

Issues and Strategies for Robotic Object Manipulation in Domestic Settings

Danica Kragic, Mårten Björkman, Henrik I Christensen and Jan-Olof Eklundh

Computer Vision and Active Perception Lab

Centre for Autonomous Systems

Numerical Analysis and Computer Science

Royal Institute of Technology, Stockholm, Sweden

Email:{danik, celle, hic, joe}@nada.kth.se

Abstract—Many robotic tasks such as autonomous navigation, human-machine collaboration, object manipulation and grasping facilitate visual information. Some of the major research and system design issues in terms of visual systems are robustness and flexibility.

In this paper, we present a number of visual strategies for robotic object manipulation tasks in natural, domestic environments. Given a complex fetch-and-carry type of tasks, the issues related to the whole *detect-approach-grasp* loop are considered. Our vision system integrates a number of algorithms using monocular and binocular cues to achieve robustness in realistic settings. The cues are considered and used in connection to both foveal and peripheral vision to provide depth information, segment the object(s) of interest in the scene, object recognition, tracking and pose estimation. One important property of the system is that the step from object recognition to pose estimation is completely automatic combining both appearance and geometric models. Rather than concentrating on the integration issues, our primary goal is to investigate the importance and effect of camera configuration, their number and type, to the choice and design of the underlying visual algorithms. Experimental evaluation is performed in a realistic indoor environment with occlusions, clutter, changing lighting and background conditions.

I. INTRODUCTION

One of the key components of a robotic systems operating in a dynamic, unstructured environment is robust perception. Our current research considers the problem of mobile manipulation in domestic setting where, in order for the robot to be able to detect and manipulate objects in the environment, robust visual feedback is of key importance. In case of humans, complex coordination between the eye and the hand is facilitated during execution of everyday activities such as pointing, grasping, reaching or catching. Each of these activities or actions require attention to different attributes in the

environment - while pointing requires only an approximate location of the object in the visual field, a reaching or grasping movement require more exact information about the object's pose.

In robotics, the use of visual feedback for coordination of a robotic arm motion is termed *visual servoing*, [Hutchinson et al. \(1996\)](#). Given a complex fetch-and-carry type of task, issues related to the whole *detect-approach-grasp* loop have to be considered. Most visual servoing systems, however, deal only with the *approach* step and forget about the problems such as *detecting* the object of interest in the scene or retrieving its 3D structure in order to perform grasping. A so called *teach-by-showing* approach is typically used where the desired camera placement with respect to the object is well defined and known before hand.

Our interest is the development of an architecture that integrates a number of modules where each module encapsulates a number of visual algorithms responsible for a particular task such as recognition or tracking. Our system is heavily based on the *active vision* paradigm, [Ballard \(1991\)](#) where, instead of passively observing the world, viewing conditions are actively changed so that the best results are obtained given a task at hand.

In our previous work, [Bjorkman & Kragic \(2004\)](#) we have presented a system that consists of two pairs of stereo cameras: a peripheral camera set and a foveal one. Recognition and pose estimation was performed using either one of these, depending on the size of and distance to an observed object. From segmentation based on binocular disparities, objects of interest were found using the peripheral camera set, which then triggered the system to perform a saccade, moving the object into the centre of foveal cameras achieving thus a combination of a large field of view and high image resolution. Compared to one of the recent systems, [S. Kim & Kweon \(2003\)](#), our system used both hard (detailed models) and soft modelling (approximate shape) for

"Pick Up ..."		WHERE (<i>location</i>)	
		known	unknown
WHAT (<i>identity</i>)	known	"This Cup"	"The Cup"
	unknown	"This Object"	"Something"

Fig. 1
ROBOTIC MANIPULATION SCENARIOS.

object segmentation. In addition, choice of binocular or monocular cues was used depending on the task.

In this paper, we

This paper is organized as follows,

II. BACKGROUND AND MOTIVATION

In our current system, the robot may be given tasks such as "Robot, bring me the raisins" or "Robot, pick up this". Depending on the prior information, i.e. task or context information, different solution strategies may be chosen. The first task of the above is well defined in that manner that the robot already has the internal representation of the object, e.g. the *identity* of the object is known. An example of such a task is shown in Figure 2: after being given a spoken command, the robot locates the object, approaches it, estimates its pose and finally performs grasping. More details related to this approach are given in Section III.

For the second task, the spoken command is commonly followed by a pointing gesture - here, the robot does not know the *identity* of the object, but it knows its approximate *location*. The approach considered in this work is presented in Section IV. Figure 1 shows different scenarios with respect to prior knowledge of object *identity* and *location*, with the above examples being shaded. A different set of underlying visual strategies is required for each of these scenarios. We have considered these two scenarios since they are the most representative examples for robotic fetch-and-carry tasks.

Techniques employed in terms of visual servoing and object manipulation in general depend on:

- Camera placement: Most visual servoing systems today use *eye-in-hand* cameras and deal mainly with the *approach* object step in a *teach-by-showing* manner, Malis *et al.* (2003). In our approach, we consider a combination of a stand-alone stereo and an eye-in-hand camera systems, Kragic & Christensen (2003b).
- Number of cameras: In order to extract metric information, e.g. sizes and distances, about objects

observed by the robot, we will show how we can benefit from binocular information. The reason for using multiple cameras in our system is the fact that it simplifies the problem of segmenting the image data into different regions representing objects in a 3D scene. This is often referred to as *figure-ground segmentation*. In cluttered environments and complex backgrounds, figure-ground segmentation is particularly important and difficult to perform and commonly the reason for experiments being performed in rather sparse, simplified environments.

- Camera type: zooming or not, combinations of foveal and peripheral, etc. Here, very little work has been reported, Benhimane & Malis (2003).

A. Experimental platform

The experimental platform is a Nomadic Technologies XR4000 and is equipped with a Puma 560 arm for manipulation (see Figure 3). The robot has sonar sensors, a SICK laser scanner, a wrist mounted force/torque sensor (JR3), and a color CCD camera mounted on the Barrett Hand gripper. The palm of the Barrett hand is covered by a VersaPad touch sensor and, on each finger, there are three Android sensors. On the robot's shoulder, there is a binocular stereo-head. This system, known as Yorick, has four mechanical degrees of freedom; neck pan and tilt, and pan for each camera in relation to the neck. The head is equipped with a pair of Sony XC999 cameras, with focal lengths of 18 mm.

