



Bekiroglu, Y., Damianou, A., Detry, R., Stork, J. A., Kragic, D., & Ek, C. H. (2016). Probabilistic consolidation of grasp experience. In 2016 IEEE International Conference on Robotics and Automation (ICRA 2016): Proceedings of a meeting held 16-21 May 2016, Stockholm, Sweden. (pp. 193-200). Institute of Electrical and Electronics Engineers. DOI: 10.1109/ICRA.2016.7487133

Peer reviewed version

Link to published version (if available):
[10.1109/ICRA.2016.7487133](https://doi.org/10.1109/ICRA.2016.7487133)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via IEEE at <http://ieeexplore.ieee.org/document/7487133/>. Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/pure/about/ebr-terms.html>

Probabilistic Consolidation of Grasp Experience

Yasemin Bekiroglu, Andreas Damianou, Renaud Detry, Johannes A. Stork, Danica Kragic and Carl Henrik Ek

Abstract—We present a probabilistic model for joint representation of several sensory modalities and action parameters in a robotic grasping scenario. Our non-linear probabilistic latent variable model encodes relationships between grasp-related parameters, learns the importance of features, and expresses confidence in estimates. The model learns associations between stable and unstable grasps that it experiences during an exploration phase. We demonstrate the applicability of the model for estimating grasp stability, correcting grasps, identifying objects based on tactile imprints and predicting tactile imprints from object-relative gripper poses. We performed experiments on a real platform with both known and novel objects, i.e., objects the robot trained with, and previously unseen objects. Grasp correction had a 75% success rate on known objects, and 73% on new objects. We compared our model to a traditional regression model that succeeded in correcting grasps in only 38% of cases.

I. INTRODUCTION

Grasping is a key building block of autonomous robots and as a result it has received much attention in the last three decades [1], [2], [3], [4]. Different approaches have been studied, e.g., analytic [2] and data-driven [4]. Moreover, different subproblems have been addressed, e.g., grasp planning [5], force control [6], stability estimation from sensory data after grasp execution [7] or grasp adaptation [8]. However, current robotic systems still have severe limitations in dealing with novelty, uncertainty and unforeseen situations. Limitations arise from multiple sources: noisy and incomplete perceptual data, insufficient experience and high dimensionality of the problem involving variables with complex relations. Problems such as *selecting* the relevant information from the environment, *merging* different sources of information to reduce uncertainty and, making use of experience (even from failures) by relating sensor data to *previous knowledge* remain open. To our knowledge, no model to date addresses these three issues in a principled manner. Addressing these three issues jointly is the aim of this paper.

Here, we focus on grasping, in a scenario with multiple sensory modalities, and we investigate a learning approach to encode grasping knowledge acquired from experience. We aim to provide robots with means of reasoning about object grasps and their probability of success, taking into account the information provided by complementary sensory

channels. We consider the integration of both visual and haptic cues, as they contribute substantially to grasp control [9].

We learn which sensory modalities contain information that correlates with other modalities or with robot action parameters, and learn the structure of these relationships. We study how to combine multiple sources of information, how previous experiences can be related to one another, and how to provide information about future events. Previous experiences come in the form of successful and unsuccessful grasp examples, represented by object-relative gripper poses, tactile information, object types, and how unsuccessful grasps have previously been corrected to become stable. Casting this as a representation problem, we consider each of the sensory modalities and the action parameters as projections (or *views*) of a single latent variable.

The proposed model learns a single factorized latent variable of all the views, a process which we will refer to as *consolidation*. This allows to pose questions such as: “*How and what portion of the tactile data is determined by the gripper position?*” or “*What do I know about the object type given the tactile data and the gripper position?*”. This means that the factorization provides a means to acquire knowledge of what the *relevant* information is, and have the facilities of *merging* several disparate views such as tactile and gripper position. Further, as we have both successful and unsuccessful variations we can exploit this and transfer or correct a grasp within the model, i.e. use *previous knowledge* about what differentiate a successful and a failed grasp.

II. RELATED WORK

Various approaches for avoiding or recovering unsuitable or potentially failing grasps have been proposed in the literature. In [10] the authors proposed to correct grasps by adapting to local geometry using the force-closure criterion. Contact positions were transferred between objects of the same functional class by surface geometry warping. Grasps were adapted by moving finger contacts onto the object’s surface to reach force-closure, or reject the grasp. Compared to our work, this system did not integrate experience from training data or feedback from grasp execution. Differently from [10], [11] and [8] included an off-line training phase based on examples demonstrated by a teacher. In our work, training also relies on human demonstration. The teacher shows the robot a set of grasps, and the robot autonomously explores more grasps in the neighborhood of those demonstrated by the teacher. The learning process is thus data-driven, based on self-exploration without human intervention. [11] used a programming by demonstration approach

Y. Bekiroglu is with the School of Mechanical Engineering, University of Birmingham, UK. A. Damianou is with the Department of Computer Science, University of Sheffield, UK. R. Detry, J. A. Stork, D. Kragic and C. H. Ek are with the Centre for Autonomous Systems, CSC, KTH, Sweden. R. Detry is also with University of Liège, Belgium. C. H. Ek is also with the Department of Computer Science, University of Bristol, UK. Email: Y.Bekiroglu@bham.ac.uk, andreas.damianou@sheffield.ac.uk, renaud.detry@ulg.ac.be, {jastork,dani}@kth.se, carlhenrik.ek@bristol.ac.uk

where a robot relies on human-to-robot grasp mapping to learn most likely hand preshapes for specific objects. Starting from the inferred preshape, a force controller that relied on joint angles and tactile sensing [12] was selected to handle position and orientation uncertainty. A grasp was chosen via control laws and experience was only used to select the hand preshapes for the given object. In our work, we do not consider the control paradigm but we correct grasps based on prior experience. The work of [12] did not consider the success of a grasp before executing the controller.

