point to this shape is easy to compute.

$$\text{dist}(\mathbf{x}, \mathcal{B}) = \left| \frac{\mathbf{x}^T B \mathbf{x}}{B_{3,3}} \right| \quad (22)$$

From the one essential matrix  $E$  with the maximal number of inliers the motion between the cameras can be computed using the SVD based method proposed by Hartley [19]. If more than one  $E$ -matrix is found with the same maximum number of inliers, the one is chosen with the best (i.e. smallest) quality measure  $q_E = \sigma_1 - \sigma_2$ , where  $\sigma_i$  is the  $i$ th singular value of the matrix  $E$ .

Out of this relative camera motion, a first estimate of the homing vector is derived. During the motion phase this homing vector is refined.

From the singular value decomposition

$$E = UDV^T, \quad (23)$$

the matrices  $R$  and  $S$  in  $E = RS$  are computed by

$$S = \pm VZV^T \quad (24)$$

$$R = UYV^T \quad \text{or} \quad UY^TV^T, \quad (25)$$

where

$$Z = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \text{and} \quad Y = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (26)$$

Based on one feature correspondence, the right choice between the 4 possible combinations can be made yielding a antisymmetric translation matrix  $S$  and a rotation matrix  $R$ .

### 7.1.2 Local feature map estimation

In order to start up the succession of tracking iterations, an estimate of the local map must be made. In our approach the local feature map contains the 3D world positions of the visual features, centred at the starting position of the visual homing operation. These 3D positions are easily computed by triangulation.

We only use two images, the first and the target image, for this triangulation. This has two reasons. Firstly, these two have the widest baseline and therefore triangulation is best conditioned. Our wide baseline matches between these two images are also more plentiful and less influenced by noise than the tracked features.

## 7.2 Motion phase

Then, the robot is put into motion in the direction of the homing vector and an image sequence is recorded. We rely on lower-level collision detection, obstacle avoidance and trajectory planning algorithms to drive safely [11, 36]. In each new incoming image the visual features are tracked. Robustness to tracking errors (caused by e.g. occlusions) is achieved by reprojecting lost features from their 3D positions back into the image. These tracking results enable the calculation of the present location and from that the homing vector towards which the robot is steered.

When the (relative) distance to the target is small enough, the entire homing procedure is repeated with the next image on the sparse visual path as target. If the path ends, the robot is stopped at a position close to the position where the last path image was taken. This yields a smooth trajectory along a sparsely defined visual path.

### 7.2.1 Feature tracking

The corresponding features found between the first image and the target image in the previous step, also have to be found in the incoming images during driving. This can be done very reliably performing every time wide baseline matching with the first or target image, or both. Although our methods are relatively fast this is still too time-consuming for a driving robot.

Because the incoming images are part of a smooth continuous sequence, a better solution is *tracking*. In the image sequence, visual features move only a little from one image to the next, which enables to find the new feature position in a small search space.

A widely used tracker is the KLT tracker of Kanade, Lucas, Shi, and Tomasi [47]. KLT starts by identifying interest points (corners), which then are tracked in a series of images. The basic principle of KLT is that the definition of corners to be tracked is exactly the one that guarantees optimal tracking. A point is selected if the matrix

$$\begin{bmatrix} g_x^2 & g_x g_y \\ g_x g_y & g_y^2 \end{bmatrix}, \quad (27)$$

containing the partial derivatives  $g_x$  and  $g_y$  of the image intensity function over an  $N \times N$  neighbourhood, has large eigenvalues. Tracking is then based on a Newton-Raphson style minimisation procedure using a purely translational model. This algorithm works surprisingly fast: we were able to track 100 feature points at 10 frames per second in  $320 \times 240$  images on a 1 GHz laptop.

Because the well trackable points are not necessarily coinciding with the anchor points of the wide baseline features to be tracked, the best trackable point in a small window around such an anchor point is selected. In the assumption of local planarity we can always find back the corresponding point in the target image via the relative reference system offered by the wide baseline feature.

### 7.2.2 Recovering lost features

The main advantage of working with this calibrated system is that we can recover features that were lost during tracking. This avoids the problem of losing all features by the end of the homing manoeuvre, a weakness of our previous approach [17]. This feature recovery technique is inspired by the work of Davison [9], but is faster because we do not work with probability ellipses.

In the initialisation phase, all features are described by a local intensity histogram, so that they can be recognised after being lost during tracking. Each time a feature is successfully tracked, this histogram is updated.

