# Latent-Class Hough Forests for 6 DoF Object Pose Estimation

Rigas Kouskouridas, Alykhan Tejani, Andreas Doumanoglou, Danhang Tang, Tae-Kyun Kim, *Member IEEE*

**Abstract**—In this paper we present *Latent-Class Hough Forests*, a method for object detection and 6 DoF pose estimation in heavily cluttered and occluded scenarios. We adapt a state of the art template matching feature into a scale-invariant patch descriptor and integrate it into a regression forest using a novel template-based split function. We train with positive samples only and we treat class distributions at the leaf nodes as latent variables. During testing we infer by iteratively updating these distributions, providing accurate estimation of background clutter and foreground occlusions and, thus, better detection rate. Furthermore, as a by-product, our *Latent-Class Hough Forests* can provide accurate occlusion aware segmentation masks, even in the multi-instance scenario. In addition to an existing public dataset, which contains only single-instance sequences with large amounts of clutter, we have collected two, more challenging, datasets for multiple-instance detection containing heavy 2D and 3D clutter as well as foreground occlusions. We provide extensive experiments on the various parameters of the framework such as patch size, number of trees and number of iterations to infer class distributions at test time. We also evaluate the *Latent-Class Hough Forests* on all datasets where we outperform state of the art methods.

**Index Terms**—Object detection, pose estimation, Hough forests, one-class training, 6 DoF Pose Estimation

✦

## 1 INTRODUCTION

AMONG the most challenging tasks in computer vision is the one of estimating the 3D pose of an object due to its practical implication and its fundamental importance applications like robotic manipulation [19] and tracking [8]. In order to efficiently fulfill the 3D object detection and pose estimation task, a computer vision method should tackle several cascading issues that hinder its effective application. Although recent emergence of consumer depth sensors provides additional cue in favour of textureless objects, background clutter, partial occlusions and large scale changes still put barriers to this problem. Template matching techniques [16] can tackle clutter and occlusion to some degree, but have inherent weakness due to their holistic nature. Point-to-Point approaches fail in cases of planar, self-similar or similar to background clutter objects [2], [10] due to the fact that similar point features vote for different pose parameters. Moreover, these methods were only evaluated with the assumption of only one instance existing in the scene. The case of multiple object instances, i.e., challenging precision-recall, is left unexplored.

Another important perspect is that, prior arts in 3D object pose estimation [2], [10], [16] utilize mesh models of target objects to generate training samples. This implies that only positive samples are used and, thus, falls into the category of one-class learning. On the other hand, to explicitly tackle the aforementioned challenges, a more traditional way in 2D
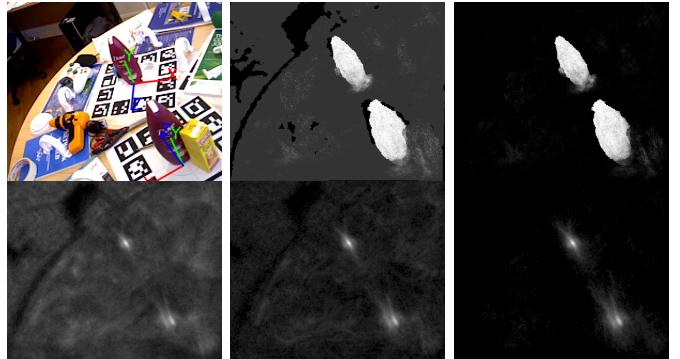


Fig. 1. An illustration of intermediate results of the iterative process. Column 1: Initial pose estimation and the corresponding hough map. Columns 2 - 3: Foreground probability masks and the respective hough maps after $\#5$ and $\#10$ iterations.

detection is to augment negative samples (clutter and occluder) during training. Due to the complexity of negative scenarios, this often results in huge amount of training samples, thus, increasing the computation burden of the system. And yet no such work can guarantee covering all cases. At present, there is a big disparity in the number of depth image datasets vs. 2D image datasets, adding a further challenging in mining for negative depth samples.

Highly motivated by these challenges, we present a novel method, called *Latent-Class Hough Forests*, for 3D object detection and pose estimation. Unlike traditional Hough Forest [12], which explicitly exploits class information, our method utilizes only the regression term during the training

---

- *R. Kouskouridas and A. Tejani contributed equally to this work.*
- *R. Kouskouridas, A. Doumanoglou, D. Tang and TK Kim are with the Department of Electrical and Electronic Engineering, Imperial College London, SW7 2AZ, UK. E-mail: {r.kouskouridas, a.doumanoglou12, d.tang11, tk.kim}@imperial.ac.uk*
- *A. Tejani is with Blippar, 1 London Bridge, London, SE1, UK. E-mail: info@alykhantejani.com*

stage. Also differing from a regression forest, *Latent-Class Hough Forests* take into account class distribution at leaf nodes. During testing, the distribution of positive / negative classes are considered as latent variables being updated iteratively, in order to provide more and more accurate voting results. As a byproduct, our system can also produce accurate occlusion-aware segmentation masks. Figure 1 demonstrates the effect of our inference algorithm. At iteration#0, the forest has no background information, thus all pixels are considered as foreground, which results in a noisy vote map. As the iterative process goes on, it is evident that background pixels are greatly suppressed and both the occlusion mask and vote map become more accurate and clean.

