

# Probabilistic Fusion of Multiple Algorithms for Object Recognition at Information Level

Matthias Lutz, Dennis Stampfer, Siegfried Hochdorfer and Christian Schlegel

University of Applied Sciences Ulm

Department of Computer Science, Prittwitzstr. 10, 89075 Ulm, Germany

{lutz, stampfer, hochdorfer, schlegel}@hs-ulm.de

**Abstract**—Reliable object recognition is a mandatory prerequisite for service robots that operate in everyday environments. Typical approaches run a single classifier for the purpose of object recognition. However, no single algorithm proved to classify across all types of objects.

We propose an approach that combines the recognition result of several methods working on different features. This reduces the effort and complexity of a single algorithm to recognize all known objects and makes the overall recognition robust. Known algorithms are extended to use a semantic output of a recognition probability for easy integration. To overcome the limitation of an algorithm to a class of objects based on their features, we introduce a probabilistic quality that defines how well an algorithm can recognize a known object type. The algorithms results are integrated using probabilistic methods to formulate a final belief.

The approach is demonstrated in practical experiments in which a service robot recognizes and grasps similar appearing objects. The experiments show that the recognition is improved by probabilistic fusion of multiple algorithms.

## I. INTRODUCTION

Service robots that navigate through and work in everyday environments must be capable to perceive the world around them. Especially in manipulation tasks, a reliable pose estimation and object classification is a mandatory prerequisite for robust real world applications.

The typical approach in object recognition for mobile manipulation is to capture a single image (in contrast to [1]) of the scene (e.g. objects on a table). Object recognition is done by executing a single algorithm on the scene image.

Impressive recognition performance has been achieved in object recognition algorithms. However, no single algorithm has proved to classify across all types of objects. They are often specialized for certain features (e.g. color, shape, texture) which limit the recognition performance to the individual objects containing these features. For example, a 3D model matcher can differ between different shapes of objects but not between two similar shapes or flavours of a product. Instead of using only a single algorithm, the combination of them uses the strengths and overcomes weaknesses of the individual ones. This leads to a more robust and reliable overall object recognition performance. For example, apple and pineapple juice may have the same shape but different color. However, a major problem when combining algorithms is the way of interpreting the results as different algorithms have different output that cannot be easily compared.

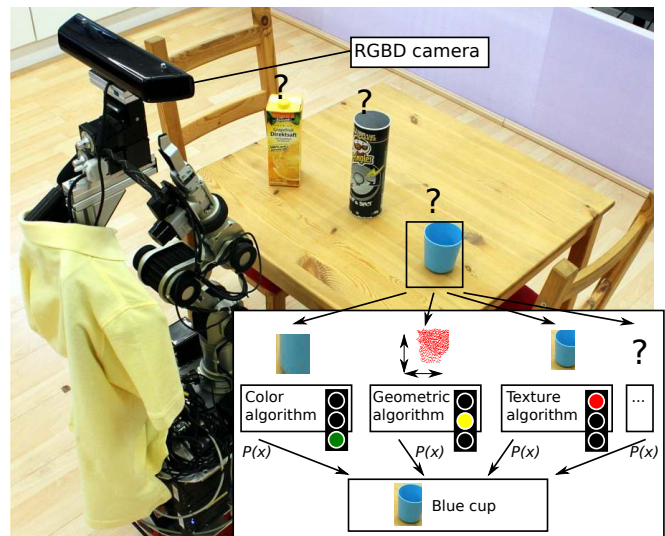


Fig. 1. A service robot in an everyday environment performing object recognition in front of a table. One object candidate (e.g. blue cup) is being recognized using several algorithms that each work on different features. A recognition probability and recognition quality (traffic light) is used to probabilistically fuse the results. The three objects are representative for three object classes in which we will recognize different types/flavours.

