

# Estimating the Relation of Perception and Action During Interaction

Roberto Martín-Martín    Arne Sieverling    Oliver Brock

**Abstract**—We assume that for perception for manipulation the relevant structure lies in the combined  $S \times A \times t$  space, the space of changing sensor signals and actions. We also assume that if we can find the right structure we achieve robustness, resilience to noise and disturbances and generalization. If our hypothesis is right and  $S \times A \times t$  is strongly structured, it should be possible to find this structure for a large class of manipulation tasks (tasks with smooth correlation between actions and changes in the sensor signal). In this work we propose a method to estimate on-line the structure of  $S \times A \times t$  in the form of a visuo-motor Jacobian. We show experimentally how the estimated Jacobian allows to solve perceptual tasks that can not be solved without knowledge of the robots actions. We also show how the Jacobian directly leads to a control law for simple manipulation tasks.

## I. INTRODUCTION

Robots are much more than just passive observers: they can interact with their environment as part of their perceptual process. These interactions reveal new sensor signals containing relevant information, and knowledge about the interactions can be exploited to interpret the sensor signal. These ideas have been explored by a family of methods called *Interactive Perception* (Bohg et al., 2016).

To be able to exploit the interactions as source of information, interactive perception methods assume an a priori specified correlation between actions and changes in the sensor signals (Barragán et al. (2014); van Hoof et al. (2012); Hausman et al. (2015)). This correlation represents the structure of the combined space  $S \times A \times t$  of sensor signals  $S$  and actions  $A$  over time  $t$  that is relevant for the perceptual task.

We assume that if we can find the right structure in  $v$  we achieve robustness and resilience to noise and disturbances. If our hypothesis is right and  $S \times A \times t$  is strongly structured, it should be possible to find this structure for a large group of manipulation tasks by interacting and observing the resulting changes in sensor signal.

In this work we propose to find the relevant structure in this combined space by estimating the (possibly dynamic) correlation of  $A$  and changes in  $S$ . We assume that the relationship between actions and changes in sensor signals is sufficiently smooth to be estimated recursively from pairs of actions and observed changes. We exploit the acquired model that relates  $A$  and changes in  $S$  to address new perceptual

All authors are with the Robotics and Biology Laboratory, Technische Universität Berlin, Germany. We gratefully acknowledge the funding provided by the Alexander von Humboldt foundation and the Federal Ministry of Education and Research (BMBF), by the European Commission (EC, SOMA, H2020-ICT-645599) and the German Research Foundation (DFG, Exploration Challenge, BR 2248/3-1).

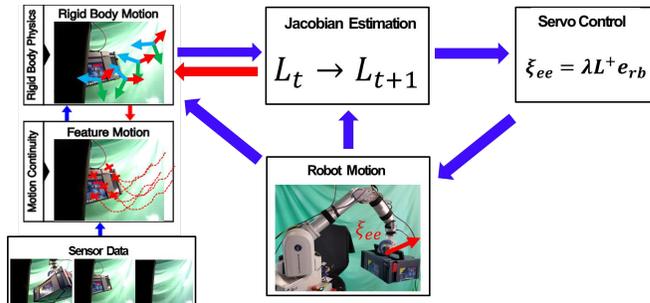


Fig. 1. An overview of the pipeline: square blocks denote recursive estimation loops and arrows the information flow; the robot learns a linear model – a Jacobian – that relates its actions and their consequences in the sensor space (in the form of rigid body motion); the model can be used to improve perception and to control the manipulation

tasks that cannot be solved passively and to achieve tasks defined as goals in  $S$ .

Methodologically, our work extends our previously presented method for *Online Interactive Perception (online-IP)* of articulated objects (Martín-Martín and Brock, 2014). The method is based on interconnected Bayesian recursive estimation processes that infer the kinematic structure of mechanisms in real-time from a continuous RGB-D sensor stream.

We propose to extend this method by integrating and exploiting knowledge about the robot interaction. To do so, the robot needs a *forward model* that relates robot actions to changes in the environment, the structure in the combined  $S \times A \times t$  space. Only then, the robot can exploit the knowledge about the interaction to complete missing sensor information and to make the perceptual process more robust. Because such *interactive forward models* are usually dynamic and task specific we would like our robot to learn them from ongoing interactions. We propose a learning method to estimate the interactive forward model recursively from pairs of actions and their caused changes in the sensor signal.

