



EU FP7 CogX

ICT-215181

May 1 2008 (52months)

DR 2.5:

Qualitative models of object behaviour, and grasping of novel objects

Renaud Detry², Michael Zillich¹, Andreas Richtsfeld¹, Johann Prankl¹, Sergio Roa³, Sebastian Zurek⁴, Yasemin Bekiroglu², Thomas Mörwald¹, Markus Vincze¹, Danica Kragic², Jeremy Wyatt⁴, Geert-Jan Kruijff³

¹TUW, Vienna

²KTH, Stockholm

³DFKI GmbH, Saarbrücken

⁴BHAM, Birmingham

<detryr@kth.se>

Due date of deliverable: May 25, 2012

Actual submission date: May 28, 2012

Lead partner: KTH

Revision: final

Dissemination level: PU

This deliverable reports work related to object manipulation. The first problem discussed in this report is the extraction of qualitative models of object behaviour. We present a learning algorithm capable of extracting a discrete representation of a sensorimotor space. We also present a method for identifying different modes of object interactions, which allows us for instance to predict whether an object will turn left, right, or not turn at all if a particular push is applied.

Grasping novel objects is the second problem discussed in this report. We present means of detecting new objects from vision. We also present a novel approach that allows an agent to plan grasps onto novel objects by matching parts of the new object to parts of previously-grasped objects. We finally extend the work done on tactile-based grasp stability estimation in the previous period to allow the robot to assess grasp stability from both touch data *and* task requirements.

1	Tasks, objectives, results	6
1.1	Qualitative models of object behaviour	6
1.1.1	Planned work	6
1.1.2	Actual work performed	6
1.1.3	Relation to the state-of-the-art	7
1.2	Grasping of novel objects	8
1.2.1	Planned work	8
1.2.2	Actual work performed	8
1.2.3	Relation to the state-of-the-art	10
2	Annexes	14
2.1	Roa et al. “Robust Vector Quantization for Inference of Substochastic Sequential Machines”	14
2.2	Zurek et al. “Identification of qualitative states from the behaviour of objects”	15
2.3	Richtsfield et al. “Towards Scene Understanding Object Segmentation Using RGBD-Images”	16
2.4	Richtsfield et al. “Segmentation of Unknown Objects in Indoor Environments”	17
2.5	Richtsfield et al. “Implementation of Gestalt Principles for Object Segmentation”	18
2.6	Balzer et al. “Isogeometric Finite-Elements Methods and Variational Reconstruction Tasks in Vision – A Perfect Match”	19
2.7	Mörwald et al. “Fitting B-Spline Curves to Complex Shaped Boundaries” .	20
2.8	Mörwald et al. “Self-Monitoring to Improve Robustness of 3D Object Tracking for Robotics”	21
2.9	Detry et al. “Generalizing Grasps Across Partly Similar Objects”	22
2.10	Bekiroglu et al. “A Probabilistic Framework for Task-Oriented Grasp Stability Assessment”	23
	 References	 24

Executive Summary

This report presents the work done in the final year of CogX on the topics of (1) modelling qualitative object behaviour and (2) grasping novel objects. Regarding qualitative behaviour models, we first developed an algorithm that extracts probabilistic finite state representations of a dynamical system. This algorithm is applicable to the extraction of qualitative states from sensorimotor data gathered during the execution of a task by a robot. Second, we developed a method for identifying different modes of object interactions, which allowed us for instance to predict whether an object will turn left, right, or not turn at all if a particular push is applied.

Regarding (2), we made three novel contributions. First we developed a method for detecting novel objects in 3D scenes. This method was further integrated in the object tracker discussed in the previous periods, to allow the tracker to trigger texture-based object detection when its belief on the object pose becomes too low. Second, we developed an agent capable of learning the shape of parts by which objects are often grasped, which subsequently allows the agent to plan grasps on partly familiar objects. Finally, we introduced a task model that includes kinematic grasp parameters and tactile signals, which allows the agent to model the stability of a grasp with respect to a given task.

The work presented in this report led to four peer-reviewed conference publications, and six more conference and journal submissions. The work presented here follows up on DR 2.4 (forward models, grasping previously unseen objects) and on DR 5.4 (learning of cross-modal concepts).

Qualitative models of object behaviour

We address the problem of finding qualitative representations of dynamical systems. The task is to infer probabilistic finite-state machines that model the interaction between a robot and an object. In this case, a robot performs pushing actions and sequences of object poses are stored to be used in the learning process. New algorithms were developed for discretisation of sensorimotor spaces and extraction of finite-state probabilistic models [59, 73] (Annexes 2.1 and 2.2). In order to evaluate them, we tested their ability to find qualitative representations of artificial dynamical systems with noisy features, i.e., from data generated by probabilistic finite-state automata where states are gaussian noise distributions.

Role of qualitative models of object behaviour in CogX

Robots need to have capabilities for introspection, abstraction and memory in order to use their acquired knowledge in future tasks. Then, the models that are obtained by the robot during the interaction with objects

can be used subsequently to plan actions, to reason, to test a theory of the behaviour of the object after some action. Additionally, given the high-dimensionality of the sensorimotor spaces, it is useful to find more coarse-grained abstractions (concepts) from these interactions, in order to be used for human-driven learning and communication tasks involving language and gestures.

Contribution to the CogX scenarios and prototypes

This work contributes to the Dexter scenario. By learning how objects move when it interacts with them, the agent is able to predict how objects would move if certain manipulation plans are executed. Moreover, the graph-based nature of the representation encodes the probabilistic transitions that lead to subsequent system states which is particularly useful in planning.

Grasping of novel objects

Grasping novel objects is the second problem addressed in this deliverable. It encompasses two sub-problems: detecting novel objects and planning grasps onto novel objects. For the former we learned 3D perceptual grouping principles to segment objects from RGBD images of cluttered scenes and describe objects as grouped surface patches [56, 57, 58] (Annexes 2.3, 2.4, 2.5). These surface patches are described as NURBS [47, 4] (Annexes 2.7 and 2.6), allowing for a flexible and compact representation of objects. Once objects are detected and object models in terms of surface patches are extracted, these can be tracked using methods presented in DR 2.4 [46]. We extended this work to include a self-assessment of the tracker regarding its current performance and object state (moving, occluded) which allows for robust tracking by optimally combining tracking and re-detection as required [48] (Annex 2.8). This work led to a conference and workshop publication at CVPR and CVWW respectively [4, 56] (Annexes 2.6 and 2.3).

Regarding grasping, we present a method that allows a robotic agent to learn prototypical parts by which objects are often grasped, from a set of grasps demonstrated by a teacher [18]. Prototypes subsequently allow the agent to grasp novel objects that contain a part that resembles one of the prototypes. This work led to an ICRA publication [18] (Annex 2.9). We also present a model of task-oriented grasp stability, and means of inferring task stability from tactile data [8] (Annex 2.10). This work was evaluated on the KTH manipulation platform (industrial arm and dexterous hand). The work on task stability partly builds on the contributions presented in DR 2.4 [7].

Role of grasping of novel objects in CogX

One of the aims of CogX is to create an agent that is able to familiarise itself with an unknown environment. While exploring its environment, the agent comes into contact with novel objects. The agent has to grasp some of these objects to fulfil its task. Even when grasping is not required by the task itself, manipulating objects represents an efficient exploration strategy, and it allows the agent to fill object-related knowledge gaps.

In this report, we present means of detecting novel objects from vision, and means of exploiting previously-acquired object knowledge to grasp novel objects. We also present a method for exploiting tactile data *and* task requirements to assess the stability of a grasp. These three contributions improve the efficiency and the robustness with which our agent familiarises itself with novel objects.

Contribution to the CogX scenarios and prototypes

This work contributes to the Dexter scenario, where the agent is required to interact with novel objects. Novel objects first need to be identified as such. This problem is solved using the novelty detection method discussed above. Part-based planning is then used to plan grasps on the objects, provided that partly similar objects have been handled previously. Novelty detection also contributes to the George scenario, where the robot learns novel objects in interaction with a tutor.

1 Tasks, objectives, results

1.1 Qualitative models of object behaviour

1.1.1 Planned work

This deliverable reports work related to Task 2.7:

Task 2.7: Extracting qualitative states. To be able to perform introspection on possible qualitative cause we require a model that has not only continuous states, but also qualitative states, with qualitative explanations for the transitions between them. We will use the notion of force-aspect graphs to devise a learning algorithm capable of partitioning the continuous configuration space of the modular motor learning predictions into sets of qualitatively similar stable states, plus their basins of attraction. (M33–M39)

This deliverable presents two contributions that address this task. The first contribution is an algorithm that learns a probabilistic finite-state representation of a dynamical system from sensorimotor data. The second contribution is an algorithm that extracts different modes of object-effector interaction in manipulative actions.

This deliverable contributes to the realisation of the sixth measurable objective:

Objective 6: Methods for perception and manipulation of objects that enable a robot to actively explore objects, to extend its manipulative skills, and its understanding of these.

Extracting qualitative models provides the agent with a compact representation of its sensorimotor space, which allows it to better understand its skills. Compact representations also simplify the decision process required by the high-level planning of actions or exploration, providing the agent with a small set of discrete choices instead of the full range of possible motions.

1.1.2 Actual work performed

We developed an algorithm [59] (Annex 2.1) that extracts probabilistic finite state representations of a dynamical system. The resulting representation has the form of a set of states, with transitions of different probabilities between the states. A dynamical system can be represented as a tuple $\langle I, O, S, P \rangle$, where I, O and S are input, output and state spaces respectively and P a set of conditional probabilities. I, O and S need to be quantised in order to extract qualitative representations of the system. Here, we use a modification of the Growing Neural Gas algorithm [24, 54, 55, 59]

for quantisation which is robust for finding the right clusters in the presence of noise (RobustGNG). In order to evaluate the quantisation ability of the algorithm, we performed a clustering task with Gaussian distributions of different types. The algorithm is successful in finding the right number of clusters, by making use of an information-theoretic method, namely the Minimum Description Length (MDL) criterion. The algorithm employs an incremental way of learning where no prior information about the maximal number of clusters or iterations is needed. Its stopping criterion is a measure of graph stability based on MDL, where each node in a graph is a discrete latent variable. To evaluate the extraction of probabilistic machines, we used Noisy Automata where states are Gaussian distributions and transitions are probabilistic. The algorithm was able to infer qualitative states and construct corresponding probabilistic machines which include quantisers for the input space, the output space and the state space, corresponding transitions functions and their probabilities.

The second contribution related to qualitative behaviour models is a method for identifying qualitative states in robot-object interactions [73] (Annex 2.2). The interaction studied here is a robot pushing an object with a single finger. We present means of clustering the sensorimotor data obtained during short exploratory pushes. The sensorimotor data are composed of the starting position and orientation of the finger with respect to the object, and the object displacement that results from the push. Our results demonstrate that the algorithm enumerates states that accord with human judgement. For instance, our system extracts discrete behaviours that correspond to the object turning left, turning right, or moving straight.

We note that the two contributions presented above are similar in spirit, and have complementary roles in this deliverable. The finite state machine model is able to capture complex interactions involving sequences of multiple states and their transition probabilities, while the second contribution focuses on short interactions, and reasons on sensorimotor data captured at the beginning and at the end of each interaction. However, while the state machine is tested on artificial data, the state clustering approach is evaluated on a concrete robot problem, and it includes heuristics that allow the agent to process high-dimensional sensorimotor data efficiently.

1.1.3 Relation to the state-of-the-art

The quantisation algorithm builds upon previous implementations of an algorithm based on Neural Gas algorithm [54, 55], which uses information-theoretic properties to decide the proper number of clusters. The new algorithm is incremental and can work in online settings. A decision to add nodes to a graph is based on an online estimation of error and nodes can also be removed depending on information-theoretic measures. We also incorporated additional efficiency improvements in the learning process.

To extract probabilistic finite-state machines from these dynamical systems, we applied the CrySSMEx algorithm [35] with some modifications to allow the quantisation based on RobustGNG and improve the learning convergence.

The clustering algorithm for extracting qualitative states was inspired by the “push-stability diagram,” introduced by Brost [12, 13]. We also used some of the discretisation ideas found in Kuipers’ work (e.g. [49]).

1.2 Grasping of novel objects

1.2.1 Planned work

This deliverable reports work related to Task 2.8:

Task 2.8: Grasping novel objects. Based on our object models, we will investigate the scalability of the system with respect to grasping novel, previously unseen objects. We will demonstrate how the system can execute tasks that involve grasping based on the extracted sensory input (both about the scene and individual objects) and taking into account its embodiment. (M27–M50)

Task 2.8 spans the second half of the project. Grasping novel objects requires (1) the ability to detect novel objects, (2) the ability to *plan* grasps onto novel objects, and (3) the ability to execute the planned grasps robustly. All three points are addressed in this report. The first point is addressed through 3D perceptual grouping. The second point is addressed with a method for planning grasps from partial object snapshots. The third point is addressed with a model of touch-based task stability.

This deliverable contributes to the realisation of the sixth measurable objective:

Objective 6: Methods for perception and manipulation of objects that enable a robot to actively explore objects, to extend its manipulative skills, and its understanding of these.

The novelty detection and grasping work presented here fulfil both the perception and manipulation objectives. The method for learning graspable parts allows the agent to understand its manipulation skills, by extracting recurrent patterns from the agent’s experience.

1.2.2 Actual work performed

A prerequisite for grasping novel objects is detection of these objects in the first place. While this is comparatively easy for simple scenes of isolated objects on a table surface, cluttered scenes containing arbitrary arrangements (such as stacks and piles) of unknown objects still poses a challenge. The

first contribution in this section then is a method to segment objects from RGBD images of cluttered scenes. After pre-segmentation of the RGBD input image based on surface normals, surface patches are estimated using a mixture of planes and NURBS (non-uniform rational B-splines) and model selection is employed to find the best representation for the given data. We then construct a graph from surface patches and relations between pairs of patches and perform graph cut to arrive at object hypotheses segmented from the scene. The energy terms for patch relations are learned from user annotated training data, where support vector machines (SVM) are trained to classify a relation as being indicative of two patches belonging to the same object. We show evaluation of the relations and results on a database of different test sets, demonstrating that the approach can segment objects of various shapes in cluttered table top scenes. This work led to a conference and workshop publication at CVPR and CVWW respectively [4, 56] (Annexes 2.6 and 2.3) and submissions to IROS, ICPR and DAGM [57, 58, 47] (Annexes 2.4, 2.5, 2.7).

The second contribution is an extension of the tracking work presented in DR 2.4 [46] with a method for self-assessment of the tracker. Real world settings in object tracking pose challenges such as automatically detecting tracking failure, real-time processing, and robustness to occlusion, illumination, and view point changes. This work presents a 3D tracking system that is capable of overcoming these difficulties using a monocular camera. We present a method of Tracking-State-Detection (TSD) that takes advantage of commercial graphics processors to map textures onto object geometry, to learn textures online, and to recover object pose in real-time. Our system is able to handle 6 DOF object motion during changing lighting conditions, partial occlusion and motion blur while maintaining an accuracy of a few millimetres. Furthermore using TSD we are able to automatically detect occlusions or whether we lost track, and can then trigger a SIFT-based recognition system that is trained during tracking to recover the pose. Evaluations are presented in relation to ground truth pose data and examples present TSD on real-world scenes presented in video sequences. This work led to a conference publication at ROBIO [48] (Annex 2.8).

The third contribution related to grasping novel objects is an agent that has the ability to identify parts by which objects are often grasped [18] (Annex 2.9). As a result, the agent is able to quickly plan grasps onto novel objects that partly resemble objects that it has grasped before. Our agent extracts experience from a set of grasps demonstrated by a teacher. Demonstrations are conducted by placing objects of various shapes and sizes within the robot hand and instructing the robot to close the hand. The final configurations of the hand with respect to the 3D object shapes are used as training data. The agent searches these data for parts that recur in the vicinity of the hand across different grasps. To this end, the agent first extracts shape segments of predefined sizes around the grasping point

of each grasp example. This process provides it with a set of prototype *candidates*. The agent then computes pairwise shape similarities between all candidates, and clusters the candidate in the space induced by the similarity measure. The agent only conserves the cluster centres, which altogether form a dictionary of grasp prototypes. By keeping the number of cluster low, we can effectively compress grasping experience into a dictionary that is orders of magnitude smaller than the original set of grasp examples. The dictionary allows the agent to plan grasps from a single partial 3D snapshot of a novel object. The agent attempts to fit all the prototypes to the snapshot, and it executes the grasp that corresponds to the best-fitting prototype. This work led to an ICRA publication [18] (Annex 2.9).

The fourth contribution related to grasping novel objects is a joint model of object grasping parameters, tactile imprints, and task stability [8] (Annex 2.10). In DR 2.4, we presented a general-purpose model of touch-based grasp stability. As noted in the report, however, stability is not an absolute property. Instead, stability largely depends on the task that the agent is performing. For instance, a grasp aimed at seizing a hammer for hitting on a nail needs to be more firm than a grasp aimed at pouring water from a bottle.

We have extended the model of DR 2.4 to include task-related information. The result is a generative model of the class of an object, grasp parameters, task, tactile imprints, and grasp stability. The joint probability of these variables is modelled with a Bayesian network. The model is learned from experiments performed both in simulation and on a real robot. The model allows our agent to reason on any of the variables listed above, given observations of the other variables. For instance, given the tactile feedback gathered after closing the hand on an bottle, the agent is able to decide whether it is safe to use the grasp to pour water off the bottle. If it is not, the agent can compute whether the grasp is good enough to simply transport the bottle. If it is the case, the agent could potentially move the bottle to another location from which it could try a grasp that is better-suited for pouring.

1.2.3 Relation to the state-of-the-art

Various approaches to segment objects either in 2D images or in point clouds exist, where early approaches aimed to formulate generic Gestalt principles to organise 2D scenes into objects. Gestalt principles are also used by Kooststra et al. in [38] and [37] where the authors developed a symmetry detector to initialise segmentation based on a Markov Random Field (MRF). Furthermore Kooststra et al. developed a quality measure based on Gestalt principles to rank segmentation results. Many state-of-the-art approaches formulate image segmentation as energy minimisation with a Markov Random Field (MRF) [9, 65, 11, 60]. In addition to an appearance model computed from

colour and texture, which is commonly used to better distinguish foreground from background, Bergstrom et al. [9] formulate an objective function where it is possible to incrementally add constraints generated through human-robot interaction. In [68] Werlberger et al. propose a variational model for interactive segmentation using a shape prior. This method is based on minimising the Geodesic Active Contour energy. The approach by Hager et al. [28] is able to segment objects from cluttered scenes in point clouds generated from stereo by using a strong prior 3D model of the scene and explicitly modelling physical constraints such as support and handles dynamic changes such as object appearance/disappearance. It is however limited to parametric models (boxes, cylinders), whereas our approach is only limited by the amount and type of training data.

In his groundbreaking paper [33], Horn shows how to approximate the Laplacian by second-order finite differences (FD) on the image grid and solve the resulting algebraic system by a fixed-point scheme. Extensions of Horn's method are too numerous to list here but let us explicitly mention the most recent ones like Harker's and O'Leary's [29, 30] as well as that due to Durou et al. [20], who describe a powerful total-variation-based algorithm capable of resolving discontinuities in the depth map without prior segmentation of the gradient field. Our work relates to the class of kernel methods [21, 50], which can be thought of as mesh-free FEMs in disguise. Similarly, Kovési applies a basis $\{b_j\}$ of shapelets to the normal adaption problem in scene space [39]. Only a few authors explicitly consider the classical, i.e., non-isogeometric FEM: Hicks employs it for integrating normal fields with three-dimensional support into a foliation of surfaces [32]. Generalisations of Horn's method applicable to such spatially varying normal fields are presented by Balzer [3] and Delaunoy and Prados [17]. None of aforementioned methods is compatible with the geometry representation of contemporary CAD packages. Higher-degree polynomial bases and a multi-scale mechanism are per se possible, at least on polygon meshes, but quite challenging to implement.

A review of the massive body of literature on B-Spline curve fitting would go far beyond the scope of this paper. We briefly give an overview of relevant work and afterwards point out the most common approaches to which we want to apply our methods. One of the most fundamental summaries on B-Splines and least-squares fitting to point-clouds was done in the well know book of Piegl et. al. [51] where they minimise a functional. This method was carefully investigated in [52] where they especially focus on the squared distance function and their approximants used for least-squares fitting. In [66, 67, 10] new point-curve distance functions are introduced to improve the convergence rate and robustness. In [70] Yang et al. propose an active implicit B-Spline model and find the zero set of a bivariate tensor-product B-Spline function using the trust region algorithm [53]. Fitting B-Spline curves to point-clouds in the presence of obstacles is introduced in [22,

23], where they minimise a functional subject to an inequality constraint. Hu et al. [34] present a method where they take advantage of both algebraic and geometric distance minimisation and therefore avoid additional constraints. Often it is necessary to modify an existing curve fitting method to apply to a problem with certain characteristics, such as noise, outliers, unknown degree of freedom (DOF) and so forth [27, 72, 5]. Our approach extends the *Squared Distance Minimisation* (SDM) of [67], i.e. we are modifying their error term for the functional to be minimised. Further we want to overcome the problem of specifying the degree of freedom manually and add control points and knots depending on the error of the curve.

Tracking the pose of an object in image sequences is a classical problem in robot vision, where current approaches aim at improving robustness in tough real-world scenarios [19, 14, 42, 64, 36]. [41] use a combination of edges and textures for tracking. Their approach extracts point features from surface texture and use them together with edges to calculate object pose. This turns out to be very fast as well as robust against occlusion. Our approach not only uses patches but the whole texture, which usually lets the pose converge very quickly to the accurate pose. Since the algorithm runs on the GPU, it is as fast as the method in [41]. More recent approaches aim to solve most of the problems of tracking, such as [63] where the authors are matching the camera image with pre-trained keyframes and then minimising the squared distance of feature points taking into account neighbouring frames. The approach described in [44] uses a modified version of the Active Appearance Model which allows for partial and self occlusion of the objects and for high accuracy and precision. In [16] the authors minimise the optical flow resulting from the projection of a textured model and the camera image. To compensate for shadows and changing lighting they apply an illumination normalisation technique. The work presented in [25] describes an approach for real-time visual servoing using a binocular camera setup to estimate the pose by triangulating a set of feature points. As in our approach [61] takes advantage of robust Monte Carlo particle filtering to determine the pose of the camera with respect to SIFT features, which are localised in 3D using epipolar geometry. Missing in all the above methods is a detection when tracking fails rather than reporting tracking trapped in a local optimum. The proposed tracking state detection (TSD) addresses this problems and we develop an approach to to work fully automatically.

Recent approaches to plan grasps onto novel objects rely on methods that learn a direct mapping from visual cues to grasp parameters. Authors have studied the association of grasping strategies to various kinds of visual cues. Grasps associated to local visual features [62, 45] have the advantage of being easily transferable across objects, as many objects share similar components. However, local features suffer from a poor geometric resolution, which makes it difficult to accurately associate them to the 6D pose of a gripper, let alone finger preshape parameters. Conversely, grasps associated to a model of a

whole object [15, 26] benefit from increased geometric robustness, but the resulting models will not apply to novel objects. Authors have explored this trade-off between transferability and robustness by associating grasps to object parts of varying size [2, 6, 31, 43, 71]. An important distinctive point of our work is that we provide the agent with means of optimising this transferability-robustness trade-off internally, by allowing it to select prototypical parts of varying size, depending on their occurrence statistics in the training database. The result is a compact dictionary of parts that lend themselves to grasping.

As argued above, task-related constraints are important for grasp planning. The geometry of a grasp (i.e., the side by which an object is grasped) is often crucial for the execution of a task. This problem has been studied for instance by Xue et al. [69], who manually encoded the expertise about task semantics provided by a human tutor. Another task-related aspect of grasping is that different tasks require different levels of robustness to external object disturbances (in term of the force a grasp is able to apply onto an object). This problem has been studied by Li et al. [40], and more recently by Aleotti et al. [1], who defined task-related grasp quality measures which combined task knowledge with analytical stability measures used in traditional grasp stability studies.

In our work, we combined supervised task learning with experience-based stability learning. This allowed stability to be assessed in a task-oriented manner. This is especially beneficial for energy-efficient control: when a task (e.g., *hand-over*) does not require strong grasping, a relatively smaller gripping force can be applied.

2 Annexes

2.1 Roa et al. “Robust Vector Quantization for Inference of Substochastic Sequential Machines”

Bibliography S. Roa and G.-J. Kruijff: “Robust Vector Quantization for Inference of Substochastic Sequential Machines”. Submitted to Journal of Neurocomputing, 2012.

Abstract The article explores the problem of discretizing the continuous evolution of a dynamical system. The article proposes an algorithm to learn a probabilistic discrete state, an input and an output space representation of the system, together with probabilistic transition functions. The method is based on the CrySSMEx algorithm for extracting substochastic finite state machines, and a new Vector quantization algorithm. We performed experiments on Vector quantization with artificial data generated using Gaussian noise distributions. The quantization algorithm is able to find the optimal number of clusters. It induces a good model of the data, avoiding overfitting. Data stemming from Noisy automata were used to test the algorithm for extracting sequential finite state machines. The induced models represent accurately the behavior of these discrete dynamical systems.

Relation to WP The interaction between a robot and an object leads to different object behaviors. These behaviors can be learned by using predicting models inferring the causal relationships in these interactions. Additionally, these models are qualitative representations in which the sensorimotor space is discretized to find meaningful abstractions. The algorithm presented here is used for obtaining those qualitative representations.