We propose an approach for object recognition that combines the recognition result of several methods that work on different features (fig. 1). In order to combine the results of object recognition algorithms, we extend them to use a semantic output interface of a recognition probability. Due to the probabilistic result as interface, algorithms are decoupled from fusion which makes the integration of new algorithms easy. The individual recognition results are integrated using established probabilistic methods to formulate a final belief. The result specifies the object type and the recognition probability. However, algorithms do not perform equally well on each object type, which is because features may be unsuitable for an object. Thus results must not be considered equally. For this purpose, we introduce a probabilistic quality measure for fusion that defines how well an algorithm can recognize a known object. For example, texture features are less useful than geometric features to recognize a textureless blue cup.

We demonstrate the approach (video at [2]) by recognizing and grasping household goods, where we focus on objects

with similar appearance. Recognition is done on a scene image captured with a Microsoft Kinect RGBD camera or closeup images of objects captured by a high-res camera. Such images can be obtained with an eye-in-hand camera as described in [1]. The implementation uses a simple 3D model matcher for the object shape, color detector, barcode algorithm [3] and two alternative text reading (OCR) algorithms [4], [5]. The experiments show that the recognition is significantly improved by probabilistically integrating the results of individual algorithms. Even almost identical objects can be distinguished.

## II. RELATED WORK

Impressive approaches like MOPED [6], Bag of Words [7] or algorithms in the Point Cloud Library [8] exist for object recognition. Instead of using them apart, we combine them in one recognition system to use their strengths.

A feature level fusion approach for object classification to detect cars and other traffic participants is proposed in [9]. Features are extracted and combined out of multiple sensor data sources applying a boosted cascade to classify the objects instead of combining the results in the end.

Work has been done in recognition using multiple sensors or input images. A framework for multi-sensor object recognition is presented in [10]. It is based on Conditional Random Fields incorporating temporal information to recognize cars using laser ranger data and vision data. The focus is on gaining information out of temporal relations of multiple observations of a car. [11] runs the MOPED algorithm on several images of different angles at the scene. Their approach combines the results from a single algorithm performed on multiple images. Similar to this, the field of active object recognition uses images/observations from several views on an object. In [12], a robot drives around a table and recognizes objects. The object recognition result is based on the combination of observations using one texture based recognition algorithm. Instead of using different input images for one algorithm, our method is based on one input image and different algorithms.

The work presented in [1] for active object recognition obtains close-up images of different views of objects by systematically inspecting them. It can be used to obtain input images for the system described in this paper.

The authors in [13] propose a recognition infrastructure built on a distributed message passing architecture to combine recognition algorithms. The architecture is divided into attention detector and pose estimator. Attention is used to find many prospective objects in a scene with one algorithm including false positives which are filtered in the detection step using a second algorithm. While the architecture allows to run several algorithms, it is unclear how the results can be fused. The architecture is only described in attention and filter steps.

The approach in [14] is the most similar one compared to ours. The authors combine the output of algorithms for color, depth and texture features for object recognition based on histograms. The common output interface of the algorithms is a distance measure of histogram comparisons. The object

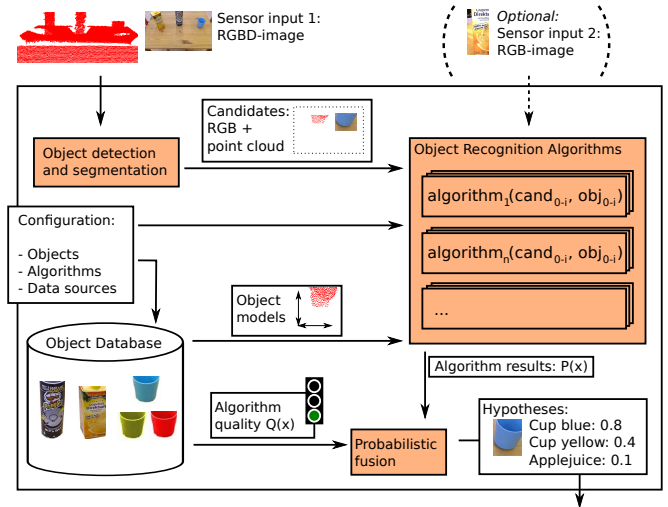


Fig. 2. Overall structure of the approach which consists of three parts: object detection/segmentation, object recognition algorithms and probabilistic fusion. In object detection and segmentation, the RGBD image is processed to generate prospective object candidates (point clouds and cropped RGB images of objects). The recognition algorithms are run on each candidate against the database. The algorithms output recognition probabilities which are fused under consideration of the quality that they can recognize a certain object. The final result is a list of object hypotheses.

identification is done using k-Nearest Neighbor but not by probabilistic fusion of algorithm results.