We demonstrate that the learned model can be applied to improve perception and manipulation in three aspects: a) It allows to predict the motion of controllable degrees of freedom even under occlusions. b) It allows to separate controllable from uncontrollable degrees of freedom in the environment. c) It allows to generate actions that fulfill a given manipulation task robustly.

In the remainder of this report we first present briefly our method to estimate recursively forward models. Then, we present preliminary results that support the improvements in perception and manipulation aforementioned.

## II. BAYESIAN RECURSIVE ESTIMATION OF INTERACTIVE FORWARD MODELS

Learning a complete interactive forward model that can predict the outcome of any robot action, without any prior assumptions, would require an unfeasible amount of data and processing power, due to the high dimensionality of the state space and the variety of tasks. Previous methods to learn interactive forward models are restricted to a specific task and generate enough experimental data using physical models (Barragán et al., 2014) or continuous executions of the task (Lenz et al., 2015; Levine et al., 2016).

In this work, we exploit prior knowledge about the interactive forward model by assuming that it changes smoothly wrt. time and configuration space. This assumption is reasonable for those tasks, where an action of the robot causes a proportional reaction in the environment. This property allows to approximate locally the forward model by a linear function  $L$  that maps changes in the sensor stream to robot motion:  $\dot{s} = L\dot{x}$ .

For visual control tasks this linear model is known as the interaction matrix (Chaumette and Hutchinson, 2006). The simplicity of this model allows for efficient on-line updates. Jägersand et al. (1997) presented an efficient update rule for visual servoing tasks.

Our method estimates the interaction matrix  $L$  recursively solely based on data of robot motion and feature change, during interaction. We formulate this estimation as a Kalman filter, analogous to the other Bayesian recursive processes used in online-IP. We rewrite  $\dot{s} = L\dot{x} = Hl$  with

$$l = (L_{1,1}, L_{1,2}, \dots, L_{1,n}, L_{2,1}, L_{2,2}, \dots, L_{m,n})^T$$

$$H = \begin{bmatrix} \dot{x}^T & 0 \\ & \ddots \\ 0 & \dot{x}^T \end{bmatrix}$$

Now  $l$  is a vector that contains all the elements of the interaction matrix and can be used as the state of a Kalman filter with process model  $l_{k+1} = l_k + Q$  and linear measurement model  $\dot{s}_{k+1} = Hl_k + R$ . Formulating this estimation process as a Kalman filter has the advantage that priors about the uncertainty of the features can be taken into account for estimation.

The previously presented recursive estimation of a linear forward model can be applied to any features defined in sensor space. We use as features the poses of the moving rigid bodies estimated by online-IP. A benefit of using online-IP is that it generates uncertainty bounds for the estimated poses, which we can directly use to balance the recursion in the estimation of the forward model.

## III. EXPERIMENTAL EVALUATION

### A. Online Interactive Perception using Interactive Forward Models

One of our objectives is to exploit knowledge about the actions to improve perception. We extend online-IP to use the estimated interactive forward model to predict the motion of

the rigid bodies based on the robot actions. Originally one of the Bayesian recursive processes in (Martín-Martín and Brock, 2014) estimates the motion of the rigid bodies using as measurement the motion of features in the visual stream obtained with an RGB-D sensor. The pose and velocity of each moving rigid body is estimated using an Extended Kalman Filter (EKF). The velocity of the rigid bodies is used to predict the next pose based on a *passive* forward model that doesn't exploit knowledge about the actions. We improve the prediction step by using the estimated interactive forward model and the actions of the robot to predict the motion of the rigid bodies in the following manner:

$$\hat{p}_{t+1} = p_t + L\dot{x}$$

Using the forward model to predict how its actions affect the environment, the robot is capable of tracking the rigid body more accurately and robustly, even when there is no visual information to correct the prediction, e.g. due to visual occlusions.

*Experiments:* To demonstrate the integration of the estimated forward models into our perceptual system we first estimate the forward model that relates motions of the robot and motions of a grasped object. The interactive forward model is not trivial because some robot motions have no effect on the pose of the object: the robot grasps the object with a *cylindrical grasp* that does not transmit rotations along the main axis of the cylinder (see Fig. 2).