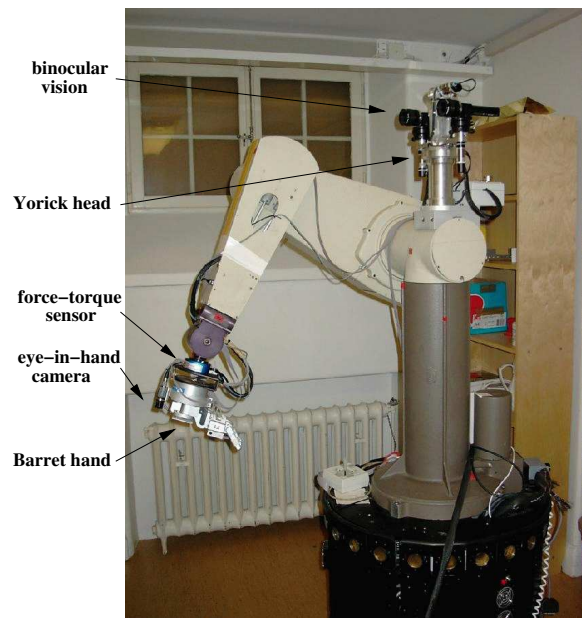


Fig. 3
EXPERIMENTAL PLATFORM.



Fig. 2

DETECT-APPROACH-GRASP EXAMPLE.

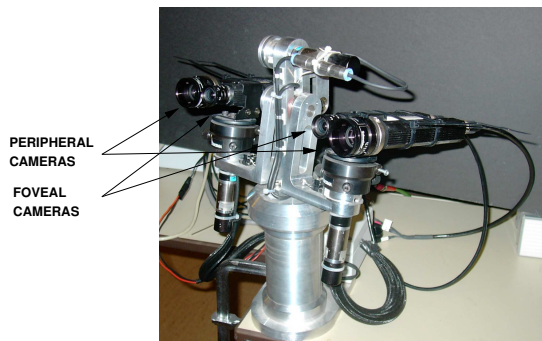


Fig. 4

THE YORICK STEREO-HEAD.

For some of the experimental results that will be presented further on, the binocular stereo-head system shown in Figure 4 was used. Here, the head is equipped with two pairs of Sony XC999 cameras, with focal lengths 28 mm and 6 mm respectively. The motivation for this combination of cameras will be explained related to the examples.

B. Stereo System Modeling - Epipolar Geometry

If a binocular set of cameras is available, differences in position between projections of 3D points onto the left and right image planes (disparities) can be used to perform figure-ground segmentation and retrieve the information about three-dimensional structure of the scene. If the relative orientation and position between cameras is known, it is possible to relate these disparities to actual metric distances. One of the commonly used settings is where the cameras are rectified and their optical axes mutually parallel, [Kragic & Christensen \(2003b\)](#). However, one of the problems arising is that the part of the scene contained in the field of view of both cameras simultaneously is quite limited.

Another approach is to estimate the epipolar geometry

continuously from image data alone, [Bjorkman \(2002\)](#). Additional reason for this may be that small disturbances such as vibrations and delays introduce significant noise to the estimation of the 3D structure. In fact, an error of just one pixel leads to depth error of several centimeters on a typical manipulation distance. Therefore, for some of the manipulation tasks, the epipolar geometry is estimated robustly using Harris' corner features, [Harris & Stephens \(1988\)](#). Such corner features are extracted and matched between the camera images using normalized cross-correlation. The vergence angle α , gaze direction t , relative tilt r_x and rotation around the optical axes r_z , are iteratively sought using

$$\begin{pmatrix} dx \\ dy \end{pmatrix} = \begin{pmatrix} (1+x^2)\alpha - yr_z \\ xy\alpha + r_y + xr_z \end{pmatrix} + \frac{1}{Z} \begin{pmatrix} 1 - xt \\ -yt \end{pmatrix}, \quad (1)$$

where Z is the unknown depth of a point at image position (x, y) . The optimization is performed using a combination of RANSAC ([Fischler & Bolles, 1981](#)) for parameter initialisation, and M-estimators ([Huber, 1981](#)) for improvements.

This optical flow model ([Longuet-Higgins, 1980](#)) is often applied to motion analysis, but has rarely been used for stereo. The reason for this is because the model is approximate and only works for relatively small displacements. In our previous work we have, however, experimentally shown that this model is more robust than the essential matrix in the case of binocular stereo heads, [Björkman & Eklundh \(2002\)](#), even if the essential matrix leads to a more exact description of the epipolar geometry, [Longuet-Higgins \(1981\)](#).

III. MANIPULATING *known* OBJECTS

If a robot is to manipulate a known object, some type of representation is typically known in advance. Such a representation may include object textural and/or geometrical properties which are sufficient for the object to be located and manipulation task to be performed.

For realistic settings, a crude information about objects location can sometimes be provided from the task level. e.g. “Bring me red cup from the dinner table.” However, if the location of the object is not provided, it is up to the robot to search the scene. The following sections give examples of how these are problems are approached in the current system.

A. Detect

If we can assume that the object is in the field of view from the beginning of the task, a monocular recognition system can be used to locate the object in the image, [Zillich et al. \(2001\)](#).

However, when a crude information about object’s current position is not available, detecting a known object is not an easy task since a large number of false positives can be expected. Candidate locations have to be analysed in sequence which may be computationally too expensive, unless the robot has an attentional system that delivers the most likely candidate locations first, using as much information about the requested object as possible.

A natural approach here is to facilitate a binocular system that provides metric information as an additional cue. Since the field of view of a typical camera is quite limited, binocular information can only be extracted from those parts of the 3D scene that are covered by both cameras’ field of view. In order to make sure that an object of interest is situated in the centre of each camera’s field of view, the head is able to actively change gaze direction and vergence angle, i.e. the difference in orientation between the two cameras. In our system, stereo based figure-ground segmentation is intended for mobile robot navigation and robot arm transportation to the vicinity of the object. More detailed information about an object’s pose is provided using a monocular model based pose estimation and tracking.