In [13], grasp corrections were synthesized by matching to a database of stable grasps based on similarity in tactile measurements. If a match similar enough to the current tactile measurements was found in the database, the current grasp was adjusted accordingly. However, an unsuccessful look-up initiated tactile-based reconstruction of local surface geometry and re-planning to adapt the grasps to the actual local object shape. In contrast to our work, the assumption was that the recorded stable grasp that resulted in the most similar tactile reading was the best correction of the current grasp. As a statistical modeling of grasp correction was not employed, novel grasps were not synthesised due to lacking a continuous mapping within a probabilistic framework. A method for grasp adaptation by learning a statistical model to adapt the hand posture based on perceived contacts were presented in [8]. Kinesthetic demonstration learning was used to train a Gaussian mixture model (GMM) for prediction of desired joint values and finger pressure from contact signatures. For this a human teacher improved robot grasps while the robot generated a database of poses and contacts. As in our work, real robot data was included in the training process. However, we explicitly include examples of failed grasps and therefore are able to infer what is shared between good and bad grasps. This allows our model not to just perform a stable grasp but do so in a manner which is conditioned on the unsuccessful grasp.

A recent GMM-based grasp adaptation approach was introduced in [14]. Based on an object-level impedance controller, a grasp stability estimator was first learned in the object frame. Once a grasp was predicted to be unstable by the stability estimator, a grasp adaptation strategy was triggered according to the similarity between the new grasp and training examples. However learning was achieved with positive data only. Compared to this work, to correct a grasp we utilize both positive and negative training samples, we select relevant features to alter the relevant factors and also model uncertainty in our estimate in a principled way.

Differently from all the approaches discussed above, our method involves learning a representation capable of parameterizing, the pose and tactile sensors associated with grasps applied to several different objects. Further, the representation parametrizes grasps that are deemed both good and unsuccessful. We argue that there is a certain underlying characteristic of a good grasp that generalizes across different objects, e.g., grasping sufficiently closer to the center of mass. By learning a representation which reveals these structures we aim to “correct” grasps to adhere to the relevant

factors while factors irrelevant to the success of the grasp can be retained. Further, as the representation is shared between the pose and the tactile domain we can hallucinate how a specific grasp will “feel” in terms of tactile feedback given a specific pose and vice versa.

III. GRASPING

We learn visual and tactile characteristics of grasps which allows to answer several important grasp-related questions. The overall system can be summarized in three steps: exploration, training and inference, as illustrated in Figure 1. During the exploration step, the robot gathers visual and tactile data from both successful and unsuccessful grasping trials. After training our model using the extracted data from the exploration step, through the inference step we can make predictions on target variables given any subset of the observed variables. In our experimental evaluation through inference using the trained model, we demonstrate three applications: correct faulty grasps, predict the shape of hand-object contacts, or recognize objects. In our analysis we focus on mainly grasp correction, i.e., what is a better grasping pose given an unstable grasp.

This section introduces the sensory modalities we consider, describes the extracted training data seen in the table in Figure 1 and explains how we merge the information provided by multiple sensors and in order to perform the inference based applications.

A. Robot Sensing: Tactile sensing and vision are two information sources that are relevant to grasping. Tactile sensing allows the robot to measure the shape of contacts between its hand/fingers and an object. Contact shapes are directly linked to grasp stability, as they relate to the net force applicable by the gripper onto the object. Contact areas also characterize the stability of a grasp by relating to the magnitude of friction forces that exist between the hand and the object. Complementarily, vision provides the robot with information about, for instance, the position, shape or identity of an object. These data are useful to put contact information into perspective – a specific contact reading may relate to a different grasp outcome depending on the size of an object, for instance.

Our experimental robot platform consists of a Kuka articulated arm equipped with a Schunk SDH hand (Figure 1). The platform implements the two sensor channels listed above. Tactile sensing is provided by capacitive pads attached to the inner faces of three fingers’ distal links. Each tactile pad measures pressure along a grid of 6×13 pressure cells, yielding a total of 234 pressure cells. Vision is provided by an RGB camera pointed towards the robot’s workspace.

Learning a visual model from raw pixel data is often prohibitive, because of the high dimensionality of an image. We reduce the dimensionality of vision data by extracting the pose of the target object. Estimating the pose of an object during a grasp is a challenging task, as the object is often largely occluded by the gripper. To overcome this problem, we continuously track the pose of known objects [15]. In this way, by the time the gripper reaches an object, an estimate of

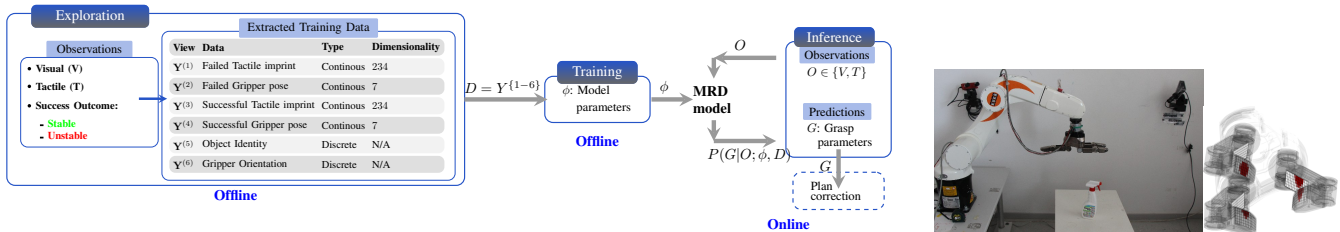


Fig. 1. The system overview (left): Our approach involves exploration, training and inference steps in order to answer grasp-related questions. During grasping trials, the robot gathers visual and tactile observations as well as success outcomes of each grasp, i.e., whether or not lifting leads to slippage or rotation. The extracted data, D , from these observations are used to train Manifold Relevance Determination (MRD) model (see Section IV) and obtain model parameters, ϕ . The MRD model is then used to infer grasp parameters, $P(G|O; \phi, D)$, given a set of observations, O , based on the obtained model with the aim of using those parameters for plan corrections. The characteristics of each view are seen on the right. Our robot platform (right): composed of a Kuka industrial arm, a three-finger Schunk Dextrous hand equipped with tactile sensing arrays (on the right side), and a monocular camera.

the object’s pose is already available, and the robot can track the object’s movement throughout the grasp even if only a fraction of the object’s surface is visible. We note that the tracker that we use in this work requires mesh and texture models of all objects, a limitation that we intend to remove in future work.