When tracking, some features are lost due to invisibility because of e.g. occlusion. Because our local map contains the 3D positions of each feature, and the last robot position in that map is known, we can reproject the 3D feature in the image. Svoboda shows that the world point  $\mathbf{X}_C$  (i.e. the point  $\mathbf{X}$  expressed in the camera reference frame) is projected on point  $\mathbf{p}$  in the image:

$$\mathbf{p} = \frac{K}{2e}(\lambda\mathbf{X}_C - \mathbf{t}_C), \quad (28)$$

wherein  $\lambda$  is the largest solution of

$$\lambda = \frac{b^2(-e)\mathbf{X}_{C3} \pm a\|\mathbf{X}_C\|}{b^2\mathbf{X}_{C3}^2 - a^2\mathbf{X}_{C1}^2 - a^2\mathbf{X}_{C2}^2}. \quad (29)$$

Based on the histogram descriptor, all trackable features in a window around the reprojected point  $\mathbf{p}$  are compared to the original feature. When the histogram distance is under a fixed threshold, the feature is found back and tracked further in the next steps.

### 7.2.3 Motion computation

When in a new image the feature positions are computed by tracking or backprojection, the camera position (and thus the robot position) in the general coordinate system can be found based on these measurements.

It is shown that the position of a camera can be computed when for three points the 3D positions and the image coordinates are known. This problem is known as the *three point perspective pose estimation problem*. An overview of the proposed algorithms to solve it is given by [18]. We chose the method of Grunert, and adapted it for our omnidirectional case.

Also in this part of the algorithm we use RANSAC to obtain a robust estimation of the camera position. Repeatedly the inliers belonging to the motion computed on a three-point sample are counted, and the motion with the greatest number of inliers is kept.

### 7.2.4 Robot motion

In subsection 7.1.1, it is explained how the position and orientation of the target can be extracted from the computed epipolar geometry. Together with the present pose results of the last subsection, a homing vector can easily be computed. This command is communicated to the locomotion subsystem. When the homing is towards the last image in a path, also the relative distance and the target orientation w.r.t. the present orientation is given, so that the locomotion subsystem can steer the robot to stop at the desired position. This is needed for e.g. docking at a table.

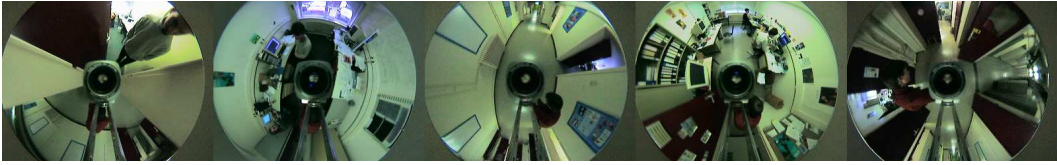


Figure 11: Some of the images which were the input for the map building experiment.

## 8 Experiments

### 8.1 Test platform

We have implemented the proposed algorithm, using our modified electric wheelchair ‘Shar-ioto’. A picture of it is shown in the left of fig. 1. It is a standard electric wheelchair that has been equipped with an omnidirectional vision sensor (consisting of a Sony firewire colour camera and a Neovision hyperbolic mirror, right in fig. 1). The image processing is performed on a 1 GHz laptop. As additional sensors for obstacle detection, 16 ultrasound sensors and a Lidar are added. A second laptop with a 840 MHz processor reads these sensors, receives visual homing vector commands, computes the necessary manoeuvres, and drives the motors via a CAN-bus. More information on this platform can also be found in [35] and [11].

### 8.2 Map building

The wheelchair was guided around in a large environment, while taking images. The environment was a large part of our office floor, containing both indoor and outdoor locations. This experiment yielded a database of 545 colour images with a resolution of  $320 \times 240$  pixels. The total distance travelled was approximately 450 *m*. During a second run 123 images were recorded to test the localisation. Some of these images are shown in fig. 11.

After place clustering with a fixed place size threshold (in our experiments 0.5), this resulted in a set of 53 clusters. Using the Dempster-Shafer based evidence collection, 6 of 41 link hypotheses were rejected, as shown in fig. 12. Fig. 13 shows the resulting 59 place prototypes along with the accepted interconnections.

Instead of keeping all the images in memory, the database is now reduced to only the descriptor sets of each prototype image. In our experiment, the memory needed for the database was reduced from 275 MB to 1.68 MB.