The paper in hand presents a novel method for object detection and 6 DoF pose estimation in real scenarios by adopting a part-based strategy into the random forest framework. The previous conference version [40] has been extended by a more detailed discussion of the key ingredients of the method and an extensive analysis of the parameters of the technique. Moreover, we propose an additional dataset inspired by industrial settings as well as reporting more experiments on three different datasets. Our main contributions can be summarized as follows:

- We propose the *Latent-Class Hough Forests*, a novel patch-based approach to 3D object detection and pose estimation; It performs one-class learning at the training stage, and iteratively infers latent class distributions at test time.
- We adapt the state of the art 3D holistic template feature, LINEMOD [14], to be a scale invariant patch descriptor and integrate it into the random forest framework via a novel template-based splitting function.
- During the inference stage, we jointly estimate objects' 3D location and pose as well as a pixel wise visibility map, which can be used as an occlusion aware figure-ground segmentation.
- We provide two new, more challenging public datasets for *multi-instance* 3D object detection and pose estimation, comprising *near and far range 2D and 3D clutter* as well as *foreground occlusions* in domestic and industrial scenarios. To the best of our knowledge, we are the first to provide a fully-annotated bin-picking dataset.

In the remainder of this paper we first discuss related work in Sec. 2 before introducing our method in Sec. 3. Following this, in Sec. 4, we provide a quantitative and qualitative analysis of our results as well as a comparison to current state of the art methods. Finally, in Sec. 5, we conclude with some final remarks and a discussion of future work.

## 2 RELATED WORK

Throughout the years several techniques for the detection and registration of objects in 3D environments have been proposed. According to the literature, three main categories can be distinguished: Template matching, learning-based methods and Point-to-Point techniques. The simplicity along with facile training sessions render template matching methods as one of the most widely used solutions for object detection tasks.

From the very known baseline techniques of LINEMOD [14] and its extension [31], to the classic implementation of Distance Transform approaches [24], template matching methods have found application in contemporary vision tasks and robotics modules, respectively. On the other hand, learning-based approaches impose upon laborious training sessions with numerous training samples with view to extract highly representative object models [25], [20]. Point-to-Point techniques build upon point pair features to construct object models based on point clouds. A representative method of this category is the one presented by Drost *et al.* [10]. In turn, simple pixel-based features have been also employed to tackle the object pose estimation problem. More recently, Brachmann *et al.* [5] introduced a new representation in form of a joint 3D object coordinate and class labelling (extended for tracking in [21]), which, however, suffers in cases of occlusions.

Moreover, in [23] a method for fine pose estimation by representing geometric and appearance information as a collection of 3D shared parts and objectness, has been presented. Song *et al.* [36] proposed a computationally expensive approach to the 6 DoF pose estimation problem that slides exemplar SVMs in the 3D space, while in [4] shape priors are learnt by soft labelling random forest for 3D object classification and pose estimation. Wu *et al.* [43] designed a model that learns the joint distribution of voxel data and category labels using a Convolutional Deep Belief Network, while the posterior distribution for classification is approximated by Gibbs sampling. From the relevant literature we could also identify the works of Aldoma *et al.* [1] and Buch *et al.* [7] that propose a final step of fine pose refinement for false positive / outlier removal. Last, Wohlhart *et al.* [42] showed how a Convolutional Neural Network can be trained to learn a 3D pose estimation-wise descriptor.

It is well understood that modeling objects as a collection of parts increases robustness to intra-class variation, pose change and even occlusion. The implicit shape model, introduced by Leibe *et al.* [22], learns, via unsupervised clustering, class-specific visual codebooks and spacial distributions for each entry. Codebook entries are then detected in the test image and used to cast probabilistic votes in the Hough space based on the learnt spatial distributions. Gall and Lempitsky showed, with the class-specific Hough forest [12], how part-based modeling can be effectively combined with generalized Hough voting for object detection under the random forest framework [6]. Tang *et al.* [37] combined Hough Forest with DOT [15] with a template matching split function, while highly efficient, requires extensive and diverse background images for training.

On the other hand, one-class training stands for the learning process that imposes upon training without negative samples. Introduced by Moya *et al.* [27], [28] and further developed by Tax [39] and Scholkopf [33], these approaches lay their foundations within the support vector framework and aim to derive an enclosing decision boundary of the training data as a whole from a few supporting samples. Other techniques such as the works of Bishop [3] and Para *et al.* [29] approached the problem in a probabilistic manner and tried to find the underlying density model for the training data as a whole.