In general, this part of the system can be seen as a visual front-end, responsible for delivering 3D data about the observed scene. Such information is extracted using a three-step process, which includes epipolar geometry estimation, image rectification and calculation of dense disparity maps. The generation of this data is done continuously at a rate of 8 Hz, independently of the task at hand and used by more high-level processes for further interpretation. Further information on this part of the system can be found in [Bjorkman \(2002\)](#). Since most methods for dense disparity estimation assume the image planes to be parallel, image rectification has to be performed using the estimated epipolar geometry before disparities can be estimated. The current system includes seven different disparity algorithms, from simple area

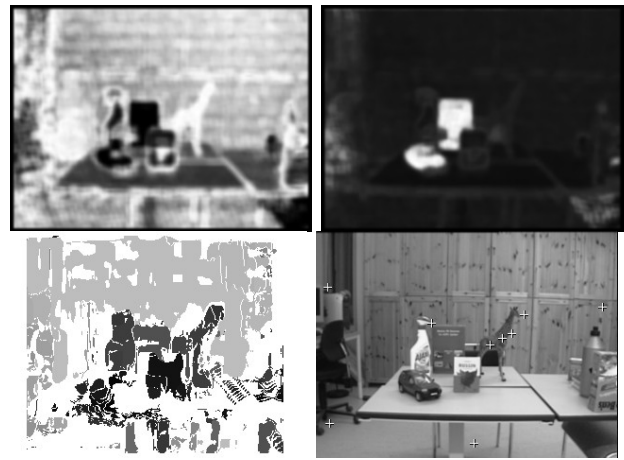


Fig. 6

OBJECT HYPOTHESIS GENERATION

correlation, [Konolige \(1997\)](#) to more complicated graph-cut methods, [Kolmogorov & Zabih \(2001\)](#). The benefit of using a more advanced global method, is the fact that they often lead to denser and more accurate results. However, even if density is important, the computational cost of these methods makes them infeasible for our particular application which means that correlation based methods are typically used in practice. The second column of Figure 5 shows two examples of disparities calculated using sums of absolute differences.

Currently, we use two kinds of visual cues for this purpose, 3D size and hue histograms. These cues were chosen since they are highly object dependent and relatively insensitive to changing lighting conditions, object pose and viewing direction.

The images in Figure 6 show an example where the giraffe in the centre of the lower-right image is requested. The upper images illustrate the saliency maps generated using the hue histograms of the giraffe (left) and a blue box (right) respectively. From the disparity map (lower-right) and the saliency map based on the giraffe, a number of candidate locations are found as shown in the lower left image. We further use recognition to verify that a requested object has indeed been found. With attention and recognition applied in a loop, the system is able to automatically search the scene for a particular object, until it has been found by the recognition system. Two recognition modules are available for this purpose: i) a feature based module based on Scale Invariant Feature Transform (SIFT) features [Lowe \(1999\)](#), and ii) an appearance based module using colour histograms.

Most recognition algorithms expect the considered object to subtend a relatively large proportion of the images. If the object is small, it has to be approached

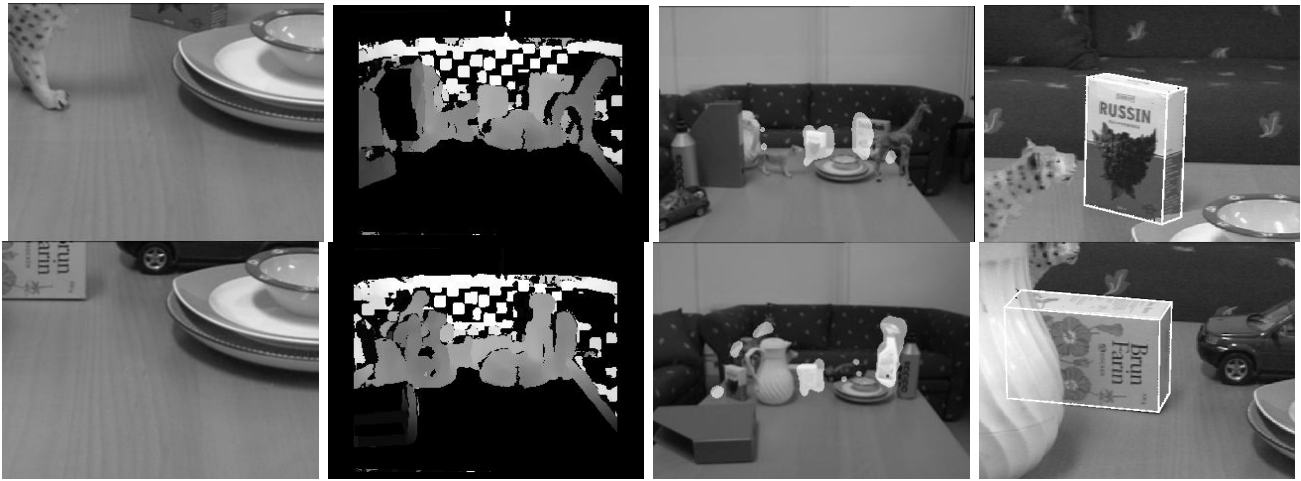


Fig. 5

AN EXAMPLE OF BINOCULAR FIGURE-GROUND SEGMENTATION AND POSE ESTIMATION. THE FIRST COLUMN SHOWS THE FOVEAL IMAGES BEFORE A SACCADIC MOVEMENT HAS BEEN ISSUED. DISPARITY MAPS CAN BE SEEN IN THE SECOND COLUMN WITH OBJECT HYPOTHESES SHOWN IN THE THIRD. THE LAST TWO IMAGES SHOW THE POSE OF RECOGNISED OBJECTS BEING CORRECTLY ESTIMATED.

before is can be detected. A more efficient solution is a system equipped with wide field as well as foveal cameras, like the stereo-head system used for the example presented here. Hypotheses are found using the wide field cameras, while recognition is done using the foveal ones. An alternative solution would be using an eye-in-hand camera and only approach the object through the manipulator, keeping the platform itself static.