### 8.3 Localisation

From this map, the motion model is computed offline as explained in section 5.2. Now, for the separate test set, the accuracy of the localisation algorithm is tested. A typical experiment is illustrated in fig. 14.

In total, for 78% of the trials the maximum of the belief function was located at the closest place at the first iteration, after the second and third belief update this percentage

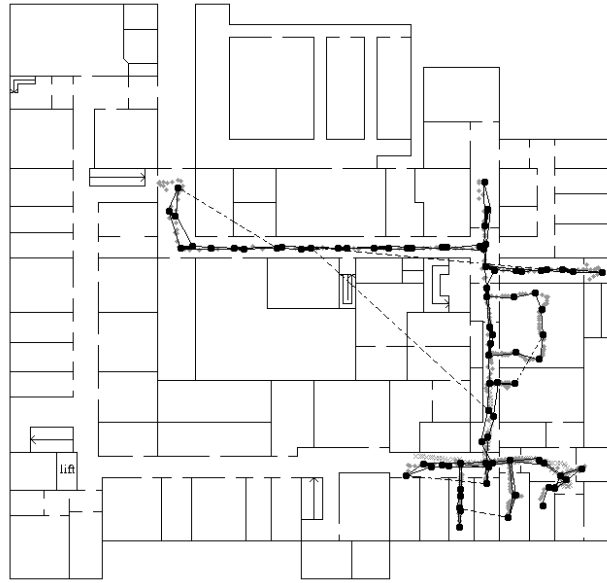


Figure 12: Topological loop closing: accepted hypotheses are shown in thick black lines, rejected in dashed thin black lines.

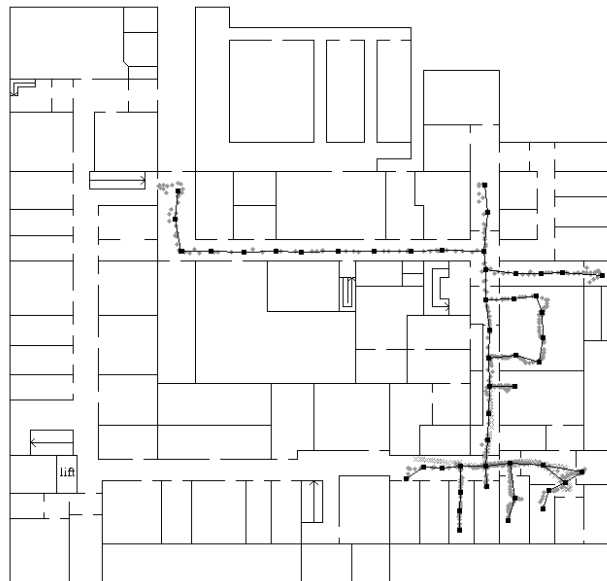


Figure 13: The resulting topological map: locations of the place prototypes with interconnections.

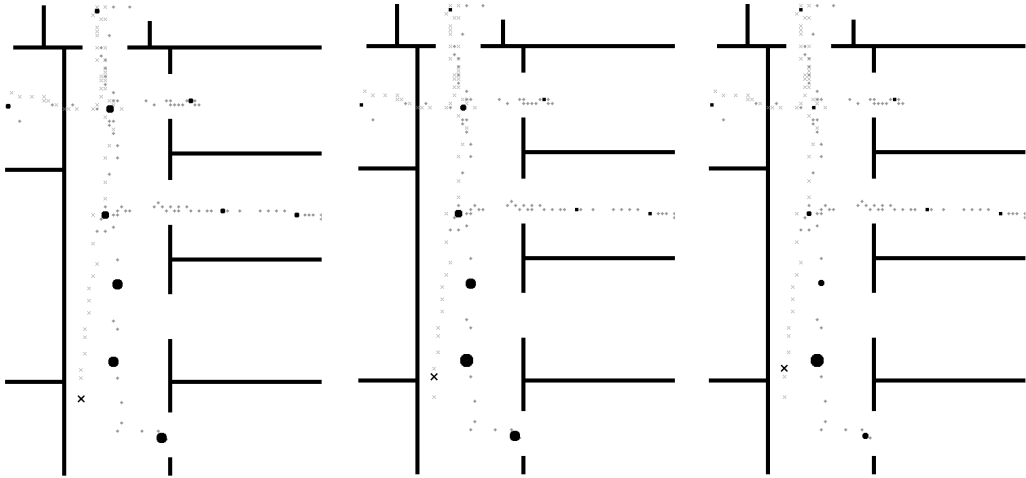


Figure 14: Three belief update cycles in a typical localisation experiment. The black  $\times$  denotes the location of the new image. Place prototypes with a higher belief value are visualised as larger black circles.

raised to 89% and 97%.