Occlusion-handling is well-related to object segmentation problem. In traditional 2D scenarios, existing methods can be categorised into tackling occlusion in training [26], [30] or inference stage [17], [18]. Utilizing depth cues is relatively new and recently Wang *et al.* [41] approached occlusion reasoning by explicitly learning clutter and oclusion scenarios during the training session, whilst our method falls into the inference category.

# 3 PROPOSED METHOD

In the field of object detection and 3D pose estimation, LINEMOD [14], a binary RGB-D feature, has demonstrated both state of the art accuracy and efficiency. However, so far it has been combined with a holistic template matching scheme, which has inherent problem with occlusion due to the nature of holism. Moreover, as a near-neighbour search, this scheme slows down linearly as the number of templates grows. And the fact that LINEMOD is not scale invariant often leads to thousands of templates per object, in order to cover multiple scales and numerous viewpoints. To be more robust to occlusion and clutter, we start off by combining a state of the art 3D feature and a part-based detector. In this work, we choose the state of the art part-based detector Hough Forest [12]. However in our case, naively combining them does not work because: a) As a binary feature, LINEMOD only considers orientation information whilst discarding magnitude. This provides efficiency but degrades the accuracy in the case of a patch-based detector. b) No negative (clutter / occlusion) information is available during training, which means the *classification* term in Hough Forest cannot be adopted. c) Moreover, not knowing the boundary between positive and negative samples leads to large amount of false positives during inference.

To address these issues, we propose *Latent-Class Hough Forests* to datamine useful clutter / occlusion information from inference stage and transfer the knowledge to the detector. In section 3.1 we describe how to modify the similarity measurement of LINEMOD and integrate it into the split function of *Latent-Class Hough Forests*. Section 3.2 depicts the inference process that jointly and iteratively updates the latent class distribution and voting results.

## 3.1 Learning

*Latent-Class Hough Forests* are an ensemble of randomized binary decision trees trained using the general random forest framework [6]. During training, each tree is built using a random subset of the complete training data. Each intermediate node in the tree is assigned a split function and threshold to optimize a measure of information gain; this test is then used to route incoming samples either left or right. This process is repeated until some stopping criteria is met, where a leaf node containing application-specific contextual information is formed. Each stage in this learning process is highly application dependent and we will discuss each in turn below. Fig. 2 illustrates the high level idea underlying our training module. A typical split function of a random forest can be formulated as below:

$$h_i(x) = \begin{cases} 0, & S(x, \rho_i) \leq \tau_i \\ 1, & S(x, \rho_i) > \tau_i \end{cases}, \qquad (1)$$

where $\rho_i$ is the parameter and $\tau_i$ is the threshold stored at node $i$. $S$ is a test function that evaluates the input sample given $\rho_i$. The process of training a tree is to decide the optimal $\rho_i$ and $\tau_i$ by measuring the information gain.

### 3.1.1 Training Data

In order to capture reasonable viewpoint coverage of the target object, we render synthetic RGB and depth images by placing a virtual camera at each vertex of a subdivided icosahedron of a fixed radius. A tree is trained from a set of patches sampled from the training images. We extract patches with size relative to the bounding box of the rendered object, while the template features are evenly spread across each patch; features capturing the image gradients are taken only from the object contours and features capturing the surface normals are taken from the body of the object. Moreover, the collection and representation of template features is the same as described in [14].

### 3.1.2 Split Function

It has been shown in the 2D detection problem [37] that the original test function (so called two-pixel test) of Hough Forest does not work well with a binary feature. Thus, doing a naive holistic patch comparison, or the two-dimenson / two-pixel tests (as used in [34], [11], [38]) can lead to test patches taking the incorrect route at split functions. To compensate the loss of magnitude information whilst keeping the efficiency, one possible way is to utilize the orientation information in a more effective way. To this end, a non-linear template matching test function is adopted as described below:

$$
\begin{aligned}
S(\mathcal{X}, \rho) = S(\mathcal{X}, \mathbb{T}) &= \sum_{r \in \mathcal{P}} g(\mathbf{ori}(\mathcal{X}, r), \mathbf{ori}(\mathcal{O}, r)) \\
&= \sum_{r \in \mathcal{P}} \Big( \max f_m\big(\mathcal{X}(r), \mathcal{O}(r)\big) \Big)
\end{aligned}
\qquad (2)
$$

where $\mathcal{X}$ is a testing patch, $\mathbb{T} = (\mathcal{O}, \mathcal{P})$ is the template of a patch $\mathcal{O}$ with a list $\mathcal{P}$ of features. $f_m$ is the dot product between the gradient orientation at location $r$ of RGB-D patch $\mathcal{X}$ and and $\mathcal{O}$, respectfully. The similarity measure $f_m$ is similar to [14], while, here we show how it can be adapted to work over patches. Combining Eq. 1 and 2 gives us a non-linear split function, which performs much better than axis-aligned and linear split functions, yet, has the complexity similar to an axis-aligned one, since it involves only bitwise operations that can SSE-accelerated.