2.2 Zurek et al. “Identification of qualitative states from the behaviour of objects”

Bibliography Sebastian Zurek, Marek Kopicki, and Jeremy Wyatt: “Identification of qualitative states from the behaviour of objects”. University of Birmingham, technical report, 2012.

Abstract For a robotic agent interacting with its environment, it is natural to represent its sensory input and motor output as continuous state spaces. This poses a challenge for controlling the behaviour of a robot, since at almost every instant it will observe a novel situation and will have an infinite choice of motor commands that it could deploy. An objective for robotics research is to devise algorithms that can extract qualitative states, in accord with human judgement. We present an algorithm that uses the behaviour of an object, when manipulated and observed by a robot, to discover the qualitative states in perception-action space. Thus we take the definition of a qualitative state to be a set of points in state space that behave similarly under a given action. The algorithm is evaluated by using data from a simulation of a robotic finger pushing an object.

Relation to WP This paper directly addresses the first topic of this deliverable by presenting a method for extracting qualitative states of object behaviours. The paper also shows the applicability of the method to a robot problem.

2.3 Richtsfeld et al. “Towards Scene Understanding Object Segmentation Using RGBD-Images”

Bibliography Richtsfeld, Andreas; Mörwald, Thomas; Prankl, Johann; Balzer, Jonathan; Zillich, Michael; Vincze, Markus: “Towards Scene Understanding Object Segmentation Using RGBD-Images”, Proceedings of the 2012 Computer Vision Winter Workshop (CVWW), 2012.

Abstract We present a framework for detecting unknown 3D objects in RGBD-images and extracting representations suitable for robotics tasks such as grasping. We address cluttered scenes with stacked and jumbled objects where simplistic plane pop-out methods are not sufficient. We start by estimating surface patches using a mixture of planes and NURBS (non-uniform rational B-splines) fitted to the 3D point cloud and employ model selection to find the best representation for the given data. We then construct a graph from surface patches and relations between patches and perform graph cut to arrive at object hypotheses segmented from the scene. The energy terms for patch relations are learned from user annotated training data, where we train a support vector machine (SVM) to classify a relation as being indicative of two patches belonging to the same object given a vector of relation features, such as proximity or color similarity. We show preliminary results demonstrating that the approach can segment objects of various shapes in cluttered table top scenes.

Relation to WP A prerequisite for grasping novel objects (Task 2.8) is the detection of novel objects, i.e. objects for which no instance or category model is available. The above work addresses this problem in a learning framework being essentially only limited by the amount and diversity of training data.

2.4 Richtsfeld et al. “Segmentation of Unknown Objects in Indoor Environments”

Bibliography Richtsfeld, Andreas; Mörwald, Thomas; Prankl, Johann; Zillich, Michael; Vincze, Markus: “Segmentation of Unknown Objects in Indoor Environments”, submitted to the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2012.

Abstract We present a framework for segmenting unknown objects in RGBD-images suitable for robotics tasks such as object search, grasping and manipulation. While handling single objects on a table is solved, handling complex scenes poses considerable problems due to clutter and occlusion. After pre-segmentation of the input image based on surface normals, surface patches are estimated using a mixture of planes and NURBS (non-uniform rational B-splines) and model selection is employed to find the best representation for the given data. We then construct a graph from surface patches and relations between pairs of patches and perform graph cut to arrive at object hypotheses segmented from the scene. The energy terms for patch relations are learned from user annotated training data, where support vector machines (SVM) are trained to classify a relation as being indicative of two patches belonging to the same object. We show evaluation of the relations and results on a database of different test sets, demonstrating that the approach can segment objects of various shapes in cluttered table top scenes.

Relation to WP A prerequisite for grasping novel objects (Task 2.8) is the detection of novel objects, i.e. objects for which no instance or category model is available. The above work addresses this problem in a learning framework being essentially only limited by the amount and diversity of training data.

2.5 Richtsfeld et al. “Implementation of Gestalt Principles for Object Segmentation”

Bibliography Richtsfeld, Andreas; Zillich, Michael; Vincze, Markus: “Implementation of Gestalt Principles for Object Segmentation”, submitted to the International Conference on Pattern Recognition (ICPR), 2012.

Abstract Gestalt principles have been studied for about a century and were used for various computer vision approaches during the last decades, but became unpopular because the many heuristics employed proved inadequate for many real world scenarios. We show a new methodology to learn relations inferred from Gestalt principles and an application to segment unknown objects, even if objects are stacked or jumbled and tackle also the problem of segmenting partially occluded objects. The relevance of the relations for object segmentation is learned with support vector machines (SVMs) during a training period. We present an evaluation of the relations and show results of the segmentation framework.

Relation to WP A prerequisite for grasping novel objects (Task 2.8) is the detection of novel objects, i.e. objects for which no instance or category model is available. The above work investigates a set of 3D Gestalt principles used in the work reported in Annexes 2.3 and 2.4, for the detection on novel objects based on learning the importance of these Gestalt principles from training examples.

2.6 Balzer et al. “Isogeometric Finite-Elements Methods and Variational Reconstruction Tasks in Vision – A Perfect Match”

Bibliography Balzer, Jonathan; Mörwald, Thomas: “Isogeometric Finite-Elements Methods and Variational Reconstruction Tasks in Vision – A Perfect Match”, IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012.

Abstract Inverse problems are abundant in vision. A common way to deal with their inherent ill-posedness is reformulating them within the framework of the calculus of variations. This always leads to partial differential equations as conditions of (local) optimality. In this paper, we propose solving such equations numerically by isogeometric analysis, a special kind of finite-elements method. We will expose its main advantages including superior computational performance, a natural ability to facilitate multi-scale reconstruction, and a high degree of compatibility with the spline geometries encountered in modern computer-aided design systems. To animate these fairly general arguments, their impact on the well-known depth-from-gradients problem is discussed, which amounts to solving a Poisson equation on the image plane. Experiments suggest that, by the isogeometry principle, reconstructions of unprecedented quality can be obtained without any prefiltering of the data.

Relation to WP Fitting of parametric surface models is one part of the work reported in Annex 2.3, where plane and NURBS models compete for an optimal representation of the underlying point cloud data in a model selection framework. As such it is an enabling technology for Task 2.8.

2.7 Mörwald et al. “Fitting B-Spline Curves to Complex Shaped Boundaries”

Bibliography Mörwald, Thomas and Prankl, Johann and Zillich, Michael and Vincze, Markus: “Fitting B-Spline Curves to Complex Shaped Boundaries”, Submitted to the Joint German/Austrian Pattern Recognition Symposium (DAGM-OAGM), 2012.

Abstract Finding the boundary of some region and computing a curve to approximate it best is a common task in computer vision and image processing. This paper describes an approach of fitting B-Splines to 2D point-clouds for robustly finding the boundary of complex shapes. The problems of common B-Spline fitting methods are discussed. New techniques to overcome this problems, namely the Asymmetric Distance Minimization, Error-Adaptive Knot Insertion and Concavity Filling are applied and considered as the main contribution of our work. We will show how our fitting approach leads to satisfying solutions, even by employing a generic initialization scheme and without knowing the required degree of freedom. All improvements are discussed and demonstrated on difficult problems from real sensor data.

Relation to WP The work in this paper is yet another sub-problem of the work reported in Annex 2.6: finding the exact boundary of the data points contributing to a model, projected onto the parametric surface. This is a requirement for constructing precise and “dense” object models with surface patches stitched seamlessly together.

2.8 Mörwald et al. “Self-Monitoring to Improve Robustness of 3D Object Tracking for Robotics”

Bibliography Mörwald, Thomas; Zillich, Michael; Prankl, Johann; Vincze, Markus: “Self-Monitoring to Improve Robustness of 3D Object Tracking for Robotics”, IEEE International Conference on Robotics and Biomimetics (ROBIO), 2012.

Abstract In robotics object tracking is needed to steer towards objects, check if grasping is successful, or investigate objects more closely by poking or handling them. While many 3D object tracking approaches have been proposed in the past, real world settings pose challenges such as automatically detecting tracking failure, real-time processing, and robustness to occlusion, illumination, and view point changes. This paper presents a 3D tracking system that is capable of overcoming these difficulties using a monocular camera. We present a method of Tracking-State-Detection (TSD) that takes advantage of commercial graphics processors to map textures onto object geometry, to learn textures online, and to recover object pose in real-time. Our system is able to handle 6 DOF object motion during changing lighting conditions, partial occlusion and motion blur while maintaining an accuracy of a few millimetres. Furthermore using TSD we are able to automatically detect occlusions or whether we lost track, and can then trigger a SIFT-based recognition system that is trained during tracking to recover the pose. Evaluations are presented in relation to ground truth pose data and examples present TSD on real-world scenes presented in video sequences.

Relation to WP While not directly related to detection/grasping of novel objects, this work extends previous work reported in deliverable DR 2.4 on tracking objects for grasping and manipulation. Reasoning about the current state of tracking is an important factor when employing tracking within a larger system that has to make informed decisions, such as aborting a grasping movement in case reliable pose estimates are no longer available.

2.9 Detry et al. “Generalizing Grasps Across Partly Similar Objects”

Bibliography Detry, Renaud; Ek, Carl Henrik; Madry, Marianna; Piater, Justus; Kragic, Danica : “Generalizing grasps across partly similar objects”, IEEE International Conference on Robotics and Automation, 2012.

Abstract The paper starts by reviewing the challenges associated to grasp planning, and previous work on robot grasping. Our review emphasizes the importance of agents that generalize grasping strategies across objects, and that are able to transfer these strategies to novel objects. In the rest of the paper, we then devise a novel approach to the grasp transfer problem, where generalization is achieved by *learning*, from a set of grasp examples, a dictionary of object parts by which objects are often grasped. We detail the application of dimensionality reduction and unsupervised clustering algorithms to the end of identifying the size and shape of parts that often predict the application of a grasp. The learned dictionary allows our agent to grasp novel objects which share a part with previously seen objects, by matching the learned parts to the current view of the new object, and selecting the grasp associated to the best-fitting part. We present and discuss a proof-of-concept experiment in which a dictionary is learned from a set of synthetic grasp examples. While prior work in this area focused primarily on shape analysis (parts identified, e.g., through visual clustering, or salient structure analysis), the key aspect of this work is the emergence of parts from *both* object shape *and* grasp examples. As a result, parts intrinsically encode the intention of executing a grasp.

Relation to WP This work is concerned with transferring grasping knowledge across known objects and to novel objects. We developed a method that allows a robot to identify parts by which objects are often grasped, thereby allowing the robot to easily grasp novel objects that contain a familiar part.

2.10 Bekiroglu et al. “A Probabilistic Framework for Task-Oriented Grasp Stability Assessment”

Bibliography Bekiroglu, Yasemin; Song, Dan; Wang, Lu; Kragic, Danica : “A Probabilistic Framework for Task-Oriented Grasp Stability Assessment”, KTH Royal Institute of Technology (Technical report), 2012.

Abstract We present a probabilistic framework for grasp modeling and stability assessment. The framework facilitates assessment of grasp success in a goal-oriented way, taking into account both geometric constraints for task affordances and stability requirements specific for a task. We integrate high-level task information introduced by a teacher in a supervised setting with low-level stability requirements acquired through a robot’s self-exploration. The conditional relations between tasks and multiple sensory streams (vision, proprioception and tactile) are modeled using Bayesian networks. The generative modeling approach both allows prediction of grasp success, and provides insights into dependencies between variables and features relevant for object grasping.

Relation to WP This work is concerned with the exploitation of touch data *and* task requirements to assess the stability of a grasp. As different tasks impose different constraints on object-gripper bonds, it is important to take tasks into account when assessing stability. The ability to assess stability from touch is particularly important when grasping novel objects, as the configuration of the grasp is less certain than when grasping a known object.

References

- [1] J. Aleotti and S. Caselli. Interactive teaching of task-oriented robot grasps. *Robotics and Autonomous Systems*, 58(5):539–550, 2010.
- [2] Jacopo Aleotti and Stefano Caselli. Part-based robot grasp planning from human demonstration. In *IEEE International Conference on Robotics and Automation*, 2011.
- [3] Jonathan Balzer. A Gauss-Newton Method for the Integration of Spatial Normal Fields in Shape Space. *J. Math. Imaging Vis.* In press.
- [4] Jonathan Balzer and Thomas Mörwald. Isogeometric finite-elements methods and variational reconstruction tasks in vision – a perfect match. In *IEEE Conference on Computer Vision and Pattern Recognition*, Providence, RI, USA, Jun 2012.
- [5] L Barbieri, F Bruno, M Muzzupappa, and J Pernot. Constrained fitting of B-spline curves based on the force density method. In *International Conference on Innovative Methods in Product Design*, volume 46, 2011.
- [6] C. Bard and J. Troccaz. Automatic preshaping for a dextrous hand from a simple description of objects. In *IEEE International Workshop on Intelligent Robots and Systems*, pages 865–872. IEEE, 1990.
- [7] Yasemin Bekiroglu, Janne Laaksonen, Jimmy Alison Jørgensen, Ville Kyrki, and Danica Kragic. Assessing grasp stability based on learning and haptic data. *IEEE Transactions on Robotics*, 2011.
- [8] Yasemin Bekiroglu, Dan Song, Lu Wang, and Danica Kragic. A probabilistic framework for task-oriented grasp stability assessment. Technical report, KTH Royal Institute of Technology, 2012.
- [9] N. Bergström, M. Björkman, and Danica Kragic. Generating object hypotheses in natural scenes through human-robot interaction. In *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, pages 827–833. IEEE, 2011.
- [10] A Blake and M Isard. *Active Contours*, volume 17. Springer, 1998.
- [11] Y Y Boykov and M.-P. Jolly. Interactive graph cuts for optimal boundary and region segmentation of objects in N-D images. In *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, volume 1, pages 105–112 vol.1, 2001.
- [12] Randy C. Brost. Planning robot grasping motions in the presence of uncertainty. Technical Report CMU-RI-TR-85-12, The Robotics Institute, Carnegie-Mellon University, Pittsburgh, PA, July 1985.

- [13] Randy C. Brost. Automatic grasp planning in the presence of uncertainty. *Int. J. Robotics Research*, 7(1):3–17, 1988.
- [14] J. Chestnutt, S. Kagami, K. Nishiwaki, J. Kuffner, and T. Kanade. Gpu-accelerated real-time 3d tracking for humanoid locomotion. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2007.
- [15] Charles de Granville, Joshua Southerland, and Andrew H. Fagg. Learning grasp affordances through human demonstration. In *IEEE International Conference on Development and Learning*, 2006.
- [16] Hans de Ruiter and Beno Benhabib. Visual-model-based, real-time 3d pose tracking for autonomous navigation: methodology and experiments. *Autonomous Robots*, 25:267–286, 2008.
- [17] A. Delaunoy and E. Prados. Gradient Flows for Optimizing Triangular Mesh-based Surfaces: Applications to 3D Reconstruction Problems Dealing with Visibility. *Int. J. Comput. Vision*, 95:100–123, 2011.
- [18] Renaud Detry, Carl Henrik Ek, Marianna Madry, Justus Piater, and Danica Kragic. Generalizing grasps across partly similar objects. In *IEEE International Conference on Robotics and Automation*, 2012.
- [19] T. Drummond and R. Cipolla. Real-time tracking of complex structures with on-line camera calibration. In *British Machine Vision Conference (BMVC'99)*, pages 574–583, 1999.
- [20] J.-D. Durou, J.-F. Aujol, and F. Courteille. Integrating the Normal Field of a Surface in the Presence of Discontinuities. *Proc. EMMCVPR*, pages 261–273, 2009.
- [21] S. Ettl, J. Kaminski, M. Knauer, and G. Häusler. Shape reconstruction from gradient data. *Appl. Optics*, 47(12):2091–2097, 2008.
- [22] S Flory. Fitting curves and surfaces to point clouds in the presence of obstacles. *Computer Aided Geometric Design*, 26(2):192–202, 2009.
- [23] S Flory and M Hofer. Constrained curve fitting on manifolds. *Computer-Aided Design*, 40(1):25–34, 2008.
- [24] Bernd Fritzke. A growing neural gas network learns topologies. In *Advances in Neural Information Processing Systems 7*, pages 625–632. MIT Press, 1995.
- [25] J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendón-Mancha. Binocular visual tracking and grasping of a moving object with a 3d trajectory. *Journal of Applied Research and Technology*, 7(03):259–274, 2009.

- [26] C. Goldfeder, M. Ciocarlie, H. Dang, and P.K. Allen. The Columbia grasp database. In *IEEE International Conference on Robotics and Automation*, 2009.
- [27] O Grove. From CT to NURBS: Contour Fitting with B-spline Curves. *Computer Aided Design And Applications*, 8(1):3–21, 2011.
- [28] Gregory D Hager and Ben Wegbreit. Scene parsing using a prior world model. *The International Journal of Robotics Research*, 2011.
- [29] M. Harker and P. O’Leary. Least squares surface reconstruction from measured gradient fields. *Proc. CVPR*, 1:1–7, 2008.
- [30] M. Harker and P. O’Leary. Least squares surface reconstruction from gradients: Direct algebraic methods with spectral, Tikhonov, and constrained regularization. *Proc. CVPR*, 1:2529–2536, 2011.
- [31] A. Herzog, P. Pastor, M. Kalakrishnan, L. Righetti, T. Asfour, and S. Schaal. Template-based learning of grasp selection. In *The PR2 Workshop (Workshop at IROS’11)*, 2011.
- [32] R. Hicks. Designing a mirror to realize a given projection. *J. Opt. Soc. Am. A*, 22(2):323–330, 2005.
- [33] B. Horn. Height and gradient from shading. *Int. J. Comput. Vision*, 5(1):37–75, 1999.
- [34] Mingxiao Hu, Jieqing Feng, and Jianmin Zheng. An additional branch free algebraic B-spline curve fitting method. *The Visual Computer*, 26(6-8):801–811, 2010.
- [35] H. Jacobsson. The crystallizing substochastic sequential machine extractor - CrySSMEx. *Neural Computation*, 18(9):2211–2255, 2006.
- [36] Georg Klein and Tom Drummond. Robust visual tracking for non-instrumented augmented reality. In *ISMAR IEEE/ACM International Symposium on Mixed and Augmented Reality*, 2003.
- [37] Gert Kootstra, Niklas Bergström, and Danica Kragic. Fast and Automatic Detection and Segmentation of Unknown Objects. In *Humanoids*, Bled, 2011.
- [38] Gert Kootstra, Niklas Bergström, and Danica Kragic. Gestalt Principles for Attention and Segmentation in Natural and Artificial Vision Systems. In *SPME*, Shanghai, 2011.
- [39] P. Kovesi. Shapelets correlated with surface normals produce surfaces. *Proc. ICCV*, 2:994–1001, 2005.

- [40] Z. Li and S. Sastry. Task oriented optimal grasping by multifingered robot hands. In *Robotics and Automation. Proceedings. 1987 IEEE International Conference on*, volume 4, pages 389–394, January 2003.
- [41] Lucie Masson, Michel Dhome, and Frederic Jurie. Robust real time tracking of 3d objects. In *International Conference on Pattern Recognition, ICPR*, 2004.
- [42] Philipp Michel, Joel Chestnutt, Satoshi Kagami, Koichi Nishiwaki, James Kuffner, and Takeo Kanade. Gpu-accelerated real-time 3d tracking for humanoid autonomy. In *JSME Robotics and Mechatronics Conference (ROBOMECH'08)*, June 2008.
- [43] A. T. Miller, S. Knoop, H. Christensen, and P. K. Allen. Automatic grasp planning using shape primitives. In *IEEE International Conference on Robotics and Automation*, volume 2, pages 1824–1829, 2003.
- [44] Pradit Mittrapiyanuruk, Guilherme N. Desouza, and Avinash C. Kak. Accurate 3d tracking of rigid objects with occlusion using active appearance models. In *WACV/MOTION*, pages 90–95, 2005.
- [45] L. Montesano and M. Lopes. Learning grasping affordances from local visual descriptors. In *IEEE International Conference on Development and Learning*, 2009.
- [46] Thomas Mörwald, Marek Kopicki, Rustam Stolkin, Jeremy Wyatt, Sebastian Zurek, Michael Zillich, and Markus Vincze. Predicting the unobservable, visual 3d tracking with a probabilistic motion model. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA11)*, May 2011.
- [47] Thomas Mörwald, Johann Prankl, Michael Zillich, and Markus Vincze. Fitting b-spline curves to complex shaped boundaries. In *Joint German/Austrian Pattern Recognition Symposium (DAGM-OAGM) (submitted)*, Graz, Austria, Aug 2012.
- [48] Thomas Mörwald, Michael Zillich, Johann Prankl, and Markus Vincze. Self-monitoring to improve robustness of 3d object tracking for robotics. In *IEEE International Conference on Robotics and Biomimetics (RO-BIO)*, Phuket, Thailand, Dec 2011.
- [49] Jonathan Mugan and Benjamin Kuipers. Autonomous learning of high-level states and actions in continuous environments. *IEEE Trans. Autonomous Mental Development*, 4(1):70–86, 2012.
- [50] H.-S. Ng, T.-P. Wu, and C.-K. Tang. Surface-from-gradients without discrete integrability enforcement: A gaussian kernel approach. *IEEE T. Pattern Anal.*, 32:2085–2099, 2010.

- [51] Les Piegl and Wayne Tiller. *The NURBS book*. Monographs in visual communication. Springer, 1996.
- [52] Helmut Pottmann and Michael Hofer. Geometry of the Squared Distance Function to Curves and Surfaces. *Visualization and mathematics III*, (90):223–244, 2003.
- [53] M.J.D. Powell. On the global convergence of trust region algorithms for unconstrained optimization. *Math. Prog.*, 29:297–303, 1984.
- [54] A.K. Qin and P.N. Suganthan. Robust growing neural gas algorithm with application in cluster analysis. *Neural Networks*, 17(8-9):1135 – 1148, 2004. New Developments in Self-Organizing Systems.
- [55] K Qin and N Suganthan. Enhanced neural gas network for prototype-based clustering. *Pattern Recognition*, 38(8):1275–1288, 2005.
- [56] Andreas Richtsfeld, Thomas Mörwald, Johann Prankl, Jonathan Balzer, Michael Zillich, and Markus Vincze. Towards scene understanding – object segmentation using rgb-d-images. In *Proceedings of the 2012 Computer Vision Winter Workshop (CVWW)*, Mala Nedelja, Slovenia, February 2012.
- [57] Andreas Richtsfeld, Thomas Mörwald, Johann Prankl, Michael Zillich, and Markus Vincze. Segmentation of unknown objects in indoor environments. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (submitted)*, Vilamoura, Algarve, Portugal, Oct 2012.
- [58] Andreas Richtsfeld, Michael Zillich, and Markus Vincze. Implementation of gestalt principles for object segmentation. In *21st International Conference on Pattern Recognition (ICPR) (submitted)*, Tsukuba, JAPAN, Nov 2012.
- [59] Sergio Roa and Geert-Jan Kruijff. Robust vector quantization for inference of substochastic sequential machines. *Neurocomputing*, 2012. submitted.
- [60] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. "Grab-Cut": interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.*, 23(3):309–314, August 2004.
- [61] J.R. Sánchez, H. Álvarez, and D. Borro. Towards real time 3d tracking and reconstruction on a gpu using monte carlo simulations. In *IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 185 –192, Oct. 2010.
- [62] A. Saxena, J. Driemeyer, and A. Y. Ng. Robotic Grasping of Novel Objects using Vision. *International Journal of Robotics Research*, 27(2):157, 2008.

- [63] L. Vacchetti, V. Lepetit, and P. Fua. Stable real-time 3d tracking using online and offline information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2004.
- [64] Luca Vacchetti, Vincent Lepetit, and Pascal Fua. Combining edge and texture information for real-time accurate 3d camera tracking. In *IEEE/ACM International Symposium on Mixed and Augmented Reality*, 2004.
- [65] S Vicente, V Kolmogorov, and C Rother. Joint optimization of segmentation and appearance models. *2009 IEEE 12th International Conference on Computer Vision*, (Iccv):755–762, 2009.
- [66] Wenping Wang, Helmut Pottmann, and Yang Liu. Fitting B-Spline Curves to Point Clouds by Squared Distance Minimization. *ACM TOG*, page 41, 2004.
- [67] Wenping Wang, Helmut Pottmann, and Yang Liu. Fitting B-spline curves to point clouds by curvature-based squared distance minimization. *ACM Transactions on Graphics*, 25(2):214–238, 2006.
- [68] Manuel Werlberger, Thomas Pock, Markus Unger, and Horst Bischof. A Variational Model for Interactive Shape Prior Segmentation and Real-Time Tracking. In *International Conference on Scale Space and Variational Methods in Computer Vision (SSVM)*, Voss, Norway, 2009.
- [69] Z. Xue, J. Zoellner, and R. Dillmann. Automatic optimal grasp planning based on found contact points. In *IEEE/ASME Int. Conf. on Advanced Intelligent Mechatronics*, pages 1053–1058, 2008.
- [70] Zhouwang Yang, Jiansong Deng, and Falai Chen. Fitting unorganized point clouds with active implicit B-spline curves. *The Visual Computer*, 21(8-10):831–839, 2005.
- [71] Li (Emma) Zhang, Matei Ciocarlie, and Kaijen Hsiao. Grasp evaluation with graspable feature matching. In *RSS Workshop on Mobile Manipulation: Learning to Manipulate*, 2011.
- [72] Xiuyang Zhaoa, Caiming Zhang, Bo Yang, and Pingping Li. Adaptive knot placement using a GMM-based continuous optimization algorithm in B-spline curve approximation. *Computer-Aided Design*, 43(6):598–604, 2011.
- [73] Sebastian Zurek, Marek Kopicki, and Jeremy Wyatt. Identification of qualitative states from the behaviour of objects. Technical report, University of Birmingham, 2012.