## 8.4 Visual servoing

### 8.4.1 Initialisation phase

During the initialisation phase of one visual homing step, correspondences between the present and target image are found and the epipolar geometry is computed. This is shown in fig. 15.

To test the correctness of the initial homing vector, we took images with the robot positioned at a grid with a cell size of 1 meter. The resulting homing vectors towards one of these images (taken at  $(6, 3)$ ) are shown in fig. 16. This proves the fact that even if the images are situated more than 6 metres apart the algorithm works thanks to the use of *wide baseline* correspondences.

### 8.4.2 Motion phase

We present a typical experiment in fig. 17. During the motion, the top of the camera system was tracked in a video sequence from a fixed camera. This video sequence, along with the homography computed from some images taken with the robot at reference positions, permits calculation of metrical robot position ground truth data.

Repeated similar experiments showed an average homing accuracy of 11 cm, with a standard deviation of 5 cm, after a homing distance of around 3 m.

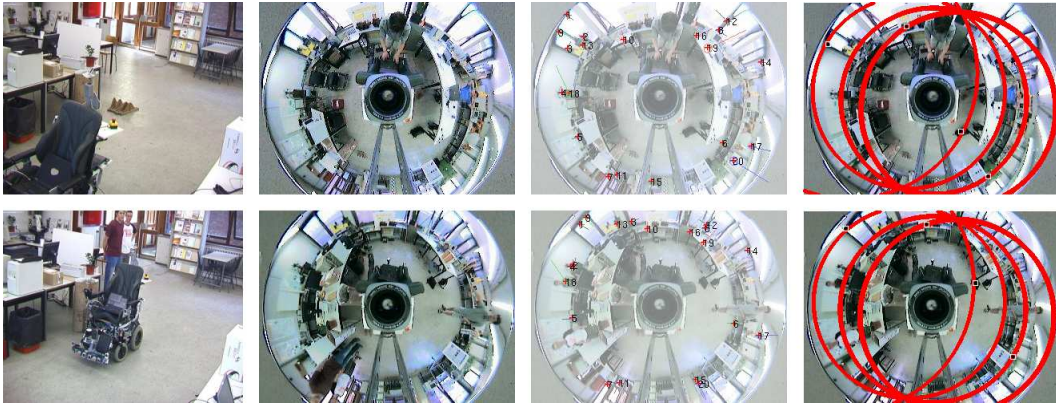


Figure 15: Results of the initialisation phase. Top row: target, bottom row: start. From left to right, the robot position, omnidirectional image, visual correspondences and epipolar geometry are shown.

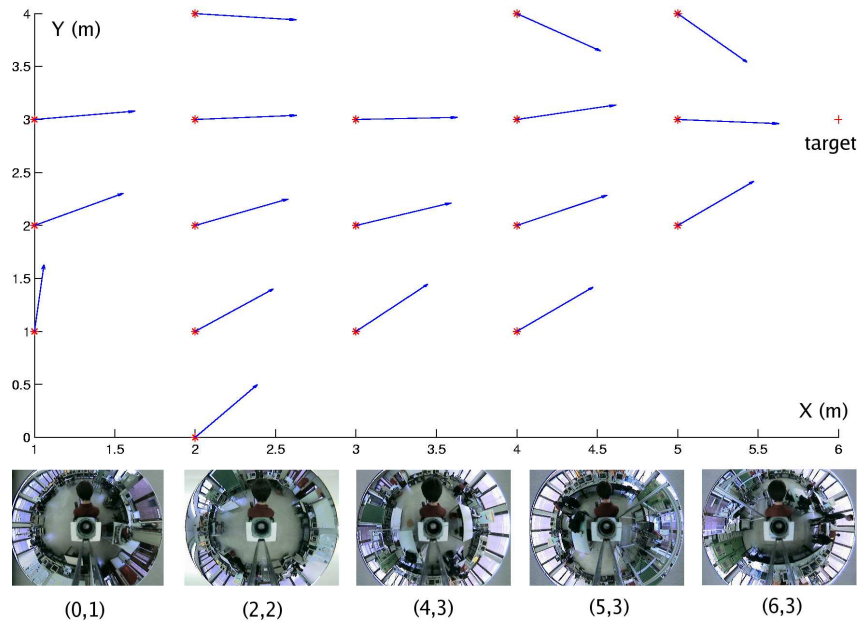


Figure 16: Homing vectors from 1-meter-grid positions and some of the images.