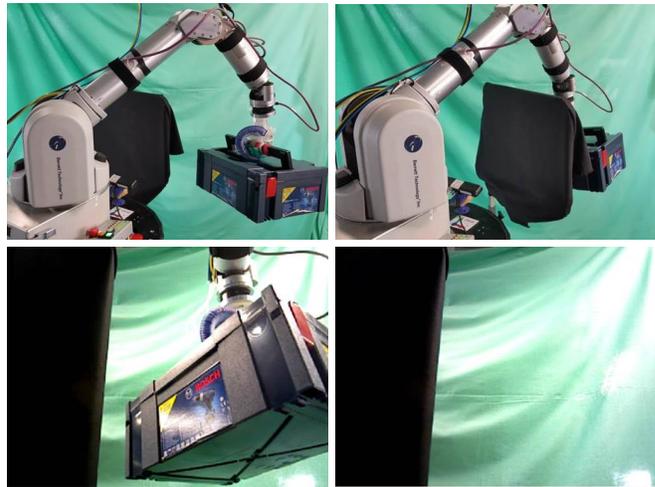


Fig. 2. The robot grasps an object with a cylindrical grasp and manipulates it; rotations along the main axis of the cylindrical grasp are not transmitted to the object; first row: external view of the experiment; second row: robot view from an RGB-D; left column: the robot observes the outcome of its actions and learns an interactive forward model; right column: the object is occluded and its pose is predicted using the forward model

After learning the interactive forward model the robot moves the object to a region where it becomes visually occluded. We compare the two forward models that predict the motion of the object, the passive forward model that uses the (last) estimated object velocity, and the interactive forward model that relates actions of the robot into changes in the environment. The results are depicted in Fig. 3.

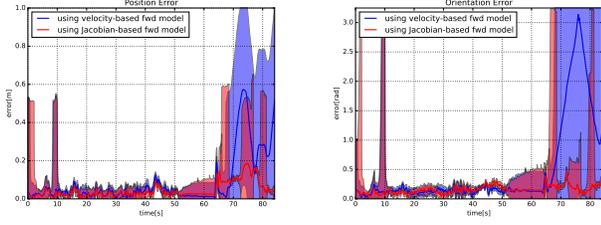


Fig. 3. Position and orientation error of the object pose using the interactive and the passive forward models; the passive forward model (blue) uses the velocity of the object to predict its pose; the interactive forward model (red) uses the on-line estimated model to predict the object pose based on the robot’s actions

After approximately 65 seconds the robot moves the object to a region where it becomes visually occluded and predicts the object pose using the forward models. The pose predicted using the forward model based on constant velocity starts to drift as soon as the object changes its velocity. Using the interactive forward model, the robot can predict the object pose correctly, but the pose uncertainty increases due to the lack of visual feedback. Approximately 6 seconds later the object reappears in the visual field. Using the interactive forward model the robot re-identifies the object and assign correctly visual features, which allows to reduce the error and the uncertainty about the pose. The process using the passive forward model cannot identify the object and continues drifting. The occlusion-appearance process is repeated two more times.

### B. Identifying Controllable Objects

For each rigid body in the scene, we estimate an interactive forward model  $L^{rb^i}$  relating the motion of object  $i$  to the motion of the robot end-effector. We express both the rigid body velocity twist  $\xi_{rb^i}^f$  and the end-effector rigid body twist  $\xi_{ee}^f$  in the same coordinate system  $f$ . Therefore  $L^{rb^i} = \frac{\partial \xi_{rb^i}^f}{\partial \xi_{ee}^f}$ .



Fig. 4. Left: external view of the experiment; right: robot view from an RGB-D at different steps of the experiment; two objects move - one controlled by the robot and the other controlled by a human

For a fully controllable object and a perfectly calibrated system,  $L^{rb^i}$  should converge to the identity matrix. Therefore we define the controllability of the box as the similarity to the unit matrix, which we measure as the maximum element of  $|L^{rb^i} - I|$ .

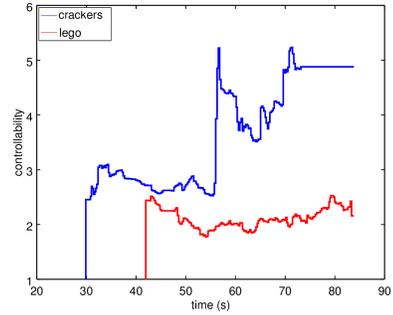


Fig. 5. Controllability of the two objects during 30 seconds of motion; the red cracker box (red line) has a lower controllability than the Lego box (blue line) throughout the whole experiment

*Experiments:* To demonstrate that we can correctly classify controllable and uncontrollable objects, we gather data of the robot interacting with an object in the presence of moving distractors. We place two objects in the world and move them. One objects is moved by the human experimenter and its motion is uncorrelated to robot actions. The other object is grasped by the robot with a suction gripper and moved (see Fig. 4).