17th Computer Vision Winter Workshop
Matej Kristan, Rok Mandeljc, Luka Čehovin (eds.)
Mala Nedelja, Slovenia, February 1-3, 2012

Towards Scene Understanding – Object Segmentation Using RGBD-Images

Andreas Richtsfeld, Thomas Mörwald, Johann Prankl, Jonathan Balzer,
Michael Zillich and Markus Vincze

Vienna University of Technology
Gusshausstrae 25-29, 1040 Vienna

Abstract. We present a framework for detecting unknown 3D objects in RGBD-images and extracting representations suitable for robotics tasks such as grasping. We address cluttered scenes with stacked and jumbled objects where simplistic plane pop-out methods are not sufficient. We start by estimating surface patches using a mixture of planes and NURBS (non-uniform rational B-splines) fitted to the 3D point cloud and employ model selection to find the best representation for the given data. We then construct a graph from surface patches and relations between patches and perform graph cut to arrive at object hypotheses segmented from the scene. The energy terms for patch relations are learned from user annotated training data, where we train a support vector machine (SVM) to classify a relation as being indicative of two patches belonging to the same object given a vector of relation features, such as proximity or color similarity. We show preliminary results demonstrating that the approach can segment objects of various shapes in cluttered table top scenes.

1. Introduction

Segmenting unknown objects from generic scenes is one of the elusive goals of computer vision and in general a very ill defined problem. Thanks to the recent introduction of cheap and powerful 3D sensors (such as the Microsoft Kinect or Asus Xtion-PRO) which deliver a dense point cloud plus color for almost any indoor scene, a renewed interest in 3D methods holds the promise to push the envelope slightly further.

In this work we aim at segmenting unknown objects of arbitrary (but reasonably compact) shape from table top scenes, where objects need not be standing isolated but can be jumbled in heaps. An example for such a scene is shown in Fig. 1. More-



Figure 1. Segmented objects from a cluttered table top scene with stacked and jumbled objects.

over we want a compact and accurate representation of object shapes, suitable in a robotics domain for various manipulation tasks.

The dense and reliable point cloud delivered by a Kinect sensor allows us to robustly fit planar surface patches to parts of the point cloud. These planes are fast to compute and capture a good range of typical man made objects. In order to also model curved objects with high accuracy we furthermore fit NURBS (non-uniform rational B-splines), replacing planes whenever NURBS provide a better fit. We use model selection [12] to find the combination of planes and NURBS optimally explaining the point cloud data.

Segmenting objects from the scene then amounts to identifying groups of surface patches that are likely to belong to the same objects. I.e. we perform perceptual grouping, but not as is more traditionally done in 2D using e.g. edges and junctions, but using 3D surface features and relations. We define several pairwise relations such as proximity or color similarity and form a relation feature vector. Each of these relations indicates to a certain degree that the respective surface patches are likely to belong to the same object, with e.g. closeness being a very good indicator and color similarity being far weaker.

We address this with a learning approach where we use human-annotated ground truth to offline train an SVM to categorize a relation feature vector as either indicating same or different object for a given

pair of patches. We then construct a graph from all the surface patches and pairwise relations and use the output of the SVM as the pairwise energy term in a graph cut based segmentation. The resulting segmentation is able to detect many typical objects as they arise in robotics tasks (books, boxes, cups, bowls), provided that single surfaces are big enough to be captured in sufficient detail by the Kinect sensor and that enough training data is provided to the SVM to capture all arising surface relations.

The key novelty in our approach lies in combining planes and NURBS with learned relations in a 3D perceptual grouping approach.

The paper is structured as follows. The next section sets the presented work into context with related work in this research field. Sec. 3 explains the approach in detail for the different components of the framework. Experimental results are shown in Sec. 4, before the work ends with a conclusion and outlook in Sec. 5.

2. Related Work

Various approaches to segment objects either in 2D images or in point clouds exist. Early approaches aimed to formulate generic Gestalt principles to organise 2D scenes into objects. For an overview of this early work in perceptual organisation we want to refer to Boyer and Sarkar [2]. More recently Zillich [22] proposed an any-time perceptual grouping framework to segment convex parts in images. Gestalt principles are also used by Koosstra et al. in [11] and [10]. They developed a symmetry detector to initialize segmentation based on a Markov Random Field (MRF). Furthermore Koosstra et al. developed a quality measure based on Gestalt principles to rank segmentation results.

Many state-of-the-art approaches formulate image segmentation as energy minimization with an MRF [1, 19, 3, 17]. In addition to an appearance model computed from colour and texture, which is commonly used to better distinguish foreground from background, Bergstrom et al. [1] formulate an objective function where it is possible to incrementally add constraints generated through human-robot interaction. In [20] Werlberger et al. propose a variational model for interactive segmentation using a shape prior. This method is based on minimizing the Geodesic Active Contour energy.

Active segmentation is proposed in Mishra et al. [13] and [14] where an image

point is fixated and the shortest path in a log polar transformed edge image is computed. In addition to the edges computed from colour and texture information the above authors propose to use the depth image from stereo cameras to improve segmentation.

The approach by Hager et al. [9] is able to segment objects from cluttered scenes in point clouds generated from stereo by using a strong prior 3D model of the scene and explicitly modelling physical constraints such as support and handles dynamic changes such as object appearance/disappearance. It is however limited to parametric models (boxes, cylinders).

The problem of fitting higher order surfaces to point clouds was already addressed by the framework of Leonardis et al. [12]. They segment range images by estimating piecewise linear surfaces, modelled with bivariate polynomials. Furthermore they developed a Model Selection framework, which is used to find the best interpretation of the range data in terms of Minimum Description Length (MDL). Instead of using bivariate polynomials we first describe the scene with simple plane models and then substitute planes with NURBS if the approximation of the point cloud is better in terms of MDL. Additionally, we cluster surface patches to objects depending on learned patch relations.

3. Approach

Our approach consists of four major parts, namely plane fitting, NURBS fitting, model selection and object segmentation. The first two parts abstract from the raw point cloud to surface patches. The model selection part determines the combination of surface patches which optimally represents the underlying point cloud. Object segmentation finally uses relations between surface patches to estimate which of the surface models belong together thus forming object hypotheses.

3.1. Plane Fitting

We chose to fit planes into the point cloud as a first abstraction step. Man-made environments contain many planar surfaces, so planes do in fact quite well describe a good portion of the scenes we are interested in from a robotics point of view. Furthermore, planes are easy to fit using robust methods, as opposed to higher-parametric models such as NURBS or superquadrics.

Plane fitting is typically done with a robust method

(e.g. some variant of RANSAC [8]) by sequentially fitting models and removing the inliers, and iterating until there are too few points remaining to support a model. Taking into account the scene layout we are typically going to encounter, we added an additional step. Even though we explicitly tackle cluttered scenes with objects lying around in heaps, often at least these heaps can be separated. To this end we iteratively fit a plane, remove its inliers and then perform clustering based on pairwise Euclidean distances on the remaining point cloud. We then perform RANSAC on the subset of points belonging to a single cluster, leading to a significantly increased inlier rate in comparison to the whole remaining point cloud. For scenes with objects lying on a supporting surface this very effectively speeds up the iterative RANSAC procedure.

For all point cloud operations we use PCL (Point Cloud Library) [18], which provides various RANSAC methods as well as clustering and other basic operations.

The resulting planes represent the point cloud not always optimally, because e.g. a cylindrical surface will be represented with a number of planar stripes. This issue is corrected by NURBS fitting and model selection described in the following subsections.

3.2. NURBS Fitting

For representing free-form surfaces there are a number of geometric models available. Most widely used in industry are NURBS (non-uniform rational B-splines). The reasons for their popularity are the convenient manipulation and the ability to represent all conic sections, i.e. circles, cylinders, ellipsoids, spheres and so forth. The possibility for refinement through knot insertion allows for adaption to local irregularities, while selecting a certain degree of freedom gives reason about the measured surface we want to fit to.

A good overview of the properties and advantages of NURBS can be found in Chapter 1.1.2 in [6]. NURBS are a generalisation of B-splines, that allow for exact representation of a wide range of objects. For simplicity we will focus on B-splines for now and will postpone the move to NURBS to future work.

3.2.1 B-splines

The mathematical concept of B-splines would go far beyond the scope of this paper. So for those who

are interested let us refer to the well known book by Piegl et al. [15] and start from their definition of B-spline surfaces in Chapter 3.4.

$$\mathbf{S}(\xi, \eta) = \sum_{i=1}^n \sum_{j=1}^m N_{i,p}(\xi) M_{j,p}(\eta) \mathbf{B}_{i,j} \quad (1)$$

The basic idea behind this formulation is to manipulate the B-spline surface $\mathbf{S} : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ of degree p , by changing the values of the *control grid*. The i, j -element of the control grid is called *control point* $\mathbf{B}_{i,j} \in \mathbb{R}^3$ which defines the B-spline surface at its region of influence determined by the *basis functions* $N_{i,p}(\xi), M_{j,p}(\eta)$. $(\xi, \eta) \in \Omega$ are called parameters defined on the domain $\Omega \subset \mathbb{R}^2$.

Refinement is established by *knot insertion*, i.e. increasing the number of control points, and therefore the degrees of freedom, without changing the surface \mathbf{S} . A detailed exposition of knot refinement is available in Chapter 5.3 in [15].

3.2.2 Point-Cloud Fitting

Given a set of points $\mathbf{q}_h \in \mathbb{R}^3$ with $h = 1 \dots k$ and $k > mn$ we want to fit a B-spline surface \mathbf{S} with $n > p, m > p$ and $p \geq 1$. Writing Eq. (1) as a linear system

$$\mathbf{s} = \mathbf{A}\mathbf{b} \quad (2)$$

where $\mathbf{s} \in \mathbb{R}^{k \times 3}$ are points on the B-spline surface. $\mathbf{A} = \mathbf{A}(\xi_h, \eta_h) \in \mathbb{R}^{k \times nm}$ contains the values of the basis functions at (ξ_h, η_h) and the vector of control points $\mathbf{b} \in \mathbb{R}^{nm \times 3}$ is the control grid $\mathbf{B} \in \mathbb{R}^{n \times m \times 3}$ written as vector. The (ξ_h, η_h) are precomputed parameters described in Sec. 3.2.3. We look for a solution of the overdetermined linear system (2) in the least-squares sense, i.e. a minimum of

$$d = \sum_{h=1}^k |\mathbf{q}_h - \mathbf{s}_h(\mathbf{b})|^2 \quad (3)$$

with respect to \mathbf{b} .

3.2.3 Initialisation

For minimizing the functional in Eq. (3) the parameters (ξ_h, η_h) for \mathbf{A} in Eq. (2) are required. We compute them by finding the closest point $\mathbf{s}_h(\xi_h, \eta_h)$ on the B-spline surface to \mathbf{q}_h using Newton's method. Therefore a B-spline surface is initialised from the front face of the camera-axis-aligned bounding box of the point-cloud (Fig. 2).

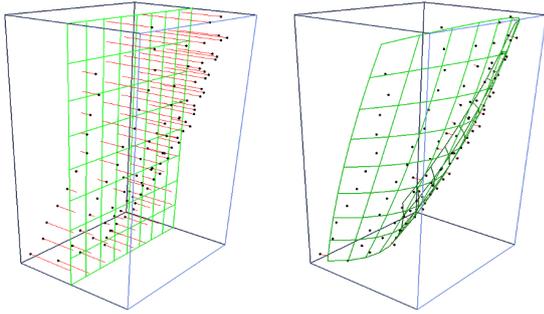


Figure 2. Fitting a B-spline surface (green) by minimizing the closest point distances (red) of a point-cloud (black). For initialisation the camera-axis aligned bounding-box (blue) is used ($m = n = 3, p = 2, w_a = 1, w_r = 0.1$).

3.2.4 Regularisation

To get a smooth surface and to avoid folding we add a regularisation term to Eq. (2) such that a control point tends to lie in the arithmetic mean of its neighbours. For control points at the interior and the boundary of the control grid we consider the 4- and 2-neighbourhood respectively. The regularisation can be written as

$$\mathbf{0} = \mathbf{R}\mathbf{b} \quad (4)$$

and is simply appended to Eq. (2).

$$\begin{bmatrix} \mathbf{s} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} w_a \mathbf{A} \\ w_r \mathbf{R} \end{bmatrix} \mathbf{b} \quad (5)$$

w_a and w_r are the weights defining the influence of the point matching and regularisation.

3.3. Model Selection

For assembling surface patches to object hypotheses, we first need to segment the point cloud into individual patches and estimate the surface model parameters. As shown in the previous section for planes, this can be done with the robust estimation method RANSAC by sequentially fitting models and filtering the inliers. A plane has only three parameters, therefore random sampling is an appropriate approach. For NURBS, where the number of parameters is three times the number of control points, an intractable number of random samples would be necessary to select inliers of a surface patch. For this reason we first explain the point cloud in terms of piecewise planar patches using the sequential RANSAC approach described in Sec. 3.1 and then greedily merge planes and substitute them with NURBS by using Model Selection and a Minimum Description

Length (MDL) criterion. In the following paragraph we briefly describe the basic mathematical tool which is introduced by Leonardis et al. [12] for the purpose of range image segmentation and we describe the proposed framework. Our formulation is most similar to the formulation by Prankl et al. [16] who use Model Selection to detect planes in image pairs.

The idea of Model Selection is that the same data point can not belong to more than one surface model. Hence an over-complete set of models is generated and the best subset in terms of an MDL criterion is selected. To select the best model, the savings for each surface hypothesis H can be expressed as

$$S_H = S_{data} - \kappa_1 S_m - \kappa_2 S_{err} \quad (6)$$

where S_{data} is the number of data points N explained by the hypothesis H , S_m stands for the cost of coding different models and S_{err} describes the cost for the error added. κ_1 and κ_2 are constants to weight the different terms. As proposed in [12] we use the number of parameters to define S_m . For the cost S_{err} experiments have shown that the Gaussian error model $\mathcal{N}(\mu_{err}, \sigma_{err}^2)$ and an approximation of the log-likelihood has a superior performance. Hence the cost of the error results in

$$S_{err} = -\log \prod_{i=1}^N p(f_i|H) = \quad (7)$$

$$\approx \sum_{i=1}^N (1 - p(f_i|H)) \quad (8)$$

and accordingly the substitution of Eq. 8 in Eq. 6 yields the savings of a model

$$S_H = \frac{N}{A_m} - \kappa_1 S_m - \frac{\kappa_2}{A_m} \sum_{i=1}^N (1 - p(f_i|H)), \quad (9)$$

where A_m is a normalization value for merging two models.

For modelling surface patches we then propose a two step algorithm, where first the savings for individual point clusters are compared and then neighbouring point clusters are greedily merged if the savings of the merged cluster

$$S_{ij} > S_i + S_j \quad (10)$$

is higher than the savings of two individual clusters. Alg. 1 summarizes the proposed surface modelling pipeline.

Algorithm 1 Modelling of surface patches

```

Detect piecewise planar surface patches
for  $i = 0 \rightarrow$  number of patches do
    Fit nurbs to patch  $i$ 
    Compute MDL savings  $S_{i,nurbs}$  and  $S_{i,plane}$ 
    if  $S_{i,nurbs} > S_{i,plane}$  then
        Substitute the model  $H_{i,plane}$  with  $H_{i,nurbs}$ 
    end if
end for
Create Euclidean neighbourhood pairs  $P_{ij}$  for surface patches
for  $k = 0 \rightarrow$  number of neighbours  $P_{ij}$  do
    Greedily fit nurbs to neighbouring patches  $P_{ij}$ 
    Compute MDL savings  $S_{ij}$  to merged patches
    if  $S_{ij} > S_i + S_j$  then
        Substitute individual models  $H_i$  and  $H_j$  with merged nurbs model  $H_{ij}$ 
    end if
end for
    
```

3.4. Object Segmentation

The previous sections explained how to find the best representation of a point cloud with different surface patch models. In the last processing step we are now interested in grouping these surface patches into object hypotheses. To this end it is first necessary to figure out which relations between surface patches contribute to the probability that these patches belong together and are part of one and the same object. We define the following relations, which are calculated for neighbouring surface patches:

- r_{ch} ... difference of patch colour
- r_{tr} ... difference of patch texture
- r_{cb} ... colour distance along patch border
- r_{di} ... distance along patch border
- r_{cu} ... curvature along patch border

The first two of the five relations describe differences of global patch properties, while the other three describe differences of local properties along the borders of the patches. The detailed implementation of the relations is explained and discussed in the result Sec. 4.

Each of these relations is defined to produce a value between 0 (same) and 1 (different) - with the exception of r_{cu} , but the degree to which that value indicates two surfaces as belonging to the same object is different for each relation. So a r_{ch} value of

0.3 will typically have a completely different meaning than a r_{tr} value of 0.3. And moreover these will be dependent on the scenes and objects encountered.

We address this with a learning approach. We define relation vectors \mathbf{r}

$$\mathbf{r} = \{r_{ch}, r_{tr}, r_{di}, r_{cb}, r_{cu}\} \quad (11)$$

and train an SVM to classify a relation vector as indicating same or different object.

For the offline training phase of the SVM we hand-annotated a set of depth images. Relation vectors between neighbouring surfaces that belong to the same object represent positive training examples, and those between neighbouring patches belonging to different objects or to an object and background represent negative examples.

We use the libsvm package [5], a free SVM software package, with a radial basis function as kernel:

$$\mathbf{K}(x_i, x_j) = e^{\gamma \|x_i - x_j\|_2} \quad (12)$$

In the online phase, the SVM is capable to provide not only a binary decision *same* or *notsame* for each \mathbf{r} , but also a probability value $p(\text{same} | \mathbf{r})$ for each decision, based on the theory introduced by Wu et al. [21].

The last processing step makes a global decision and answers the question, which groups of patches form objects. To this end we define a graph, where patches represent nodes and relations represent edges. The graph is not fully connected (which would be computationally prohibitive), as we only define relations between surface patches which are close neighbours. We then employ graph-cut segmentation, introduced by Felzenszwalb et al. [7], using the above probability values as the pairwise energy terms.

4. Experiments

Each learning approach is only as good as its training data, in our case training images for the SVM. The training images must be complex enough, so that the trained SVM can later distinguish between objects, which are e.g. next to each other or stacked. On the other side the images must also contain simple examples to learn the typical relations between the surface patches on the same object.

We created a training set of 27 images together with annotations. All of the images show a tabletop scene, 17 of the images show several boxes and

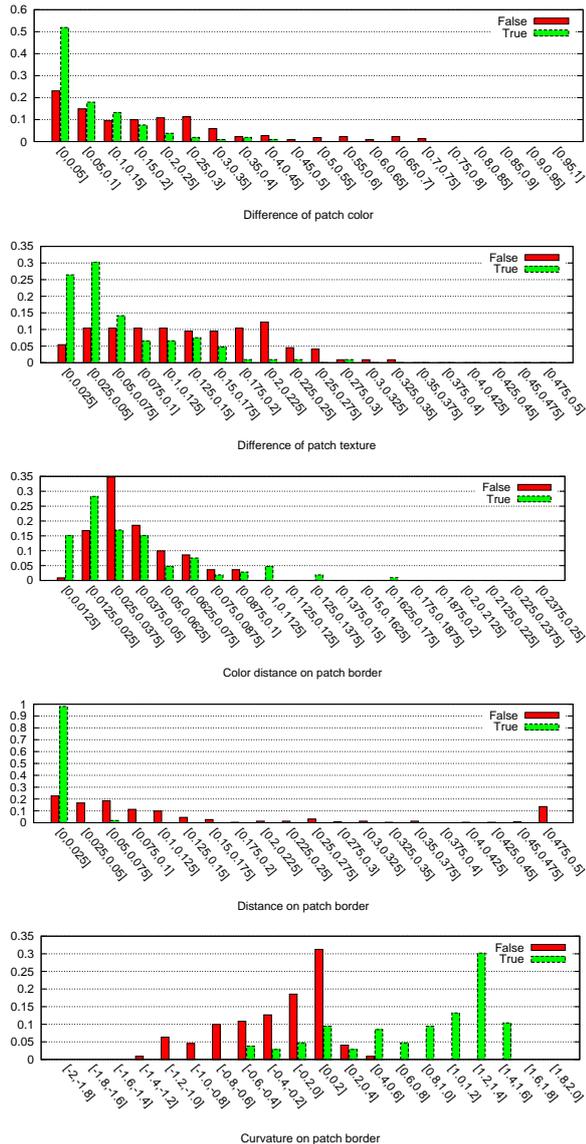


Figure 3. Histograms of relation values for positive and negative examples for relations r_{ch} , r_{tr} , r_{cb} , r_{di} , r_{cu} .

the other ten images show cylindric objects, such as cups, cookie packaging or other kitchenware. The first two rows of Fig. 4 show some of the objects, used for the training period. To give an intuition into the meaning of the different relation values, Fig. 3 shows the distribution of relation values for positive and negative training examples.

The difference of the color histogram is calculated as the Fidelity distance a.k.a. Bhattacharyya coefficient of the UV components in the YUV color space. The comparison of the UV components is less susceptible to brightness changes in the scene than a comparison in the RGB color space. The first histogram in Fig. 3 shows for rising equality (i.e. values close to 0) also a rising number of true examples,

while the negative examples are widely distributed over the histogram.

The texture rate of a surface patch is the rate between the number of canny edge pixels on the surface patch and the sum of all surface patch pixels. The distribution shows that surfaces with similar texture rate tend to belong together, with again the negative examples more spread than the positive ones. Shadows on uniform surfaces are a problem, because they can cause edges and thus fake texture.

The third histogram shows again color distance, but now as mean value between neighbouring points (in the image space) on the border of the patches. The distributions of the positive and negative examples are nearly similar. One reason for this is that a lot of our training objects have different colours on neighbouring surfaces, such as the boxes in the top left image of Fig. 4. Another reason for the weak performance of colour is the wrong assignment of colour to 3D points at occlusion boundaries, which is an artefact of the Kinect sensor we were using.

The distance value along patch borders is calculated as mean distance between neighbouring border points. As expected, the positive examples show a high peak for very small distances. Note that for now we do not learn that an object bisected by an occluder should be treated as one object.

The last histogram of Fig. 3 shows the curvature relation, expressed as the mean angle between neighbouring points along patch borders. It shows that positive (convex) curvature typically indicates same object (e.g. two sides of a box joining) whereas negative (concave) curvature typically indicates two different objects.

The results of the evaluation of the SVM and the graph-cut algorithm of four test sets with 50 images is shown in Tab. 1. Examples of the results are shown in Fig. 4. The first four rows show successful examples, the last two image rows show nearly successful segmentation at the systems limit.

5. Conclusion

We presented a framework for segmenting unknown objects in RGBD-images of cluttered table top scenes, by first approximating surfaces with a combination of planes and NURBS and then segmenting the scene based on learned relations between surface patches. One of the problems we are still facing is the combinatorial explosion when trying to replace several planes (think of a couple of stripes ap-

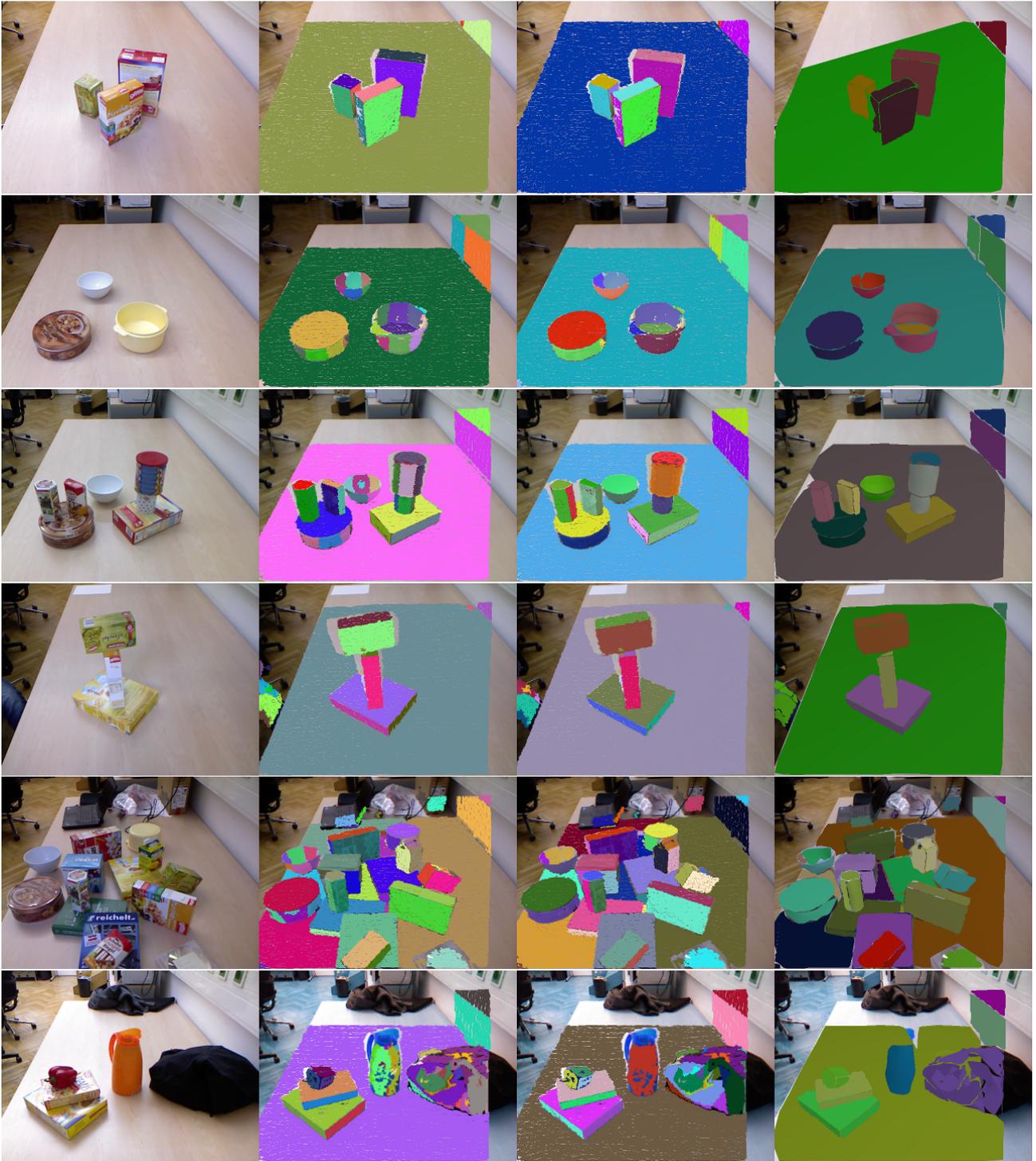


Figure 4. Selected examples of the proposed approach. The first column shows the color image, the second one the estimated planes, the third the plane patches and NURBS after model selection, still as point cloud. The last column shows the segmented object models.

proximating a cylinder) simultaneously with a single NURB, as we would go from pairs to n -tuples with possibly large n . We intend to address this by following not a plane-first-then-NURBS approach, but by first identifying clusters of points forming a sin-

gle smooth surface area, based on curvature derived from normal vectors. To this end we will employ an improved iterative normal estimation scheme [4] and then trying to locally find the optimal combination of planes and NURBS (again using model selection)

Set	Nr.	SVM acc. (#)	u.seg	o.seg
Boxes	14	100% (158)	0.2%	0.6%
St. Boxes	16	91.5% (224)	1.2%	12.2%
Cylinders	11	91.8% (135)	1.8%	9.1%
Mixed	9	84.7% (406)	6.9%	39.0%

Table 1. Results of object segmentation for different test sets with boxes, stacked boxes, cylinders and mixed objects. Columns: Type of set, number of images, accuracy of SVM prediction (number of relation vectors), under-segmentation and over-segmentation of objects.

for that particular surface area rather than globally for the whole scene.