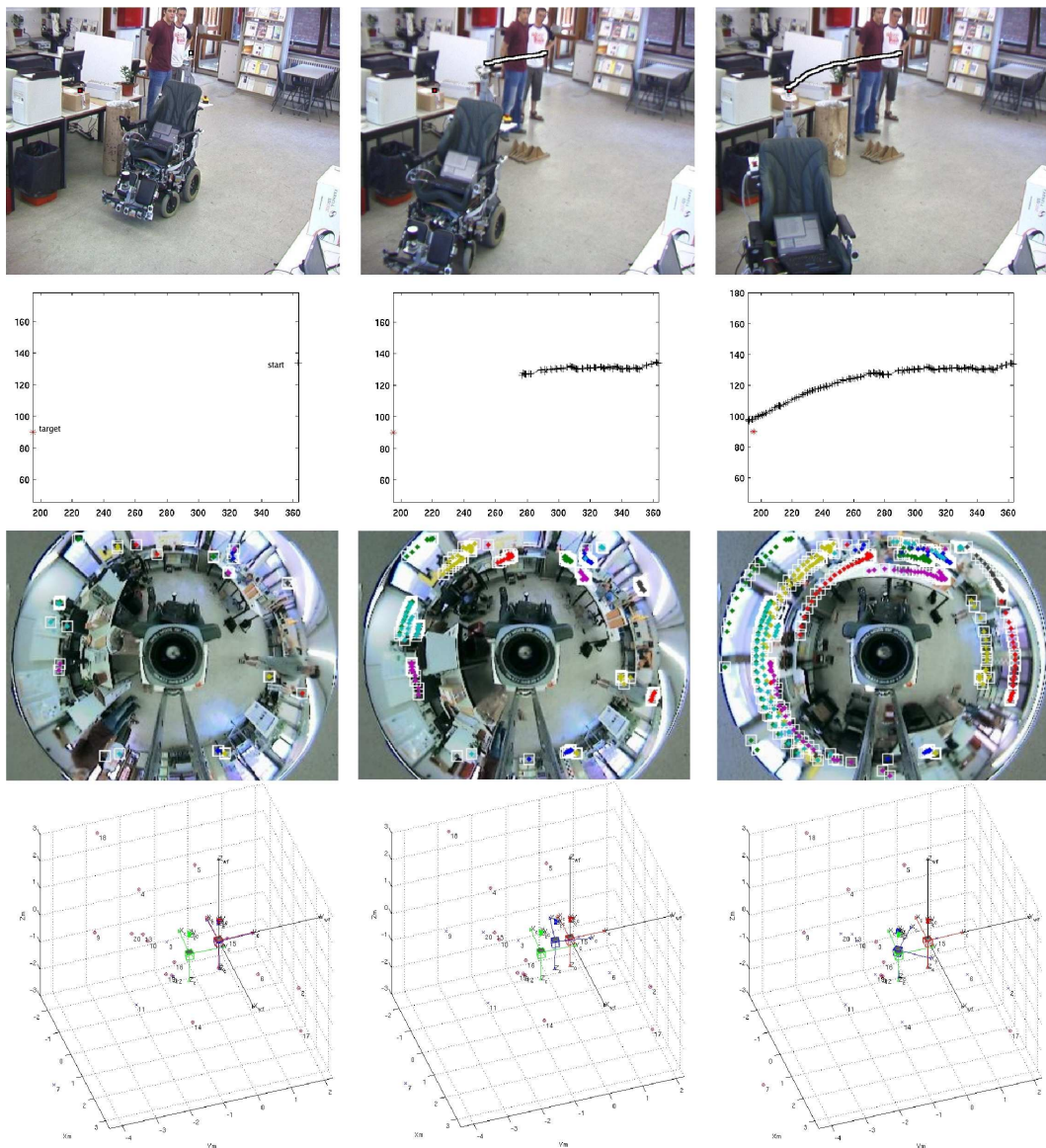


Figure 17: Three snapshots during the motion phase: in the beginning (left), half (centre) and at the end (right) of the homing motion. The first row shows the external camera image with tracked robot position. The second row shows the computed world robot positions [cm]. The third row shows the colour-coded feature tracks. The bottom row shows the sparse 3D feature map (encircled features are not lost).