## III. OBJECT RECOGNITION METHOD

### A. Structure / Architecture

The overall structure of the object recognition process is illustrated in fig. 2. It can work on two sources of data: RGB-images with depth information (RGBD) and RGB-images. RGBD images are first segmented in a detection step to separate the objects in the scene. The results are object candidates which consist of a cropped RGB image and a segmented 3D point cloud. As an alternative, the object recognition can be configured at runtime to take RGB data from a different camera to input high-resolution images of one individual object as e.g. obtained with an active recognition approach like [1].

The object recognition runs several algorithms that identify objects, estimate their poses or both. They work on different features, e.g. color, texture or geometrical properties. For each candidate, each algorithm is run against each object type in the database. The algorithm's result is a probability that indicates the quality of recognition, i.e. the belief that the object from the database matches the candidate. In contrast to combining features at the algorithmic level, the combination of the recognition results with a recognition probability brings a stable and semantic interface. It allows the integration of existing algorithms since these can be extended to return a recognition probability. It is therefore not necessary to combine algorithms at the feature level.

The object recognition can be configured to only recognize objects which are expected in the scene given the current

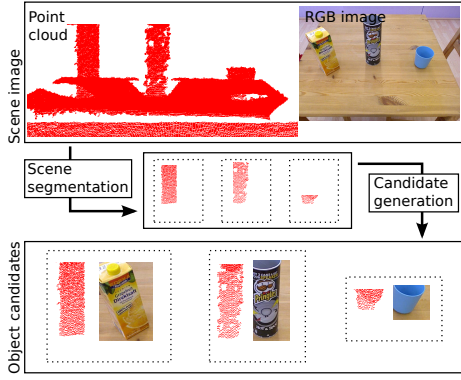


Fig. 3. The object detection and segmentation process. The scene image is segmented based on the point cloud. The resulting point clouds are used to crop the scene images which generates object candidates consisting of point cloud and RGB image.

application. This increases both speed and recognition performance as the set of objects to recognize is smaller, e.g. no hammer is expected or required for making coffee. The set of algorithms to run may further be configured for speed, e.g. not to execute time-consuming OCR algorithms if quick recognition is needed.

In a last step, the results are probabilistically fused to a final hypothesis. The fusion considers an algorithm quality which states how well an algorithm is able to identify an object. Objects which are not recognized sufficiently well (recognition probability below threshold) are considered and reported as obstacles.

### B. Scene Segmentation

The step-by-step process of segmentation is illustrated in fig. 3. The current implementation uses the basic method described in [15] for scene segmentation. In a first step, a plane is searched and removed from the 3D point cloud of the scene. The remaining points are segmented. In addition to that, we crop the RGB scene image with respect to the candidate point cloud to remove the surrounding scene from the candidate. This relieves the algorithms from finding objects in the full scene image and simplifies the object recognition from recognition in a cluttered scene to recognition of a single object in full view.

### C. Object Recognition Algorithms

The object classification itself is done in separated algorithms using different information they generate from the input data. The algorithm interface described in fig. 4 shows the input and output as well as the execution of the algorithms.

