

Seeing with Your Hands: a Better Way to Obtain Perception Capabilities with a Personal Robot

Christian Potthast and Gaurav S. Sukhatme

I. INTRODUCTION

For centuries humans have been dreaming of robotic systems assisting them in their lives. Slowly but steadily robots get involved or take over human tasks. Today we can find a variety of different robots, deployed and used in all kinds of different environments for a variety of tasks and purposes. These mechanical helpers are deployed in application sites like the manufacturing industry, search and rescue scenarios or as assistive robots in our daily live, to just name a few. Benefits to utilizing a robotic system which either assists or takes over human tasks ranges from economical benefits over saving human lives to improve human quality of life. In factories robots are used to optimize highly repetitive tasks, increasing production margins. In hazardous environments robots are put to use instead of humans, preventing death or life threatening injuries. In our homes or gardens, we find robots vacuuming or mowing the lawn.

Today's robots however, are still very limited in tasks they can perform. They are either restricted to a very specific or known environment, remote operated, simplified in their task execution or not fully autonomous. These limitations arise from deficits in perception, manipulation and state estimation algorithm, to name only a few reasons why operating a fully autonomous general purpose robot in an unstructured environments is still impossible to realize. Incentives to overcome these limitations and utilize more sophisticated and smarter robotic systems are plenty available.

Surely everyone of us has dreamed of a robotic servant that takes care of all the unloved duties in our household. A robot that could handle such duties when we are at work or pursue our hobbies. Tasks like cleaning, washing the dishes, doing laundry, setting the table and so many more things could be performed by a personal robot butler, Fig. 1. Compared to household robots, which only benefit individuals, robots in healthcare could potentially reward entire societies. With healthcare costs rising for several years robots could help lowering or stabilizing these costs. They could be operated in hospitals, senior homes or other healthcare facilities. Tasks like patient transport within the hospital, food delivery, errands and so on, could be performed by a robot in the future. Furthermore they could assist elderly in their day to day living, like in senior homes or facilitate home care. These two applications and their possibilities are just a small glimpse of the endless possibilities how robots could help.

The Authors are with the Robotic Embedded Systems Laboratory, University of Southern California (USC), Los Angeles, CA 90089, USA {potthast, gaurav}@usc.edu



Fig. 1. PR2 robot in a household environment.

Without any doubt such robots will change the world and our way of living in available for the public.

To make such visions reality, robots have to be able to operate in human environments. They have to be capable of interacting with the environment, designed and built for human comfort. As a consequence such robots have to mimic human task execution, which includes environment or object manipulation, navigation, object recognition, reasoning and so on. Tasks which are natural and trivial for us humans have to be performed and sufficiently executed by a robotic system. For all of these tasks a prerequisite is sensing, the ability to gather environment information. Environment information is important because the quality of the gathered information determines in many cases the success or feasibility of executing a specific task. In general, the more information we have available prior to executing a task the better the execution and performance. How well the robot is able to observe the environment depends mainly on two details: The kind of sensor and the sensor position.

Today the main sensors for visual feedback used on robots, which operate in human environments, are laser scanners, mono- and stereo-cameras, time-of-flight cameras, depth sensors and RGB-D sensors. Each sensor type has different sensing properties, leading to different observations depending on the environment. Choosing the right sensor can be important for a robot to be able to accomplish a task successful. However even more important as selecting the right sensor is to control what the sensor observes. The sensor observation is mainly influenced by two things, the sensor pose and the objects in the environment. Position and orientation of the sensor influences the field of view, therefore controlling what the sensor perceives. Objects on the other hand influence the line-of-sight of the sensor. As a consequence everything which lies not in line-of-sight of the

sensor is obstructed and therefore not visible at this moment. The non-visible portion is generally referred to as occlusion.