### 8.4.3 Timing results

The algorithm runs surprisingly fast on the rather slow hardware we used: the initialisation for a new target lasts only 958 ms, while afterwards every 387 ms a new homing vector is computed. For a wheelchair driving at a cautious speed, it is possible to keep on driving while initialising a new target. This resulted in a smooth trajectory without stops or sudden velocity changes.

## 9 Conclusion

This chapter describes and demonstrates a novel approach for a service robot to navigate autonomously in a large, natural complex environment. The only sensor is an omnidirectional colour camera. As environment representation, a topological map is chosen. This is more flexible and less memory demanding than metric 3D maps. Moreover, it does not show error build-up and enables fast path planning. As natural landmarks, we use two kinds of fast wide baseline features which we developed and adapted for this task. Because these features can be recognised even if the viewpoint is substantially different, a limited number of images suffice to describe a large environment.

Experiments show that our system is able to build autonomously a map of a natural environment it drives through. The localisation ability, with and without knowledge of previous locations, is demonstrated. With this map, a path towards each desired location can be computed efficiently. Experiments with a robotic wheelchair show the feasibility of executing such a path as a succession of visual servoing steps.

## Acknowledgments

This work is partially supported by the Inter-University Attraction Poles, Office of the Prime Minister (IUAP-AMS), the Flemish Institute for the Advancement of Science in Industry (IWT), and the Fund for Scientific Research Flanders (FWO-Vlaanderen, Belgium).

## References

- [1] D. Aliaga, "Accurate catadioptric calibration for real-time pose estimation in room-size environments," in Proc. IEEE Int. Conf. Computer Vision (ICCV), Vancouver, pp. 127–134, 2001.
- [2] A. Baumberg, "Reliable feature matching across widely separated views," Computer Vision and Pattern Recognition, Hilton Head, South Carolina, pp. 774-781, 2000.
- [3] R. Basri, E. Rivlin, and I. Shimshoni, "Visual homing: Surfing on the epipoles," in IEEE International Conference on Computer Vision ICCV'98, pp. 863-869, Bombay, 1998.
- [4] K. Beevers and W. Huang, "Loop Closing in Topological Maps", ICRA, 2005.

- 
- [5] H. Bischof, H. Wildenauer, and A. Leonardis, "Illumination insensitive eigenspaces," In Proc. ICCV01, pages 233-238. IEEE Computer Society, 2001.
- [6] H. van Veen, H.K. Distler, S.J. Braun and H.H. Bülthoff, "Navigating through a virtual city: Using virtual reality technology to study human action and perception," *Future Generation Computer Systems* 14, 231-242 (1998).
- [7] B. Cartwright, T. Collett, "Landmark Maps for Honeybees," *Biol. Cybernetics*, 57, pp. 85-93, 1987.
- [8] Cheng Chen and Han Wang, "Appearance-based topological Bayesian inference for loop-closing detection in cross-country environment", *IROS*, pp. 322-327, Edmonton, 2005.
- [9] A. Davison, "Real-time simultaneous localisation and mapping with a single camera," *Intl. Conf. on Computer Vision, Nice*, 2003.
- [10] A. P. Dempster, "Upper and Lower Probabilities Induced by a Multivalued Mapping", *The Annals of Statistics* (28), pp. 325-339, 1967.
- [11] E. Demeester, M. Nuttin, D. Vanhooydonck, and H. Van Brussel, "Fine Motion Planning for Shared Wheelchair Control: Requirements and Preliminary Experiments," *International Conference on Advanced Robotics, Coimbra*, pp. 1278-1283, 2003.
- [12] E. W. Dijkstra, "A note on two problems in connection with graphs", *Numerische Mathematik*, 1: 269-271, 1959.
- [13] M. Fischler, R. Bolles, "Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis," *Comm. of the ACM*, Vol 24, pp 381-395, 1981.
- [14] M. Franz, B. Schölkopf, H. Mallot, and H. Bülthoff, "Where did I take that snapshot? Scene-based homing by image matching," *Biological Cybernetics*, 79, pp. 191-202, 1998.
- [15] C. Geyer and K. Daniilidis, "Mirrors in motion: Epipolar geometry and motion estimation," *Proc. ICCV*, p. 766, Nice, 2003.
- [16] T. Goedemé, T. Tuytelaars, and L. Van Gool, "Fast Wide Baseline Matching with Constrained Camera Position," *Computer Vision and Pattern Recognition, Washington, DC*, pp. 24-29, 2004.
- [17] T. Goedemé, T. Tuytelaars, G. Vanacker, M. Nuttin and L. Van Gool, "Feature Based Omnidirectional Sparse Visual Path Following," *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2005, Edmonton*, 2005.
- [18] R. Haralick, C. Lee, K. Ottenberg, and M. Nölle, "Review and analysis of Solutions of the Three Point Perspective Pose Estimation Problem," *International Journal of Computer Vision*, 13, 3, pp. 331-356, 1994.
- [19] R. Hartley, "Estimation of relative camera positions for uncalibrated cameras," *2nd European Conference on Computer Vision*, pp. 579-587, Springer-Verlag, LNCS 588, 1992.
- [20] A. Jain, "Fundamentals of digital image processing," Prentice Hall, 1989.