The test function of Eq. 2 performs well within the object boundary, but poorly around the edge. Mainly because in inference stage, clutter and occluder around target object make the similarity measurement between patches to fail. See Fig. 3 for an illustration of this issue. To tackle this, we modify the test function by incorporating an efficient $z$-value check:

$$
\begin{cases}
S(\mathcal{X}, \mathbb{T}) &= \sum\limits_{r \in \mathcal{P}} f(\mathcal{X}, \mathcal{O}, c, r) g(\mathbf{ori}(\mathcal{X}, r), \mathbf{ori}(\mathcal{O}, r)), \\
f(\mathcal{X}, \mathcal{O}, c, r) &= \delta(|(D(\mathcal{X}, c) - D(\mathcal{X}, r)) - (D(\mathcal{O}, c) - D(\mathcal{O}, r))| < \tau)
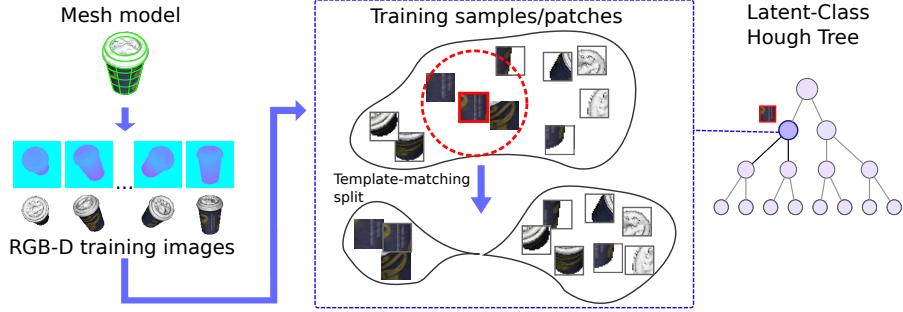\end{cases}
$$
$$(3)$$

Fig. 2. During training, a random patch $\mathbb{T}$ is selected (red frame) as a template. The similarity between it and all other patches is measured and splitted based on a threshold $\tau$ (dotted circle). This process is repeated until the optimal $\mathbb{T}$ and $\tau$ are found and stored in the current node.
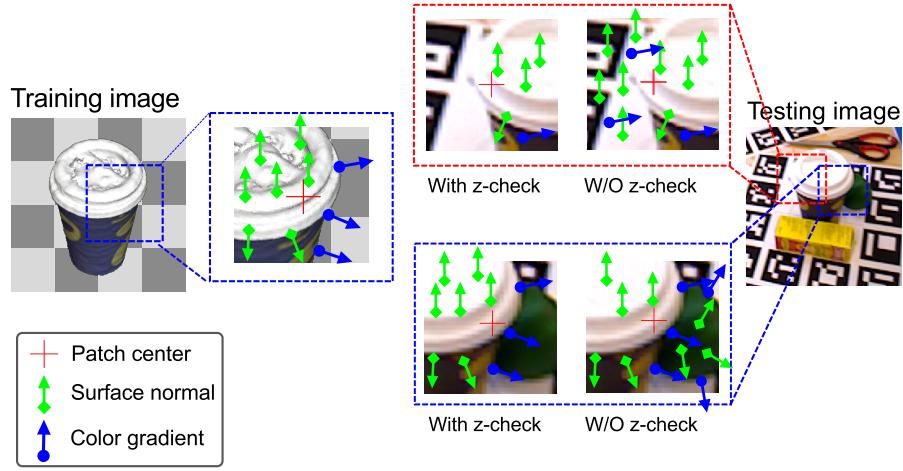


Fig. 3. $z$-value check enhance the robustness against clutter and occluders. Blue patches indicate a true positive match and red ones a false positives one. Without the z-check, features of the planar background of the false patch could match the training patch and become a false positive.

where, for a patch centered at position $c$, $D(a, b)$ retreives $z$-value from position $b$ of patch $a$ and $f$ is an indicator function that checks if the depth difference between two patches is larger than a threshold $\tau$.

The original LINEMOD is not scale invariant, which results in a large amount of templates sampled from scale space. Inspired by [34], we achieve scale invariance by normalizing feature offsets by its depth value.

$$\begin{cases} S(\mathcal{X}, \mathbb{T}) & = \sum_{r \in \mathcal{P}} f(\mathcal{X}, \mathcal{O}, c, r) g(\mathbf{ori}(\mathcal{X}, \frac{r}{D(\mathcal{X},c)}), \mathbf{ori}(\mathcal{O}, \frac{r}{D(\mathcal{O},c)})), \\ f(\mathcal{X}, \mathcal{O}, c, r) & = \delta(|(D(\mathcal{X}, c) - D(\mathcal{X}, \frac{r}{D(\mathcal{X},c)})) - (D(\mathcal{O}, c) - D(\mathcal{O}, \frac{r}{D(\mathcal{O},c)}))| < \tau) \end{cases}$$
(4)

During training, at each split node, we randomly choose a patch as a template $\mathbb{T}$ and measure its similarity with all other patches. Those patches with similarity larger than a threshold $\tau$ will go to one child node, whilst the rest will go to the other. This split is measured with information gain and repeated multiple times until an optimal one is found.