Acknowledgements

The research leading to these results has received funding from the European Communitys Seventh Framework Programme [FP7/2007-2013] under grant agreement No. 215181, CogX.

References

- [1] N. Bergström, M. Björkman, and D. Kragic. Generating object hypotheses in natural scenes through human-robot interaction. In *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, pages 827–833. IEEE, 2011. 2
- [2] K. L. Boyer and S. Sarkar. Perceptual organization in computer vision: status, challenges, and potential. *Comput. Vis. Image Underst.*, 76(1):1–5, 1999. 2
- [3] Y. Y. Boykov and M.-P. Jolly. Interactive graph cuts for optimal boundary amp; region segmentation of objects in N-D images. In *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, volume 1, pages 105 –112 vol.1, 2001. 2
- [4] F. Calderon, U. Ruiz, and M. Rivera. Surface Normal Estimation with Neighborhood Reorganization for 3D Reconstruction. *Progress in Pattern Recognition Image Analysis and Applications*, 4756:321–330, 2007. 7
- [5] C.-c. Chang and C.-j. Lin. LIBSVM : A Library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27:1—27:27, 2011. 5
- [6] J. A. Cottrell, T. J. R. Hughes, and Y. Bazilevs. Iso-geometric Analysis. *Continuum*, page 355, 2010. 3
- [7] P. F. Felzenszwalb and D. P. Huttenlocher. Efficient Graph-Based Image Segmentation. *International Journal of Computer Vision*, 59(2):167–181, Sept. 2004. 5
- [8] M. A. Fischler and R. C. Bolles. Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cortography. *Communications of the ACM*, 24(6), 1981. 3
- [9] G. D. Hager and B. Wegbreit. Scene parsing using a prior world model. *The International Journal of Robotics Research*, 2011. 2
- [10] G. Kootstra, N. Bergström, and D. Kragic. Fast and Automatic Detection and Segmentation of Unknown Objects. In *Humanoids*, Bled, 2011. 2
- [11] G. Kootstra, N. Bergström, and D. Kragic. Gestalt Principles for Attention and Segmentation in Natural and Artificial Vision Systems. In *SPME*, Shanghai, 2011. 2
- [12] A. Leonardis, A. Gupta, and R. Bajcsy. Segmentation of range images as the search for geometric parametric models. *International Journal of Computer Vision*, 14(3):253–277, Apr. 1995. 1, 2, 4
- [13] A. K. Mishra and Y. Aloimonos. Active Segmentation. *I. J. Humanoid Robotics*, 6(3):361–386, 2009. 2
- [14] A. K. Mishra, Y. Aloimonos, L. F. Cheong, and A. Kassim. Active Segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 6(3):361–386, Aug. 2011. 2
- [15] L. Piegl and W. Tiller. *The NURBS book*. Monographs in visual communication. Springer, 1997. 3
- [16] J. Prankl, M. Zillich, B. Leibe, and M. Vincze. Incremental Model Selection for Detection and Tracking of Planar Surfaces. In *Proceedings of the British Machine Vision Conference*, pages 87.1—87.12. BMVA Press, 2010. 4
- [17] C. Rother, V. Kolmogorov, and A. Blake. "Grab-Cut": interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.*, 23(3):309–314, Aug. 2004. 2
- [18] R. Rusu and S. Cousins. 3d is here: Point cloud library (pcl). In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 1–4. IEEE, 2011. 3
- [19] S. Vicente, V. Kolmogorov, and C. Rother. Joint optimization of segmentation and appearance models. *2009 IEEE 12th International Conference on Computer Vision, (Iccv)*:755–762, 2009. 2
- [20] M. Werlberger, T. Pock, M. Unger, and H. Bischof. A Variational Model for Interactive Shape Prior Segmentation and Real-Time Tracking. In *International Conference on Scale Space and Variational Methods in Computer Vision (SSVM)*, Voss, Norway, 2009. 2
- [21] T. Wu and C. Lin. Probability estimates for multi-class classification by pairwise coupling. *The Journal of Machine Learning Research*, 5:975–1005, 2004. 5
- [22] M. Zillich. Incremental Indexing for Parameter-Free Perceptual Grouping. In *31st Workshop of the Austrian Association for Pattern Recognition*, pages 25–32, 2007. 2

Isogeometric Finite-Elements Methods and Variational Reconstruction Tasks in Vision – A Perfect Match

Jonathan Balzer
University of California
Los Angeles, CA 90095
USA
balzer@cd.ucla.edu

Thomas Mörwald
Vienna University of Technology
1040 Vienna
Austria
moerwald@acin.tuwien.ac.at

Abstract

Inverse problems are abundant in vision. A common way to deal with their inherent ill-posedness is reformulating them within the framework of the calculus of variations. This always leads to partial differential equations as conditions of (local) optimality. In this paper, we propose solving such equations numerically by isogeometric analysis, a special kind of finite-elements method. We will expose its main advantages including superior computational performance, a natural ability to facilitate multi-scale reconstruction, and a high degree of compatibility with the spline geometries encountered in modern computer-aided design systems. To animate these fairly general arguments, their impact on the well-known depth-from-gradients problem is discussed, which amounts to solving a Poisson equation on the image plane. Experiments suggest that, by the isogeometry principle, reconstructions of unprecedented quality can be obtained without any prefiltering of the data.

1. Introduction

1.1. Motivation

Vision is dominated by inverse problems in the sense that from an observation, one wishes to make inferences about its cause, e.g., shape from shading aims at computing the shape of a Lambertian object from the gray values it induces on the image plane. Inverse problems often struggle with ill-posedness, meaning that they admit no solution at all, or if they do, it is either ambiguous or does not depend continuously on the input data. A general strategy to deal with this issue is to turn away from classical or strict solutions to ones that are merely optimal with respect to an application-dependent cost. Indeed, we assess that energy minimization is ubiquitous in vision. It can be categorized roughly into two classes: A major line of work, following direct discretization and quantization of the objective func-

tion, combinatorializes the optimization problem and solves it by one of many graph-based algorithms, among which graph cuts have proven to be particularly successful in recent years [7].

A quite different strategy is restricting all considerations to a *continuous* version of the optimization problem. In that case, the energy takes the form of a functional on an infinite-dimensional linear space of functions. The *calculus of variations* is mainly concerned with deriving a condition of (local) optimality for such functionals, their so-called Euler-Lagrange equation [8]. Generally, the latter is a, possibly nonlinear, partial differential equation (PDE) of arbitrary order and as such rarely solvable by analytical means. In most vision-related works, the remedy of choice is approximating the occurring differential operators by weighted differences of the function values on neighboring grid points, thereby transforming the original PDE into a system of algebraic equations [14].

Digital images and the regular arrangement of their pixels provide an ideal computational grid for the finite-differences method (FDM). This is probably why, in the

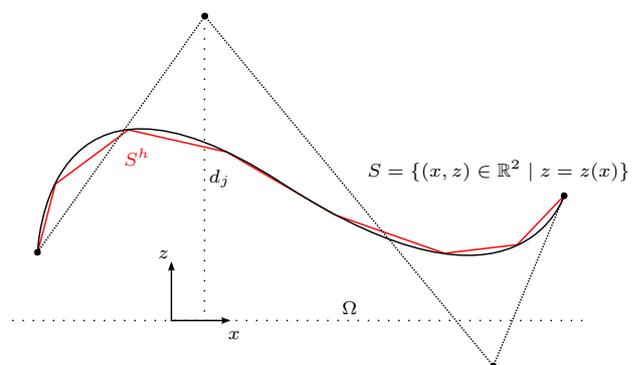


Figure 1. Isogeometric analysis is a finite-elements method to solve partial differential equations defined on a spline surface S (black) with parametric domain Ω . Classical FEMs operate on polygonal meshes such as S^h (red).

realm of variational reconstruction, the finite-elements method (FEM) has been widely neglected, and if not so, it is conducted on unstructured meshes. In this paper, we advocate *isogeometric analysis*, a particular variant of the FEM that operates on B-spline patches. We will argue in Section 2 why it addresses the needs of vision extremely well. To further support our claims and to show the method in action, in Section 3, it is applied to a well-known problem in visual reconstruction, the integration of a gradient field into a depth map. Although just a case study, the proposed algorithm outperforms state-of-the-art ones in terms of efficiency, accuracy, and robustness.

2. Finite-element methods for vision

2.1. Conventional approach

To foster a better understanding of what follows let us briefly summarize the fundamental principles guiding the finite-elements method. Its central idea is to approximate the *candidate solutions* $z : S \rightarrow \mathbb{R}$ of a PDE posed on a domain $S \subset \mathbb{R}^3$, say a surface, by representing them as linear combinations

$$z(x) = \sum_{j=1}^n d_j b_j(x) \quad (1)$$

of some pre-defined basis functions $b_j : S \rightarrow \mathbb{R}$. The coefficients d_j are sometimes called the *degrees of freedom* (DOFs). This way, the search for a minimum of the underlying energy functional is confined to a vector of coefficients $\mathbf{d} = (d_j) \in \mathbb{R}^n$ in a tractable n -dimensional subspace of the original infinite-dimensional function space. In other words, the underlying functional is discretized directly, not its Euler equation. This bears some resemblance with graph-based optimization, only with the difference that images of solutions, mostly the set of real numbers, remain un-quantized.

The FEM is preferred over the FDM especially when geometry and/or topology of S are more complex than e.g. the image plane's because it inevitably couples digital geometry representation with the shape of the b_j . In fact, as illustrated in Figure 1, S is usually replaced by a polygonal approximation S^h , a collection of *finite elements* (FEs) or *mesh*. The functions¹ $b_j : S^h \rightarrow \mathbb{R}$ are then constructed as to model the local behavior of the overall solution (1) by supporting bivariate polynomials locally around elementary geometric entities, like the vertices of the mesh. FEMs admit two modes of refinement: An *h-refinement* increases the local element density. The smoothness of the solution can be controlled by raising the polynomial degree of the basis (*p-refinement*). Let us record in preparation of the next

¹Note that with slight abuse of notation, we do not distinguish the basis functions on S and its approximation S^h .

paragraph that any FEM is characterized by 1. the choice of geometric model for S , and 2. the type of basis functions b_j . For a comprehensive introduction to the classical FEM, we refer to [14].

2.2. Isogeometric analysis and its implications

Contemporary product development consists of a design phase in which the geometry of a part is specified by the engineer, and a phase in which its physical properties (e.g., stress resistance, electrical/thermal conductivity, and others that can be modeled by a PDE) are validated at hand of FE analysis. Often, the two phases are carried out iteratively. All the worse, that at each such cycle, a conversion becomes necessary because – for historical reasons – computer-aided design (CAD) systems and conventional FEMs found on totally different representations of the computational domain. This conversion is an error-prone, costly, and time-consuming process, as it almost certainly requires manual assistance by a trained user. As a remedy, Hughes *et al.* [12] propose to refrain from tessellating the CAD-native spline models into polygonal meshes, the input format presumed by current FE solvers. Their insight is that every CAD model already possesses its own set of basis functions to do FE analysis with. The concept truly deserves the predicate *isogeometric*: Translated from greek, it means that *one and the same* representation, a spline basis, is used for specifying the geometry and solving PDEs on it. To our best knowledge, Elguedj *et al.* [5] are the only authors so far to consider isogeometric analysis (IGA) in a vision-related scenario, specifically for optical flow estimation. However, with their background in material testing of metal sheets, they do not elaborate on its significance for reconstruction and vision in general, the study of which will be the first contribution of this paper. In fact, isogeometric methods convince by the following advantages:

1. **Natural parametrization:** Unlike polygonal meshes, splines carry a natural parametrization, i.e., a bijective mapping from the domain S to some subset Ω of Euclidean space. Consequently, not only the FE analysis itself but all other operations that work on planar images extend to curved surfaces in a straightforward manner, cf. Figure 2(a).
2. **Projective invariance:** Projecting the control points of a NURBS surface first and then evaluating it yields the same result as proceeding in reverse order.
3. **Meaningful priors:** A majority of man-made objects have been designed on a computer, and IGA provides a simple mechanism to embrace this prior knowledge directly into the reconstruction process. Smoothing is not only an avoidable preprocessing step. But it also necessitates scale selection, i.e., the choice of the

“correct” smoothing parameter, which, given a class of CAD surfaces to look for, is almost trivial within the IGA framework, cf. Figure 3.

4. **Direct reverse engineering:** Why initially use splines for synthesis (design) and later point clouds, triangle meshes, or similar for analysis (reconstruction)? The reverse engineering pipeline from optical measurements of a physical prototype back to a digital model can be shortened. IGA makes conversions between different representations obsolete, cf. Figure 2(b) and the video included in the supplemental material.
5. **Multi-scale reconstruction:** Polynomial bases of high degree are realized effortlessly and at any desired scale by parametric splines, whereas constructing quadratic basis functions on polygons is already far from trivial, even on such simple entities as triangles. This is significant in view of the fact that the analysis of images at multiple scales is a fundamental and often-needed technique in image processing and vision.
6. **Accuracy:** The geometric approximation power of polygon meshes is limited. Spline surfaces on the other hand are “continuous” in the sense that, although defined by a finite number of DOFs, they acquire the very shape envisioned by the designer. Furthermore, they can be evaluated *exactly* and at *arbitrary* locations. The same holds true for the results of IGA, which, by definition, are expanded in the same basis as the geometry, thus enabling subpixel resolutions without spending extra interpolation efforts. This becomes particularly important in visualization which unfortunately still requires polygonal representations of both, geometry and solutions. Note, however, that hardware-accelerated direct rendering of splines is on the way, and yet, it is much easier to transform a spline patch into a mesh than vice-versa.
7. **Efficiency and robustness:** IGA features both mechanisms of refinement known from classical FEMs on polygonal surfaces. But opposed to a mesh, a spline surface remains faithful to the original geometry when refined. Loosely speaking, a mesh could need a lot more DOFs to approximate the *geometry* “well” enough than the *solution*. This may result in a substantial computational overhead. The refinement of splines, however, is geometry-independent. Hence, it is possible in principle to tailor the resolution to the requirements of the physical process being modelled by the underlying PDE *not* the geometric model. This has two far-reaching consequences: The content of natural images is known to be concentrated on a much smaller

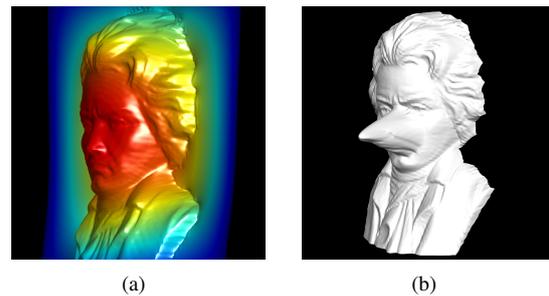


Figure 2. Advantages of the isogeometric approach: (a) PDEs can be solved on patches of arbitrary shape, all with the same framework. This example shows the color-coded solution to a Poisson problem with constant source term under homogenous Dirichlet conditions. (b) The reconstruction is available as a spline patch, readable and editable by all common CAD programs. Beethoven has undergone trimming as well as plastic surgery.

set than the collection of its gray values. Reconstruction based on the FDM always utilizes as many DOFs as there are pixels and is thereby highly redundant. By allocating resources tuned to the part of the data of actual interest, IGA leads to very efficient algorithms, which underlines its potential in real-time applications. This immediately implies the superior robustness of isogeometric methods. As an example, consider the linear case, in which, as we will see later in our case study (Section 3.4), the PDE is transformed into a sparse linear system. The smaller the system matrix, the faster it is to invert, and the better is its condition and hence the numerical stability.

3. A case study: depth from gradients

Numerous computer vision/optical metrology techniques such as photometric stereo, shape from shading, or deflectometry acquire the surface *slope* at points on an unknown object rather than their spatial locations directly. The essence of reconstruction is integrating measured normal or gradient data into a visual surface representation.

3.1. Continuous variational model

A widespread variational formulation of this problem is the following: For the sake of simplicity, we assume that the region of interest on the surface is parameterizable by an orthographic depth map $z : \Omega \rightarrow \mathbb{R}$ over the (image) plane $\Omega \subset \mathbb{R}^2$, see Figure 1. Given the measurement of a vector field $\mathbf{g}_m : \Omega \rightarrow \mathbb{R}^2$, we wish to find the function z whose gradient ∇z is closest to the data \mathbf{g}_m in terms of the squared L^2 -norm, i.e., the function minimizing the Dirichlet-type energy

$$E(z) = \int_{\Omega} \frac{1}{2} \|\nabla z - \mathbf{g}_m\|^2 dx. \quad (2)$$

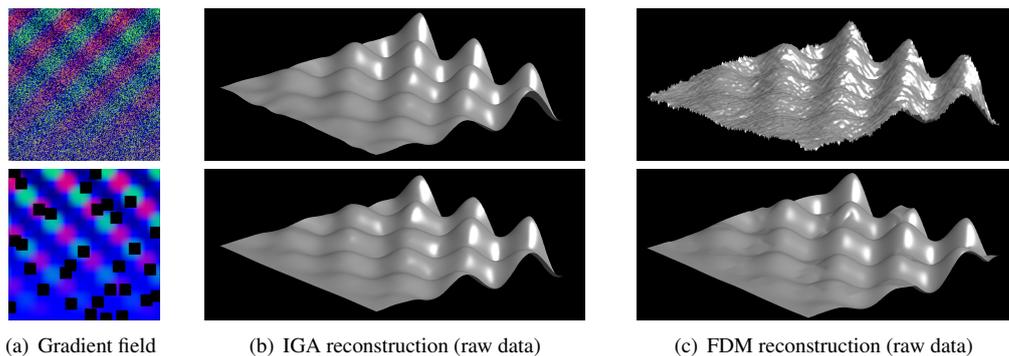


Figure 3. An example from the case study in Section 3: By isogeometric analysis, one can put effective and meaningful constraints on the class of admissible shapes and thereby alleviate ramifications of (a) noisy and incomplete data. Of course, results comparable to (b) ours could also be obtained by (c) the FDM. This would come at the cost of preprocessing the data by suitable smoothing/interpolation.

Even if \mathbf{g}_m is not integrable and no classical solution exists for which equality $\nabla z = \mathbf{g}_m$ holds, the energy (2) implies a projection² onto the curl-free portion of \mathbf{g}_m . This is right in the spirit of Section 1, where we motivated the use of variational models to tackle ill-posed problems. The Euler-Lagrange equation of (2), i.e., a necessary – in the present case even sufficient – condition for optimality, is

$$\Delta z = \operatorname{div} \mathbf{g}_m, \quad (3a)$$

in which Δ denotes the Laplace operator, and div the vector field divergence of \mathbf{g}_m . The unique infinite-dimensional least-squares- or L^2 -solution can be deduced from (3a) if and only if complemented by the condition

$$\langle \nabla z, \hat{\mathbf{o}} \rangle = \langle \mathbf{g}_m, \hat{\mathbf{o}} \rangle \quad (3b)$$

on the image boundary $\partial\Omega$ with outer unit normal $\hat{\mathbf{o}}$. This is the natural boundary condition arising from the variational principle [4]. It affords that z can move freely above $\partial\Omega$ and thus, the surface adapts optimally to the gradient field there³. Note, however, that a solution to above Neumann problem, i.e., Poisson’s equation plus aforementioned boundary condition, is unique only up to a scalar: If some $z(x)$ fulfills (3), obviously, the same holds true for $z(x) + c$ with $c \in \mathbb{R}$ because the constant vanishes under differentiation by Δ on the interior of Ω and the directional derivative $\langle \nabla(z + c), \hat{\mathbf{o}} \rangle$ on the boundary. Since the integration constant c has no influence on the shape of the surface S parametrized by z , the ambiguity can be resolved either by prescribing the distance of a single point in the computational domain or restricting the search for a solution to all depth maps which are mean-free:

$$\int_{\Omega} z(x) dx = 0. \quad (4)$$

²or *Helmholtz-Hodge decomposition* of \mathbf{g}_m

³Opposed to a *Dirichlet* or *essential* boundary condition which fixes the boundary by specifying depth values known in advance (or sometimes even guessed, causing a significant bias, cf. Figure 5(c)).

See [1, Sec. 4.1] for a detailed discussion of both, the boundary and mean-value condition.

3.2. Related work

In his groundbreaking paper [11], Horn shows how to approximate the Laplacian in (3a) by second-order finite differences (FD) on the image grid and solve the resulting algebraic system by a fixed-point scheme. Extensions of Horn’s method are too numerous to list here but let us explicitly mention the most recent ones like Harker’s and O’Leary’s [9] as well as that due to Durou *et al.* [3], who describe a powerful total-variation-based algorithm capable of resolving discontinuities in the depth map without prior segmentation of the gradient field. The inferiority of FD-based methods should be apparent from the discussion in Sections 2.2 and 3.5. Our work relates to the class of kernel methods [6, 15], which can be thought of as mesh-free FEMs in disguise. Similarly, Kovessi applies a basis $\{b_j\}$ of shapelets to the normal adaption problem in scene space [13]. Only a few authors explicitly consider the classical, i.e., non-isogeometric FEM: Hicks employs it for integrating normal fields with three-dimensional support into a foliation of surfaces [10]. Generalizations of Horn’s method applicable to such spatially varying normal fields are presented by Balzer [1] and Delaunoy and Prados [2]. None of aforementioned methods is compatible with the geometry representation of contemporary CAD packages. Higher-degree polynomial bases and a multi-scale mechanism are per se possible, at least on polygon meshes, but quite challenging to implement.