We parametrize grasps with the pose of the robot’s hand. A grasp is executed by bringing the hand to the given pose while fully open, and closing the fingers until contact.

We know that a grasp is not dependent on the positioning of an object. To make our framework invariant to object position, we instead model object-relative gripper poses (or, equivalently, gripper-relative object poses). Yet, considering only the relative pose of the object and the gripper leaves out an important bit of information: grasp stability does not only depend on the relative object-gripper configuration, but also on the orientation of the gripper. When an elongated object lies on a flat surface and the robot attempts a grasp in a top-down fashion, it is generally better to aim for the center of mass of the object. By contrast, if the object is standing and the robot attempts a sideways grasp, grasping near the object’s tip is just as acceptable as grasping near its center of mass. To account for this observation, we model both the pose of the gripper with respect to the object and the orientation of the gripper in world frame. To limit the dimensionality of the data the gripper orientation is discretized into either top-down grasp, or sideways grasp.

We wish to perform grasp correction in a manner such that the corrected grasp reflects the characteristics of the unstable. Our model allows for doing just this as it learns a single representation which allows for conditional transfer between the views. To that end, we consider all the different types of information that we have as views: tactile data, object-relative gripper pose, object identity, and gripper orientation. Furthermore, in order to allow for correcting gripper poses from tactile data, and for predicting successful and unsuccessful tactile feedback given a gripper pose, we separate the tactile data and object-relative gripper poses into separate views based on the success of the corresponding grasp. In other words, our model is composed of the following views which are also listed in Figure 1:

- Tactile imprints characterizing *unsuccessful* grasps. As explained above, tactile imprints are parametrized by a vector of 234 measurements.

- Object-relative gripper poses characterizing *unsuccessful* grasps, parametrized by a translation (3-vector) and a rotation (quaternion), seven values in total.
- Tactile imprints characterizing *successful* grasps, also parametrized by a vector of 234 elements.
- Object-relative gripper poses characterizing *successful* grasps, parametrized by seven values also.
- Object identity, consisting in a discrete object label.
- Gripper orientation, either top-down or sideways.

We note that object identity and gripper orientation are not separated into successful/unsuccessful views because the ability to correct object identity or switch from top-down grasp to sideways grasp, is not useful for our purpose.

We will now proceed to show how the problem of encoding grasping knowledge can be cast as a multi-view representation learning task. In specific, observing several modalities associated with grasping we learn a single representation that consolidates these disparate sources of information. The key characteristic of our approach is that the latent variable has a specific factorized structure which is essential to perform grasp correction.

IV. MODEL

In this paper we apply a model called Manifold Relevance Determination (MRD) [16] which is a Bayesian formulation of the multi-view Gaussian Process Latent Variable Model (GP-LVM) [17]. The MRD learns a single latent variable consolidating a set of views. Each view consists of general vector valued observations and in our specific application they are the hand pose and the tactile sensing for both the successful and the unsuccessful grasp, the object orientation and the object type, resulting in six different views in total. In Fig. 3 the characteristics of each view is given. The 6 modalities $\mathcal{Y} = \{\mathbf{Y}^{(k)}\}_{k=1}^6$ are aligned such that $\mathbf{y}_n^{(i)} \in \mathbf{Y}^{(i)}$ and $\mathbf{y}_n^{(j)} \in \mathbf{Y}^{(j)}$ are considered corresponding, from which the model learns a single latent representation $\mathbf{X} \in \mathbb{R}^d$. In Section V we describe the details of how this alignment is acquired directly from data in an unsupervised manner.

In the MRD model each observation $\mathbf{y}_n^{(i)}$ is seen as generated from a latent variable \mathbf{x}_n through a mapping $f^{(i)}$ with additive gaussian noise $\mathbf{y}_n^{(i)} = f^{(i)}(\mathbf{x}_n) + \epsilon$ where ϵ is normally distributed with a spherical covariance. By assuming that each view is independent given the latent space and placing a \mathcal{GP} -prior over the mappings we can formulate

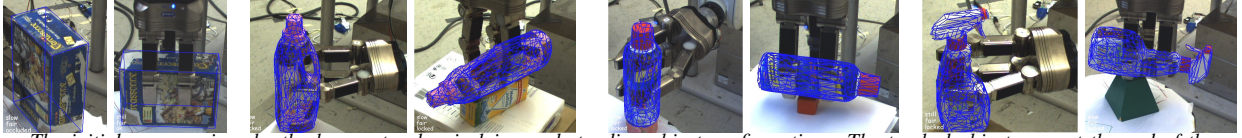


Fig. 2. The initial grasps, given by the human teacher in lying and standing object configurations: The tracked object pose at the end of the grasp is depicted with the blue wire-frame.