B. Approach

Transporting the arm to the vicinity of the object, considering a closed-loop control system, requires registration or computation of spatial relationship between two or more images. Although this problem has been studied extensively in the computer vision society, it has rarely been fully integrated in robotic systems for unknown objects. One reason for this is that high real-time demand makes the problem of tracking more difficult than when processing image sequences off-line. For cases where the object is initially far away from the robot, a simple tracking techniques can be used to keep the object in the field of view while approaching it. For this purpose we have developed and evaluated methods based on correlation and optical flow, [D. Kragic & Allen \(2001\)](#) as well as those based on integration of cues such as texture, colour and motion, [Kragic & Christensen \(2002\)](#). The latter approach is currently facilitated for tracking.

Performing final approach towards a known object depends also on the number of cameras and their placement. For eye-in-hand configuration we have adopted

a *teach-by-showing* approach, where a stored image taken from the reference position is used to move the manipulator so that the current camera view is gradually changed to match the stored reference view. Accomplishing this for general scenes is difficult, but a robust system can be made under the assumption that the objects are piecewise planar. In our system, a wide baseline matching algorithm is employed to establish point correspondences between the current and the reference image, [Kragic & Christensen \(2002\)](#). The point correspondences enable the computation of a homography relating the two views, which is then used for 2 1/2D visual servoing.

Another method uses a CAD model of the object, which in our case also includes a set of SIFT features, for full 6D pose estimation and tracking. After the object has been localised in the image, its pose is automatically initiated using SIFT features from the foveal camera image, fitting a plane to the data. Thus it is assumed that there is a dominating plane that can be mapped to the model. The process is further improved searching for straight edges around this plane. An example of this approach is shown in [Figure 7](#).

IV. MANIPULATING *unknown* OBJECTS

For general setting, manipulation of unknown objects has rarely been pursued. The primary reason is likely to be that the shape of an object has to be determined in order to successfully grasp it. Another reason is that, even if the location is given by a pointing gesture, the size also has to be known and the object segmented from its background.

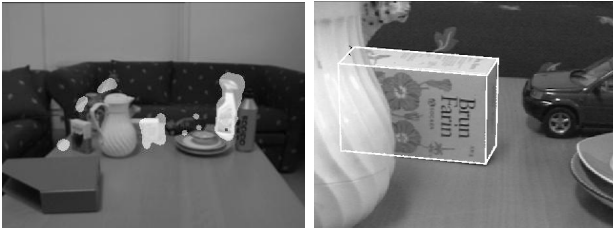


Fig. 7

POSE ESTIMATION USING SIFT FEATURES

A. Detect

Numerous methods exist for segmentation of objects in cluttered scenes. However, from monocular cues only this is very difficult, unless the object has a colour or texture distinct from its surrounding. Unfortunately, these cues are sensitive to lighting as well as pose variations. Thus, for the system to be robust, one has to rely on information such as binocular disparities or optical flow. A binocular setting is recommended, since the motion that needs to be induced should preferably be parallel to the image plane, complicating the process of approaching the object.

In our current system, binocular disparities are used for segmentation with the foveal camera set. We use this set since the focal lengths have to be relatively large in order to get the accuracy required for grasping. When the resolution in depth increases, so does the range of possible disparities. If only a fraction of these disparities are tested, e.g. the range in which the object is located, a large number of outliers can be expected, such as in the lower-left image of Figure 8. We apply a Mean-Shift algorithm, [Comaniciu et al. \(2000\)](#) to prune the data, using the fact that the points representing the object are located in a relatively small part of 3D space and the centre of these points is approximately known. After applying a sequence of morphological operation a mask is found as shown in the lower-right image.

B. Approach

Approaching an unknown object can be done either using the stereo-head or with an eye-in-hand camera. Without knowing the identity of the object the latter case is hardly feasible. It would be possible to take a sequence of images, while approaching the object, and from these estimate a disparity map, but this map would hardly be as accurate as using the disparities available from the foveal camera set.

If the stereo-head is used instead, it is essential that the robot gripper itself can be located in disparity space.

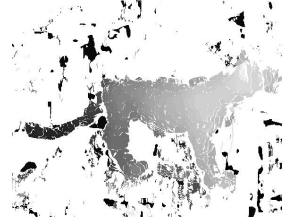


Fig. 8

FIGURE-GROUND SEGMENTATION

Using the mask derived in Section IV-A, the elongation and orientation of the object can be determine and the fingers of the gripper be placed on either side of the object. In general we will not be able, from one stereo view only, to retrieve the full 3D shape of the object. In particular, if the extension in depth is significant, it will be difficult to guarantee that the full closing grasp can be performed. This problem can be solved by moving the stereo-head to another location. This is a topic we intend to investigate further in the future.

V. GRASP

For active grasping, visual sensing will in general not suffice. One of the problems closely related to eye-in-hand configurations is the fact that when the *approach* step is finished, the object is very close to the camera, commonly covering the whole field of view. To retrieve features necessary for grasp planning is impossible. One solution to this problem is to use a wide field eye-in-hand camera, together with a stand-alone mono- or stereo vision system. Our previous work has integrated visual information with tactile and force-torque sensing for object grasping, [Kragic & Christensen \(2003a\)](#). We have, however, realised that there is a need for a system that is able to monitor the grasping process and track the pose of the object during execution. We have shown that in this way, even if the robot moves the object, grasping can successfully be performed without the need to reinitiate the whole process. This can be done even for unknown objects where the Mean-Shift strategy suggested in Section IV-A is applied on consecutive images.

VI. THE SYSTEM

Figure 9 shows a schematic overview of the basic building blocks of the system. These blocks do not necessarily correspond to the actual software components, but are shown in order to illustrate the flow of information through the system. For example, the visual front end consists of a several components, some of which are running in parallel and others hierarchically. On the other hand, action generation, such as initiating 2D or 3D tracking, is distributed and performed across multiple components.