- [21] M. Jogan, and A. Leonardis, "Panoramic eigenimages for spatial localisation," In Solina, Franc and Leonardis, Ales, Eds. Proceedings 8th International Conference on Computer Analysis of Images and Patterns(1689), pages 558-567, Ljubljana, 1999.
- [22] M. Jogan, A. Leonardis, "Robust localization using panoramic view-based recogniton," In Proceedings 15th International Conference on Pattern Recognition ICPR'00, pages 136-139, Barcelona, Spain.
- [23] S. Koenig and R. Simmons, "Unsupervised learning of probabilistic models for robot navigation", ICRA, 1996.
- [24] J. Košecká, X. Yang, "Global Localization and Relative Pose Estimation Based on Scale-Invariant Features," ICPR (4), pp. 319-322, 2004.
- [25] B. Kröse, J. Porta, A. van Breemen, K. Crucq, M. Nuttin, and E. Demeester, "Lino, the User-Interface Robot", First European Symposium on Ambient Intelligence (EUSAI 2003), pp. 264-274, Veldhoven, The Netherlands, 2003.
- [26] L. Ledwich and S. Williams, "Reduced SIFT Features For Image Retrieval and Indoor Localisation," Australasian Conf. on Robotics and Automation ACRA, Canberra, 2004.
- [27] C. Loeffler, A. Ligtenberg and G. Moschytz, "Practical Fast 1-D DCT Algorithmh with 11 Multiplications," ICASSP, pp. 988-991, 1989.
- [28] D. Lowe, "Object Recognition from Local Scale-Invariant Features," International Conference on Computer Vision, pp. 1150-1157, 1999.
- [29] J. Matas, O. Chum, M. Urban and T. Pajdla, "Robust wide baseline stereo from maximally stable extremal regions," British Machine Vision Conference, Cardiff, Wales, pp. 384-396, 2002.
- [30] G. Mariottini, E. Alunno, J. Piazzzi, and D. Prattichizzo, "Epipole-Based Visual Servoing with Central Catadioptric Camera," IEEE ICRA05, Barcelona, 2005.
- [31] K. Mikolajczyk and C. Schmid, "An affine invariant interest point detector," ECCV, vol. 1, 128-142, 2002.
- [32] F. Mindru, T. Moons, and L. Van Gool, "Recognizing color patters irrespective of viewpoint and illumination," Computer Vision and Pattern Recognition, vol. 1, pp. 368-373, 1999.
- [33] S. Arya, D. Mount, N. Netanyahu, R. Silverman, and A. Wu, "An optimal algorithm for approximate nearest neighbor searching," J. of the ACM, vol. 45, pp. 891-923, 1998,
- [34] D. Nistér, O. Naroditsky, J. Bergen, "Visual Odometry," Conference on Computer Vision and Pattern Recognition, Washington, DC, 2004.
- [35] M. Nuttin, E. Demeester, D. Vanhooydonck, and H. Van Brussel, "Shared autonomy for wheelchair control: Attempts to assess the user's autonomy," in *Autonome Mobile Systeme*, pp. 127-133, 2001.
- [36] M. Nuttin, D. Vanhooydonck, E. Demeester, H. Van Brussel, K. Buijsse, L. Desimpelaere, P. Ramon, and T. Verschelden, "A Robotic Assistant for Ambient Intelligent Meeting Rooms", First European Symposium on Ambient Intelligence (EUSAI 2003), Veldhoven, The Netherlands, pp. 304-317, 2003.