### 3.1.3 Constructing Leaf Nodes

The training data is recursively split by this process until the tree has reached a maximum depth or the number of samples arriving at a node fall below a threshold. When either one of these criteria is met a leaf node is formed from the patches reaching it. As to the information gain, since no negative information is available during training, thus, we cannot use a classification term but only the regression one for measuring information gain. However, differing from a regression forest, *Latent-Class Hough Forests* still store the class distribution in leaf nodes. Following the approach of Girshick *et al.*[13] we only store the indexes of the training patches that reached at the each leaf node and the modes of the distribution, which we efficiently calculate via the MeanShift algorithm. We create a class distribution at the leaf, however, as no background information reaches the leaves during training this distribution is initialized to $p_{fg}^l = 1$ and $p_{bg}^l = 0$ for the foreground and background probabilities, respectively.

### 3.2 Inference

After learning, we have a *Latent-Class Hough Forest* trained with positive samples only. In inference stage, we propose an iterative algorithm to datamine the negative information.
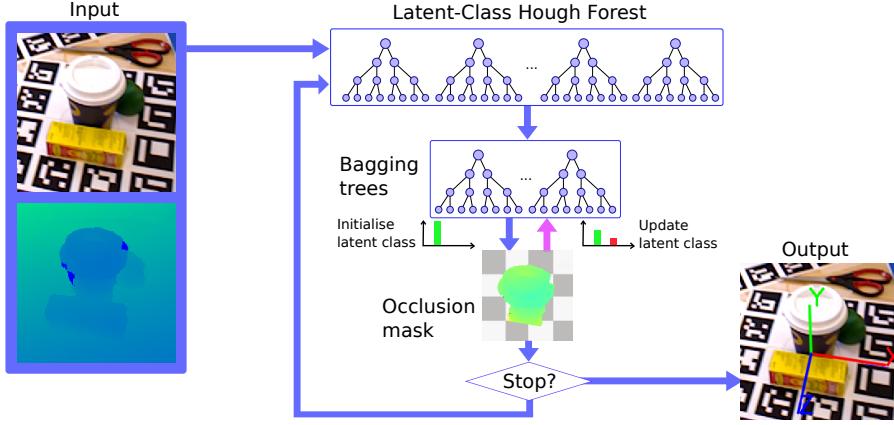
Fig. 4. Inference process: Input LINEMOD patches are extracted from RGB-D images. For each iteration, a subset of trees are drawn from the forest with bagging.

### 3.2.1 Hough Voting

Let $E(\theta)$ represent the probability of the random event that the target object exists in the scene under the 6D pose $\theta = (\theta_x, \theta_y, \theta_z, \theta_{roll}, \theta_{pitch}, \theta_{yaw})$. The goal of the voting process is to aggregate the conditional probabilities $p(E(\theta)|\mathcal{P})$ for each patch $\mathcal{P}$, given that each vote originates from a foreground patch ($p_{fg} = 1$). For a patch $P$ reaching leaf node $l$ of tree $T$, the conditional probability is formalized as follows:

$$
\begin{aligned}
p\left(E\left(\theta\right)|\mathcal{P}; \mathcal{T}\right) &= p\left(E\left(\theta\right), p^l_{fg} = 1|\mathcal{P}\right) \\
&= p\left(E\left(\theta\right)|p^l_{fg} = 1, \mathcal{P}\right) \cdot p\left(p^l_{fg} = 1|\mathcal{P}\right)
\end{aligned}
\tag{5}
$$

where $p^l_{fg}$ is the foreground probability at the leaf node, $l$. Finally, for a forest, $\mathcal{F}$, we simply average the probabilities over all trees:

$$
p\left(E\left(\theta\right)|\mathcal{P}; \mathcal{F}\right) = \frac{1}{|\mathcal{F}|} \sum_t^{|\mathcal{F}|} p\left(E\left(\theta\right)|\mathcal{P}; \mathcal{T}_t\right)
\tag{6}
$$

The first factor, $p\left(E\left(\theta\right)|p_{fg} = 1, \mathcal{P}\right)$, can be estimated by passing each patch down the forest and accumulating the votes stored at the leaf, in which votes from multiple trees can be combined in an additive manner, this gives us the same probabilities as in Eq. 6 up to a constant factor. The estimation is then deferred to the ability of locating local maxima in this aggregated space. Traditionally there are two different methods of locating targets in vote space. One is to aggregate all votes in the same position of vote space, and return the peaks with non-maximum suppression [12]. The other is to treat each vote as a data point, and then use MeanShift to locate the mode. In this case an assumption of only one instance in the scene is made to avoid local minimum [11]. The former is efficient for locating but less accurate than the latter, especially when the votes are sparsely distributed.