The most important building blocks can be summarized as follows:

- The Visual Front-End is responsible for the extraction of visual information needed for figure-ground segmentation and other higher level processes.
- Hypotheses Generation produces a number of hypotheses about the objects in the scene that may be relevant to the task at hand. The computations are moved from being distributed across the whole image to particular regions of activation.
- Recognition is performed on selected regions, using either corner features or color histograms, to determine the relevancy of observed objects.
- Action Generation triggers actions, such as visual tracking and pose estimation, depending on the outcome of the recognition and current task specification.

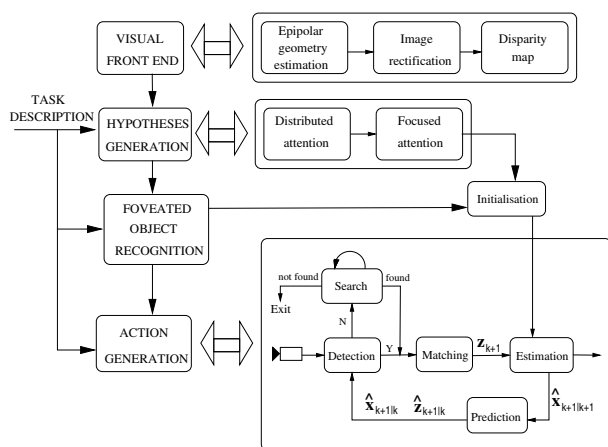


Fig. 9

BASIC BUILDING BLOCKS OF THE SYSTEM.

Due to the complexity of the software system, it was partitioned into a number of smaller modules that communicate through a framework built on an inter-process communication standard called CORBA (Common Object Request Broker Architecture), [Vinoski \(1997\)](#).

The current version of the system consists of about ten such modules, each running at a different frame rate. The lowest level frame grabbing module works at a frequency of 25 Hz, while the recognition modules is activated only upon request. In order to consume processing power, modules are shut down temporarily when not been accessed by any other module within a time frame of 10 seconds.

The experimental evaluation has been performed on a 1.2 GHz dual Athlon MP computer running under the Linux operating system. For frame grabbing, a Leutron PicProdigy board was used. This board is able to simultaneously grab images from all four cameras at full frame-rate.

VII. HYPOTHESES GENERATION

The purpose of this component is to derive qualified guesses of *where* the object of interest is located in the current scene. As mentioned earlier, this step is performed using the peripheral cameras while the recognition module uses the foveal ones. This requires a transfer from peripheral to foveal vision, or from distributed to focused attention [Palmer \(1999\)](#).

A. Distributed attention

Unlike focused attention, distributed attention works on the whole image instead of being concentrated to a particular image region. Using the available visual cues a target region, that might represent an object of interest, is identified. Even if the current system is limited to binocular disparities, it is straightforward to add additional cues, such as in the model of Itti, Koch and Niebur [L. Itti & Niebur \(1998\)](#). Here, we have concentrated on disparities because they contain valuable information about object sizes and shapes. This is especially important in a manipulation task, where the color of an object might be irrelevant, whereas the size is not.

The only top-down information needed for hypotheses generation is the expected size of an object of interest and the approximate distance from the camera set. A binary map is created containing those points that are located within a specified distance range. The third column of Figure 5 shows two such maps overlaid on-top of the corresponding left peripheral images. Initial hypotheses positions are then generated from the results of a difference of Gaussian filter applied to the binary map. The scale of this filter is set so as to maximize the response of image blobs representing objects of the requested size and distance.

B. Focused attention

From the generated hypotheses, a target region is selected so that the gaze can be redirected and recognition performed using the foveal cameras. To make the experimental analysis easier and safer, this is currently done manually by an external user who points at the hypothesis of choice. It is straightforward to make this fully automatic by searching through the view space until the requested object is found.

Since hypotheses are described in the peripheral cameras frame and recognition is performed using the foveal ones, the relative transformations have to be known. These are found applying a similarity model to a set of Harris' corner features similar to those used for epipolar geometry estimation in Section II-B. The relative rotations, translations and scales are continuously updated at a rate of about 2 Hz. Knowing the transformations, it is possible to translate the hypotheses positions into the foveal camera frames.

Before a saccade is finally executed, fixating the foveal cameras onto the selected hypothesis region, the target position is refined in 3D. During a couple of image frames, a high-resolution disparity map is calculated locally around the target area. A mean shift algorithm, [Comaniciu et al. \(2000\)](#), is run iteratively updating the position from the cluster of 3D points around the target position, represented by the disparity map. The maximum size of this cluster is specified using the top-down information mentioned above. The first two images of Figure 10 show these clusters highlighted in the left peripheral images before and after a saccade. The foveal images after the saccade can be seen to the right.

VIII. EXPERIMENTAL EVALUATION

As mentioned in Section VI, our system is built on a number of independently running, but communicating, modules. Since most methods used within these modules have been analysed elsewhere, we will concentrate on the integrated system as a whole, rather than analysing each individual method in isolation. The system should be considered as an integrated unit and its performance measured based on the behaviour of the complete system. The failure of one particular module does not necessarily mean that the whole system fails. For example, figure-ground segmentation might well fail to separate two nearby objects located on a similar distance, but the system might still be able to initiate pose estimation after recognition.

The following properties of the system have been evaluated, as will be described in more detail in the sections below:

- Combined figure-ground segmentation based on binocular disparities and monocular pose estimation,
- Combined monocular CCH based object recognition and monocular pose estimation,
- Robustness of figure-ground segmentation,
- Robustness towards occlusions using SIFT features,
- Robustness of pose initialisation towards rotations.

For recognition, a set of objects shown in Figure 11 was used. A database was created consisting of object models based on SIFT features and CCHs. Only one view per object was used for the SIFT models, while the CCHs were based on multiple views. Pose estimation was only considered for the first three box-like objects, automatically starting as one of these objects are recognised. For this purpose, the width, height and thickness of these objects were measured and recorded in the database.

Since the observed matching scores did not significantly differ from those already published in [Lowe \(1999\)](#), [Mikolajczyk & Schmid \(2001\)](#) and [S. Ekvall & Kragic \(2003\)](#), we have chosen not to include any additional quantitative results. A few observations have lead us to believe that recognition would benefit from CCHs and SIFT features being used in conjunction. For example, the blue car is rarely recognized properly using SIFT, since the most salient features are due to specularities. However, the distinct colour makes it particularly suitable for CCHs, which on the other hand have a tendency of mixing up the tiger and the giraf, unlike to recognition module based on SIFT features.