3.3. B-splines

Here, as a prototypical application of IGA in vision and as the second contribution of this paper, we present the first FE method for gradient field integration based on *B-splines*. In order not to cloud the key ideas by a complicated index calculus, let us assume for the remainder of

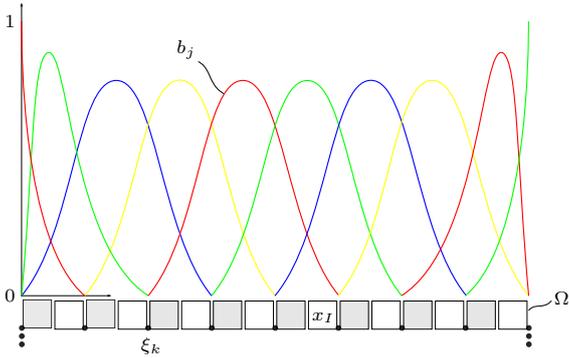


Figure 4. The knots ξ_k , indicated by black dots, uniquely define the B-spline basis functions b_j on the image plane Ω , consisting of pixels x_I (white and gray squares).

the theoretical discussion that, with slight abuse of earlier notation, $\Omega \subset \mathbb{R}$. Everything applies mutatis mutandis to the two-dimensional case. A B-spline is a scalar-valued function of the form (1) with the elements $b_j(x)$ of the basis being compactly supported polynomials defined by the so-called *Cox-de Boor recursion*. An indispensable ingredient of the recursion formula is the fixed *knot vector* $\xi = (\xi_1, \dots, \xi_k, \dots, \xi_{n+p+1}) \in \mathbb{R}^{n+p+1}$. It may be understood as a set of real parameters in calculating the value $b_j(x)$, which we could (but for the sake of notational simplicity will not) indicate by writing $x \mapsto b_j(x; \xi)$. The knot vector tiles the parameter domain Ω into smaller intervals (the *finite elements* of IGA), restricted to which, the spline is a polynomial of degree p . Repetition of knots is allowed and leads to a decrease in smoothness. In particular, the spline is interpolating at some ξ_k if and only if ξ_k appears $p+1$ times in ξ . Throughout the paper, we assume so-called *open* knot vectors which are interpolating at both ends ξ_1 and ξ_{n+p+1} . Multivariate basis functions can be constructed from tensor products of B-splines in a single variable, the respective Cartesian coordinate of the domain Ω . To extend the image of (1) to higher dimensions, say $d \in \mathbb{N}$, one simply chooses coefficient *vectors* $\mathbf{d}_j \in \mathbb{R}^d$. At modeling curves, surfaces, or solids in \mathbb{R}^3 , these are commonly referred to as *control points*. The knot vector is not to be confused with the much finer pixel grid $\mathbf{x}_h = (x_1, \dots, x_I, \dots, x_N) \in \mathbb{R}^N$, which is staggered with respect to ξ and contains no redundant abscissae. The geometric relationship between both is illustrated in Figure 4, alongside with the second-degree B-spline basis induced by the shown knot vector. *h*-refinement is established by knot insertion. Note that, since the length of ξ is by definition $n + p + 1$, during this process, the dimension n of the approximation space must grow with p remaining constant. The effect of degree elevation upon n and the knot vector is less obvious. A detailed exposition of *p*-refinement, B-splines, and their non-uniform rational generalization NURBS can be found in [16].

3.4. Isogeometric discretization

Multiplication with a test function φ and integration over the domain Ω brings (3) into a variational form, the starting point of any FEM. The equation is said to hold in a weak sense if

$$\int_{\Omega} \Delta z \varphi dx = \int_{\Omega} \operatorname{div} \mathbf{g}_m \varphi dx \quad (5)$$

for arbitrary φ in the Sobolev space $H_0^1(\Omega)$, and a fortiori, for all elements of the spline basis i.e. $\varphi \in \{b_1, \dots, b_i, \dots, b_n\}$. Integration by parts and application of the Gauss theorem conveniently relax the differentiability assumptions on both, the solution z and the measured gradient field \mathbf{g}_m , by one order:

$$\int_{\Omega} \langle \nabla z, \nabla b_i \rangle dx = \int_{\Omega} \langle \mathbf{g}_m, \nabla b_i \rangle dx + \int_{\partial\Omega} \langle \nabla z - \mathbf{g}_m, \hat{\mathbf{o}} \rangle b_i dx \quad (6)$$

for all $i \in \{1, \dots, n\}$. Recall that in the strong form (3a), the divergence operator is applied to the data, and that this differentiation amplifies contained noise. Suppose that the boundary condition (3b) holds⁴, then the last integral vanishes. The actual discretization is performed by inserting the representation (1) of a candidate solution. The unknown function $z(x)$ is reduced to a finite sequence of unknown coefficients d_j . Thanks to linearity, integration can be restricted to known functions only so that we obtain a linear system $\mathbf{K}\mathbf{d} = \mathbf{f}$ coupling with each other *stiffness matrix* $\mathbf{K} = ((K_{ij}))$, *displacement* $\mathbf{d} = (d_j)$, and *force vector* $\mathbf{f} = (f_i)$. The terminology originates from linear elasticity being the primary application of early FEMs and is appropriate regardless of the physical background. In view of (6), we get

$$K_{ij} = \int_{\Omega} \langle \nabla b_i, \nabla b_j \rangle dx, \quad f_i = \int_{\Omega} \langle \mathbf{g}_m, \nabla b_i \rangle dx. \quad (7)$$

These integrals are typically computed by Gauss quadrature. In particular, we found the midpoint rule to be sufficient here, which yields, for any two pixels being unit length apart, a simple summation over the grid points x_I in \mathbf{x}_h :

$$K_{ij} \approx \sum_{I=1}^N \langle \nabla b_i(x_I), \nabla b_j(x_I) \rangle, \quad (8a)$$

$$f_i \approx \sum_{I=1}^N \langle \mathbf{g}_m(x_I), \nabla b_i(x_I) \rangle. \quad (8b)$$

⁴As a caveat, this must be taken care of later by appropriate modification of the stiffness matrix.

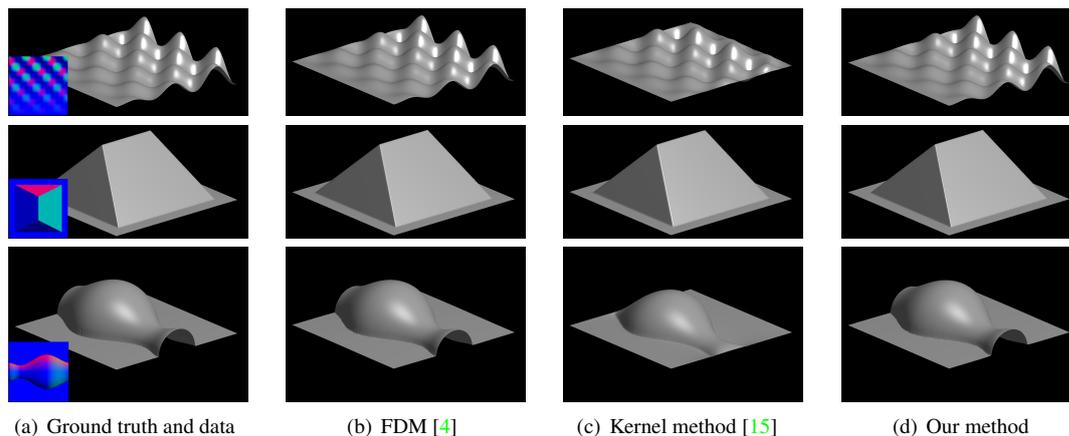


Figure 5. (a) Ground truth surfaces and gradient fields \mathbf{g}_m in false-color representation: The components of corresponding unit normals directly map to the values of the three RGB channels. (b)-(d) Reconstruction results for different numerical methods to solve the Euler-Lagrange equation (3).

	N	s	p	#DOFs	RMSE [px]
Waves	256	6	3	4489	$7.6 \cdot 10^{-5}$
Tent	256	8	1	66049	$6.2 \cdot 10^{-4}$
Vase	256	8	1	66049	$3.3 \cdot 10^{-2}$
Paraboloid	256	2	2	16	$1.3 \cdot 10^{-14}$
Beethoven	256	7	2	16900	$1.3 \cdot 10^{-1}$

Table 1. Parameters of conducted experiments: N is the size of the input data in one dimension, s the scale, and p the polynomial degree of the spline patch. The number of unknowns is denoted by #DOFs. RMSE stands for the root mean square error between the true and reconstructed gradient fields.

The field \mathbf{g}_m is given per pixel x_I . The b_i , b_j , and all of their derivatives can be evaluated at arbitrary locations by the Cox-de Boor formula. Because of their compact support, \mathbf{K} is sparse and thus efficient to invert.

Note that we actually solve the Poisson problem (3) on a planar (but not necessarily degree one) B-spline $\tilde{S} = \Omega$, initially coinciding with the image plane. This simplifies the integrals (7) because the coordinate transformation from Ω to \tilde{S} is just the identity map. A solution \mathbf{d} of $\mathbf{K}\mathbf{d} = \mathbf{f}$ then defines the soughtafter shape S as follows: Since \tilde{S} and S share a common parameter domain Ω , we can simply move the control points of the original patch \tilde{S} by the entries d_j of \mathbf{d} in z -direction, cf. Figures 1 and 6.

3.5. Numerical tests

We created a MATLAB implementation of the numerical approach presented in the previous section. In our effort to support reproducibility of research, all data sets and code will be made available online upon publication. The configurations of our IGA solver during various experiments are summarized in Table 1: All images used in the study were

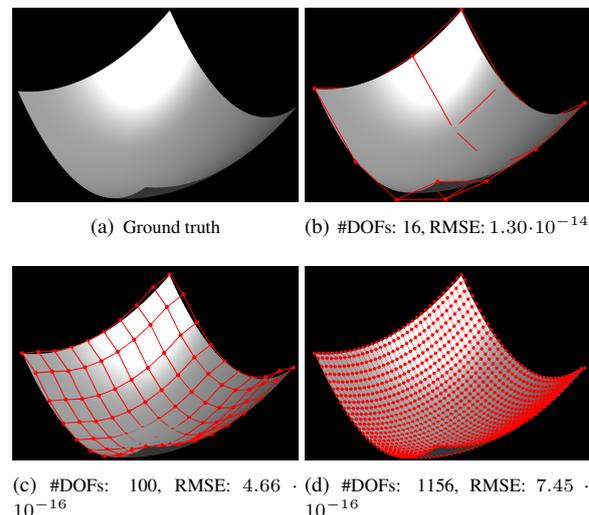


Figure 6. Scale and polynomial degree can be optimally adapted to the data in the isogeometric approach. The reconstruction in (b) remains close to ground truth utilizing as little as 16 unknown coefficients.

of size $N \times N$. We define the scale as the integer $s \in \mathbb{N}$ such that the knot vector consists of $2^{-s} \cdot N$ intervals, e.g., we have $N = 16$ and $s = 1$ in Figure 4. The root mean square error (RMSE) essentially equals (2) normalized with the area of Ω .

Besides the standard FDM, we fed our data into the kernel method (KM) described in [15], which we chose, because it has been published recently, possesses an FEM-flavor, and its implementation is available for download. The results on a first set of examples is shown in Figure 5: Coordinate axes are suppressed because by (4), we are most interested in the shape but less in the exact size of the object. The z -direction should be obvious from each of the

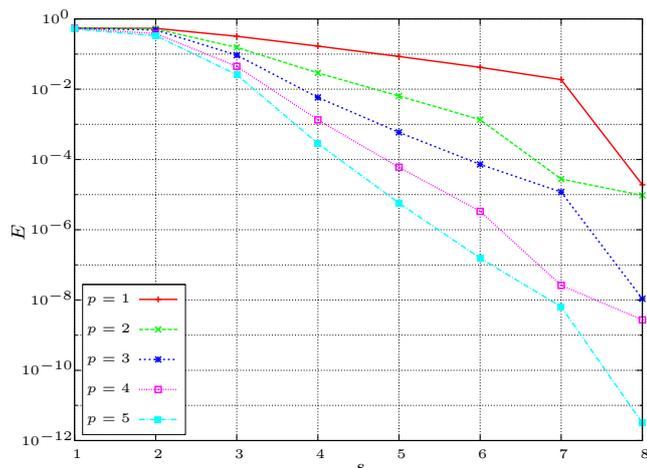


Figure 7. Reconstruction error in dependence on scale s and polynomial degree p for the waves data set.

		t_{Assembly} [s]	t_{Solve} [s]	#DOFs	RMSE [px]
Waves	FDM	189.39	3.94	256×256	$2.8 \cdot 10^{-3}$
	KM	562.63	338.65	221952	$2.9 \cdot 10^{-1}$
	IGA	22.74	0.82	4489	$7.6 \cdot 10^{-5}$
Tent	FDM	189.39	4.23	256×256	$5.6 \cdot 10^{-2}$
	KM	559.92	209.58	221952	$7.7 \cdot 10^{-2}$
	IGA	40.92	8.72	66049	$1.3 \cdot 10^{-1}$
Vase	FDM	189.39	3.84	256×256	$5.3 \cdot 10^{-1}$
	KM	557.81	234.85	221952	$2.7 \cdot 10^{-1}$
	IGA	42.12	10.19	66049	$3.3 \cdot 10^{-2}$
Paraboloid	FDM	189.39	4.34	256×256	$3.5 \cdot 10^{-4}$
	KM	564.23	357.15	221952	$5 \cdot 10^{-1}$
	IGA	17.35	0.01	16	$1.3 \cdot 10^{-14}$
Beethoven	FDM	189.39	4.18	256×256	$5.5 \cdot 10^{-1}$
	KM	557.83	245.4	221952	$3.7 \cdot 10^{-1}$
	IGA	27.17	3.53	16900	$1.3 \cdot 10^{-1}$

Table 2. Quantitative comparison of algorithms: t_{Assembly} is the time needed for equation assembly, t_{Solve} the time for solving the resulting linear system. The value of t_{Assembly} for our implementation is rather pessimistic as assembling the stiffness matrix involves many loops which are handled slowly by the MATLAB kernel.

original surfaces in Figure 5(a). The sinusoidal wave example clearly confirms that unlike previous works, we do not require the unknown surface to be periodic by enforcing the correct boundary condition (3b). It is evident from Figure 5(c) that the KM fails to do so. The visible impression is that all three algorithms resolve the discontinuities in the tent gradient field quite well. A quantitative assessment shows that ours, probably due to its non-local nature, performs slightly worse, see Table 2. The overall reconstruction time $t_{\text{Assembly}} + t_{\text{Solve}}$ does not admit general statements about the performance, it does however allow a direct comparison of algorithms under equal environmental conditions. All ground truth depth maps and their derivatives

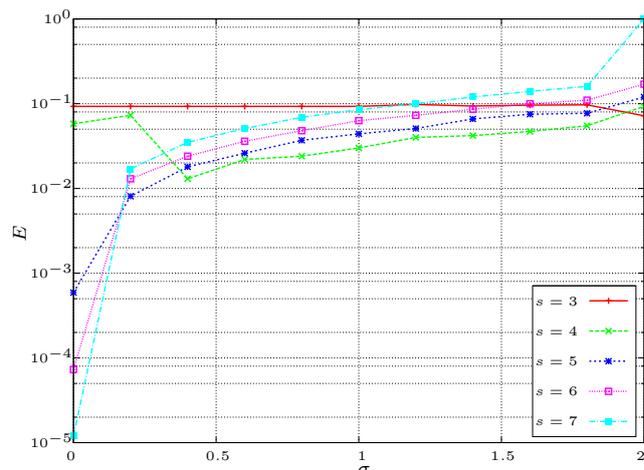


Figure 8. Residual of the energy (2) over the variance σ of noise imposed on the gradient field. These curves were obtained for reconstructions of the waves data set at scales ranging from 3 to 7.

arose from analytical formulas, except for the synthetic vase, which exhibits infinite gradients unless these are estimated numerically. In fact, we assume that there are no occlusions because we believe that they should be avoided in the measurement process anyway. Recovering information that is just not present in the data by the numerical reconstruction algorithm seems at least somewhat questionable. Still, the isogeometric discretization principle naturally extends to discontinuous depth maps. Only energy functional and optimization method have to be adapted. Section 2.3 of the technical report attached as supplemental material outlines how to proceed.

The outcome of additional experiments back up the claims of Section 2.2: Growing computational demands of the FDM can only be met by naive downsampling of the data. Consider on the other hand the paraboloid in Figure 6, which is known to have degree 2. Thus, in theory, one should be able to express it by only a few quadratic B-splines. Consequently, the number of DOFs for computation by IGA can be reduced to the minimum *without ever affecting the data*. The p/s -diagram in Figure 7 substantiates this further: Since the depth map of the sinusoidal wave can be decomposed into an infinite power series, the reconstruction quality can only grow with s and p . Furthermore, we investigated the influence of two common disturbances of g_m upon the quality of reconstruction, Gaussian noise and clutter, see Figure 3(a). As predicted earlier, opposed to the FDM, our method turns out to be highly robust with respect to both, provided that scale and polynomial degree have been set appropriately. How this can be done in tune with the knowledge about the scene that the image portrays is beyond the scope of this paper. Also, it is fair to say in regard to Figure 3, that competitive results are not impossible to obtain by other methods. The point is that this would

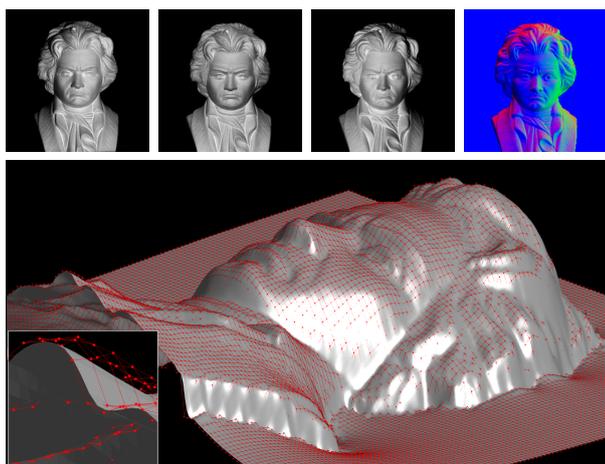


Figure 9. The Beethoven data set (first three images in the top row) was taken from [17]. Gradients (top right) were obtained exploiting the photometric stereo effect and then integrated into a quadratic B-spline patch (bottom). The magnified excerpt shows that the control point lattice is in general not interpolating.

require more or less expensive prefiltering, whereby again, the data cannot remain untouched. Note that the curves in Figure 8 at $\sigma = 0$ and $\sigma = 2.0$ are in reverse order with respect to the error measure E . A simple explanation is that, while for flawless data, higher-order polynomials lead to better approximations, the converse is true in the presence of high-frequency noise which will also be integrated at finer resolutions, cf. top row of Figure 3(c). The result of a final experiment on real-world data is depicted in Figure 9.

4. Open problems

In this paper, we identified isogeometric analysis a general strategy to deal with Euler-Lagrange equations in computer vision and demonstrated its advantages at hand of the depth-from-gradients problem. Although a very promising numerical tool, we would like to point out in closing two particular limitations it entails: First, the construction of a multivariate basis from tensor products of univariate B-splines prohibits a purely local refinement, for inserting a knot in one direction always alters all orthogonal translates of the modified basis function. So-called *T-splines* have been developed to address this issue and found their way into commercially available CAD software. Second, in order to remove the assumption that the unknown surface is parametrizable by a depth map and thus warrant reconstruction in scene space, we have to be able to conduct IGA on surfaces that consist not of one but rather a whole series of patches, not necessarily watertight. This problem has not seen a satisfactory solution, yet, but in our opinion, mathematical techniques such as domain decomposition are quite promising in this respect.

Acknowledgements

This work was conducted while the first author worked at King Abdullah University of Science and Technology. The research of T. Mörwald leading to these results has received funding from the European Community's Seventh Framework Programme [FP7/2007-2013] under grant agreement No. 215181, CogX.

References

- [1] J. Balzer. A Gauss-Newton Method for the Integration of Spatial Normal Fields in Shape Space. *J. Math. Imaging Vis.* In press. 4
- [2] A. Delaunoy and E. Prados. Gradient Flows for Optimizing Triangular Mesh-based Surfaces: Applications to 3D Reconstruction Problems Dealing with Visibility. *Int. J. Comput. Vision*, 95:100–123, 2011. 4
- [3] J.-D. Durou, J.-F. Aujol, and F. Courteille. Integrating the Normal Field of a Surface in the Presence of Discontinuities. *Proc. EMMCVPR*, pages 261–273, 2009. 4
- [4] J.-D. Durou and F. Courteille. Integration of a Normal Field without Boundary Condition. *Proc. PACV*, 1, 2007. 4, 6
- [5] T. Elguedj, J. Réthoré, and A. Buteri. Isogeometric analysis for strain field measurements. *Comput. Methods Appl. Mech. Engrg.*, 200:40–56, 2010. 2
- [6] S. Ettl, J. Kaminski, M. Knauer, and G. Häusler. Shape reconstruction from gradient data. *Appl. Optics*, 47(12):2091–2097, 2008. 4
- [7] P. Felzenszwalb and R. Zabih. Dynamic Programming and Graph Algorithms in Computer Vision. *IEEE T. Pattern Anal.*, 33:721–740, 2011. 1
- [8] I. Gelfand and S. Fomin. *Calculus of Variations*. Dover, 2003. 1
- [9] M. Harker and P. O’Leary. Least squares surface reconstruction from gradients: Direct algebraic methods with spectral, Tikhonov, and constrained regularization. *Proc. CVPR*, 1:2529–2536, 2011. 4
- [10] R. Hicks. Designing a mirror to realize a given projection. *J. Opt. Soc. Am. A*, 22(2):323–330, 2005. 4
- [11] B. Horn. Height and gradient from shading. *Int. J. Comput. Vision*, 5(1):37–75, 1999. 4
- [12] T. Hughes, J. Cottrell, and Y. Bazilevs. Isogeometric analysis: CAD, finite elements, NURBS, exact geometry and mesh refinement. *Comput. Methods Appl. Mech. Engrg.*, 194(39-41):4135 – 4195, 2005. 2
- [13] P. Kovési. Shapelets correlated with surface normals produce surfaces. *Proc. ICCV*, 2:994–1001, 2005. 4
- [14] K. Morton and D. Mayers. *Numerical Solution of Partial Differential Equations*. Oxford University Press, 2005. 1, 2
- [15] H.-S. Ng, T.-P. Wu, and C.-K. Tang. Surface-from-gradients without discrete integrability enforcement: A gaussian kernel approach. *IEEE T. Pattern Anal.*, 32:2085–2099, 2010. 4, 6
- [16] L. Piegl and W. Tiller. *The NURBS Book*. Springer, 1996. 5
- [17] W. Snyder. NC State University Image Analysis Laboratory Database, 2002. 8

Self-Monitoring to Improve Robustness of 3D Object Tracking for Robotics

Thomas Mörwald, Michael Zillich, Johann Prankl and Markus Vincze

Automation and Control Institute, Vienna University of Technology, AT

Abstract—In robotics object tracking is needed to steer towards objects, check if grasping is successful, or investigate objects more closely by poking or handling them. While many 3D object tracking approaches have been proposed in the past, real world settings pose challenges such as automatically detecting tracking failure, real-time processing, and robustness to occlusion, illumination, and view point changes. This paper presents a 3D tracking system that is capable of overcoming these difficulties using a monocular camera. We present a method of Tracking-State-Detection (TSD) that takes advantage of commercial graphics processors to map textures onto object geometry, to learn textures online, and to recover object pose in real-time. Our system is able to handle 6 DOF object motion during changing lighting conditions, partial occlusion and motion blur while maintaining an accuracy of a few millimetres. Furthermore using TSD we are able to automatically detect occlusions or whether we lost track, and can then trigger a SIFT-based recognition system that is trained during tracking to recover the pose. Evaluations are presented in relation to ground truth pose data and examples present TSD on real-world scenes presented in video sequences.

I. INTRODUCTION

Robotic object grasping requires to determine the object's pose and to track the object during the approach with sufficient accuracy (Figure 1, left). After grasping it is necessary to confirm if the grasp was successful and is stable [1], i.e. the object moves together with the end effector without slipping (Figure 1, right). For learning about physical behaviour of objects as in [2–5] the robot has to observe its motion (Figure 1, middle). Again accuracy of tracking but also detecting whether the tracked object has been lost, for example after toppling over, are important to decide whether a certain trajectory should be taken into account for learning. For a robot operating in a complex unpredictable environment, the challenge is to develop a tracking method that is robust to different lighting conditions, partial occlusion, and motion blur.

Today this is achieved best by model-based tracking of objects and numerous solutions using different feature types, models and mathematical frameworks have been developed, where the today's computational power allows for several real-time solutions. However, practical application of these methods is often limited for various reasons. For example, some methods report good results, without giving actual numbers on accuracy [6–9]. Others are capable of handling partial occlusion or changing lighting conditions [9–12] but can not differentiate between deteriorating tracking conditions and lost tracks. Some methods are restricted in their degrees of freedom, e.g. 140 degrees of rotation as in [11],

require off-line learning [10] or are limited to either textured [13, 14] or low-textured objects [15]. Also recovery from lost tracks is rarely handled with a few exceptions [13, 14], which are tracking-by-detection approaches.

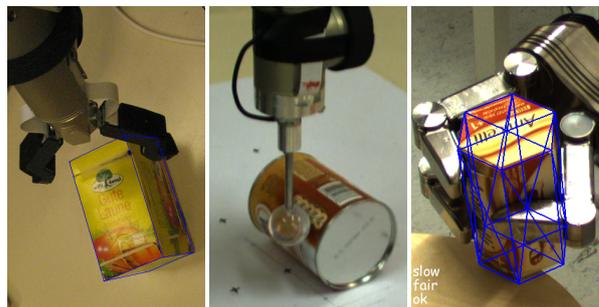


Fig. 1. Tracking for robotic applications. Left: grasping; middle: learning about object motion; right: grasp stability.

Another requirement in robotics is computational efficiency to react to observed situations in time. Consider again the grasping scenario, where we want to use visual servoing to adapt the grasping movement on-line. Hence, we require real-time performance, i.e. processing time within the frame rate of a typical camera (25-50 Hz).

To meet all these requirements we propose to tackle the core problem of detecting tracking failure and take advantage of this supervisory knowledge to achieve automatic object tracking using texture mapping, pose recovery and online learning. Hence, the approach is based on the following methods:

- *Tracking-State-Detection (TSD)*: To know whether we are tracking correctly, whether the object is occluded or whether we lost track we employ our novel TSD method. The knowledge of the tracking state, including speed and confidence of tracking, allows for triggering online learning or pose recovery.
- *Texture mapping*: We take advantage of texture, if available, to boost robustness of tracking, especially in cluttered scenes.
- *Pose recovery*: To initialise tracking and recover lost tracks we use distinctive features placed on the object's surface.
- *Online learning*: We learn these feature points and surface texture of the object automatically while tracking.

Our main contribution is the TSD, since it is the key to use the other methods automatically.