the marginal likelihood of the model,

$$p(\mathcal{Y}|\mathbf{X}, \theta) = \prod_{k=1}^K \int p(\mathbf{Y}^{(k)}|\mathbf{F}^{(k)})p(\mathbf{F}^{(k)}|\mathbf{X}, \theta)d\mathbf{F}^{(k)}, \quad (1)$$

where $\mathbf{F}^{(i)}$ is the realization of the mapping and θ are hyper-parameters defining the form of the prior. Different shared GP-LVM approaches when learning \mathbf{X} from Eq. 1 add different types of constraints on \mathbf{X} . What makes the MRD different from other shared GP-LVM models is that it allows for feature selection on the latent space. This means that when generating view (i) the mapping $f^{(i)}$ can select the dimensions of \mathbf{X} which are relevant for encoding the variations in $\mathbf{Y}^{(i)}$. This is referred to as a *factorized* latent representation where a latent dimension can be responsible for encoding any combination of views. We will refer to a dimension which generates several views as shared between those views and a dimension that is only responsible for a specific view to be private to that view. This feature selection is implemented by using an Automatic Relevance Determination (ARD) [18] prior for the mappings by associating each view with a weight vector $\mathbf{W}^{(i)} \in \mathbb{R}^d$ such that $w_m^{(i)}$ determines the relevance of dimension m for generating view i , e.g., if $w_m^{(i)} = 0$ dimension m is independent of $\mathbf{Y}^{(i)}$. We say that the set of weight vectors $\mathcal{W} = \{\mathbf{W}^{(k)}\}_{k=1}^K$ factorizes the latent space. The MRD model is shown in Figure 3.

The ARD prior introduces several new parameters into the model. This makes the model much less constrained making training challenging. This is addressed by variational approach to approximately marginalize out the latent space from the model. This means that both the ARD weights and the latent space can be learned from data, thus providing a natural factorization.

A. Gaussian Process Predictions: A trained model implies that we have learned the latent representation \mathbf{X} , the hyper-parameters defining the characteristics of each generating mapping $\theta^{(i)}$ and the weights $\mathbf{w}^{(i)}$ which select the relevant generating parameters of \mathbf{X} for each view. Prediction of view (i) from a previously unknown latent point \mathbf{x}_* can be made by conditioning on the training data $p(\mathbf{y}_*^{(i)}|\mathbf{x}_*, \mathbf{Y}^{(i)}, \mathbf{X}, \theta^{(i)})$. The \mathcal{GP} specifies that all instantiations of the function is jointly Gaussian. As a Gaussian is self-conjugate this means that also this conditional distribution is a Gaussian. This leads to the predictive equations of the mean and the variance,

$$\begin{aligned} \mu(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}\mathbf{Y}^{(i)} \\ \sigma^2(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T k(\mathbf{X}, \mathbf{X})^{-1}k(\mathbf{x}_*, \mathbf{X}), \end{aligned} \quad (2)$$

where $k(\cdot, \cdot)$ is the kernel or covariance function of the process, in the experimental section we use a linear combination

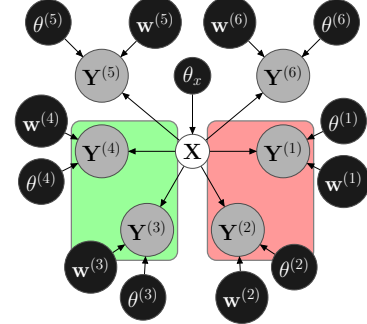


Fig. 3. Manifold Relevance Determination model for grasping.

of a radial basis and a white noise kernel. To present the quantitative results, we report the mean prediction as a point estimate and use the variance as an indication of the certainty in the prediction.

B. Multi-view Inference for Grasping: View Transfer:

Once we have learned the model we wish to perform inference by transferring information between one (or several) views. For instance, given the gripper pose of an unstable grasp $\mathbf{y}_*^{(2)}$ we can infer a gripper pose which would correspond to a stable execution, i.e. $\mathbf{y}_*^{(4)}$. A trained model in combination with the predictive equations in Eq 2 allows to generate new points in any of the views $\mathbf{y}_n^{(i)}$ given the corresponding latent location \mathbf{x}_n . Our aim in this paper is to *transfer* the information from one view to another in a conditional manner. Observing an instance of view 2, $\mathbf{y}_*^{(2)}$ we want to use the model to alter the *relevant* parameters in order to make it match the corresponding view 4 i.e., we want to infer or transfer one modality from another. This is done using a three-step process.

In the first step, we determine the latent representation \mathbf{x}_* corresponding to $\mathbf{y}_*^{(2)}$. This is done by using an approximation of the true posterior $p(\mathbf{x}_*|\mathbf{y}_*^{(2)}, \mathbf{Y}^{(2)})$. This posterior is not informative of the dimensions of \mathbf{X} which are independent of $\mathbf{Y}^{(2)}$. Due to the factorization of \mathbf{X} , we will only be able to determine the latent location of the relevant dimensions for view 2, i.e. the dimensions of \mathbf{X} which have non-zero ARD weight in $\mathbf{w}^{(2)}$. In order to determine the projection onto view 4, we need to also find the dimensions, if any, that are relevant for view 4 but irrelevant for view 2 i.e. the non-shared dimensions which are the ones where $\mathbf{w}_m^{(2)} = 0$ and $\mathbf{w}_m^{(4)} \neq 0$. We determine the remaining dimensions of \mathbf{x}_* using a nearest neighbour approach and can then use the generative mapping to map to the corresponding gripper position $\mathbf{y}_*^{(4)}$. We perform a nearest neighbor search in the subspace which is shared between the views to recover the closest point in the training data to \mathbf{x}_* . In the last step, we replace the non-shared dimensions of \mathbf{x}_* with those of the nearest neighbor. This is the same process as

described in the work of [16]. Importantly, this means that we let view 2 constrain view 4 by the factors which they share and then alter, through the nearest neighbor search, the factors required to make it match view 4. The same procedure can be performed to transfer information between any set of views. If we are given observations from more than one view simultaneously we can find \mathbf{x}_* through the joint posterior. As different modalities are likely to provide different information this should increase the number of latent dimensions which constrains the output further. As an example if we want to infer the stable gripper pose $\mathbf{y}_*^{(4)}$ from the stable tactile $\mathbf{y}_*^{(3)}$ and the object identity $\mathbf{y}_*^{(5)}$, we just alter the first step of the inference and use the corresponding posterior $p(\mathbf{x}_* | \mathbf{y}_*^{(3)}, \mathbf{Y}^{(3)}, \mathbf{y}_*^{(5)}, \mathbf{Y}^{(5)})$ to determine the observation dependent latent location. In this manner we can perform a large range of different inference problems within the same model.