A. Binocular Segmentation and Pose Estimation

The first presented experiments illustrate the typical behaviour of the system with binocular disparity based figure-ground segmentation and SIFT based recognition. Results from these experiments can be seen in Figure 5. The first column shows the left foveal camera images prior to the experiments. It is clear that a requested object would be hard to find, without peripheral vision controlling a change in gaze direction. However, from the disparity maps in the second column the system is able to locate a number of object hypotheses, which can be shown as white blobs overlaid on-top of the left peripheral camera image in the third column of the figure.

The matching scores of the recognition module for these two examples were 66% and 70% respectively, measured as the fraction of SIFT features being matched to one particular model. Once an object has been recognised, pose estimation is automatically initiated. This is

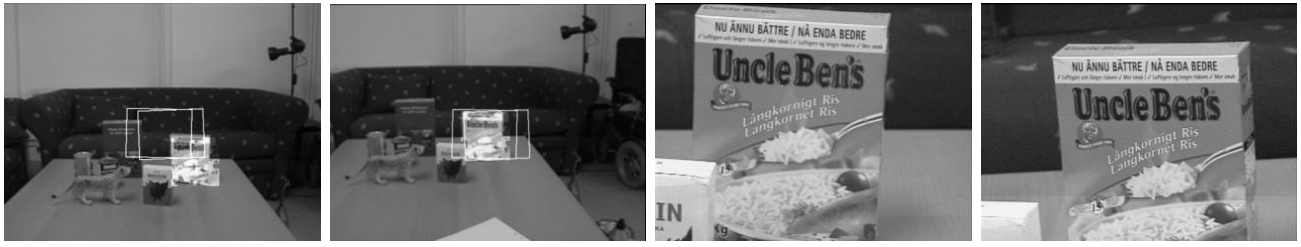


Fig. 10

THE FIRST TWO IMAGES SHOW A TARGET REGION BEFORE AND AFTER A SACCAD (THE RECTANGLES SHOW THE FOVEAL REGIONS WITHIN THE LEFT PERIPHERAL CAMERA IMAGE) AND THE FOVEAL CAMERA IMAGES AFTER EXECUTING A SACCAD ARE SHOWN IN THE LAST TWO IMAGES.

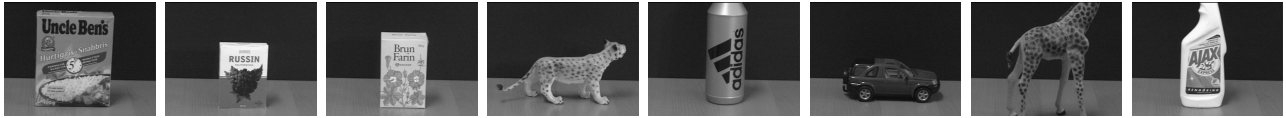


Fig. 11

OBJECTS USED FOR EXPERIMENTAL EVALUATION.

done using SIFT features from the left and right foveal camera images, fitting a plane to the data. Thus, it is assumed that there is a dominating plane that can be mapped to the model. The process is further improved searching for straight edges around this plane. The last two images in the fourth columns show an example of this being done in practice.

B. Monocular CCH Recognition and Pose Estimation

Figure 12 shows two examples of recognition and pose estimation based on monocular CCH. Here, object recognition and rotation estimation serve as the initial values for the model based pose estimation and tracking modules. With the incomplete pose calculated in the recognition (first figure from the left) and orientation estimation step, the initial full pose is estimated (second figure from the left). After that, a local fitting method matches lines in the image with edges of the projected object model. The images obtained after convergence of the tracking scheme is shown on the right. It is important to note, that even under the incorrect initialization of the two other rotation angles as zero, our approach is able to cope with significant deviations from this assumption. This is strongly visible in the second example where the angle around camera's Z-axis is more than 20° .

C. Robustness of disparity based figure-ground segmentation

As mentioned in Section VII, object location hypotheses are found slicing up the disparities into a binary

map of pixels located within a given depth range. There are some evident disadvantages associated with such a procedure. First of all, an object might be tilted and extend beyond this range. This can be seen in the upper left image in Figure 13 - but it does not occur in the second image on the same row. However, since a more accurate localization is found through the focused attention process, a saccade is issued to the approximately same location. This is shown in the last two images on the upper row.

Another challenge occurs if two nearby objects are placed on almost the same distance, especially if the background lacks sufficient texture. Then the objects might merge into a single hypothesis, which is shown on the second row of Figure 13. In our experiments this seemed more common when a global disparity method [Kolmogorov & Zabih \(2001\)](#) was used and is the reason why we normally use simple area correlation. The global optimisation methods tend to fill in the space between the two objects, falsely assuming that rapid changes in disparities are unlikely and thus should be suppressed. In practice, it is preferable if the textureless area between the objects are left unassigned. The right two images on the last row show that pose estimation might still be possible, even when hypotheses are merged.

D. Robustness of SIFT based recognition towards occlusions

In a cluttered environment, a larger fraction of objects are likely to be occluded. These occlusions affect most