Each algorithm is supplied with the image data and point cloud data separated by the preceding segmentation step, dependent on its needs. Second, the algorithm needs reference data to compare the measurements against the models of the object types. The information required by each algorithm is stored in the object database. Each algorithm has to deliver an object type-probability pair as the most important result. The probability  $P(x)$  describes how well the candidate is

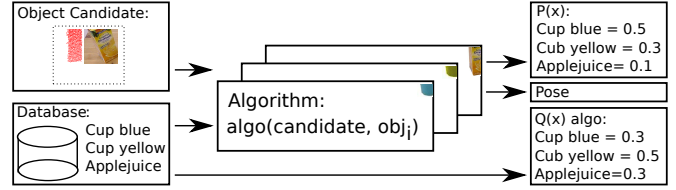


Fig. 4. Overview of an algorithm's input and output: each algorithm is executed for each object candidate against each object in the database. The algorithm's output is a list with the probability that the candidate matches the searched object and an object pose. The recognition quality of the algorithm is taken from the database and used for fusion.

recognized as the object type  $x$ . Beside the classification, each algorithm may deliver a pose estimation of the object.

Each time an algorithm is executed, it is configured to search for a single object type only. The results of the single recognition runs are combined to a list per candidate (fig. 4). To enhance the overall recognition performance, the features of the algorithms could be calculated and cached with the candidate to be used in later runs searching for other object types on the same candidate.

To cover a wide range of everyday objects, multiple algorithms using different features are applied to the image and the point cloud. A simple custom 3D model matcher algorithm uses a point cloud for object recognition and pose estimation. In case of similar objects (different flavors of the same item), the color of the packaging could deliver the required information to distinguish the objects. Therefore a color histogram algorithm is used. Objects labeled with barcodes or text are recognized by two alternative OCR and one barcode algorithm.

The data required by a new algorithm is added to the object database. The other algorithms and their data stays unaffected. It is even possible to implement specialized classifiers for sole object detection, e.g. door handles.

### D. Recognition probability $P(x)$

Probabilistic fusion requires each algorithm to calculate a recognition probability. It states how well the candidate is recognized. The algorithms are configured and executed to search for a specific sole object type, therefore the result  $P(x)$  is the probability how well the reference data matches the candidates data.

Typical algorithms do not calculate a probabilistic result but very often they do calculate a somehow scaled score value (e.g. MOPED score [6]). The calculation of the probability has to be done for each algorithm individually. In some cases this can be done very easily, like for barcode algorithms. Since EAN or UPC barcodes on products are equipped with checksums, one can assume a detection probability near to one if a barcode is recognized. For other algorithms like the simple custom model matcher, the recognition probability has to be determined based on the model data. The simple model matcher is a case of a feature based algorithm and calculates its recognition probability using the distribution of the features in a feature space. The algorithm compares the input data to the models



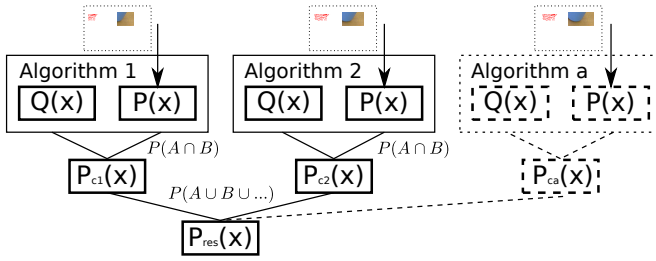


Fig. 5. Fusion of algorithm results: recognition quality  $Q(x)$  is included in recognition probability  $P(x)$ . These are probabilistically combined to a final belief  $P_{res}$ .

from the database. The models contain reference features and their normal distribution. The Mahalanobis distance from the detected feature towards the corresponding feature in the object type model is calculated. Since the Mahalanobis distance follows the Chi-Square distribution, the probability can be calculated using the cumulative distribution function.

The calculation  $P(x)$  for the other algorithms, like the color histogram or the OCR algorithm, are done in a similar way.

#### E. Algorithm Quality $Q(x)$

The results of the algorithms heavily depend on the type of object to recognize. Not all algorithms perform equally well on all object types. It is obvious that a 3D model matcher cannot distinguish between a red and a blue cup if they are of the same shape. Thus, we define a quality  $Q(x)$  for each algorithm and object type, which tells the probability that an algorithm is able to identify an object of type  $x$ .