Fig. 5 shows how the controllability for both boxes changes over time during the interaction. The red box has consistently lower value and thus can clearly be identified as the object under robot’s control.

### C. Visual Servoing

Estimating an interactive forward model allows to formulate a servo control law that purposefully controls the motion of the manipulated object without prior knowledge about the contact and the camera configuration. Similar to classical visual servoing control, we use a pseudo-inverse of the estimated interactive forward model  $\hat{L}^+$  and a proportional gain  $\lambda$ , to define a control velocity twist  $\xi_{cmd}$  that minimizes the error  $e_\xi$  between the objects current and desired goal pose:

$$\xi_{cmd} = -\lambda \hat{L}^+ e_\xi$$

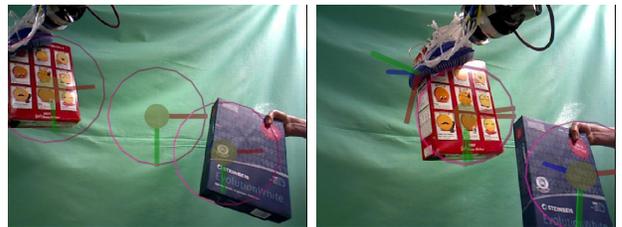


Fig. 6. Robot view of two objects moving, one controlled by itself the other controlled by an experimenter; left: the robot identifies the controllable object, estimates the interactive forward model and uses it to bring the controllable object to the goal (circle at the center of the image); right: the pose of the controllable object converges to the goal

*Experiments:* Two objects move in front of the robot, one grasped by the robot and the other moved by an experimenter (Fig. 6). In the initial phase, the robot moves randomly the

object and observes the motion of both object to identify the controllable object and estimate its interactive forward model. In the second phase (approximately after 20 seconds) the object is confident enough to use the interactive forward model to move the controllable object to the desired goal pose.

#### REFERENCES

- J. Bohg, K. Hausman, B. Sankaran, O. Brock, D. Kragic, S. Schaal, and G. S. Sukhatme, “Interactive perception: Leveraging action in perception and perception in action,” *CoRR*, vol. abs/1604.03670, 2016. [Online]. Available: <http://arxiv.org/abs/1604.03670>
- P. R. Barragán, L. P. Kaelbling, and T. Lozano-Prez, “Interactive Bayesian Identification of Kinematic Mechanisms,” in *International Conference on Robotics and Automation*, 2014, pp. 2013–2020.
- H. van Hoof, O. Kroemer, H. Ben Amor, and J. Peters, “Maximally informative interaction learning for scene exploration,” in *International Conference on Intelligent Robots and Systems*, 2012, pp. 5152–5158.
- K. Hausman, S. Niekum, S. Osentoski, and G. S. Sukhatme, “Active Articulation Model Estimation through Interactive Perception,” in *International Conference on Robotics and Automation*, 2015.
- R. Martín-Martín and O. Brock, “Online interactive perception of articulated objects with multi-level recursive estimation based on task-specific priors,” in *Intelligent Robots and Systems, IEEE/RSJ International Conference on*. IEEE, 2014, pp. 2494–2501.
- I. Lenz, R. Knepper, and A. Saxena, “Deepmpc: Learning deep latent features for model predictive control,” in *Robotics Science and Systems (RSS)*, 2015.
- S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *Journal of Machine Learning Research*, vol. 17, no. 39, pp. 1–40, 2016.
- F. Chaumette and S. Hutchinson, “Visual servo control. I. basic approaches,” *Robotics & Automation Magazine, IEEE*, vol. 13, no. 4, pp. 82–90, 2006.
- M. Jägersand, O. Fuentes, and R. Nelson, “Experimental evaluation of uncalibrated visual servoing for precision manipulation,” in *Robotics and Automation, 1997 IEEE International Conference on*, vol. 4. IEEE, 1997, pp. 2874–2880.