To accommodate both efficiency and accuracy, we propose a 3-stage localization technique in which we first aggregate all votes in 2D image space, use this as a score function and locate the top $N$ vote peaks as valid hypotheses; and then use them as initialization for MeanShift to locate modes in 3D translation, $(\theta_x, \theta_y, \theta_z)$; finally find the mode in rotation, $(\theta_{roll}, \theta_{pitch}, \theta_{yaw})$ given translation.

### 3.2.2 Update Class Distribution

As mentioned in Section 3.1, class distributions are initialized as all foreground. For each pixel $\mathbf{x} = (\mathcal{P}, l)$, located at the center position $l$ of patch $\mathcal{P}$, the initial probability of being foreground is $p^{\mathbf{x}}_{fg} = 1$. Inspired by [22], for each valid hypothesis $\theta$, we backproject to obtain a consensus patch set, i.e. peak of the probability distribution in the Hough space where multiple patches cast their vote, is considered as valid only in cases where its score is above a threshold. This threshold value is different for each object while, our three stage localization technique includes also thresholds per stage. All the thresholds are object-specific and are estimated via trial and error. All consensus voting patches are considered to be foreground and the rest background.

$$
p^{\mathbf{x}}_{fg} = \delta((\sum_\theta (p^{\mathbf{x}}_{fg}|\theta) p(\theta)) > 0),
\tag{7}
$$

in which $p(\theta)$ indicates us whether this hypothesis is valid or not and $(p^{\mathbf{x}}_{fg}|\theta)$ suggests whether $\mathbf{x}$ has voted for $\theta$. As long as $\mathbf{x}$ has voted for one valid hypothesis, it is, then, considered as a foreground pixel ($\delta(\cdot) = 1$). In other words, after a valid hypothesis is found via patch voting, we move down to pixel level, where patches that voted for this hypothesis are back-projected to the scene. Since we have an estimation of the object's center in the scene and we know also its diameter, pixels (of the voting patches) that are not spatial consistent (distance to the object centre larger than the diameter) are considered to belong to the background and the rest to the foreground. With this we can update the latent class distribution of each leaf node $l$ by:

$$
p^{\mathbf{l}}_{fg} = \frac{\sum_{\mathbf{x}} p^{\mathbf{x}}_{fg}}{|\mathbf{x}|}, \mathbf{x} \in l,
\tag{8}
$$

which can be translated to the calculation of the normalized portion of foreground patches that arrives at this leaf node.

### 3.2.3 Iterative Process

The estimation obtained with previous step is inaccurate and uncertain. However, with an iterative process, we update the probability distribution $p_{fg}^l$ for every iteration. To avoid error accumulation of each round, we propose to draw a subset of trees with bagging, such that not all the trees are updated by previous iteration. The whole process is illustrated in Figure 4 and described in Algorithm 1.

### 3.2.4 Final segmentation

With the estimated hypotheses $\{\theta\}$ and a foreground probability mask $\mathcal{Z}$ generated with $p_{fg}^{\mathbf{x}}$, we can easily obtain a final occlusion-aware segmentation by

$$\mathcal{M} = \mathbb{B}(\theta)) \cap \mathcal{Z}, \tag{9}$$

where $\mathbb{B}$ is a bounding box centered at $\theta$. This is helpful for further refinement step such as ICP, which does not work well when occlusion presents.

---

**Algorithm 1** Inference process

---
**Require:** An input image $\mathcal{I}$; A Hough Forest $\mathcal{F}$
1: **repeat**
2:     Draw a subset of trees $\mathcal{F}^*$ with bagging
3:     Randomly sample a set of patches $P$ from $I$
4:     Propagate $P$ down $\mathcal{F}^*$ and vote for $\theta$ (Eq. 6)
5:     Backproject to obtain a foreground mask (Eq. 7)
6:     Partition $P$ into positive and negative subsets by foreground mask

$$P^{fg} = \{\mathbf{x}|\mathbf{x} \in Z\}$$
$$P^{bg} = P \backslash P^{fg}$$

7:     Update the probabilities at leaf nodes (Eq. 8) given $P^{fg}$ and $P^{bg}$.
8: **until** Maximum iteration