To determine the quality, the algorithms performance is evaluated with respect to a single object type. Given the model data and the algorithms configuration, the algorithm is executed on several instances of the object type in different poses and places. The resulting list of recognition probabilities calculated by the algorithm is the foundation for the evaluation of the algorithms quality. To enforce a Bernoulli process, the recognition results are thresholded. The rate between correctly and falsely recognized objects is taken as the quality of the algorithm in relation to the object type  $Q(x)$ .

The recognition threshold to determine  $Q(x)$  is calculated maximizing the recognition probability. Therefore the algorithm is executed against a large set of segmented image or point cloud data containing various objects and obstacles. Each data set is recognized against all object types, the results are compared against ground truth. The optimal threshold is then used to determine the  $Q(x)$  as described above.

#### F. Fusion

The results of each recognition algorithm  $P(x_{c_a})$  performed on a candidate  $c$  are fused probabilistically (fig. 5). Since the performance of an algorithm varies depending on the object type to recognize, the quality of an algorithm  $Q(x)$  damps the recognition result. Assuming independence of the values, this is done using  $P_{c_a}(X_{c_a} \cap Q_a)$ . By again assuming independent measurements, the now damped results are combined using

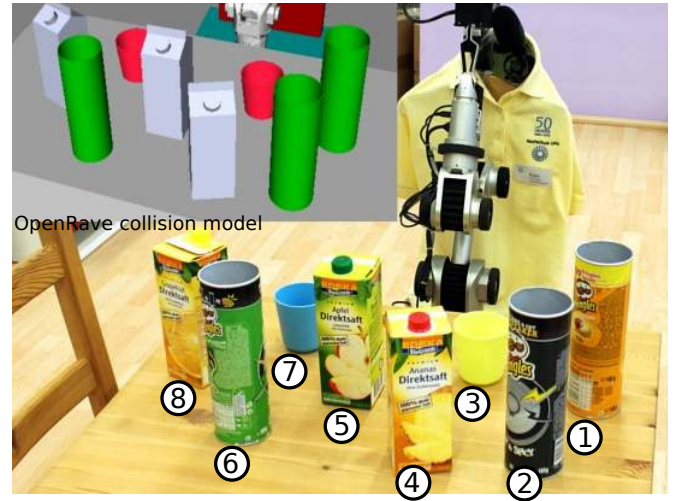


Fig. 6. Scene image of experiment 1 with eight objects of three classes. Upper left: The collision model of OpenRave used for manipulation planning to grasp objects after recognition.

$$P_{res}(X_{c_{a=1}} \cup X_{c_{a=2}} \cup \dots) = P\left(\bigcup_{a=1}^n X_{c_a}\right)$$

where  $x_{c_a}$  is the result of one object recognition algorithm for one object type. The fusion increases the recognition probability monotonically. For example, when an algorithm working on a barcode does not recognize one, the probability is not decreased. Instead, the probability of another object type may be increased if its barcode is recognized. This does not limit the recognition performance as the probability for the correct object type will outrun the others.

### IV. EXPERIMENTS AND RESULTS

The experiments are run on the service robot “Kate” which is placed in front of a table in a dining-room environment. The purpose is to recognize eight objects in two experiments. The experiments focus on the recognition of few but similar objects which are representative for household goods. The recognition performance is discussed in detail. The estimated object pose is evaluated by grasping the objects, which is the only requirement for the pose certainty.

#### A. Experimental Setup

Kate is based on a Pioneer 3-DX platform and equipped with a Microsoft Kinect RGBD camera mounted on a pan-tilt-unit. A small high-resolution (2560x1920 pixel) RGB iDS imaging uEye camera is mounted near the tool center point of the Neuronics Katana Manipulator and can be used to take close-up pictures of the objects.

The object recognition uses algorithms for color histograms, a simple 3D model matcher, two different OCR algorithms [4], [5] and one barcode algorithm [3].

Grasping of objects is done using OpenRave [16] (fig. 6). This enables safe and collision free manipulation in the scene taking into account objects and obstacles.



Fig. 7. RGB images of the candidates in experiment 1. All but object 8 were identified. The wrong classification is due to very similar appearance in both color and shape of objects 8 and 4.