In our daily life we usually deal with two types of occlusions, self occlusion and object occlusion. Although both types limit observability of the environment they are fundamentally different. Self occlusion prevents observing non line-of-sight surfaces on an object, due to material properties, but allow a partial view of the object. In contrast object occlusion has a more drastic effect because objects which are within an occluded space could potentially remain entirely unobserved. How much and where occlusion occurs depends on the sensor pose, number of objects as well as object shape. In general as the number of objects in the environment increases, so does occlusion. Consequently an increase in objects simultaneously increases the effort and difficulty of observing the environment. This becomes especially noticeable in household environments, as our living environments tend to be cluttered, filled with all kinds of objects, walls or furnitures.

To fully observe and gather as much information about the environment as possible we need to observe the occluded portion. Observing the occluded portion can be done in two different ways. By either moving and reorient objects to a new location, establishing line-of-sight to the occluded portion as well as to the previous occluded object surface or we can reposition the sensor to a new location. The first approach however is not always suitable or desirable. Some objects are simply too big to move and every object relocation takes time, due to the difficulty of object grasping itself. In general for efficiency it is better to move the sensor to a new location to observe the occluded portion of the environment.

Observing occluded space by positioning the sensor to a new location consists of two problems. New suitable viewing poses need to be found in which the occluded space lies in line-of-sight of the sensor, enabling the robot to observe the space. Finding such new poses is treated as a sensor placement problem, also known as the *Next-Best-View* (NBV) problem. The NBV approach seeks a single additional sensor placement in order to increase the knowledge of the environment. The second problem is more fundamental because the robot has to have the ability to move the sensor to the new desired pose. Sensor placement on the robot as well as motion constraints due to obstacles in the environment oftentimes hinder the robot to move the sensor into the new pose. Hence, in general the set of possible NBV poses is constrained by the robot's abilities.

Many robots expected to operate in household environments are human-like designed, with the main sensors usually placed in some kind of sensor head. As a consequence only a limited set of possible sensor pose movements are available. Thus repositioning the sensor into a new pose is therefore only possible if the entire robot moves to a new position. In kitchen environments this is however most likely not possible as generally kitchen countertops are only accessible from one side. Thus, in many scenarios the robot is unable to fully observe the environment because suitable new sensor poses can not be reached. In Fig. 2 the cluttered tabletop

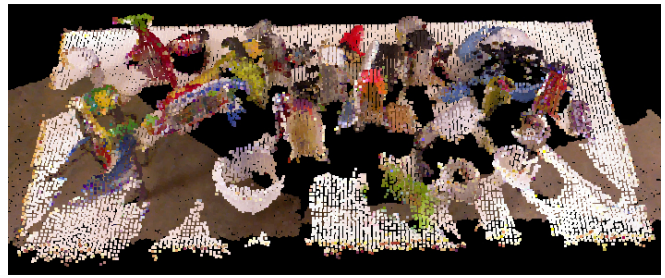


Fig. 2. A tabletop is observed using a head mounted Kinect sensor. When keeping the robot stationary the sensor motion is highly constrained, results in many missing data points. Each missing data point refers to space which is still unobserved.

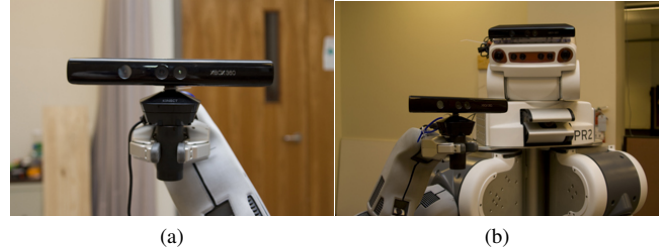


Fig. 3. The PR2 robot is holding a Microsoft Kinect sensor in its manipulator gripper, to increase the possible sensor poses.

environment could only be partially observed due to sensor pose constraints.