Importantly through this procedure we have implemented the three aspects of reasoning we aimed for; by factorizing the data we only need to alter the factors that are *relevant* to make the grasp stable, by using several views we can constrain the prediction further implementing a natural approach of *merging* different information and finally due to the alignment of the data we can exploit *previous knowledge*.

V. EXPERIMENTS

In order to acquire training data we let our robot explore grasping configurations on four different objects. We used the four home-environment objects shown in Figure 5 d: a box, and three bottles of different shapes that we call the *spray bottle*, the *cylindrical bottle* and the *oval bottle*. We chose these objects because of several reasons: They were large enough for the Schunk hand to grasp with three fingers, similar in weight and had different geometrical features and deformation properties. For example, the spray bottle had a less regular shape compared to the other objects, the cylindrical bottle was the most deformable object and the box was the least deformable among all the objects.

Training the MRD model requires examples of successful and unsuccessful grasps. To maximize the accuracy of the model, a dense sampling of the space of gripper-object poses is desirable, as it will lead to a finer tactile model, and a greater ability to associate stable grasps to the unstable ones, and thus to correct gripper poses. For the purpose of this paper, we opted for a dense exploration of the space of object-gripper poses around pre-defined grasping points. Grasps were generated by randomly sampling an isotropic distribution $P(g_r | g_i)$ centered around an initial grasp g_i . $P(g_r | g_i)$ was defined as the product of a position and an orientation distribution. Let us denote the decomposition of a pose g into position and orientation by p and o respectively. We define $P(g_1 | g_2)$ with

$$P(g_1 | g_2) = \mathcal{N}(p_1; p_2, \sigma_p^2 \mathbf{I}) \frac{e^{\sigma_o o_1^T o_2} + e^{-\sigma_o o_1^T o_2}}{2} \quad (3)$$

where \mathcal{N} is a isotropic Gaussian kernel, and the fraction corresponds to a pair of antipodal von-Mises Fisher distri-

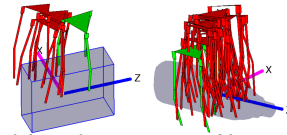


Fig. 4. Examples of data alignment: Unstable grasps (red) are associated with the example stable grasp (green). The object-relative hand grasping configurations were plotted by using a simplified hand model.

butions (a Gaussian-like distribution on the rotation group [19], [20]). The resulting grasps were distributed a few centimeters/degrees away from g_i . Executing multiple grasps in the neighborhood of g_i allowed the robot to learn the effect of small disturbances in hand positioning. In a real-world scenario it is important for the robot to learn the relations between the perceptions and the stability outcome in a region of an object rather than in a single location, because it is not reasonable to expect that the robot will always be able to grasp an object exactly at the same location.

Initial grasps were defined on the middle parts of the objects as shown in Figure 2. These initial grasps were parametrized by the pose of the hand with respect to the object. Each grasping experiment during the exploration of the objects in the neighborhood of the initial grasps was then run in the following way: An object was placed at an arbitrary position reachable by the robot. The standing/lying configuration of the objects also varied. The robot estimated the pose of the object and executed a random grasp \hat{g}_r in the neighborhood of g_i . In our experiments, after preshaping the hand, grasping is run by simultaneously closing the fingers and applying constant closing torques on all joints. Once the hand had stopped closing the fingers, the robot recorded the tactile imprints and the pose of the object. At the end of the grasp executions, the robot obtained the object-relative hand poses by comparing the vision-based object pose estimates to the known hand poses. The robot then attempted to lift the object by 5 cm. If the object slipped or rotated in the hand while being lifted, the grasp was marked as unstable. If lifting could be achieved robustly without any slippage or rotation, the grasp was marked as stable.

In total 584 grasps were executed, 134 for the box, 170 for the oval bottle, 138 for the cylindrical bottle and 142 for the spray bottle, during the exploration process explained above. Half of these grasps were stable and the other half were unstable.

A. Data Alignment: In order to learn how to alter an unstable grasp to reach a stable grasp the views need to be aligned in the data presented to the model at training time. This means that failed grasps are aligned with grasps that we want them to be corrected to. We use Self Organizing Maps (SOMs) [21] to perform this alignment.

We train SOMs using pose data (3D position and 4D orientation) from both stable and unstable grasp samples for each object and standing/lying configuration. When learning the associations the SOM chooses only a small subset of the stable grasps retaining. This suggests that we can choose fewer alternative solutions rather than trying to establish one-to-one correspondences which would force the learning to be