### B. Experiment 1

The first experiment (fig. 6) shows eight different objects, some of them have the same shape and differ only in their color or text (different flavors of a product). The cropped candidate RGB images as generated by the segmentation are illustrated in fig. 7. Fig. 8 shows the candidates 1-8 on the x-axis, with the probabilistic results of the algorithms as well as the resulting fused results (black bars). Some of the objects are distinguishable by their shape e.g. the pringles (candidate 1, 2, 6) from the cups (candidate 3, 7). The different flavors of the pringles (candidate 1, 2, 6) are examples for objects that were separated using the color histogram algorithm.

The results in fig. 8 show the best hypotheses of each candidate only. The separation of the different object types is not visible in this figure. An example for the separation is shown in fig. 9, all object type hypotheses for candidate seven are visible. It is a good example for two objects which are not distinguishable by a single algorithm. The probabilities for the yellow and the blue cup (fig. 9 object type hypothesis 1 and 2, red bars) are the same for the model matcher (same shape reference data). The correct classification is achieved by fusing with the results of the color histogram algorithm (fig. 9 object type hypothesis 1 and 2, white bar).

Candidate 8 and 4 are an example for objects that are nearly impossible to distinguish given the Kinect data. The shapes and colors are the same and the texture is too less detailed given the resolution and the distance of the sensor.

The recognized scene represented in the OpenRave collision model (fig. 6, upper left) contains the complete 3D models of the objects to overcome the 2.5D Kinect sensor data. All objects except the tetra pack in candidate 8, which is too similar to another object type (same type as candidate 4), are recognized correctly. All objects were successfully grasped.

### C. Experiment 2

The purpose of the second experiment is to recognize three similar objects (fig. 10) using closeup image data. All three objects are very similar, one of them was not recognized correctly in the first experiment (candidate 8). The two objects on the left (fig. 10, pineapple juice and grapefruit juice) were not distinguished by color as in experiment 1, too. Again, the recognition probabilities of the color histogram (fig. 11, yellow bars) show very small values for candidate one and two. The small  $P(x)$  comes from the low color histogram recognition quality  $Q(x)$  for the two similar colored (yellow) tetra packs. They damp the recognition results  $P(x)$ . The

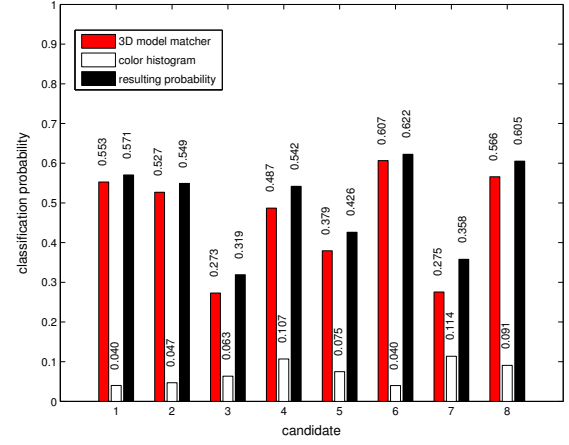


Fig. 8. Rounded classification probabilities for all candidates in experiment 1. Only the best object type hypotheses for each candidate is shown. The color of the bars show both results of the recognition algorithms (including  $Q(x)$ ) as well as the resulting fused probability.

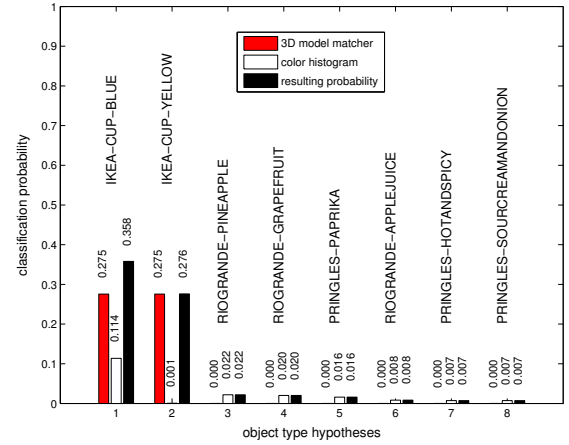


Fig. 9. Rounded classification probabilities of candidate three for all object types to recognize in experiment 1. The color of the bars show both results of the recognition algorithms (including  $Q(x)$ ) as well as the resulting probability.

grapefruit juice (candidate three) is correctly recognized with a recognition probability of 0.14, later damped by  $Q(x) = 0.2$  to  $P(x) = 0.029$ .  $Q(x)$  is found empirically, the low value can be explained by two different objects with the same color (grapefruit and pineapple juice).