---

## 4 EXPERIMENTS

We perform experiments on three 3D pose estimation datasets. The first is the publicly available dataset of of Hinterstoisser *et al.* [16], which contains 13 distinct objects each associated with an individual test sequence comprising of over 1,100 images with close and far range 2D and 3D clutter. Each test image is annotated with ground truth position and 3D pose. We also introduce our two new datasets, called *Domestic Environments Dataset* and *Bin-picking Dataset* to further evaluate the efficiency of our method in real scenarios. In all tests we use the metric defined in [16] to determine if an estimation is correct. More formally, for a 3D model $\mathcal{M}$, with ground truth rotation $\boldsymbol{R} = (\theta_{pitch}, \theta_{yaw}, \theta_{roll})$ and translation $\boldsymbol{T} = (\theta_x, \theta_y, \theta_z)$, given an estimated rotation, $\hat{\boldsymbol{R}} = (\hat{\theta}_{pitch}, \hat{\theta}_{yaw}, \hat{\theta}_{roll})$ and translation, $\hat{\boldsymbol{T}} = (\hat{\theta}_x, \hat{\theta}_y, \hat{\theta}_z)$, the matching score is defined as

$$m = \underset{\boldsymbol{x} \in \mathcal{M}}{avg} ||(\boldsymbol{R}\boldsymbol{x} + \boldsymbol{T}) - (\hat{\boldsymbol{R}}\boldsymbol{x} + \hat{\boldsymbol{T}})|| \tag{10}$$

for non-symmetric objects and

$$m = \underset{\boldsymbol{x_1} \in \mathcal{M}}{avg} \underset{\boldsymbol{x_2} \in \mathcal{M}}{\min} ||(\boldsymbol{R}\boldsymbol{x_1} + \boldsymbol{T}) - (\hat{\boldsymbol{R}}\boldsymbol{x_2} + \hat{\boldsymbol{T}})|| \tag{11}$$

for symmetric objects. An estimation is deemed correct if $m \leq k_m * d$, where $k_m$ is a chosen coefficient and $d$ is the diameter of $\mathcal{M}$. We would like to note that the symmetric version (Eq. 11) of the matching score assumes complete symmetry around an axis, like bowls. However, other objects that with, for example, square shape have an advantage if assumed symmetric.

Unlike [16], in which only the top N detections from each image are selected, we also compute precision-recall curves and present the F1-Score which is the harmonic mean of precision and recall. We argue that this is a more accurate form of comparison, as directly comparing detections is inaccurate as some images may be more challenging than others and the number of target objects may be unknown (as is the case in our new datasets).

In the proposed method a forest is currently trained for one object while there are several works (e.g. [5], [9]) that utilize a single forest for multiple objects. Our method could be extended to multiple objects per forest by adding three different objective functions that would be used for measuring information gain. During training, similar to classic Hough Forests [12], we could randomly select to perform either entropy minimization of the class distribution or entropy minimization of the of the parameters. Practically this means that we could perform multi-class classification and 6D object pose estimation, addressing, thus, one shortcoming of our method that assumes at least one object instance to always be present in the scene.

In Sec. 4.1 we present extensive experiments on the various parameters of the framework such as as patch size, number of trees and number of iterations to infer class distributions at test time. In Sec. 4.2.1 we perform self comparison tests highlighting the benefits of adding scale-invariance to the template similarity measure (Eq. 3) and using co-training to update the latent class distributions (Algorithm 1). Following this, in Sec. 4.2.2 we compare the performance of our method on the famous dataset of Hinterstoisser *et al.* [16] against the state of the art works of LINEMOD [14], Drost *et al.* [10], Rios Cabrera *et al.* [31] and Brachmann *et al.* [5]. Moreover, in Sec. 4.2.3 and 4.2.4 we present our new datasets and additional comparative analysis.

### 4.1 Framework Parameters

Parameter optimization was performed on a validation dataset that was created by randomly selecting a subset of our own datasets. For each object class, we train a *Latent-Class Hough Forest* with varying number of trees and patch sizes. Moreover, during inference and for the co-training stage we experiment with different number of iterations, while the number of hypothesis to be backprojected per iteration is set $N = 10$. We choose 10 as it is greater than the number of instances present in all datasets, however this number is not fixed and can be adapted based on the application. Furthermore, in all experiments for parameter optimization the coefficient $k_m$ is

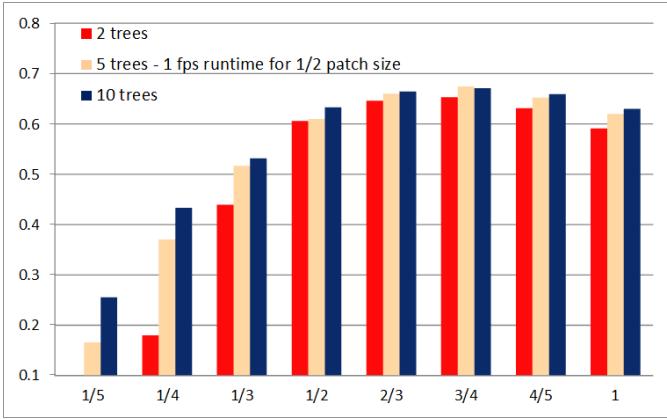set to the value of 0.15, the results with this coefficient are also found to be visually correct.



Fig. 5. F1-Scores for the 6 different patch sizes and 3 different forests in our validation dataset.