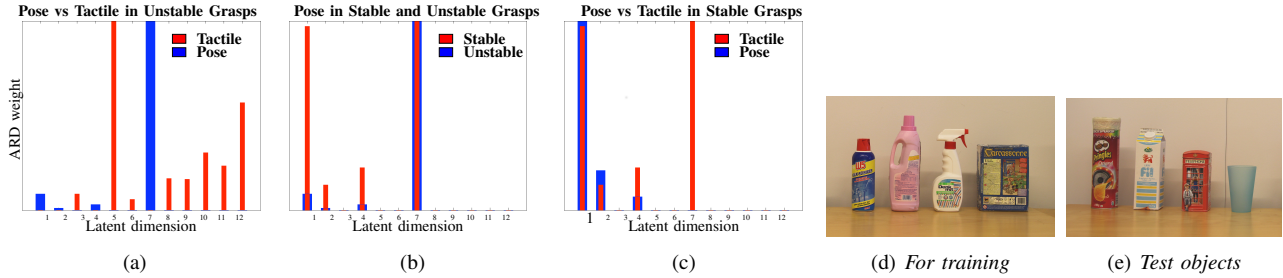


Fig. 5. (a-c): the ARD weights learned for the different modalities. The y-axis shows the value of the weight and the x-axis indicates the corresponding latent dimension. The left most plot shows the unsuccessful pose weights in blue and the weights for the associated tactile data as red. The two modalities share very little information, as there is little overlap across different weight components for tactile and pose data. This supports the expectation that there should not be a specific structure between the unstable grasp and their tactile imprints, as there are many different possibilities to do a bad grasp. The middle pane illustrates the scales for the stable and unstable grasps showing a significant sharing of information (the seventh dimension) even though different factors have different importance. The right most pane shows the scales for the tactile and the gripper position for stable grasps. As can be seen there is information in the tactile sensory data that is not represented by the gripper position, since using different objects results in different imprints. (d) For training objects used in the experiments: the cylindrical bottle, oval bottle, spray bottle and the box (from left to right), (e) The test objects: pringles, milk bottle, money box and mug (from left to right)

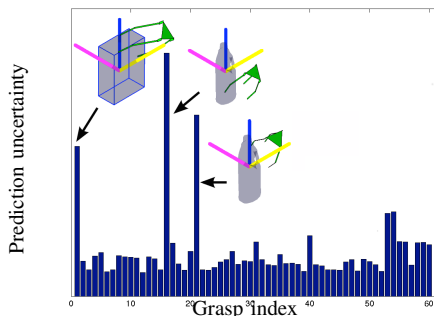


Fig. 6. Predicted grasps with associated model uncertainty: From left to right the uncertainty in the model is shown for each prediction in the test set. The generated grasps for which the model is most uncertain are also seen. Predictions with high uncertainties are not executed.

too specific and not likely to generalize beyond the simple case. Resulting associations can be intuitively interpreted. For instance, Figure 4 shows that for unstable grasps (as on the box) which are around one corner of the object when the object is lying on its elongated edge a better alternative is to move towards the center of the edge. If the unstable grasps (as on the oval bottle) are already sufficiently close to the center, the hand should be moved towards the object. In summary, in the data we have 104 stable grasps associated with 292 unstable grasps. In our experiments, we use randomly chosen 60 unstable grasps as a test set to evaluate the trained model.

B. Grasping: Correction and Prediction:

The MRD model learns a factorization of the observations which represents which factors are responsible for generating each view. From these weights we can infer the amount of information that is shared between two views and how much of their variations are independent. This is important, as a factor shared between two views can be inferred from *either* of the views and does therefore represent information that can be easily transferred from one view to another. In Figure 5a-c the weights recovered after training for some of the views are shown. We evaluate the trained model with four main experiments. We compare our approach to a standard regression method and produce inference on grasp pose given

unstable pose and the object identity. Our approach allows to perform inference given any subset of the variables without training new models which is the case for the standard regression method. We also evaluate our model using test objects seen in Figure 5e. As a different application we present prediction results on the tactile parameters. As a final experiment we present another application of this property and infer the object identity given tactile observations.

1) Correcting unstable grasps: In the first set of experiments we compare our MRD model with a non-linear regression model trained with unstable grasps as input and stable versions as output according to the association provided by the SOM. In specific, we use a \mathcal{GP} -prior to model a functional relationship between unstable and stable grasps. The results of both the MRD and the \mathcal{GP} are shown in Figure 7. As can be seen the MRD model performs significantly better compared to the regression approach, both with respect to removing non-applicable grasps and correcting the unstable ones. Comparing the predicted grasps, the MRD model is more structured in a manner that is consistent with the object having significantly smaller number of implausible grasps. This indicates that compared to the regression baseline, the MRD model has uncovered and corrected just the specific factors needed to correct the grasps, while the regression model focuses on erroneous portion of the unstable grasps to represent the stable grasp. Furthermore, the MRD will also be able to handle scenarios where an unstable grasp can be corrected in several different ways due to its multi-modal structure. This is not possible in a regression model which is unimodal and will model the response by the mean of the stable grasps for which there is no guarantee that it will be stable.

The results shown in Figure 7 for both the baseline and the MRD are based on the mean prediction of the \mathcal{GP} . The average success rate for the MRD and the \mathcal{GP} is 75% and 38% respectively. However, we can also use the variance of the prediction as a measure of the models uncertainty. In Figure 6 we show the uncertainty and the associated executed grasps. As can be seen the predictions on which the model is

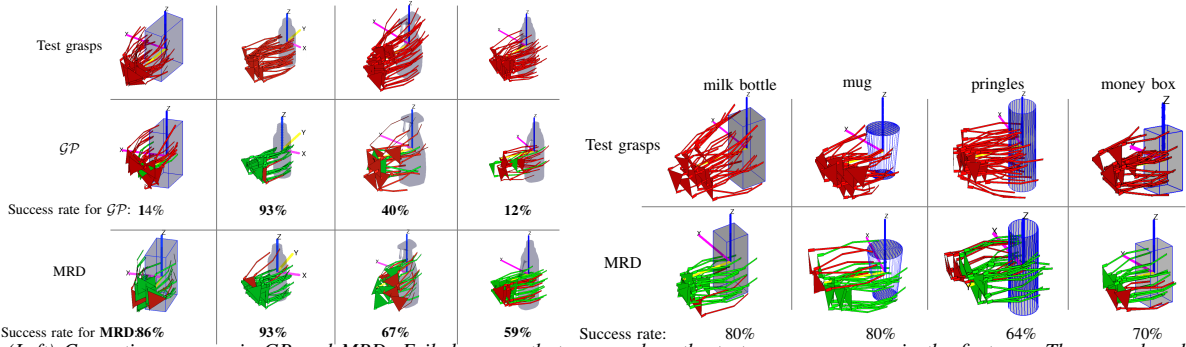


Fig. 7. (Left) Correcting grasps via GP and MRD: Failed grasps that are used as the test grasps are seen in the first row. The second and the third rows show the same grasps corrected via GP and MRD respectively. Green grasps are corrected grasps that have succeeded, red ones have failed. Success rates are also shown in the second and third rows. Predictions use both the object pose and identity. See text for more details. (Right) Correcting grasps onto new objects: Test grasps on new objects that have failed are seen in the first row. The same grasps corrected via MRD are seen in the second row. Green grasps are corrected grasps that have succeeded, red ones have failed. Success rates are seen in the last row. The average success rate over the test objects is 73% .