To overcome this problem and facilitate the robot with more observation flexibility we utilize the high degree of freedom manipulators, to observe the environment by attaching a sensor to it, shown in Fig. 3. While it is inconceivable for humans to have eyes in our hands, a robot has no biological boundaries, which makes it a purely technical problem. Attaching a sensor to a manipulator has huge advantages since the possible sensor locations are increased by the entire workspace of the manipulator. As a consequence the robot can potentially observe areas which would have been otherwise inaccessible. If needed the robot can use the sensor on the manipulator to increase the knowledge of the scene. However the sensor is not permanently attached since it could potentially interfere with manipulation tasks. Instead the robot picks up the sensor when needed and places it back after the environment is sufficiently observed.

To show conceptual proof and validate our system we used the PR2 robot from Willow Garage. The PR2 robot is a research and development platform operated by the robot operating system (ROS) backed by a large user community contributing algorithms and functionality. The goal of the PR2 is to push robotics into general household environments. To be able to handle these environments the robot is equipped with a mobile base, two high degree of freedom manipulators as well as several high-end sensors. For data acquisition of the environment we are using a Microsoft Kinect sensor, mounted on the sensor head and a second Kinect grasped by the robot gripper.

In the following section we introduce an approach to cope with the difficulties of using an eye-in-hand camera to observe the environment. Finally we illustrate and discuss advantages of this observation technique.

II. APPROACH

Our approach consists of two main phases that are iterated until a scene has been sufficiently observed. In the first phase we seek a pose which increases our environment knowledge maximally. An increase in knowledge is equivalent to reducing the unobserved space in the environment. The pose with the maximum reduction of unknown space is chosen to be the next best view. In the second phase we are optimizing data acquisition by enabling the robot to collect data while the sensor is in motion. This is achieved by planning a trajectory to the next best view pose with intermediate waypoints to acquire additional information. The waypoints are selected such that the path to the NBV pose has maximal reduction in occluded space. In the following we will describe our representation of the environment as well as each phase in more detail.

Occupancy grid

Commonly the data structure of choice representing Kinect sensor measurements is a point cloud. A point cloud is a set of vertices representing points in Cartesian coordinates as well as other information like color, normals, and so on. A point cloud however, has no information stored about the space where no measurements have been returned. Thus there is no distinction in point clouds between space which is free and space which could not be measured due to occlusion or range limitations. A distinction between the two types of space however is necessary, since unknown space arising from occlusion is the space we want to observe. Observing and reducing the unknown space in the environment provides the robot with new information. The unknown space is estimated by using an occupancy grid representation in combination with a ray traversal algorithm. An occupancy grid is a discretized representation of the environment dividing the space in equally sized voxels. Each voxel has a probability distribution over the two possible states free and occupied, where a probability of 0.5 indicates no knowledge about the state of the voxel. Given new measurements the voxels are updated using a sequential Bayesian update formulation. A ray traversal algorithm can efficiently estimate space which is free and space which is not observed yet and therefore unknown. To distinguish between the space, therefore the probability of the voxel, rays are shot into the occupancy grid evaluating whether line-of-sight is established.

Next Best View Estimation

To estimate the pose which maximizes the knowledge increase about the environment the most, we have to solve the NBV problem. In the NBV setting we are given a partially observed scene in which we want to find the pose which maximizes the gain in information by reducing the unknown space. From an information theoretical point of view observing a voxel to be occupied is equally valuable as detecting it as free. Voxel state changes in the occupancy grid are therefore a measure of gain in information about the environment. Intuitively, we aim to find a new viewing pose which has the largest gain in information when compared to

all other possible poses. Finding the next best pose is a non trivial problem since a guess has to be made about space which is unknown and therefore no information is available. This means that we have to predict how the occupancy grid changes if we were to take a snapshot from every available pose. Thus, all poses have to be evaluated by conducting 'virtual snapshots'. The pose with the highest gain in information is the next best view pose.