The paper proceeds as follows: In Section III-A we formulate tracking as particle filtering using a modified version of the Sequential Importance Resampling (SIR) filter and show how to draw observations by projecting the model into image space. Section III-B describes how to evaluate the particle weights from observations. Section III-C introduces TSD to give evidence of the current tracking quality, speed and whether tracking has lost the object. In Section III-D we show how surface texture of a tracked object can be captured online from the camera image. Section IV briefly explains what methods we use for initialisation and re-detection of the object. In Section V we evaluate our approach with respect to the requirements established above.

Additionally to the results presented in the paper we provide a video in Section V-E to demonstrate robustness and especially our novel Tracking-State-Detection.

II. RELATED WORK

Tracking the pose of an object by analysing a stream of TV images in real-time goes back to the early eighties [16, 17]. One of the first successful approaches of tracking objects based on edges was the RAPiD system [18]. It used points on model edges and searched for corresponding image edges the edge gradient. Subsequent approaches aimed at improving robustness in tough real-world scenarios [6, 7, 10, 19, 20]. Approaches based on globally matching model primitives with primitives extracted from the camera image [21–25] have been used for applications such as robot and car tracking, but were later replaced by improved versions of the RAPiD type.

[9] also use edges and textures for tracking. Their approach extracts point features from surface texture and use them together with edges to calculate object pose. This turns out to be very fast as well as robust against occlusion. Our approach not only uses patches but the whole texture, which usually lets the pose converge very quickly to the accurate pose. Since the algorithm runs on the GPU, it is as fast as the method in [9]. The work presented in [15] uses edge features to track but does not take into account texture information. This makes it less robust against occlusion. Since the search area in that approach is very small, it is also less robust against fast movement and gets caught in local minima.

More recent approaches aim to solve most of the problems of tracking, such as [12] where the authors are matching the camera image with pre-trained keyframes and then minimizing the squared distance of feature points taking into account neighbouring frames. The approach described in [11] uses a modified version of the Active Appearance Model which allows for partial and self occlusion of the objects and for high accuracy and precision. In [26] the authors minimize the optical flow resulting from the projection of a textured model and the camera image. To compensate for shadows and changing lighting they apply an illumination normalization technique.

In [27] the authors introduce real-time tracking to robotic manipulation. They are using the method proposed in [28], where they project the CAD model into image space, and try

to minimize a cost functional for the distance to image edges found along the gradients of the edges of the model. The work presented in [29] describes an approach for real-time visual servoing using a binocular camera setup to estimate the pose by triangulating a set of feature points. As in our approach [13] takes advantage of robust Monte Carlo particle filtering to determine the pose of the camera with respect to SIFT features, which are localized in 3D using epipolar geometry.

Missing in all methods is to detect when tracking fails rather than reporting tracking trapped in a local optimum. The proposed TSD proposes to solve this and we develop the approach to make it work automatically.

III. TRACKING

The work in this paper identifies the object by using colour and edge information from shape and texture. We project a model, typically consisting of triangles or quads with attached texture, into image space and compare it against the camera image. The pose is estimated using a modified version of the Sequential Importance Resampling (SIR) particle filter [30]. Image processing methods such as Gaussian smoothing and edge extraction as well as pixel-wise comparison of the projected model is accelerated using a typical graphics processing unit (GPU). We first introduce pose estimation and the measure to obtain confidence values from the image data before we explain TSD.

A. Pose estimation

Visual observation of the trajectory of the object is the problem of finding the transformations \mathbf{T}_t given a sequence of images I_t , sampled over the time $t = [1 \dots t_e]$. The transformations \mathbf{T}_t are represented as

$$\begin{aligned} \mathbf{T}_t(\mathbf{x}_t) &= \begin{bmatrix} \mathbf{R}_t & \mathbf{p}_t \\ \mathbf{0} & 1 \end{bmatrix} \\ \mathbf{R}_t &= \mathbf{R}_t(\alpha, \beta, \gamma) \\ \mathbf{p}_t &= \mathbf{p}_t(x, y, z) \end{aligned}$$

where $\mathbf{R}_t(\alpha, \beta, \gamma)^T$ are rotation matrices and $\mathbf{p}_t = [x, y, z]^T$ translations respectively. This results in a state vector $\mathbf{x}_t = [x, y, z, \alpha, \beta, \gamma]^T$ of 6 DOF. Note that we actually use quaternions to avoid the problems of rotations in Euclidean space.

A particle filter, such as the SIR (Sequential Importance Resampling), explained in [31] and more detailed in [30], estimates the current state \mathbf{x}_t based on the previous state \mathbf{x}_{t-1} and the current observation \mathbf{y}_t . Starting from the Bootstrap Filter in [30], Algorithm 1 describes our modified version.

The first modification lies in in step 2a, where we adjust system noise Ω according to the confidence of the previous tracking step c_{t-1} . This means that as the confidence of the particles increases, their degree of distribution decreases, leading to faster convergence and less jitter. Note that from Equation (1) it follows that we do not use a physical motion model. Given the requirements for tracking accuracy and speed for a typical table top scenario we chose a basic

Algorithm 1 Bootstrap Filter, modified with respect to importance sampling.

- 1) Initialisation
 - a) For $i = 1, \dots, N$, sample $\mathbf{x}_0^i \sim p(\mathbf{x}_0)$ and set $t = 1$.
- 2) Importance sampling
 - a) For $i = 1, \dots, N$, sample $\tilde{\mathbf{x}}_t^i \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}^i, c_{t-1})$ with

$$p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i, c_{t-1}) \sim \begin{cases} \Omega(\mathbf{x}_{t-1}, \sigma_{t-1}^2) \\ \sigma_t = (1 - c_t)\sigma_0 \end{cases} \quad (1)$$
 - b) For $i = 1, \dots, N$ of $\tilde{\mathbf{x}}_t^i$, evaluate the confidences

$$c_t^i \sim p(\mathbf{y}_t | \tilde{\mathbf{x}}_t^i) \quad (2)$$
 using Equation (6).
 - c) Normalize the confidence values for the importance weights

$$w_t^i = \frac{c_t^i}{\sum_{i=0}^N c_t^i} \quad (3)$$
- 3) Selection step
 - a) Resample with replacement N particles \mathbf{x}_t^i from the set $\tilde{\mathbf{x}}_t^i$ according to the importance weights.
 - b) Set $t = t + 1$ and go to step 2

standard deviation $\sigma_{0,p}$ of 0.03 m for the translational and $\sigma_{0,\theta} = 0.5$ rad for the rotational degrees of freedom.

The second modification, as proposed already in our previous works [32] and [33], is to use iterative particle filtering for increased responsiveness to rapid pose changes. This means that we perform steps 2 and 3 of Algorithm 1 several times on the same image. Figure 2 shows the improvement over conventional particle filtering when using $k = 8$ iterations with $N = 100$ particles each vs. 1 iteration with 800 particles. It can be seen that the iterative version follows the motion much faster.

To initialise the pose \mathbf{x}_0 we use the method described in Section IV.

B. Image Processing and Matching

At time-step t for each particle i , we project the geometric model of the object, described by vertices, faces and textures, into the image space using the transformation \mathbf{T}^i and standard techniques of computer graphics such as perspective transformation and texture mapping. In image space we compute the edges of the model \mathbf{g}_M^i and of the image captured by the camera \mathbf{g}_I^i .

For each point (u, v) on the model M in image space we can compute the deviation of the gradients by superimposing the projected model over the image. The match m^i of a particle is defined as the sum of the differences of the gradients, and s_i is a normalising constant given by the sum

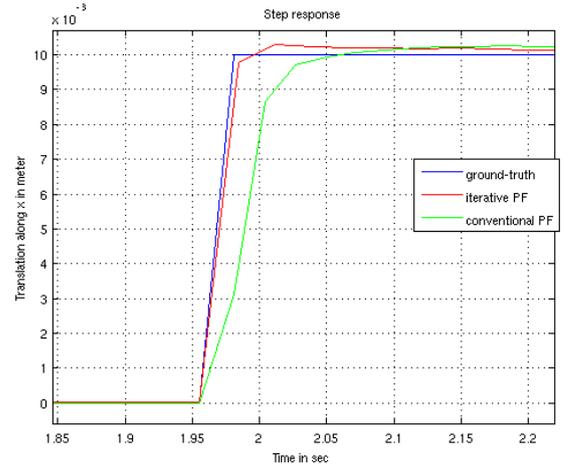


Fig. 2. Step response (1 cm) of conventional and iterative particle filtering using the same amount of particles within one frame, 1x800 and 8x100 respectively.

over all model gradients.

$$\begin{aligned} m^i &= \sum_{(u,v) \in M} |\mathbf{g}_M^i(u, v) - \mathbf{g}_I^i(u, v)| \\ s^i &= \sum_{(u,v) \in M} |\mathbf{g}_M^i(u, v)| \end{aligned} \quad (4)$$

Instead of computing the difference of gradients, the difference of the colour with respect to the hue in HSV (Hue, Saturation, Value) colour space is used:

$$\begin{aligned} m^i &= \sum_{(u,v) \in M} |h_M^i(u, v) - h_I^i(u, v)| \\ s^i &= \sum_{(u,v) \in M} |h_M^i(u, v)| \end{aligned} \quad (5)$$

where h_M^i and h_I^i are the hue values of the projected model and the image respectively. The advantage of using colour based tracking is increased robustness against edge based clutter. Of course it is less robust against changing lighting but the combination of both kinds of cues can significantly improve the overall performance.

We now define the confidence c^i of a particle \mathbf{x}^i as

$$c^i = \frac{1}{2} \left(\frac{m^i}{s^i} + \frac{m^i}{\frac{1}{N} \sum_{j=1}^N s^j} \right) \quad (6)$$

where the first term is simply the match normalised with respect to s_i per particle and the second term is normalised with respect to the mean over all particles, de-weighting particles with a low number of pixels. This prevents the system from getting stuck in poses with a small number of pixels.

The overall confidence of the current observation t is then calculated by simply taking the mean of the confidences

$$c_t = \frac{1}{N} \sum_{i=1}^N c^i \quad (7)$$

C. Tracking-State-Detection (TSD)

As outlined above observing the current state of the tracker is important for assessing the validity of the output as well as allowing to trigger recovery from lost tracks. TSD is a mechanism that gives evidence of tracking speed, quality and overall state in a qualitative and quantitative manner.

1) *Speed*: The velocity is calculated as the first derivative of the translation \mathbf{p}_t and rotation θ_t of the pose respectively.

$$\begin{aligned} \dot{\mathbf{p}}_t &= \begin{bmatrix} \frac{dx}{dt}, \frac{dy}{dt}, \frac{dz}{dt} \\ \frac{d\alpha}{dt}, \frac{d\beta}{dt}, \frac{d\gamma}{dt} \end{bmatrix} \\ \dot{\theta}_t &= \begin{bmatrix} \frac{d\alpha}{dt}, \frac{d\beta}{dt}, \frac{d\gamma}{dt} \end{bmatrix} \end{aligned} \quad (8)$$

We first apply a low-pass filter to remove noise and then normalise $\dot{\mathbf{p}}_t$ and $\dot{\theta}_t$

$$\begin{aligned} \tilde{v} &= \max(v, \omega) \\ v &= \frac{1}{fk\sigma_{0,p}} \text{lpf}(|\dot{\mathbf{p}}_t|) \\ \omega &= \frac{1}{fk\sigma_{0,\theta}} \text{lpf}(|\dot{\theta}_t|) \end{aligned} \quad (9)$$

where $fk\sigma_0$ is the maximum velocity the particle filter allows with respect to frame rate f , the number of iterations k and the maximum possible standard deviation $\sigma_t = \sigma_0$. Then we qualify the output by applying thresholds c_{v*} that indicate whether the object is *still* or moving *slow* or *fast* ($v < c_{v1} = 0.01$, $c_{v1} \leq v \leq c_{v2} = 0.1$ and $v > c_{v2}$ respectively).

2) *Quality*: To give a statement about the quality of the current pose we use Equation (7) which corresponds to the match of a pose hypothesis to the image evidence. Again we classify this measure to obtain qualitative statements by applying thresholds to distinguish if tracking is *good*, *fair* or *bad* ($c_t > c_{q1} = 0.5$, $c_{q1} \geq c_t \geq c_{q2} = 0.3$ and $c_t < c_{q2}$ respectively).

3) *State*: We decide on the overall tracking state *occluded*, *lost* or tracked *ok* based on confidence and speed, modelling the fact that confidence can decrease as a result of occlusion or motion. To this end we introduce a visibility flag $b_{visible}$ by comparing confidence against a dynamic threshold c_d derived from the current speed.

$$\begin{aligned} b_{visible} &= (\text{lpf}(\max(c_d - c_t, 0)) \leq c_{th,lost}) \\ c_d &= \max(c_{max} - \tilde{v}, c_{min}) \end{aligned} \quad (10)$$

where $c_{max} = 0.5$ and $c_{min} = 0.3$ define the range of c_d and $c_{th,lost} = 0.1$ defines the limit to be reached to declare an object to be lost or occluded. We now define an object to be not visible if the low-pass filtered confidence is low with respect to c_d . This dynamic threshold compensates for low confidence during moderate movement. Table I shows how we decide on the tracking state $_t$, based on visibility, speed, quality and previous tracking state $_{t-1}$. (\neg means logical not)

TABLE I
 DECISION ON TRACKING STATE

$b_{visible}$	speed	quality	state $_{t-1}$	state $_t$
false	still	-	ok	occluded
false	\neg still	-	ok	lost
false	\neg still	-	occluded	lost
true	still	good	occluded	ok
true	still	good	lost	ok

If the tracked pose does not move (speed = *still*) and $b_{visible} = \text{false}$, the assumption is that the object is temporarily occluded. If the tracked pose however does move while $b_{visible} = \text{false}$, this means the tracker drifts off, chasing a wrong local maximum. This follows from the fact that the tracker essentially always follows the local maximum. This

definition does not allow detecting occlusion while the object is moving.

To recover from the state where the object is occluded or lost, speed and quality have to be *still* and *good* respectively.

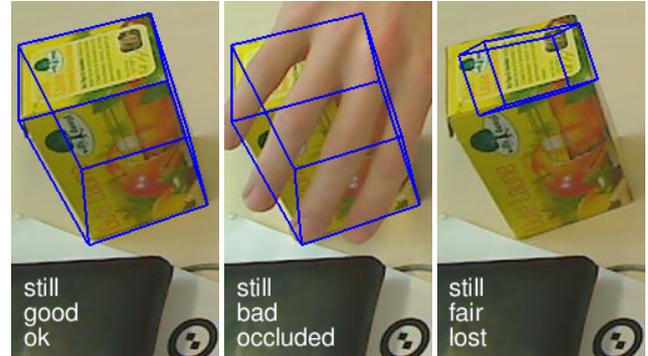


Fig. 3. Tracking-State-Detection: From left to right: ok, occluded and lost tracking.

Note that as we normalise velocity and the confidence the above thresholds apply to a large range of objects and situations.

Figure 4 shows the different values of Equation (10). Between $t = 1.0$ and $t = 2.0$ the object moves, with the confidence going down, but the tracker does not lose it. At about $t = 4.5$ we removed the object from the field of view which was detected at $t = 6.3$. Note how the confidence rises immediately after the object was removed, because the particle filter converged on a false local maximum. This suggests that confidence alone is a bad indicator for successful tracking. Furthermore it shows the importance of a dynamic threshold c_d . At the peak at $t = 1.5$ where the object moves without being lost the confidence is equally low compared to $t = 5$. However, only in the latter case tracking really failed.

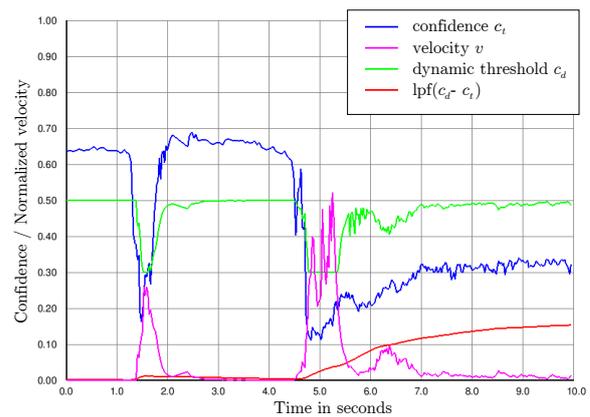


Fig. 4. Tracking-State-Detection: The peak on the left indicates movement without losing the object, whereas on the right side the tracker detects that it lost the object as the red line rises above the lost-threshold $c_{lost} = 0.1$

D. Texture Mapping

Tracking is based on a CAD model which (initially) does not include surface texture. This is sufficient for non-textured

objects, where all we can observe are edges resulting from occlusion and surface discontinuity. For textured objects the additional edges provided by texture on surfaces significantly improve performance and especially robustness. The question is how to get the texture on the model faces. One possibility is to use a 3D editor to generate the model together with the texture mapping. However, we found that the textures mapped onto the geometry often do not align properly when compared to the real object. So we instead propose a mechanism that grabs the texture from the live camera image. Starting from tracking the wire-frame of the CAD model, culling away occluded edges, we successively capture the colour map from the camera and store it as texture together with the projection matrix that maps a point in model space (i.e. vertices) to image space (i.e. pixel coordinates of the texture). Mapping can be done either manually or automatically.

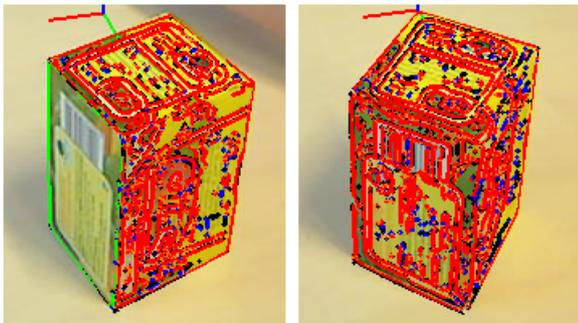


Fig. 5. Successively learning the texture of an object providing correct alignment. (red: matching edges from a textured face, green: matching edges from a non-textured face).

1) *Manual mapping*: To yield accurately aligned textures without blurring, distortions or occlusion it is best to move the object by hand to a proper position (correct alignment, small angle between view vector and negative face normal, no motion), and trigger texture capturing manually.

2) *Automatic mapping*: Of course it is more convenient to capture the textures automatically while moving the object, or by moving the camera around the object, e.g. by using a robotic arm with a camera attached to the end effector. To tell whether the object is in a good position with respect to the camera we calculate the angle between the normal of each face of the geometric model and the current view vector with respect to the centre of the object. We apply TSD (Section III-C) to check if the object is neither moving nor occluded and if tracking quality is good.

3) *Capture the Texture*: Now that we obtained a good pose for a specific face we can copy the image from the camera, cut out the respective image region and generate the UV-coordinates for the vertices. We project the 3D coordinates from model to image space (u, v) using the transformation \mathbf{T} provided by the tracker and the camera projection. Then we compute the bounding rectangle (BR) for all the vertices of the face and re-scale the UV-coordinates to the range $[0 \dots 1]$ with respect to the BR and store this area of the image.

IV. INITIALISATION AND RECOVERY

Object detection is used for pose initialisation of learned models and also triggered if the object is lost during tracking as defined by the TSD.

While edges are well suited for fast tracking we use highly discriminating SIFT features for detection, following a fairly standard learning and recognition scheme. During the learning phase SIFT features (again we use a GPU implementation [34]) are detected in keyframes and mapped to the surface model using the known 3D pose from the tracker. SIFT features falling outside the object boundary are discarded. Keyframes are indicated either manually by button press or automatically using the TSD as described in Section III-D. To speed up recognition SIFT features are represented using a codebook (one per object). SIFT descriptors are clustered using an incremental mean-shift procedure and each 3D location on the object surface is assigned to the according codebook entry.

In the recognition phase SIFT features are detected in the current image and matched to the codebook. To robustly estimate the 6D object pose we use the OpenCV pose estimation procedure in a RANSAC [35] scheme. More details and experimental results are given in [36].

V. RESULTS

All experiments were performed on a PC with an Intel Core 2 Quad (Q6600, 2.4 GHz) CPU, a NVIDIA GeForce GTX 285 GPU and a Logitech Webcam Pro 9000 run at a resolution of 640x480 pixels. We evaluated the approach by using virtual rendered image sequences with known ground truth as well as live sequences where we obtain ground truth from a calibration pattern rigidly attached to the object.

A. Evaluation of the Tracking Error

For a measure of the error we used the scheme proposed in Section IV-B in [2], where a large number, $n = [1 \dots N]$, of randomly chosen points $\mathbf{q}^{1,n}$ are rigidly attached to the object surface at the ground-truth pose and compared to the corresponding points of the tracked pose $\mathbf{q}^{2,n}$.

$$E_d = \frac{1}{N} \sum_{n=1}^N |\mathbf{q}_d^{2,n} - \mathbf{q}_d^{1,n}| \quad (11)$$

with $d \in \{x, y, z\}$.

Before evaluating our method in terms of the above error metric, let us briefly consider the possible sources of errors in our system, such as errors from calibration, geometric modelling, image quantisation and finally the tracking algorithm itself. Concretely we identify the following sources of errors:

- *Mechanical Error*: Positioning the calibration pattern rigidly on the object introduces a small unknown error which can safely be considered to be in sub-millimetre range.
- *Camera Error*: The pose of the calibration pattern is detected with a standard DLT algorithm, followed by a non-linear optimisation of the pose using the sparse bundle adjustment implementation by Lourakis [37].

- *Quantisation Error*: Depending on image resolution a digital camera introduces a pixel quantisation error. In our evaluation we use a resolution of 640x480 with a focal length of ~ 500 in pixel-related units. This leads to an error of about 0.5-1.5 mm when tracking at a distance of 0.5-1.5 m parallel to the image plane. This error is even higher for the orthogonal direction, which shows up in Table II.
- *Modelling Error*: For modelling we measured the main dimensions of the objects used, but we used simplified models that do not account for deviations like small details, chamfers or slightly bulging cardboard surfaces. Unfortunately we do not have a measure for the *Modelling Error* but since we mainly used basic shapes, where correct modelling is simple, we assume this error to be negligible.
- *Texturing Error*: We found that textures added during modelling phase do not align properly and therefore introduced the methods described in Section III-D. Manually capturing textures triggered by pressing a button incorporates less error than automatic capturing based on tracking-state-detection.
- *Tracking Error*: The failure of the tracker to accurately locate the local maximum, depending on the challenges posed by current viewing conditions.

B. Accuracy and Precision

Accuracy is defined to be the closeness of a quantity to its actual value, which in our case is measured using Equation (11), where the pose of tracking \mathbf{p}_t^2 is compared to the pose of the virtual object and the pose detected by the calibration software respectively. We evaluated the mean accuracy with respect to the poses of several trajectories using

$$E_{acc} = \frac{1}{J t_e} \sum_{j=1}^J \sum_{t=1}^{t_e} E_t \quad (12)$$

where $j = [1 \dots J]$ are the trajectories of poses $t = [1 \dots t_e]$ under unchanged conditions, i.e. tracking J times on a sequence of t_e images.

Precision, also called repeatability, is the degree of deviation of a quantity under unchanged conditions, which is measured using Equation (11), where the pose of tracking $\mathbf{q}_{t,j}^{2,n}$ is compared to its own mean with respect to J , the number of repetitions:

$$\mathbf{q}_t^{1,n} = \frac{1}{J} \sum_{j=1}^J \mathbf{q}_{t,j}^{1,n} \quad (13)$$

Table II shows the results of the accuracy and precision evaluation, where *static* refers to the pose after convergence and *dynamic* is the mean error of trajectories j , both over a set of trials J . For evaluation we used box shaped and cylindrical objects. The virtual objects give indication about the *Tracking Error* and *Quantisation Error* (all other errors being ruled out), whereas the difference between virtual and real objects are due to *Mechanical*, *Camera*, *Modelling* and

Texture Error, where we assume the *Modelling* and *Texture Error* to play the main roles.

TABLE II
 ACCURACY AND PRECISION

Target Object	Accuracy [mm]				Precision [mm]			
	static		dynamic		static		dynamic	
	x,y	z	x,y	z	x,y	z	x,y	z
box (virt.)	0.4	2.3	1.5	5.6	0.2	1.1	0.7	3.2
box (real)	2.0	5.5	2.6	7.7	1.1	2.9	1.6	4.9
cylinder (virt.)	0.9	4.4	2.4	10.0	0.4	1.9	1.3	5.7
cylinder (real)	3.0	16.5	3.9	21.9	0.5	2.5	1.6	8.8

We evaluated the dynamic errors under the following conditions:

- Linearly moving objects with different velocities
- Rotating objects
- Arbitrary moving objects (i.e. toppling, rolling)
- Partially occluded objects
- Changing illumination

Table II indicates that curved objects are typically harder to track than box-shaped objects.

A typical trajectory for arbitrary movement is shown in Figure 6 where the tracked pose is compared to the virtual with respect to translations, rotations and the error measured by Equation (11). For generating the virtual poses the poses from the pattern detection and bundle adjustment were used and low-pass filtered to remove jitter.

C. Robustness

We tested our approach against various situations including

- fast movement introducing motion blur,
- occlusion,
- changes in lighting,
- large distances, small objects,
- different objects (high resolution, curved surfaces, low texture, ...).

Since robustness is hard to put in numbers the reader is referred to a video to get an impression how these various challenges are handled.

D. Performance

Processing time during tracking depends, along with computational power, on the complexity of the model as well as on the number of particles used for tracking.