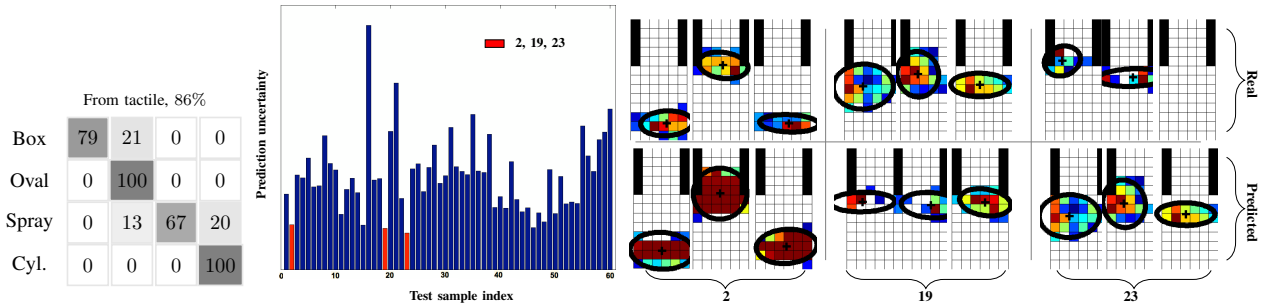


Fig. 8. Confusion matrix and the tactile modality associated with the corrected grasp predicted under the model: Left: the confusion matrix when we use the MRD model to infer the object identity from tactile data. As can be seen the model is capable of encoding the relationship between the modalities well, as given the tactile information of a grasp, it can with 86% average accuracy determine the object identity. Middle: The uncertainty in the predictions in the model. Right: The tactile data from the real execution and the model predictions for the test grasps. The uncertainty levels of these predictions are marked in red.

uncertain are also the ones that fail. Modeling uncertainty is crucial in robotics, as it allows the agent to identify the axes along which it needs to be particularly accurate, and axes along which it can afford to relax its movement to comply with other sources of constraints, such as reachability or task-related constraints.

In Figure 7 using our model we also provide results on test objects with varying shapes and sizes (seen in Figure 5e). The test set with unstable grasps include 14 grasps from the pringles bottle and 10 from each of the other three objects for lying and standing object configurations (44 test grasps in total). Grasps were collected following the data collection protocol discussed above. Given relative hand poses, the predictions for stable hand poses based on the trained model resulted in 9, 8, 8, 7 stable grasps for the pringles bottle, milk bottle, mug and the money box respectively yielding in average 73% success rate, which is similar to the performance when objects from the training set were used. These novel objects share similarities in shape, size and weight with the objects from the training set to some extent. The resulting grasps indicate that the model has learnt what to change, i.e., how to move the hand, to improve an unstable grasp and can apply it on previously unseen and similar objects.

2) Predicting Tactile Data: Being able to generate predictions on tactile readings is a useful ability in terms of successful manipulation. We show that we can generate

predictions on tactile signals in order to have expectation and therefore detect unexpected events. In Figure 8 we show examples of the predicted tactile sensory data for the corrected grasps. For most grasps the model predicts the center and orientation of the contact region well even though it is less reliable in terms of the actual pressure values. Tactile information is likely to be very challenging to model as very small changes in the gripper position might lead to losing contact resulting in a drastic change on the tactile imprints. This can be seen in the right most panel in Figure 8 where the model is certain about the prediction while the prediction is significantly different from the true estimate. This is not surprising as the gripper lost contact with one finger during execution which lead to a significant change in the tactile data.

3) Predicting Object Identity: Finally, we demonstrate another application of our model, i.e., predicting object identity from tactile observations. We use all the test grasps from the training objects. The confusion matrix obtained from classifying object identities based on tactile data can be seen in Figure 8. Although the average classification rate is 86%, predictions based on solely tactile data will inevitably fail in cases where similar readings are obtained from different parts of the objects. For example, around its center, the oval bottle has almost two parallel surfaces like a box and therefore the box is sometimes wrongly predicted as being the oval bottle based on the tactile data. The spray

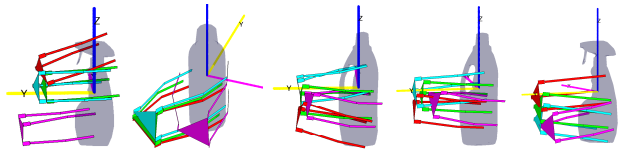


Fig. 9. Examples of corrected grasps: The red gripper indicates the failed grasp to be corrected and the green gripper depicts the output from our model. The cyan and the magenta gripper show, the closest and the furthest stable grasp in the training data to the red gripper according to the euclidean distance in gripper pose. With this figure we exemplify that our model is capable of generalizing and through its conditional approach generate better corrections compared to a data-base approach. If we have a sufficient amount of training data in the region of the failed grasp the cyan and the green gripper are very similar as can be seen in the two left-most images. The three remaining panels show grasps where the proposal under our model (green) corresponds to a smaller displacement of the failed grasp compared to the nearest neighbor (cyan) thereby generating a more “respectful” correction. Using a distance measure to mimic the conditional prediction of our model is very challenging. As can be seen in the right three panels it is not immediately obvious which is the better stable correction of the failed (red) grasp, the nearest (cyan) or the furthest (magenta) neighbor.

bottle has similar local geometry compared to the other two bottles which causes the spray bottle to be misclassified as these bottles.