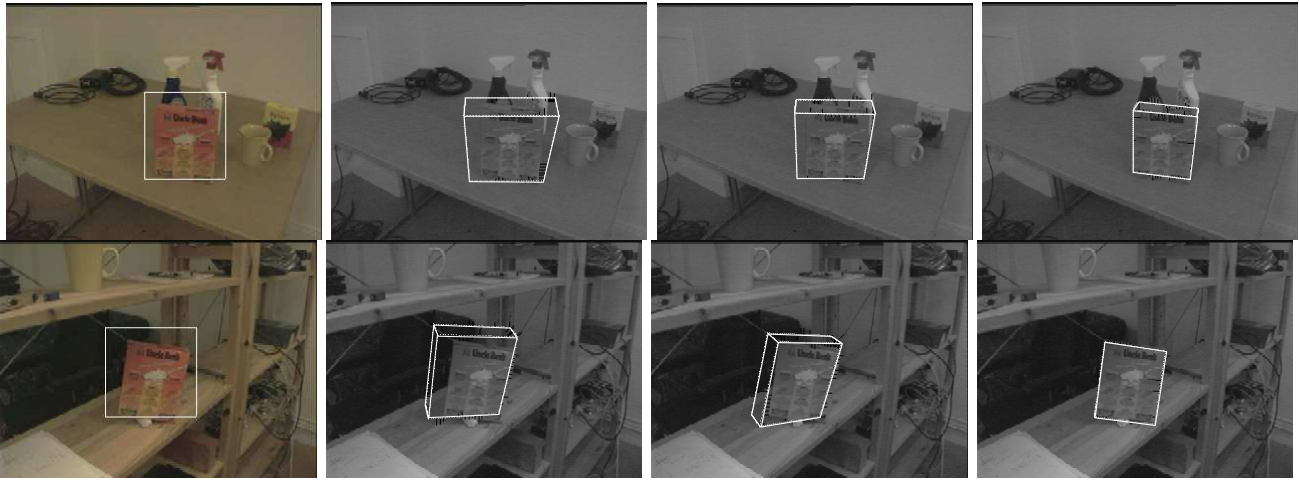


Fig. 12

FROM OBJECT RECOGNITION TO POSE ESTIMATION, (FROM LEFT): I) THE OUTPUT OF THE RECOGNITION, II) INITIAL POSE ESTIMATION, III) AFTER FEW FITTING ITERATIONS, IV) THE ESTIMATED POSE OF THE OBJECT.



Fig. 13

THE EFFECT OF IMPERFECT SEGMENTATION ON OBJECT LOCALISATION.

involved processes, in particular those of recognition and pose estimation. The first two images in Figure 14 show a scene in which the sugar box is partially occluded behind a bottle. In the first case, the recognition fails because not enough foveal features are available, while successful recognition and pose estimation is possible in the second case as shown in the third image. However, even if recognition is successful, the pose initialization might still fail when not enough edges are clearly visible. This can be seen in the last two images of Figure 14. As it is apparent from the fourth image that a failure does not necessarily mean that the results are useless, since the location of the object in 3D space is still available.

E. Robustness of pose initialisation towards rotations

Since, in SIFT based recognition, only one view was available for each object, the sensitivity of the system to rotations was expected to be high. It is already known that for efficient recognition using these features, the relative orientation between query image and object model ought to be less than about 30° . Likely because our model set only consisted of eight objects, our study indicated that slightly larger angles were in fact possible. In the three columns of Figure 15 an object was rotated about 20° , 40° and 60° respectively. The rise package was correctly recognized at a score higher than 70%. However, the break-point turned out to be highly object dependent. For example, for an object like the tiger, the

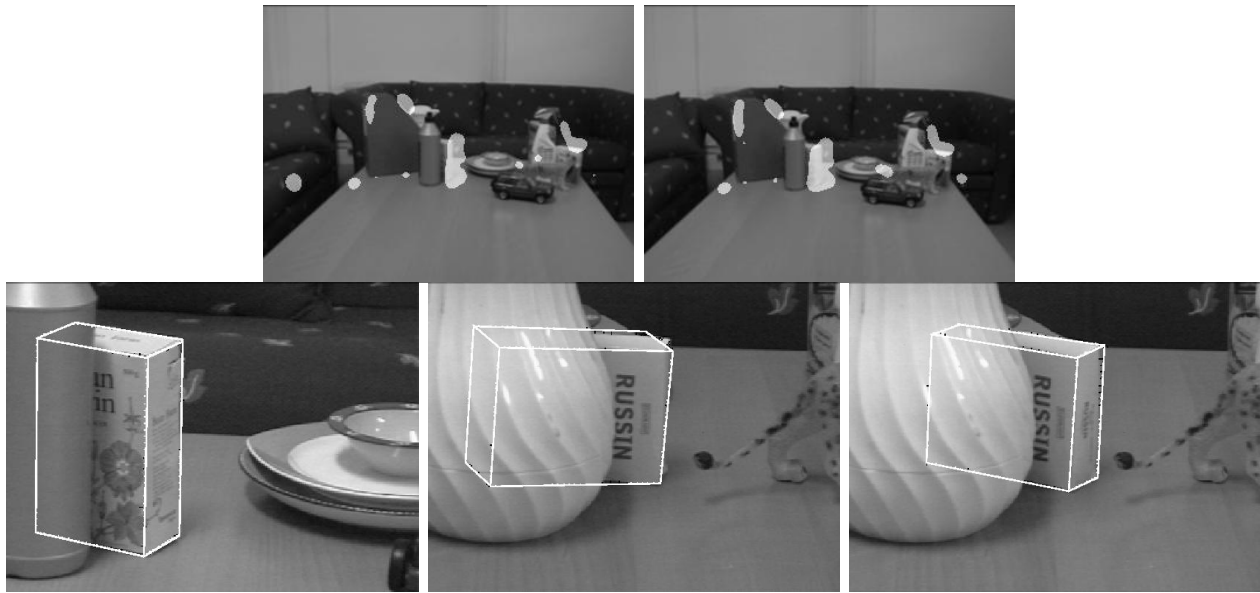


Fig. 14

THE EFFECT OF OCCLUSIONS ON SEGMENTATION AND POSE ESTIMATION.

breakpoint was as low as 20%. For a more thorough analysis on the SIFT recognition performance we refer to Lowe (1999).

As can be seen in the last two images on the upper row of Figure 15, larger rotations tends to be underestimated when the pose is initialised. However, these errors are still below what is required for the pose estimation to finally converge. The lower row shows the estimated pose after a few initial iterations. Even at an angle of 60° the process will converge, but at a somewhat slower rate. For 40° and below convergence is reach within a few frames.

IX. CONCLUSIONS

In this paper, different visual strategies necessary for robotic hand-eye coordination and object grasping tasks, have been presented. The importance of camera placement and their number have been discussed and their effect on the design and choice of visual algorithms. For realistic, domestic settings we are interested in designing robots that are able to manipulate both known and unknown objects and it is therefore important to develop methods for both cases. We have shown what are our current strategies for both of these cases.