Formally, given a set $S = \{s_1, s_2, \dots, s_n\}$ of possible viewing positions we want to find the next best viewing pose $\hat{s} \in S$. To estimate \hat{s} we define a function $F(s_n)$ which associates the quantity of how many unknown voxels can be observed from s_n . The function $F(s_n)$ is evaluated by conducting 'virtual snapshots' using a ray traversal algorithm from each pose S . The difficulty in evaluating $F(s_n)$ is to accurately predict the change of the occupancy grid from pose s_n . Simply assuming that every voxel can be seen unless the ray to the voxel under question is blocked by an occupied voxel leads to inaccurate estimations of how the occupancy grid would change. Objects in the unknown space, which have not yet been observed could potentially block the view to unknown voxels and therefore limit the actual gain in information of a viewing position s_n . Obviously this will lead to wrong estimates over the unknown space and the chosen next best pose \hat{s} is not necessarily the optimal choice. The more accurate we can predict the outcome of a potential new data acquisition pose s_n the better our choice \hat{s} . To improve accurate prediction we assign every unknown voxel in the occupancy grid an observation probability when estimating the grid change from a potential new pose. The observation probability states how likely a voxel can be observed from a position s_n . This results in the ability to compute the next best view pose \hat{s} by choosing the pose s_n with the maximum expected number of observable voxels.

The observability probability of any unknown voxel depends on the voxels in the ray shot to an unknown voxel. If the ray contains only likely to be free voxels a high observability probability is assigned to the target unknown voxel. In the case of the ray containing at least one likely to be occupied voxels we assign a very low observability probability to the voxel in question. In both previous cases the observation probability of an unknown voxel can be precisely determined, since partial information about the voxels containing in a ray are available. This is different however if the ray contains only unknown voxels. In this case we have to make a guess about how likely the unknown voxel can be observed. The probability in such a case is computed by using a decaying kernel function, depending on the number of unknown penetrated voxels. The more unknown voxels the ray needs to penetrate, the smaller is the observability probability of that target voxel. Intuitively this expresses the fact that the more unknown space a ray penetrates, the more likely it is that the ray will be blocked by an object which has not yet been seen. The NBV pose \hat{s} is estimated by choosing the pose s_n with the maximum expected number of observed voxels, evaluated by $\mathbb{E}[F(s_n)]$.

Motion Planning

To optimize data acquisition we want to be able to collect data while the sensor is repositioned to a new pose. The ideal solution would be to estimate an entire trajectory which is optimal in terms of gain in information. However, planning and estimating such a trajectory is very difficult. Planning of a trajectory by representing the sensor motion as a number of waypoints simplifies the problem drastically since only suitable waypoints have to be estimated. However the problem with finding favorable waypoints is that a geometrical approach, solving backwards for the sensor pose given an unknown region, is not possible in a cluttered environment. This makes it necessary to randomly sample sensor poses in the high dimensional sensor space and to evaluate the poses for their gain in information. To efficiently select waypoints in a high dimensional space forming a sensor motion from a given start pose to a goal pose we are using a *Probabilistic Roadmap* (PRM) planner approach. A PRM planner is a sampling based method, divided into a planning and query phase, to find a path from start to goal pose.

In the planning phase of the PRM we build a graph with nodes representing the Kinect sensor poses and edges as connections between nodes. The 6D sensor poses $s = \mathbb{R}^6$, with $S = \{s_1, s_2, \dots, s_n\}$ are randomly sampled in the Cartesian manipulator workspace and validated for their feasibility in terms of reachability and joint constraints of the manipulator using inverse kinematics. In an iterative process each newly sampled pose is inserted into the graph by connecting the new node to its k-nearest neighbors, determined with the euclidean distance as metric. Once the graph is constructed we can evaluate each of the sampled poses for their potential gain in information and determine the next best view pose.