Using closeup images, a detection of text features using OCR and a barcode detection is possible. Since the performance of the used OCR algorithms varies both of them are used concurrently. No barcode was found for candidate three (fig. 10) (apple juice), but the object is recognized by OCR. This is a good example why the quality of the barcode algorithm  $Q(X) = 0.6$  is not nearly as good as one would expect. If a barcode is recognized the object is classified very reliably. In some cases the barcode is however visible (fig. 10 candidate 3) but not recognized, for example in the case of reflections or an blurry image. Overall, the three candidates were recognized correctly, even though they are

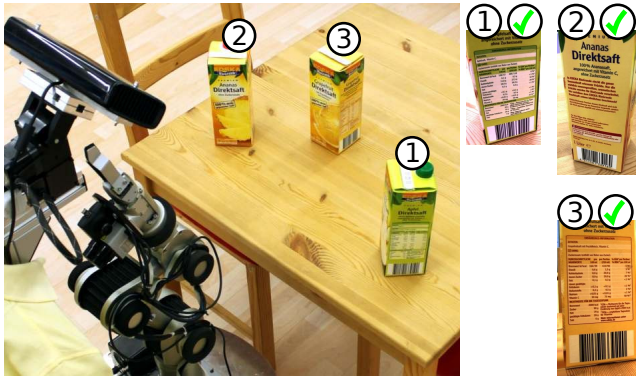


Fig. 10. The scene of experiment 2 (left). All objects were identified using color, text and barcode features on closeup images (right).

hard to distinguish. All objects were successfully grasped.

#### D. Results

The experiments show that the object recognition probability is improved by integrating the results of different algorithms. None of the algorithms alone would have been able to identify the objects reliably enough.

It is shown that even almost identical objects can be recognized reliably. The separation is achieved by applying algorithms capable of distinguishing suchlike objects. Sometimes this results in a small gap of the resulting recognition probability only. Nevertheless, this gap is enough to make the recognition successful. The recognition probability proved to be a good measure to combine the recognition result. The integration of the quality of an algorithm  $Q(x)$  also proved to be a meaningful method to incorporate the expressiveness of the recognition probability of the algorithms.

All objects were successfully grasped which shows that the pose is sufficient enough for practical manipulation applications. A video is available at [2].

#### V. CONCLUSION AND FURTHER WORK

The proposed method for object recognition combines the results of multiple recognition algorithms at information level. Fusing the results at an information level using probabilities enables the usage of probabilistic calculations/theorems and decouples the algorithms making the integration of new algorithms easy. The experiments demonstrate the performance of the approach in everyday environments. It is possible to recognize objects that are hard to distinguish. The proposed object recognition system is used in several real world mobile manipulation scenarios, see videos in [2].

Further work will evaluate more algorithms to enhance the recognition of similar objects. More work in the object pose estimation will be done. Till now the pose is estimated by one algorithm only.