Table III shows the frame rates for different numbers of faces and particles. 2x50, 3x100 and 4x300 indicates 2, 3 and 4 iterations using 50, 100 and 300 particles for each iteration respectively. Figure 7 shows the frames per second on different GPUs with respect to the total number of particles used for tracking.

E. Video

The video¹ shows how we learn texture and feature points during tracking. Then tracking identifies whether if an object

¹<http://www.youtube.com/watch?v=r5xUcDmTY3E>

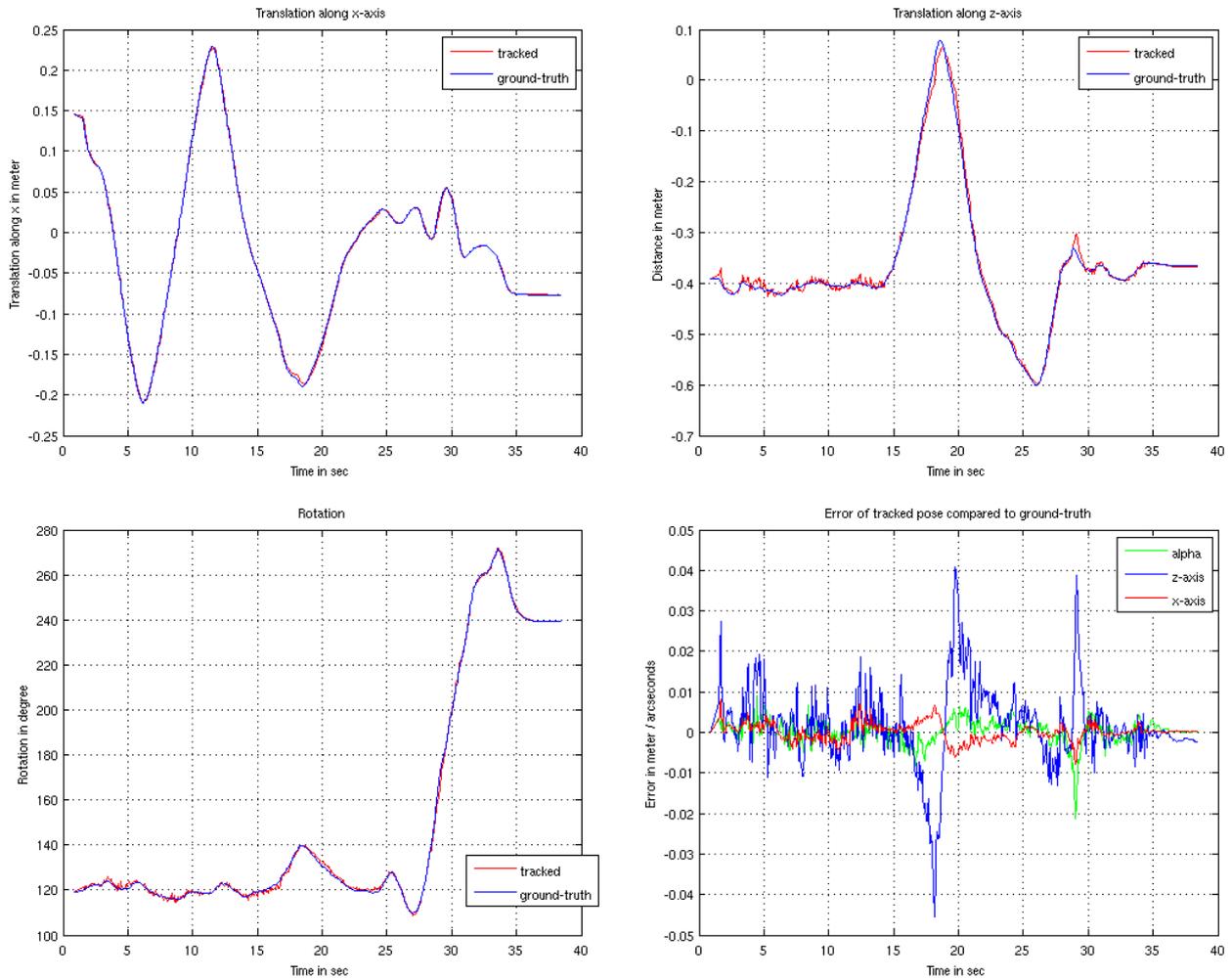


Fig. 6. Trajectory of a tracked virtual object with 45 cm x-translation followed by a 70 cm z-translation and a rotation about the objects y-axis. The lower right figure shows the pose deviations respectively. Note that the jitter results from different location of convergence of the particle filter, due to the errors mentioned in Section V-A.

TABLE III
 FRAME RATES

Example Objects	Faces	Frames per Second		
		2x50	3x100	4x300
Box	6	240	100	33
Cylinder (low)	24	220	95	30
Cylinder (mid)	96	210	90	28
Cylinder (high)	384	190	80	25

is occluded or if tracking fails, in which case pose recovery is triggered automatically. Since pose recovery also takes advantage of GPU computing the tracker slows down at this particular moments. Note that we do not interfere with the tracking system via the keyboard other than for changing the display modes.

VI. CONCLUSION AND FUTURE WORK

A. Conclusion

We presented a model-based tracking system to accurately follow the pose of an object in real-time. We developed a

novel method of Tracking-State-Detection (TSD) that analyses tracking performance to reason about tracking quality, speed and whether the object of interest is occluded or lost. This allows triggering a feature based pose recovery system, texture mapping and online learning. Accuracy, precision and performance of our approach are evaluated carefully to provide a maximum of applicability and give very good results compared to ground truth. Note that we showed tracking results for simple shapes only, which can be easily measured to compare to ground truth. However, our approach is not limited to specific shapes.

B. Future work

There are several improvements for our tracking approach. While the system has no problems in tracking both textured as well as non-textured objects pose recovery only works for the former as it depends on SIFT features. Another open issue is how to combine colour and edges as described in Equation (4) and (5). And while our novel approach for Tracking-State-Detection shows good results for a wide

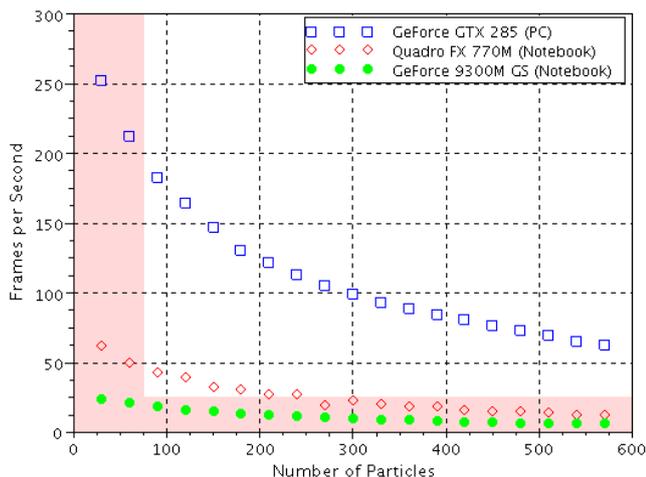


Fig. 7. Frame rate with respect to the number of particles

range of situations, the mathematical formulation in Section III-C is not yet satisfying with respect to generalisability to other tracking methods.

VII. ACKNOWLEDGEMENTS

The research leading to these results has received funding from the European Community's Seventh Framework Programme [FP7/2007-2013] under grant agreement No. 215181, CogX.

REFERENCES

- Y. Bekiroglu, D. Kragic, and V. Kyrki, "Learning grasp stability based on tactile data and hms," in *RO-MAN, 2010 IEEE*, pp. 132–137, sept. 2010.
- M. Kopicki, R. Stolkin, S. Zurek, T. Mörwald, and J. Wyatt, "Predicting workpiece motions under pushing manipulations using the principle of minimum energy," in *RSS workshop*, 2009.
- T. Mörwald, M. Kopicki, R. Stolkin, J. Wyatt, S. Zurek, M. Zillich, and M. Vincze, "Predicting the unobservable, visual 3d tracking with a probabilistic motion model," in *IEEE International Conference on Robotics and Automation, ICRA11*, May 2011.
- M. Kopicki, S. Zurek, R. Stolkin, T. Mörwald, and J. Wyatt, "Learning to predict how rigid objects behave under simple manipulation," in *IEEE International Conference on Robotics and Automation, ICRA11*, May 2011.
- D. J. Duff, T. Mörwald, R. Stolkin, and J. Wyatt, "Physical simulation for monocular 3d model based tracking," in *IEEE International Conference on Robotics and Automation (ICRA11)*, May 2011.
- J. Chestnutt, S. Kagami, K. Nishiwaki, J. Kuffner, and T. Kanade, "Gpu-accelerated real-time 3d tracking for humanoid locomotion," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2007.
- P. Michel, J. Chestnutt, S. Kagami, K. Nishiwaki, J. Kuffner, and T. Kanade, "Gpu-accelerated real-time 3d tracking for humanoid autonomy," in *JSME Robotics and Mechatronics Conference (ROBOMECH'08)*, June 2008.
- G. Klein and D. Murray, "Full-3d edge tracking with a particle filter," in *British Machine Vision Conference*, 2006.
- L. Masson, M. Dhome, and F. Jurie, "Robust real time tracking of 3d objects," in *International Conference on Pattern Recognition, ICPR*, 2004.
- L. Vacchetti, V. Lepetit, and P. Fua, "Combining edge and texture information for real-time accurate 3d camera tracking," in *IEEE/ACM International Symposium on Mixed and Augmented Reality*, 2004.
- P. Mittrapiyanuruk, G. N. Desouza, and A. C. Kak, "Accurate 3d tracking of rigid objects with occlusion using active appearance models," in *WACV/MOTION*, pp. 90–95, 2005.
- L. Vacchetti, V. Lepetit, and P. Fua, "Stable real-time 3d tracking using online and offline information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2004.
- J. Sánchez, H. Álvarez, and D. Borro, "Towards real time 3d tracking and reconstruction on a gpu using monte carlo simulations," in *IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pp. 185–192, Oct. 2010.
- M. Özuysal, M. Calonder, V. Lepetit, and P. Fua, "Fast keypoint recognition using random ferns," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009.
- M. Vincze, M. Ayromlou, W. Ponweiser, and M. Zillich, "Edge-projected integration of image and model cues for robust model-based object tracking," *The International Journal of Robotics Research*, 2001.
- H. H. Nagel, "Representation of moving rigid objects based on visual observations," *Computer*, vol. 14, pp. 29–39, August 1981.
- R. Eskenazi and R. T. Cunningham, "Real-time tracking of moving objects in tv images," *IEEE Workshop Pattern Recognition and Artificial Intelligence*, pp. 4–6, April 1978.
- C. Harris and A. Blake ed., "Tracking with rigid bodies," *Active Vision*, pp. 59–73, 1992.
- T. Drummond and R. Cipolla, "Real-time tracking of complex structures with on-line camera calibration," in *British Machine Vision Conference (BMVC99)*, pp. 574–583, 1999.
- G. Klein and T. Drummond, "Robust visual tracking for non-instrumented augmented reality," in *ISMAR IEEE/ACM International Symposium on Mixed and Augmented Reality*, 2003.
- D. Lowe, "Robust model-based motion tracking through the integration of search and estimation," *International Journal of Computer Vision*, pp. 113–122, 1992.
- D. Gennery, "Visual tracking of known three-dimensional object," *International Journal of Computer Vision*, 1992.
- D. Koller, K. Daniilidis, and H.-H. Nagel, "Model-based object tracking in monocular image sequences of road traffic scenes," *International Journal of Computer Vision*, 1993.
- A. Kosaka and G. Nakazawa, "Vision-based motion tracking of rigid objects using prediction of uncertainties," *International Conference on Robotics and Automation*, 1995.
- A. Ruf, M. Tonko, R. Horaud, and H.-H. Nagel, "Visual tracking by adaptive kinematic prediction," in *International Conference on Intelligent Robots and Systems*, 1997.
- H. de Ruiter and B. Benhabib, "Visual-model-based, real-time 3d pose tracking for autonomous navigation: methodology and experiments," *Autonomous Robots*, vol. 25, pp. 267–286, 2008.
- D. Kragic, A. T. Miller, and P. K. Allen, "Real-time tracking meets on-line grasp planning," in *IEEE Intl. Conf. on Robotics and Automation*, pp. 2460–2465, 2001.
- E. Marchand and P. Bouthemy, "A 2d-3d model-based approach to real-time visual tracking," *IVC*, vol. 19, pp. 941–955, 2001.
- J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendon-Mancha, "Binocular visual tracking and grasping of a moving object with a 3d trajectory," *Journal of Applied Research and Technology*, vol. 7, no. 03, pp. 259–274, 2009.
- A. Doucet, N. De Freitas, and N. Gordon, eds., *Sequential Monte Carlo methods in practice*. Springer, 2001.
- A. Doucet, S. Godsill, and C. Andrieu, "On sequential monte carlo sampling methods for bayesian filtering," *Statistics and Computing*, vol. 10, pp. 197–208, 2000.
- T. Mörwald, M. Zillich, and M. Vincze, "Edge tracking of textured objects with a recursive particle filter," in *Graphicon 2009*, (Moscow, Russia), 2009.
- A. Richtsfeld, T. Mörwald, M. Zillich, and M. Vincze, "Taking in shape: Detection and tracking of basic 3d shapes in a robotics context," in *Computer Vision Winter Workshop*, pp. 91–98, 2010.
- C. Wu, "<http://www.cs.unc.edu/~cewu/siftgpu/>."
- M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Comm. of the ACM*, vol. 24, pp. 381–395, 1981.
- T. Mörwald, J. Prankl, A. Richtsfeld, M. Zillich, and M. Vincze, "Blort - the blocks world robotic vision toolbox," in *ICRA workshop*, 2010.
- M. A. Lourakis and A. Argyros, "SBA: A Software Package for Generic Sparse Bundle Adjustment," *ACM Trans. Math. Software*, vol. 36, no. 1, pp. 1–30, 2009.

Generalizing Grasps Across Partly Similar Objects

Renaud Detry Carl Henrik Ek Marianna Madry Justus Piater Danica Kragic

Abstract—The paper starts by reviewing the challenges associated to grasp planning, and previous work on robot grasping. Our review emphasizes the importance of agents that generalize grasping strategies across objects, and that are able to transfer these strategies to novel objects. In the rest of the paper, we then devise a novel approach to the grasp transfer problem, where generalization is achieved by *learning*, from a set of grasp examples, a dictionary of object parts by which objects are often grasped. We detail the application of dimensionality reduction and unsupervised clustering algorithms to the end of identifying the size and shape of parts that often predict the application of a grasp. The learned dictionary allows our agent to grasp novel objects which share a part with previously seen objects, by matching the learned parts to the current view of the new object, and selecting the grasp associated to the best-fitting part. We present and discuss a proof-of-concept experiment in which a dictionary is learned from a set of synthetic grasp examples. While prior work in this area focused primarily on shape analysis (parts identified, e.g., through visual clustering, or salient structure analysis), the key aspect of this work is the emergence of parts from *both* object shape *and* grasp examples. As a result, parts intrinsically encode the intention of executing a grasp.

I. INTRODUCTION: CHALLENGES IN GRASP PLANNING

This paper studies the planning of grasping actions, or, in other words, the problem of exploiting perceptual data to select a wrist position and finger configuration to which a hand can be transported in order to grasp an object. The wrist position (or *grasping point*) corresponds to the region of the object towards which the hand will move. The finger configuration (or *hand preshape*) corresponds to the angles to which finger joints are set prior to coming in contact with the object.

Grasp planning is a complex problem. A grasp must bind a hand to an object, and prevent the object from subsequently slipping or escaping. Configurations which lead to a collision between the hand and the object or other obstacles must be avoided, and task-related constraints must be verified (certain tasks restrain the number of possible grasps, as a knife should be held specifically by its handle when the task is to cut something). Perceptual data, usually provided by vision, are noisy and often limited to a single viewpoint. For dexterous

grasping, the space of action parameters (hand positions and configurations) quickly becomes high-dimensional (a human hand has twenty-five degrees of freedom – six for the wrist position and orientation, and nineteen for the finger joint angles). Yet, despite the complexity of the problem, the frequent recurrence of grasping in everyday tasks imposes an ability to plan grasps quickly.

In robotics, grasp planning traditionally relies on contact-force analysis [3], [34]. Force analysis bases planning on a reconstruction of the geometry and physical properties of the objects that surround the agent. Provided that such a reconstruction is available, the agent searches the space of hand configurations for the configuration that best verifies grasping constraints (binding configuration, no collisions, task compatibility). In practice, the applicability of force analysis is limited by the difficulty of obtaining accurate models of object geometry, mass, and friction characteristics. Also, as the space of hand configurations is high dimensional, the optimization procedure underlying force analysis is computationally expensive. These shortcomings motivated the community to rethink the planning problem, leading for instance Borst et al. [5] to demonstrate that finding the globally optimal grasp is often not strictly worth the computational effort, as for many tasks an average grasp (in the force-analysis sense) is acceptable. The bigger leap however came with a class of methods that parted drastically from the traditional planning philosophy. Instead of searching for a grasp that optimally satisfies the various (vision-dependent) grasping constraints, these methods extract, from the agent's experience, a function that directly maps visual perceptions to grasp parameters, with the advantage of *implicitly* capturing the object's physical properties, and avoiding a costly search through the high-dimensional space of hand configurations [7], [21], [25], [30], [36].

Numerous behavioral studies tend to support the existence of similar processes in the human grasping system. It has been shown for instance that humans often grasp objects by preshaping their hand during its transportation towards the object [18], then compliantly refining the grip upon contact [19]. Concurrently, neurophysiological studies suggested that, in monkeys, the cortex encodes a set of prototype grasps, which are selectively triggered by visual stimuli [26]. It thus seems plausible, as proposed, for instance, by Johansson et al. [19], that the human grasping system relies on a set of prototypical motor programs that are selected and parametrized by visual input, therefore acting as a direct mapping from vision to action. Humans arguably possess the most sophisticated grasping system known today, being able to plan complicated grasps in just a few hundreds of

R. Detry, C. H. Ek, M. Madry, and D. Kragic are with the Centre for Autonomous Systems and the Computer Vision and Active Perception Lab, CSC, KTH Royal Institute of Technology, Stockholm, Sweden. Justus Piater is with the Intelligent and Interactive Systems group, Institute of Computer Science, University of Innsbruck, Austria. Email: {detryr, chek, madry, danik}@csc.kth.se, justus.piater@uibk.ac.at

This work was supported by the Swedish Foundation for Strategic Research, the Belgian National Fund for Scientific Research (FNRS), and the EU projects CogX (FP7-IP-027657) and TOMSY (IST-FP7-Collaborative Project-270436).

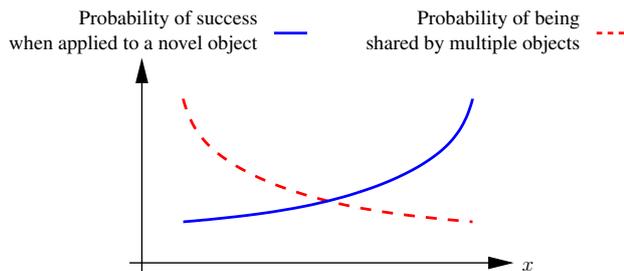


Fig. 1: Robustness-transferability trade-off in feature-based grasp planning. The x axis corresponds to the amount of information encoded by a part. Highly informative parts allow for a robust grasp application. However, these are less likely to be shared across objects.

milliseconds [17]. We believe that the possibility that such an efficient system be based on a direct vision-action mapping is a strong argument for researching vision-action mappings for robotics.

To learn a vision-to-grasp mapping for one specific object, an agent usually collects a set of grasp examples, and lets machine-learning algorithms construct a grasp predictor from these. Such a model allows the agent to quickly produce grasping plans for the object on which it trained. However, collecting grasp examples is an expensive, time-consuming process. A major focus in grasp learning is to develop methods that produce useful manipulation models from as few data as possible. A natural means of limiting the need for examples is to try and adapt memories of previous objects to the planning of a grasp onto a novel object. Many objects share similarities in shape, and similarities in grasp affordances, and both are often correlated. When a novel object appears, instead of starting to learn from scratch, an agent may instead attempt to apply to it the strategies it has acquired for partly similar objects. To this end, means of linking grasps to certain object *features* have been researched, in the hope of transferring grasps across objects that share the same features. The challenge of this task is to decide which visual cues should be captured by the features. Intuitively, a feature should capture no more no less than the specific cues that predict the applicability of a grasp. If a feature misses important cues, it risks predicting faulty grasps. If a feature includes cues that are not directly related to grasping, its transferability to other objects will be impeded. Designing a feature for grasp generalization thus involves a robustness-transferability trade-off, as illustrated in Fig. 1.

A number of methods for vision-based grasping learn a mapping from image features, such as local gradients or SIFT, to grasp parameters [24], [25], [30]. One advantage of these methods is their conceptual elegance:

- 1) Extract features from images of a set of objects.
- 2) Label these features as good or bad grasping point, either with the help of a teacher [30] or through autonomous exploration [24].
- 3) Learn a grasp classifier.

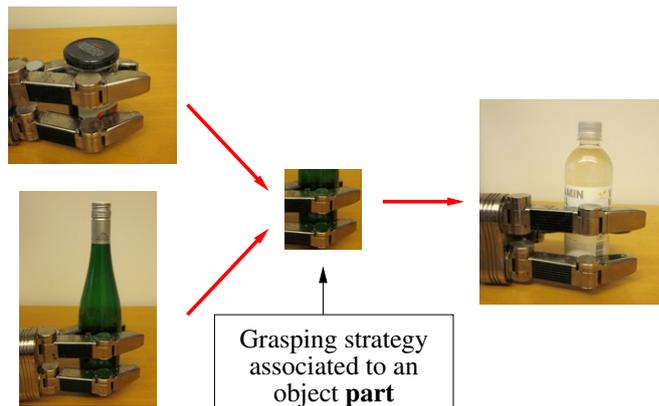


Fig. 2: Learning part-grasp associations. The agent will identify, within its visuomotor experience, recurrent associations of object *parts* and successfully executed grasps. These grasps will then be applicable to novel objects that share the same part.

- 4) Transfer grasps by classifying features obtained from images of novel objects.

Unfortunately, these methods also come with their shortcomings. From a practical viewpoint, the geometric information provided by a local feature detector is generally poor. As grasping is an intrinsically 3D interaction, it largely relies on 3D object properties, such as shape, which are only partly captured by 2D image features. It is thus difficult to link, for example, a 3D gripper orientation to an image feature.

Across the range of visual cues that have been used for designing grasp planners, 3D shape has lead to particularly good results. By contrast to methods based on image features, methods that link grasp parameters to a shape model [1], [9], [11], [14], [23] benefit from an increased geometric robustness, which makes it easier to preshape the hand to approximate object shapes, and accurately position and orient the wrist and fingers with respect to the object. Mapping grasps to 3D cues is supported by behavioral and neurophysiological studies. Behavioral studies have demonstrated the reliance of human grasping on 3D shape [16], while neurophysiologists have observed a mapping from 3D shape to action prototypes in monkeys [27].

II. LEARNING SHAPE PROTOTYPES FOR GENERALIZING GRASPS

In the rest of the paper, we present an adaptive grasp planner that learns a mapping from object shape to grasp parameters.

A. From Part to Grasp

Linking grasp parameters to the shape of the whole body of an object limits the applicability of the model to that particular object. In order to transfer grasps across objects, we instead explore the linking of grasp parameters to object parts. In order to allow the agent to generalize its acquired knowledge to novel objects, we propose to provide it with

means of identifying, within its visuomotor experience, recurrent associations of object *parts* and successfully executed grasps. For instance, the agent may have successfully transported objects such as bottles, cans, and jars, which have different sizes, but which can be seized by applying the same power grasp to their side. We propose to provide the agent with means of understanding, from a set of such examples, that any object that presents a cylindrical part can be grasped sideways with a wide-palm grasp (Fig. 2).

B. Previous Work on Part-based Grasping

Part-grasp associations have been previously suggested and studied by several research groups [2], [23], [36]. In the earlier work, the definition of parts was often either hard-coded [23], or driven by shape analysis [1], [2], [36]. There is however an increasing interest for defining parts based on grasping experience [10], [12], [15], [22], [37]. For instance, Herzog et al. [15] and Zhang et al. [37] presented two exciting data-driven approaches where a part describes an object’s shape in a fixed-size region around a grasping point. These approaches are further discussed below.

C. Method

Our work aims at learning, from a set of grasp examples, a dictionary of prototypical parts by which objects are often grasped. A key property that we wish to allow our agent to extract from experience is the spatial extent of grasp-predicting parts. For instance, in the case presented in Fig. 2, we wish our agent to learn that the relevant part is a 10cm-high cylinder. The the tap of the jar or the conic upper part of the bottle should be ignored, as they are not shared by the two objects.

Training data are provided to the agent in the form of a set of grasps demonstrated onto objects known to the agent. (The agent has previously acquired 3D point clouds that model the shape of the objects.) A grasp is parametrized by the 6D pose of the wrist (3D position and 3D orientation), and by the 6D pose of the object. Our method works as follows: First, the agent generates, from the grasp examples, a large number of part candidates of varying sizes (Section III). Most of the candidates will not generalize well. However, it is our hope that for every set of objects that share a graspable part, each object will yield one candidate that approximately captures that part. The candidates that recur across objects are identified by clustering part candidates (Section IV). Dense clusters will contain parts by which objects are often grasped, which are thus promising for grasping novel objects.