In the query phase we want to estimate a path from a given start pose P_s and an estimated goal pose P_g . The start pose P_s is the current Kinect sensor pose and the goal pose P_g is the estimated next best view pose \hat{s} hence $P_g = \hat{s}$. The intermediate waypoints are chosen in a way that the estimated path, a combination of start, end and intermediate waypoints, maximizes the gain in information. In other words we want to find the path that maximizes the number of observable unknown voxels.

To find the best path we have to evaluate all paths which connects the start pose and the goal pose with each other. However for efficiency and more natural manipulator motions we evaluate only a subset of all possible paths for their potential gain in information. The subset $r = \{r_1, r_2, \dots, r_n\}$ is chosen in a way that all paths have to strictly monotonically decrease the goal distance. This reduces the number of possible path drastically and with that the search time. To evaluate the path r we calculate the gain in information of each waypoint in r with respect to the previous waypoint. An evaluated path is therefore the relative combined gain in information of all poses in the path. The best path r_b is defined as the path with the maximal gain in information of all paths in r .

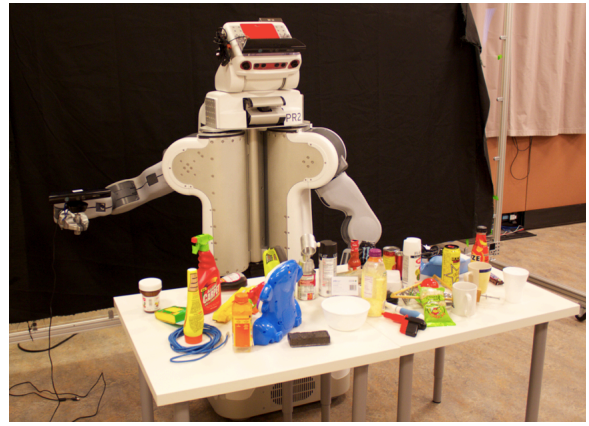


Fig. 4. To test our method we use the PR2 robot and position it in front of a highly cluttered tabletop scene. The goal is to maximally observe the scene using an in hand Kinect sensor.

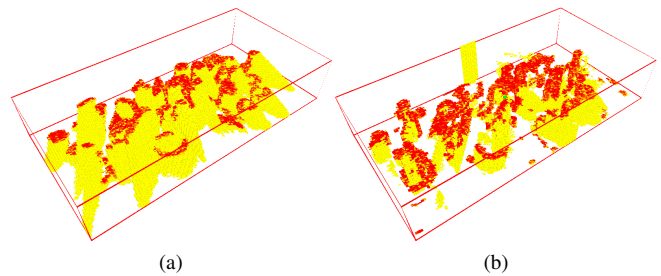


Fig. 5. Illustration of the observed environment represented as an occupancy grid (yellow unknown space, red occupied space) after a number of measurements have been taken. In (a) the robot was able to observe 87% of the environment using only the head mounted Kinect. In (b) using the in hand sensor 97% of the scene is observe.

III. RESULTS

For validation and to visualize the advantage of using an in hand Kinect sensor we are using the PR2 robot to observe a cluttered tabletop environment, as shown in Fig. 4. We assume the robot has to remain stationary, since positions around the table are not accessible to the robot. To show the advantage of an in hand camera we conduct two observation attempts, one with only the head mounted Kinect sensor and the other with the Kinect sensor held in the robots right gripper. The motion of the head mounted Kinect sensor is restricted by a pan/tilt joint, while the in hand sensor can be positioned in the entire manipulator workspace. The result of the observed environment, represented as an occupancy grid, using only the head mounted Kinect can be seen in Fig. 5a. The head mounted Kinect was able to observe 87% of the scene, however substantial parts of the space are still unobserved. In comparison, by using the in hand sensor the robot is able to observe 97% of the environment. We can see that the in hand camera increases observation ability of the robot, which in turn decreases the unobserved space by 10%. Objects previously unobserved and hidden are visible due to the increase in possible sensor positions which is especially critical for manipulation tasks, collision avoidance, object search and so on.