#### ACKNOWLEDGMENT

This work has been conducted within the ZAFH Servicerobotik (<http://www.zafh-servicerobotik.de/>). The authors gratefully acknowledge the research grants of the state of

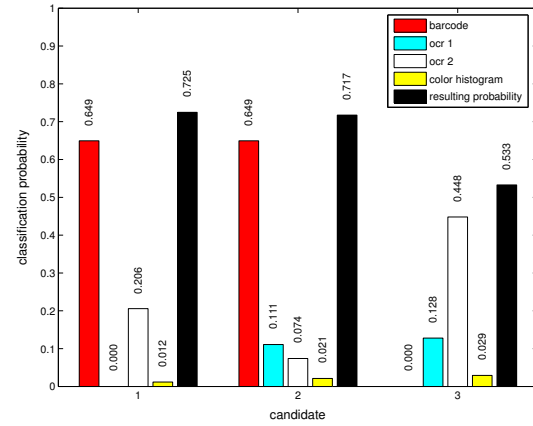


Fig. 11. Rounded classification probabilities for all candidate in experiment 2. Only the best object type hypotheses for each candidate is shown. The color of the bars show the results of the recognition algorithms (including  $Q(x)$ ) as well as the resulting probability.

Baden-Württemberg and the European Union. We thank our project partners in Ravensburg-Weingarten for their valuable feedback and discussions in the early conceptual stages.

#### REFERENCES

- [1] D. Stampfer, M. Lutz, and C. Schlegel, "Information Driven Sensor Placement for Robust Active Object Recognition based on Multiple Views," in *IEEE Int. Conf. on Technologies for Practical Robot Applications*, Woburn, MA, USA, 2012.
- [2] YouTube: Robotics@HS-Ulm, <http://www.youtube.com/roboticsathulm>.
- [3] ZBar bar code reader, <http://zbar.sf.net>, visited: Nov. 11th 2011.
- [4] R. Smith, "An Overview of the Tesseract OCR Engine," in *Int. Conf. on Document Analysis and Recognition*, vol. 2, Sept. 2007, pp. 629–633.
- [5] ABBYY OCR, <http://ocr4linux.com>, visited: Nov. 11th 2011.
- [6] A. Collet, M. Martinez, and S. S. Srinivasa, "The MOPED framework: Object Recognition and Pose Estimation for Manipulation," *The International Journal of Robotics Research*, 2011.
- [7] J. R. R. Uijlings, A. W. M. Smeulders, and R. J. H. Scha, "Real-Time Bag of Words, Approximately," in *ACM Int. Conf. on Image and Video Retrieval*, 2009.
- [8] Point Cloud Library, <http://pointclouds.org>, visited: Nov. 28th 2011.
- [9] S. Wender and K. Dietmayer, "A Feature Level Fusion Approach for Object Classification," in *IEEE Intelligent Vehicles Symposium*, Istanbul, Turkey, June 2007, pp. 1132–1137.
- [10] B. Douillard, D. Fox, and F. Ramos, "A Spatio-Temporal Probabilistic Model for Multi-Sensor Object Recognition," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, San Diego, USA, Nov. 2007, pp. 2402–2408.
- [11] A. Collet and S. S. Srinivasa, "Efficient Multi-View Object Recognition and Full Pose Estimation," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, Anchorage, Alaska, USA, May 2010, pp. 2050–2055.
- [12] R. Eidenberger and J. Scharinger, "Active Perception and Scene Modeling by Planning with Probabilistic 6D Object Poses," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, Oct. 2010, pp. 1036–1043.
- [13] M. Muja, R. B. Rusu, G. Bradski, and D. G. Lowe, "REIN - A Fast, Robust, Scalable REcognition INfrastructure," in *IEEE Int. Conf. on Robotics and Automation*, Shanghai, China, May 2011, pp. 2939–2946.
- [14] M. Attamimi, A. Mizutani, T. Nakamura, T. Nagai, K. Funakoshi, and M. Nakano, "Real-Time 3D Visual Sensor for Robust Object Recognition," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, Taipei, Taiwan, Oct. 2010, pp. 4560–4565.
- [15] M. Wopfner, J. Brich, S. Hochdorfer, and C. Schlegel, "Mobile Manipulation in Service Robotics: Scene and Object Recognition with Manipulator-Mounted Laser Ranger," in *ISR/ROBOTIK 2010*, Munich, Germany, 2010.
- [16] R. Diankov and J. Kuffner, "OpenRAVE: A Planning Architecture for Autonomous Robotics," Robotics Institute, Pittsburgh, PA, Tech. Rep. CMU-RI-TR-08-34, July 2008.