- 
- [37] T. Okuma, K. Sakaue, H. Takemura, and N. Yokoya, "Real-time camera parameter estimation from images for a Mixed Reality system," Proc. ICPR, Barcelona, Spain, 2000.
- [38] F. Pérer-Fontán, B. Sanmartín, A. Steingäß, A. Lehner, E. Kubista, M.J. Martín, B. Arbesser-Rastburg, "Measurements and modelling of the satellite-to-indoor channel for Galileo," European Navigation Conference GNSS, Rotterdam, 2004.
- [39] M. Pollefeys, L. Van Gool, M. Vergauwen, F. Verbiest, K. Cornelis, J. Tops, R. Koch, "Visual modeling with a hand-held camera," International Journal of Computer Vision 59(3), 207-232, 2004.
- [40] A. Ranganathan, E. Menegatti and Frank Dellaert, "Bayesian Inference in the Space of Topological Maps", IEEE Transactions on Robotics", 2005
- [41] J. Rekimoto and Y. Ayatsuka, "CyberCode: Designing Augmented Reality Environments with Visual Tags," Designing Augmented Reality Environments (DARE 2000), 2000.
- [42] C. Sagüés, J. Guerrero, "Visual correction for mobile robot homing," Robotics and Autonomous Systems, Vol. 50, no. 1, pp 41-49, 2005.
- [43] C. Schmid, R. Mohr, C. Bauckhage, "Local Grey-value Invariants for Image Retrieval," International Journal on Pattern Analysis and Machine Intelligence, Vol. 19, no. 5, pp. 872-877, 1997.
- [44] S. Se, D. Lowe, and J. Little, "Local and Global Localization for Mobile Robots using Visual Landmarks," In proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '01), Hawaii, USA, 2001.
- [45] G. Shafer, "A Mathematical Theory of Evidence", Princeton University Press, 1976.
- [46] H. Shatkay and L. P. Kaelbling, "Learning Topological Maps with Weak Local Odometric Information", IJCAI (2), pp. 920-929, 1997.
- [47] J. Shi and C. Tomasi, "Good Features to Track," Computer Vision and Pattern Recognition, Seattle, pp. 593-600, 1994.
- [48] D. Stricker, P. Daehne, F. Seibert, I.T. Christou, L. Almeida, R. Carlucci, and N.I. Ioannidis, "Design and development issues for archeoguide: An augmented reality based cultural heritage on-site guide," In International Conference on Augmented, Virtual Environments and Three-Dimensional Imaging, Mykonos, Greece, May 28-30 2001. IEEE.
- [49] T. Svoboda, T. Pajdla, and V. Hlaváč, "Motion Estimation Using Panoramic Cameras," Conf. on Intelligent Vehicles, Stuttgart, pp. 335-340, 1998.
- [50] T. Svoboda, "Central Panoramic Cameras, Design, Geometry, Egomotion," PhD Thesis, Czech Technical University.
- [51] A. Tapus and R. Siegwart, "Incremental Robot Mapping with Fingerprints of Places", IROS, Edmonton, 2005.
- [52] T. Tuytelaars, L. Van Gool, L. D'haene, and R. Koch, "Matching of Affinely Invariant Regions for Visual Servoing," Intl. Conf. on Robotics and Automation, pp. 1601-1606, 1999.

- [53] T. Tuytelaars and L. Van Gool, "Wide baseline stereo based on local, affinely invariant regions," British Machine Vision Conference, Bristol, UK, pp. 412-422, 2000.
- [54] I. Ulrich and I. Nourbakhsh, "Appearance-Based Place Recognition for Topological Localisation," IEEE International Conference on Robotics and Automatisation, San Francisco, CA, April 2000, pp. 1023-1029
- [55] A. Vale, M. Isabel Ribeiro, "Environment Mapping as a Topological Representation," Proceedings of the 11th International Conference on Advanced Robotics - ICAR2003 Universidade de Coimbra, Portugal, June 30 - July 1- 3, 2003.
- [56] Z. Zivkovic, B. Bakker and B. Kröse, "Hierarchical Map Building Using Visual Landmarks and Geometric Constraints", IROS, Edmonton, pp. 7-12, 2005.