The central parts of all clusters will form the dictionary used by the agent to grasp novel objects. An important aspect of our work appears at this point. As the dictionary of parts is only formed from cluster centers, it is allowed to be orders of magnitude smaller than the set of grasp examples initially provided to the agent. In the data-driven approaches discussed above [15], [37], each grasp example yields a part. By contrast, in our work, a grasp example only “votes” for the potential inclusion of a part into the dictionary, which provides us with a means of controlling the size of the

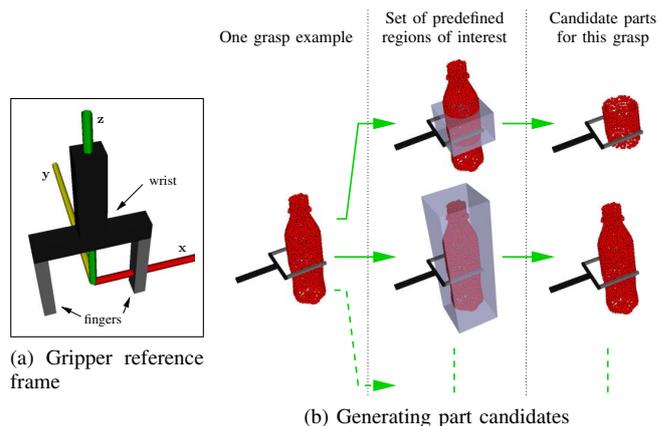


Fig. 3: Generating part candidates. The black and grey renderings on each image represent the pose of the gripper set for a sideways grasp on the soda bottle. Parts of varying sizes are generated by defining several box-shaped regions of interest centered on the gripper.

dictionary in order to keep the computational cost of planning a grasp onto a novel object reasonably low.

Also, in our work, parts emerge from both object shape and grasp examples. A key result is our ability to optimize the robustness-transferability trade-off discussed above. Not only the shape, but also the spatial extent (or size) of the parts that form the dictionary depend on the available grasp data. Our approach involves an explicit search for recurrent patterns within the agent’s visuomotor experience, which leads to the identification of parts that directly predict grasp applicability.

III. GENERATING PART CANDIDATES

Part candidates are generated by extracting object surface segments of varying size in the vicinity of grasps demonstrated by a teacher. Parts are thus represented, as the object from which they are extracted, by point clouds. This process is illustrated for a soda bottle in Fig. 3. Surface segments are extracted using a set of predefined regions of interest (ROI). These regions are centered on the gripper, as the applicability of a grasp is largely conditioned by the shape of the surface in the direct vicinity of the grasping point. ROI sizes should *a priori* vary in all directions. However, the preshape of the gripper at the time of the grasp can limit the number of regions that are interesting to look at. For instance, in the case shown in Fig. 3, it is reasonable to limit the ROI width along the x axis of the gripper to the distance that separates both fingers, as the object will usually not be larger than this gap. With more sophisticated hands, grasp preshapes can further constrain the definition of ROIs.

IV. EXTRACTING DENSE CLUSTERS OF PARTS

Graspable parts that generalize are discovered by clustering part candidates. Dense groups of similarly-shaped candidates correspond to shapes onto which grasps can be applied in order to seize several different objects. These shapes are thus likely to predict grasp applicability for novel objects.

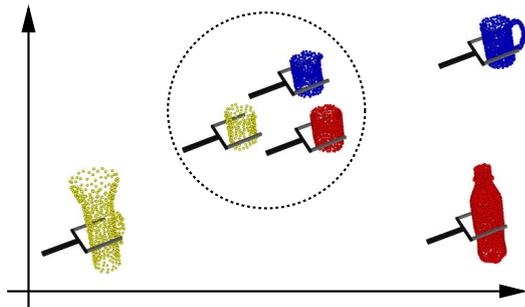


Fig. 4: Finding parts that allow for transferring grasps to a novel object. The three outer “parts” (which correspond to entire objects), will not generalize well. By contrast, the three center parts, which represent a piece of the flashlight, cup, and soda bottle, are very similar to each other. As there exist a shape similarity across these three parts extracted from different objects of the training database, it seems reasonable to assume that the grasps related to these parts are potentially applicable to novel objects.

In Fig. 4, none of three outer parts would be applicable to other objects. The three middle parts, by contrast, encode a shape-grasp relation that would be applicable to an object that has a cylindrical part of a similar diameter.

Clustering part candidates requires the definition of a measure of shape (dis)similarity. This measure is defined in the next section. Section IV-B details the clustering algorithm.

A. Measuring Part Dissimilarity

This section defines a measure part dissimilarity. We note that, as we ultimately aim at using parts for predicting gripper poses, we must measure the (dis)similarity of *grasper-relative* shapes. In other words, a cylindrical part grasped from the side should *not* be similar to the same cylindrical part grasped from the bottom.

In this work, a part is represented by a point cloud defined in a reference frame that corresponds to the 6D pose of the grasp associated to that part. Let $P = \{x_i\}_{i \in [0, n]}$ and $Q = \{y_i\}_{i \in [0, m]}$ denote the point-cloud representations of two parts, with all x_i 's and y_i 's belonging to \mathbb{R}^3 . Let us then denote by d^* an asymmetric measure of dissimilarity of P and Q , with

$$d^*(P, Q) = \sum_{i=0}^n \min_{j \in [0, m]} f(x_i, y_j), \quad (1)$$

where

$$f(x, y) = \begin{cases} \frac{\|x-y\|}{T} & \text{if } \|x-y\| \leq T, \\ 1 & \text{if } \|x-y\| > T. \end{cases} \quad (2)$$

The dissimilarity d^* is often used as error function for point-cloud alignment. In our experiments, the threshold T is set to two centimeters.

We define the dissimilarity of two parts P and Q as

$$d(P, Q) = d^*(P, Q) + d^*(Q, P). \quad (3)$$

The dissimilarity d is symmetric in its arguments. It amounts to the sum of the Euclidean distances between the points of P and their nearest neighbor in Q , and the points of Q and their nearest neighbor in P .

B. Clustering Parts

The dissimilarity measure defined in the previous section provides us with a qualitative tool for reasoning on the recurrence of shape-gripper associations across grasp examples. As expressed in the conceptual illustration of Fig. 4, we wish to find a geometric configuration with dense clusters of parts induced by our similarity measure. Dense clusters will correspond to parts that frequently occur within our database. These parts are therefore likely to be useful for grasping novel objects.

The measure described in IV-A provides a global dissimilarity measure between each item in the database from which we can generate a distance matrix

$$D_{ij} = d(P_i, P_j) \quad (4)$$

for all the entries in the database. In order to interpret the data we wish to find a geometrical configuration of the datapoints where the Euclidean distance corresponds to the dissimilarity measure we defined. One possibility is to directly apply classical multi-dimensional scaling [8] to the distance matrix. However, in this paper we are interested in finding a geometrical configuration which suits interpreting the data in terms of clusters. In order to do so we introduce additional flexibility by first interpreting the distance matrix in terms of an inner-product of Gram matrix. Distance matrices and Gram matrices can be interchanged [29] as data inducing representations. Dependent on applications there are benefits associated with each view-point. Here the use of a Gram matrix allows us to view the matrix as a covariance matrix; this approach is well known as the “kernel-trick” [4]. To that end, we use a squared exponential function to apply a non-linear transform of the space that the dissimilarity measure induces,

$$k(P, Q) = e^{-\frac{d(P, Q)^2}{\sigma}}. \quad (5)$$

The squared exponential function induces a geometrical space well-suited for clustering as it will push points that are close together closer and move points far apart even further apart. The parameter σ controls the strength of this transformation.

Discovering part clusters could be achieved directly on the distances defined above (4). However, in order to facilitate the illustration of our method in the experiments presented below, we first recover a low-dimensional approximation of the data, then cluster the data in this low-dimensional space. We recover a d dimensional approximation of the data by solving the following minimization problem,

$$\hat{\mathbf{C}} = \operatorname{argmin}_{\mathbf{C}} \|\mathbf{K} - \mathbf{C}\|_{\mathbb{F}}^2, \quad (6)$$

where \mathbf{K} is the Gram matrix whose elements are defined by $k(P_i, P_j)$ for all the entries in the database, and the rank

of \mathbf{C} is constrained to be at most d . The solution can be found in close form through an eigenvalue problem and is well-known as kernel principal component analysis [31].

Having resolved a geometrical representation of the data, we wish to partition the space in such a manner that we can discover atomic classes of grasps independent of object type. We proceed through a two-stage process. First, we want to group each point in the database into a small number of classes. Secondly, we wish to explain each class by a single representative grasp. Underpinning our approach is the notion that the dissimilarity measure contains this desired structure. This assumption implies that the grouping can be cast as a clustering problem. Clustering is a well-studied problem within computer science and datamining. It has been used extensively to create compact representations of data using mixture models [35] or for application scenarios where a significant amount of prior information about the partitioning is available [6].

The dissimilarity measure $d(\cdot, \cdot)$ is defined between each point in the database. This allows us to construct a graph $G \in \{\mathcal{V}, \mathcal{E}\}$ where each grasp is represented by a node $v_i \in \mathcal{V}$ with edges $e_{ij} \in \mathcal{E}$ connecting associated nodes. We wish to find a partitioning that respects the dissimilarity measure $d(\cdot, \cdot)$. To that end, we construct a fully connected graph. The edge weights are $e_{ij} = \mathbf{C}_{ij}$, *i.e.*, inversely proportional to the dissimilarity between the grasps according to our measure. In order to partition the space, it now remains to cut the graph into disjoint regions each representing a cluster.

In this paper we employ the normalized cuts [33] approach to partition the graph. The cut $(\mathcal{A}, \mathcal{B})$ of a graph G into two sets of disjoint nodes \mathcal{A} and \mathcal{B} is defined as,

$$\text{cut}(\mathcal{A}, \mathcal{B}) = \sum_{i \in \mathcal{A}, j \in \mathcal{B}} e_{ij}. \quad (7)$$

The normalized cuts algorithm finds the partitioning of the graph that minimizes the following objective function,

$$\text{cut}_{\text{normalized}}(\mathcal{A}, \mathcal{B}) = \frac{\text{cut}(\mathcal{A}, \mathcal{B})}{\text{assoc}(\mathcal{A}, \mathcal{V})} + \frac{\text{cut}(\mathcal{A}, \mathcal{B})}{\text{assoc}(\mathcal{B}, \mathcal{V})}, \quad (8)$$

$$\text{assoc}(\mathcal{A}, \mathcal{V}) = \sum_{i \in \mathcal{A}, j \in \mathcal{V}} e_{ij}. \quad (9)$$

The denominator grows with increasing node sets which works to penalize creating very small clusters.

V. PROOF OF CONCEPT

We now present a proof-of-concept experiment which illustrates the method suggested above. The experiment is realized on synthetic data consisting of seven two-finger grasps demonstrated on four objects (see Fig. 5a and Fig. 5b).

Three sets of regions of interest were defined for the three grasp types present in the database. Three ROIs were defined for ‘‘cylindrical’’ grasps, which correspond to the grasps number 1, 2 and 3 of Fig. 5b. Four ROIs were defined for the parallel grasps (4, 5, 6), and six ROIs for the pinch grasp (7). We note that, in the case of the synthetic data studied in this paper, considering cylindrical, parallel and pinch grasps is purely anecdotal. However, in



Fig. 6: Cylindrical grasp preshape. The finger-surface normals at the contact points are 120° apart.

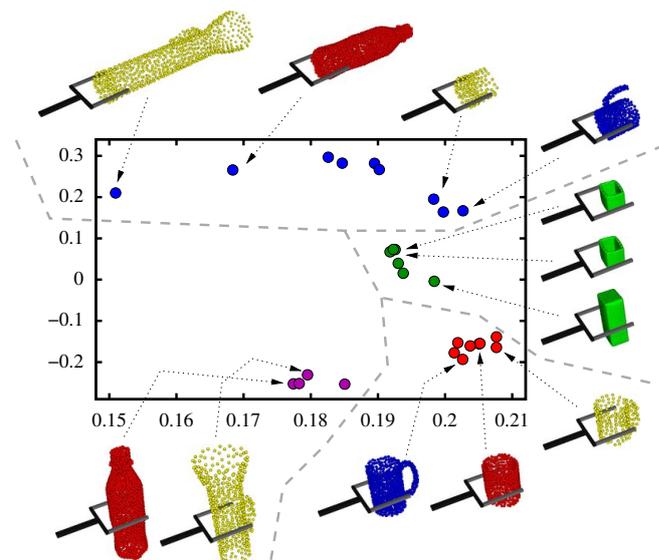


Fig. 7: Two-dimensional approximation of candidates’ geometric configuration, computed from the dissimilarity measure of Section IV-A. Dot colors indicate the data cluster to which a datapoint (part candidate) belongs (see text for details). The colors of the dots within the plot and the colors of the parts surrounding the plot are unrelated. We note that the vertical and horizontal axes are not equally scaled.

a real-case scenario, the hand preshape used for a given grasp would allow us to limit the number of parts that need to be considered as candidates. For instance, with a cylindrical grasp (Fig. 6), generating ROIs that differ in size in a direction perpendicular to the palm of the hand is more important than considering variations along directions parallel to the palm. With a parallel grasp (for instance, Fig. 2), ROIs of various lengths in a direction parallel to the palm are necessary. These observations motivated the definition of different sets of ROIs for the different types of grasps shown in Fig. 5. The part candidates generated with these ROIs are shown in Fig. 5c.

As explained in Section IV-B, kernel PCA provides us with low-dimensional approximations of our data. A two-dimensional approximation is shown in Fig. 7. This plot shows that the dissimilarity measure of Section IV-A properly separates candidate parts in groups of similarly-shaped parts. These groups can be correctly identified by the clustering algorithm of Section IV-B, as reported by the colors associ-

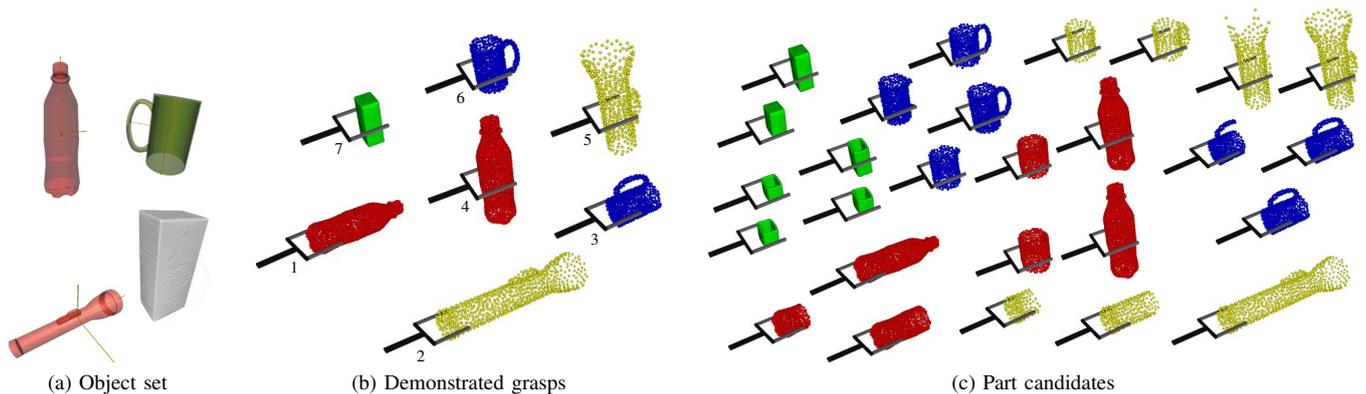


Fig. 5: Experimental data. Three of the objects are cylinders of different sizes, and one is a box. Seven grasps are synthetically demonstrated to the agent. For the cylinders, both sideways and top-down grasps are demonstrated. Fig. (c) shows the candidate parts computed from the grasps of Fig. (b). Part colors indicate which object a part is segmented from.

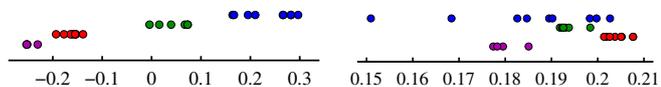


Fig. 8: Projection of the data (candidate parts) onto the first (left) and second (right) principal components of the data. Colors indicate the data cluster to which a datapoint belongs (see text for details). The elevation of the datapoints above horizontal axes is meant to help identifying clusters.



Fig. 9: Prototype parts. These parts correspond to the centers of the clusters of Fig. 7.

ated to the datapoints. In this paper, the number of clusters was determined by inspection. However, BIC-like criterions that compute an optimal number of clusters could be used instead [32]. We note that the two axes of this plot are not equally scaled. The data shows a larger variance along the vertical axis than along the horizontal axis. Fig. 8 shows the projection of the data onto its first and second principal components (which correspond to the vertical and horizontal axes of Fig. 7, respectively). Fig. 8 indicates that the first component contains enough information to identify most of the clusters computed from the dissimilarity measure. The second component leads to a clear separation of the purple and red clusters.

Despite the modest number of data, computing the central point of each cluster allows us to identify a set of prototypical graspable parts. These parts are shown in Fig. 9. We emphasize that despite its reliance on complete object shape models for learning prototypical parts, the method presented above is applicable to predicting grasps onto novel objects perceived through a single 3D snapshot. Fig. 10 illustrate the application of the first and last prototypes of Fig. 9 to a novel object. The right side of Fig. 10 shows the point-cloud



Fig. 10: Grasping a novel object using a dictionary of parts. The rightmost image shows the grasps suggested by the first and last prototypes of Fig. 9, respectively approaching the object from the side and from the top.

representation of the scene (captured by a depth sensor), and the two grasps suggested by the prototypes. The parts are aligned to the object using the pose estimation method of Detry et al. [13].

VI. DISCUSSION

The dissimilarity measure of Section IV-A provides a direct channel for injecting expert knowledge into to the method presented above. By choosing suitable dissimilarities, one can let a variety of desirable visuomotor strategies emerge from data clustering. For instance, one may argue that similarly-shaped parts may predict similar grasps despite a scale difference. Basing a similarity measure on a mix of local shape features (Spin images [20], or FPFH [28]) and global shape features (for instance, the first few moments of a point cloud) has the potential of robustly representing shape while being invariant, to some extent, to scale. Such a measure would allow an agent to understand that cylinders of different radii can be grasped in similar ways. Simultaneously, the distance matrix of Eq. 4 would be much simpler to compute from a set of compact shape features than from the original point-cloud representations. Using shape features would effectively move some of the computational effort out of the distance-matrix computation (quadratic in the number

of candidate parts), into a process linear in the number of candidate parts.

Grasp preshapes were discussed in the previous section, albeit remaining of anecdotal use. In a real-world scenario involving a dexterous hand, preshape is an essential grasping property. In such a scenario, a dissimilarity measure would benefit from the availability of preshape parameters, as it would provide an additional cue for separating unrelated parts.

VII. CONCLUSION

We reviewed the challenges associated to robotic grasping and the importance of devising means of transferring grasping strategies across objects. We then depicted a method that allows an agent to identify, within its visuomotor experience, graspable parts that generalize across objects. Part candidates are first generated by extracting object surface segments in the vicinity of grasps demonstrated by a human. Candidates are then clustered by means of nonlinear dimensionality reduction and unsupervised learning algorithms. The central elements of the resulting clusters are selected to form a dictionary of prototypical parts that can then be used for grasping novel objects. As the dictionary of parts is only formed from cluster centers, it is allowed to be orders of magnitude smaller than the set of grasp examples initially provided to the agent. A grasp example only “votes” for the potential inclusion of a part into the dictionary, which provides us with a means of controlling the size of the dictionary in order to keep the computational cost of planning a grasp onto a novel object reasonably low. Finally, not only the shape, but also the spatial extent (or size) of the parts that form the dictionary depend on the available grasp data. Prototypical parts are selected based on their recurrence across experienced grasps, which leads to the identification of parts that strongly predict grasp applicability.

REFERENCES

- [1] J. Aleotti and S. Caselli. Part-based robot grasp planning from human demonstration. In *IEEE International Conference on Robotics and Automation*, 2011.
- [2] C. Bard and J. Troccaz. Automatic preshaping for a dextrous hand from a simple description of objects. In *IEEE International Workshop on Intelligent Robots and Systems*, pages 865–872. IEEE, 1990.
- [3] A. Bicchi and V. Kumar. Robotic grasping and contact: a review. In *IEEE International Conference on Robotics and Automation*, 2000.
- [4] C. M. Bishop. Pattern recognition and machine learning, 2006.
- [5] C. Borst, M. Fischer, and G. Hirzinger. Grasping the dice by dicing the grasp. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, volume 4, pages 3692–3697, 2003.
- [6] Y. Boykov and M.-P. Jolly. Interactive Graph Cuts for Optimal Boundary & Region Segmentation of Objects in N-D Images. In *International Conference on Computer Vision*, 2005.
- [7] J. Coelho, J. Piater, and R. Grupen. Developing haptic and visual perceptual categories for reaching and grasping with a humanoid robot. In *Robotics and Autonomous Systems*, volume 37, pages 7–8, 2000.
- [8] M. Cox and T. Cox. Multidimensional scaling. *Handbook of data visualization*, Jan. 2008.
- [9] C. de Granville, J. Southerland, and A. H. Fagg. Learning grasp affordances through human demonstration. In *IEEE International Conference on Development and Learning*, 2006.
- [10] R. Detry. *Learning of Multi-Dimensional, Multi-Modal Features for Robotic Grasping*. PhD thesis, University of Liège, 2010. Supervisor: Justus Piater.
- [11] R. Detry, D. Kraft, O. Kroemer, L. Bodenhagen, J. Peters, N. Krüger, and J. Piater. Learning grasp affordance densities. *Paladyn. Journal of Behavioral Robotics*, 2(1):1–17, 2011.
- [12] R. Detry and J. Piater. Grasp generalization via predictive parts. In *Austrian Robotics Workshop*, 2011.
- [13] R. Detry, N. Pugeault, and J. Piater. A probabilistic framework for 3D visual object representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(10):1790–1803, 2009.
- [14] C. Goldfeder, M. Ciocarlie, H. Dang, and P. Allen. The Columbia grasp database. In *IEEE International Conference on Robotics and Automation*, 2009.
- [15] A. Herzog, P. Pastor, M. Kalakrishnan, L. Righetti, T. Asfour, and S. Schaal. Template-based learning of grasp selection. In *The PR2 Workshop (Workshop at IROS’11)*, 2011.
- [16] Y. Hu, R. Eagleson, and M. A. Goodale. Human visual servoing for reaching and grasping: The role of 3-d geometric features. In *IEEE International Conference on Robotics and Automation*, 1999.
- [17] L. S. Jakobson and M. A. Goodale. Factors affecting higher-order movement planning: a kinematic analysis of human prehension. *Experimental Brain Research*, 86(1):199–208, 1991.
- [18] M. Jeannerod. The timing of natural prehension movements. *Journal of Motor Behavior*, 1984.
- [19] R. S. Johansson. Sensory input and control of grip. In *Novartis Foundation Symposium*, pages 45–58, 1998.
- [20] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3D scenes. *IEEE Trans. Pattern Anal. Mach. Intell.*, 21:433–449, 1999.
- [21] I. Kamon, T. Flash, and S. Edelman. Learning to grasp using visual information. In *IEEE International Conference on Robotics and Automation*, volume 3, pages 2470–2476, 1996.
- [22] J. Kim. M. eng. Master’s thesis, Massachusetts Institute of Technology, 2007.
- [23] A. T. Miller, S. Knoop, H. Christensen, and P. K. Allen. Automatic grasp planning using shape primitives. In *IEEE International Conference on Robotics and Automation*, volume 2, pages 1824–1829, 2003.
- [24] L. Montesano and M. Lopes. Learning grasping affordances from local visual descriptors. In *IEEE International Conference on Development and Learning*, 2009.
- [25] A. Morales, E. Chinellato, A. H. Fagg, and A. P. del Pobil. Using experience for assessing grasp reliability. *International Journal of Humanoid Robotics*, 1(4):671–691, 2004.
- [26] G. Rizzolatti, R. Camarda, L. Fogassi, M. Gentilucci, G. Luppino, and M. Matelli. Functional organization of inferior area 6 in the macaque monkey. *Experimental Brain Research*, 71(3):491–507, 1988.
- [27] G. Rizzolatti and G. Luppino. The cortical motor system. *Neuron*, 31(6):889–901, 2001.
- [28] R. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (FPFH) for 3D registration. In *IEEE International Conference on Robotics and Automation*, 2009.
- [29] M. Saric, C. H. Ek, and D. Kragic. Dimensionality Reduction via Euclidean Distance Embeddings. Technical report, KTH, Royal Institute of Technology, Stockholm, 2011.
- [30] A. Saxena, J. Driemeyer, and A. Y. Ng. Robotic Grasping of Novel Objects using Vision. *International Journal of Robotics Research*, 27(2):157, 2008.
- [31] B. Schölkopf and A. Smola. Kernel principal component analysis. *Artificial Neural Networks—ICANN’97*, 1997.
- [32] G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- [33] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [34] K. B. Shimoga. Robot grasp synthesis algorithms: A survey. *The International Journal of Robotics Research*, 15(3):230, 1996.
- [35] D. Song, C. H. Ek, K. Huebner, and D. Kragic. Multivariate Discretization for Bayesian Network Structure Learning in Robot Grasping. In *International Conference on Robotics and Automation*, pages 1–8. Royal Institute of Technology, 2011.
- [36] J. D. Sweeney and R. Grupen. A model of shared grasp affordances from demonstration. In *International Conference on Humanoid Robots*, 2007.
- [37] L. E. Zhang, M. Ciocarlie, and K. Hsiao. Grasp evaluation with graspable feature matching. In *RSS Workshop on Mobile Manipulation: Learning to Manipulate*, 2011.