C. Discussion: Results in this section show that our model can be used to encode grasping knowledge and to solve different grasping-related problems. We focused on grasp correction, i.e., suggesting better grasping poses in the neighborhood of a failed grasp. Another approach to achieve that is to use a data-base and simply select a stable grasp that has previously been executed. Given repeatable conditions such an approach is likely to have a very high success rate, however it cannot generalize beyond the database as it does not naturally take the failed grasp into consideration. Our approach, on the other hand, providing a probabilistic mapping can generate new stable grasps that are not present in the training data through inference conditioned on the failed grasp. In Figure 9 we compare the predictions of our model with examples from the database and we exemplify that our model is capable of generalizing and through its conditional approach generate better corrections compared to a data-base approach. To evaluate generalization capabilities we have also conducted experiments with test objects, which resulted in similar success rates. The resulting grasps indicate that the model has learned what to change, i.e., how to move the hand, from the associated positive and negative examples in the training data, in order to improve an unstable grasp and can apply it on previously unseen objects.

VI. CONCLUSION

We have presented an approach to learn how to perform robotic grasping based on multiple sensory modalities. A model of the relationships between grasp-related parameters has been learned using both successful and failed grasping examples. We have demonstrated applications such as grasp correction, object recognition and estimating tactile imprints. The key characteristic which facilitates this is the use of a factorized representation which separately models the information that is shared between the views. The factorization allows to perform efficient inference in ambiguous scenarios

where the observations are not sufficient to discriminate the desired output. Experimental results demonstrated that the proposed learning method is capable of generating stable grasping configurations given object identity and unstable grasps. We achieved 75% success rate on 60 unstable grasps while a traditional regression approach could only achieve 38% success rate. Experiments with test objects yielded 73% success rate showing that our approach is general and capable to scale beyond the training data.

Acknowledgment This work was supported in part by EU H2020 RoMaNS, 645582, and EPSRC EP/M026477/1.

REFERENCES

- [1] V. D. Nguyen, “Constructing force-closure grasps,” in *IEEE Int. Conf. Robotics and Automation (ICRA)*, vol. 3, Apr 1986, pp. 1368–1373.
- [2] A. Bicchi and V. Kumar, “Robotic grasping and contact: A review,” in *IEEE Int. Conf. on Robotics and Automation*, 2000, pp. 348–353.
- [3] A. Sabhani, S. El-Khoury, and P. Bidaud, “An overview of 3d object grasp synthesis algorithms,” *Robotics and Autonomous Systems*, vol. 60, no. 3, pp. 326 – 336, 2012.
- [4] J. Bohg, A. Morales, T. Asfour, and D. Kragic, “Data-driven grasp synthesis – a survey,” *Robotics, IEEE Transactions on*, vol. 30, no. 2, pp. 289–309, April 2014.
- [5] M. Przybylski, T. Asfour, and R. Dillmann, “Planning grasps for robotic hands using a novel object representation based on the medial axis transform,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2011, pp. 1781–1788.
- [6] J. Romano, K. Hsiao, G. Niemeyer, S. Chitta, and K. Kuchenbecker, “Human-inspired robotic grasp control with tactile sensing,” *Robotics, IEEE Transactions on*, vol. 27, no. 6, pp. 1067–1079, Dec 2011.
- [7] Y. Bekiroglu, R. Detry, and D. Kragic, “Learning tactile characterizations of object- and pose-specific grasps,” in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 2011, pp. 1554–1560.
- [8] E. L. Sauser, B. D. Argall, G. Metta, and A. G. Billard, “Iterative learning of grasp adaptation through human corrections,” *Robotics and Autonomous Systems*, vol. 60, no. 1, pp. 55–71, 2012.
- [9] R. Johansson, “Sensory input and control of grip,” in *Novartis Foundation Symposium*, 1998, pp. 45–59.
- [10] U. Hillenbrand and M. Roa, “Transferring functional grasps through contact warping and local replanning,” in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, 2012, pp. 2963–2970.
- [11] J. Tegin, S. Ekvall, D. Kragic, J. Wikander, and B. Iliev, “Demonstration-based learning and control for automatic grasping,” *Intelligent Service Robotics*, vol. 2, no. 1, pp. 23–30, 2009.
- [12] J. Tegin and J. Wikander, “A framework for grasp simulation and control in domestic environments,” in *Mechatronic Systems*, 2006, pp. 490–495.
- [13] H. Dang and P. K. Allen, “Grasp adjustment on novel objects using tactile experience from similar local geometry,” in *IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS)*, 2013.
- [14] M. Li, Y. Bekiroglu, D. Kragic, and A. Billard, “Learning of grasp adaptation through experience and tactile sensing,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2014.
- [15] T. Mörwald, J. Prankl, A. Richtsfeld, M. Zillich, and M. Vincze, “BLORT—the blocks world robotic vision toolbox,” *Best Practice in 3D Perception and Modeling for Mobile Manipulation (Workshop at ICRA 2010)*, 2010.
- [16] A. C. Damianou, C. H. Ek, M. Titsias, and N. D. Lawrence, “Manifold Relevance Determination,” in *International Conference on Machine Learning*, June 2012, pp. 145–152.
- [17] C. H. Ek, “Shared Gaussian Process Latent Variable Models,” Ph.D. dissertation, Oxford Brookes University, Oxford, 2009.
- [18] R. M. Neal, *Bayesian Learning for Neural Networks*. New York: Springer-Verlag, 1996, vol. 8.
- [19] R. A. Fisher, “Dispersion on a sphere,” in *Proc. Roy. Soc. London Ser. A.*, 1953.
- [20] E. B. Sudderth, “Graphical models for visual object recognition and tracking,” Ph.D. dissertation, Massachusetts Institute of Technology, Cambridge, MA, USA, 2006.
- [21] J. Vesanto, J. Himberg, E. Alhoniemi, and J. Parhankangas, “Som toolbox for matlab 5,” Helsinki University of Technology, Tech. Rep., 2000.