Reflecting back to Figure 1, different scenarios can be arranged in a hierarchy depending on prior information. Even if a particular task is given, it is possible to shift between different scenarios and therefore, the underlying strategies used. For example, if the command “Pick Up This Cup” is given, but the system fails to verify the

existence of the cup, the execution may still continue as if “Pick up The Cup” was given. A vice-versa example is if the command “Pick Up This Object” was given and the system realises that the object is, in fact, a known box of raisins. Then, the system automatically changes the task to “Pick Up The Raisins”. In the future, we want to develop a more formal description for the above, in order to design a visual system framework for robotic manipulation in general.

REFERENCES

- Ballard, D. H. 1991. Animate Vision. *Artificial Intelligence*, **48**(1), 57–86.
- Benhimane, S., & Malis, E. 2003. Vision-based control with respect to planar and non-planar objects using a zooming camera. *Pages 991–996 of: IEEE IEEE International Conference on Advanced Robotics*, vol. 2.
- Bjorkman, M. 2002. *Real-Time Motion and Stereo Cues for Active Visual Observers*. Stockholm, Sweden: Doctoral dissertation, Computational Vision and Active Perception Laboratory (CVAP), Royal Inst. of Technology.
- Bjorkman, M., & Kragic, D. 2004. Combination of foveal and peripheral vision for object recognition and pose estimation. *Proceedings. IEEE International Conference on Robotics and Automation, ICRA'04*, **5**(April), 5135 – 5140.
- Björkman, M., & Eklundh, J-O. 2002. Real-Time Epipolar Geometry Estimation of Binocular Stereo



Fig. 15

THE EFFECT OF LARGE ROTATIONS ON POSE INITIALISATION.

- Heads. *IEEE Trans. Pattern Analysis and Machine Intelligence*, **24**(3), 425–432.
- Comaniciu, D., Ramesh, V., & Meer, P. 2000. Real-time Tracking of Non-Rigid Objects Using Mean Shift. *Pages 142–151 of: Proc. IEEE Conf. Computer Vision and Pattern Recognition, CVPR 2000.*
- D. Kragic, A. Miller, & Allen, P. 2001. Real-time tracking meets online grasp planning. *Proceedings. IEEE International Conference on Robotics and Automation, ICRA'01*, **3**(May), 2460 – 2465.
- Fischler, M., & Bolles, R. 1981. Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography. *Pages 381–395 of: Communications of the ACM*, vol. 24.
- Harris, C., & Stephens, M. 1988. A Combined Corner and Edge Detector. *Pages 147–151 of: Proc. Alvey Vision Conference.*
- Huber, P. J. 1981. *Robust Statistics*. John Wiley and Sons.
- Hutchinson, S., Hager, G.D., & Corke, P. 1996. A tutorial on visual servo control. *IEEE Transactions on Robotics and Automation*, **12**(5), 651–670.
- Kolmogorov, V., & Zabih, R. 2001. Computing Visual Correspondence With Occlusions Using Graph Cuts. *Pages II:508–515 of: Proc. IEEE Intl. Conf. Computer Vision.*
- Konolige, K. 1997. Small Vision Systems: Hardware and Implementation. *Pages 203–212 of: Intl. Symp. Robotics Research.*
- Kragic, D., & Christensen, H.I. 2002. Weak models and cue integration for real-time tracking. *Proceedings. IEEE International Conference on Robotics and Automation, ICRA'02*, **3**(May), 3044 – 3049.
- Kragic, D., & Christensen, H.I. 2003a. Confluence of parameters in model based tracking. *Proceedings. IEEE International Conference on Robotics and Automation, ICRA'03*, **3**(September), 3485 – 3490.
- Kragic, D., & Christensen, H.I. 2003b. A Framework for Visual Servoing. *Proceedings of the International Conference on Computer Vision Systems , ICVS 2003*, April, 345 – 354.
- L. Itti, C. Koch, & Niebur, E. 1998. A Model of Saliency-Based Visual Attention for Rapid Scene Analysis. *IEEE Trans. Pattern Analysis and Machine Intelligence*, **20**(11), 1254–1259.
- Longuet-Higgins, H. 1980. The Interpretation of a Moving Retinal Image. *Pages 385–397 of: Philosophical Trans. Royal Society of London, B-208.*
- Longuet-Higgins, H. 1981. A Computer Algorithm For Reconstructing a Scene From Two Projections. *Nature*, 133–135.
- Lowe, D. G. 1999 (Sep.). Object Recognition From Local Scale-Invariant Features. *Pages 1150–1157 of: Proc. IEEE Int'l Conf. Computer Vision (ICCV 99).*
- Malis, E., Chesi, G., & Cipolla, R. 2003. 2 1/2 D Visual Servoing with Respect to Planar Contours Having Complex and Unknown Shapes. *The International*

- Journal of robotics Research*, **22**(10-11), 841–854.
- Mikolajczyk, K., & Schmid, C. 2001 (July). Indexing Based on Scale Invariant Interest Points. *Pages 525–531 of: Proc. IEEE International Conference Computer Vision, ICCV'01.*
- Palmer, S. E. 1999. *Vision Science: Photons to Phenomenology*. Cambridge, MA: MIT Press.
- S. Ekvall, F. Hoffmann, & Kragic, D. 2003. Object Recognition and Pose Estimation for Robotic Manipulation using Color Cooccurrence Histograms. *In: Proc. IEEE/RSJ International Conference Intelligent Robots and Systems, IROS'03.*
- S. Kim, I. Kim, & Kweon, I. 2003 (October). Robust Model-based 3D Object Recognition by Combining Feature Matching with Tracking. *Pages 2123–2128 of: Proc. IEEE International Conference on Robotics and Automation, ICRA'03.*
- Vinoski, S. 1997. CORBA: Integrating Diverse Applications Within Distributed Heterogeneous Environments. *IEEE Communications Magazine*, **14**(2).
- Zillich, M., Roobaert, D., & Eklundh, J O. 2001 (December). A Pure Learning Approach to Background-Invariant Object Recognition using Pedagogical Support Vector Learning. *In: CVPR-2001*. IEEE, Kauai.