### 4.1.1 Patch Size and Number of Trees

Selecting the size of the patch is of paramount importance since large patches tend to match the disadvantages of a holistic template (i.e. sensitive clutter, occlusions etc.) while small ones are prone to noise. The size of the patch depends on the object's identity and it is relative to its size. We extract the bounding box of the target object by rendering its 3D mesh model and modify the patch size to be proportional to the size of the bounding box. Experimenting with 6 different patch sizes, as shown in Fig. 5, revealed that increasing the size of the patch has a direct effect on the F1-Score, while saturation occurs for patches with size equal to 2/3 of the bounding box and higher. Patches resembling holistic ones (i.e. 4/5 and 1) are proven to be very prone to occlusions.

We have additionally experimented with the number of the trees of our *Latent-Class Hough Forests*. As seen in Fig. 5, given a patch size of 2/3 and larger, then selecting more than ten trees is meaningless. It is apparent that, using patches with relatively large size and forests with more than five trees, puts additional computational burden to our system. For instance, we can achieve real-time execution (1 fps) if we use 1/2 as patch size, 2 trees and zero iterations at the co-training stage. However, selecting larger patches and forests results in a respective drop of the execution time.

Regarding the balance of the trees, the learned thresholds $\tau$ expand through the whole range of the values of the split function in a certain node. Usually, the nodes at the top of the tree tend to be more balanced, since there is a variety of different patches, whereas nodes at lower depths are less balanced. However, as the tree grows, each node contains less patches and balancing depends on how well the available patches can be split in two sets, according to our objective function. In fact, producing balanced nodes is not always the best way to split the samples, since it depends on the feature space the patches lay on. The forest optimizes splitting according to the semantic appearance of the patches, and balanced splitting does not always satisfy this objective.

However, we have noticed that for objects Coffee Cup, Juice Carton and Joystick the percentage of training patches that move to the child with the lowest number of samples is around 43%, implying well balanced trees, while for the rest of the objects this percentage is around 35%. Another way to measure the balance of the trees was proposed in [32] where instead of taking one patch, one could also take two or more patches and assign the patch to closest one.
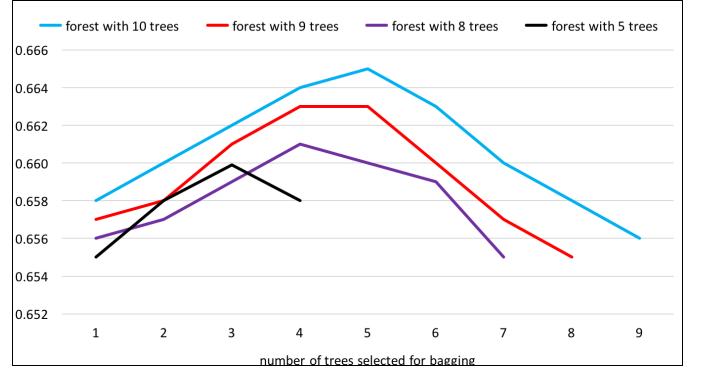


Fig. 6. The impact of bagging in our iterative process. Accuracy (F1-Scores) vs. number of trees selected for bagging under four different forest sizes.
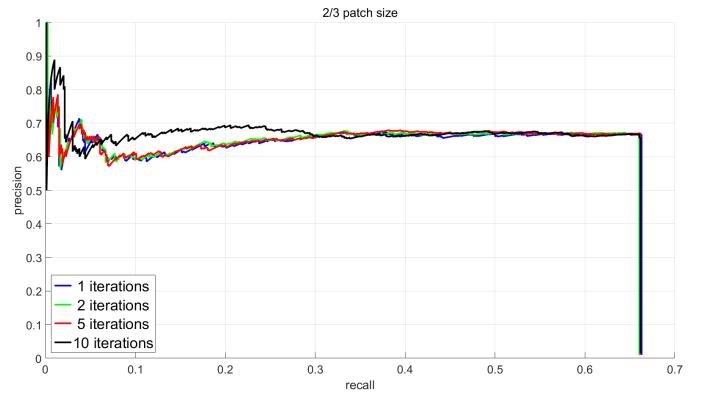


Fig. 7. The impact of the number of iterations on the performance of our system for our validation dataset on a fixed patch size of 2/3 and a forest with 10 trees.

### 4.1.2 Bagging and Number of Iterations

During inference we randomly draw a subset of trees with bagging and iteratively update class distributions in order to datamine the negative information. The impact of bagging is shown in Fig. 6. Given four different sizes of forests (i.e. 10, 9, 8 and 5 trees) we randomly select a subset of trees to obtain an initial object hypotheses set that is then used to construct a consensus pixel test. To investigate the effect of the number of selected trees, we compared the performance levels of our method when drawing a) 1 to 9 trees from a 10 tree forest, b) 1 to 8 trees from a 9 tree forest, c) 1 to 7 trees from an 8 tree forest and d) 1 to 4 from a 5 tree forest. According to the conducted experiments (see Fig. 6), our method preforms
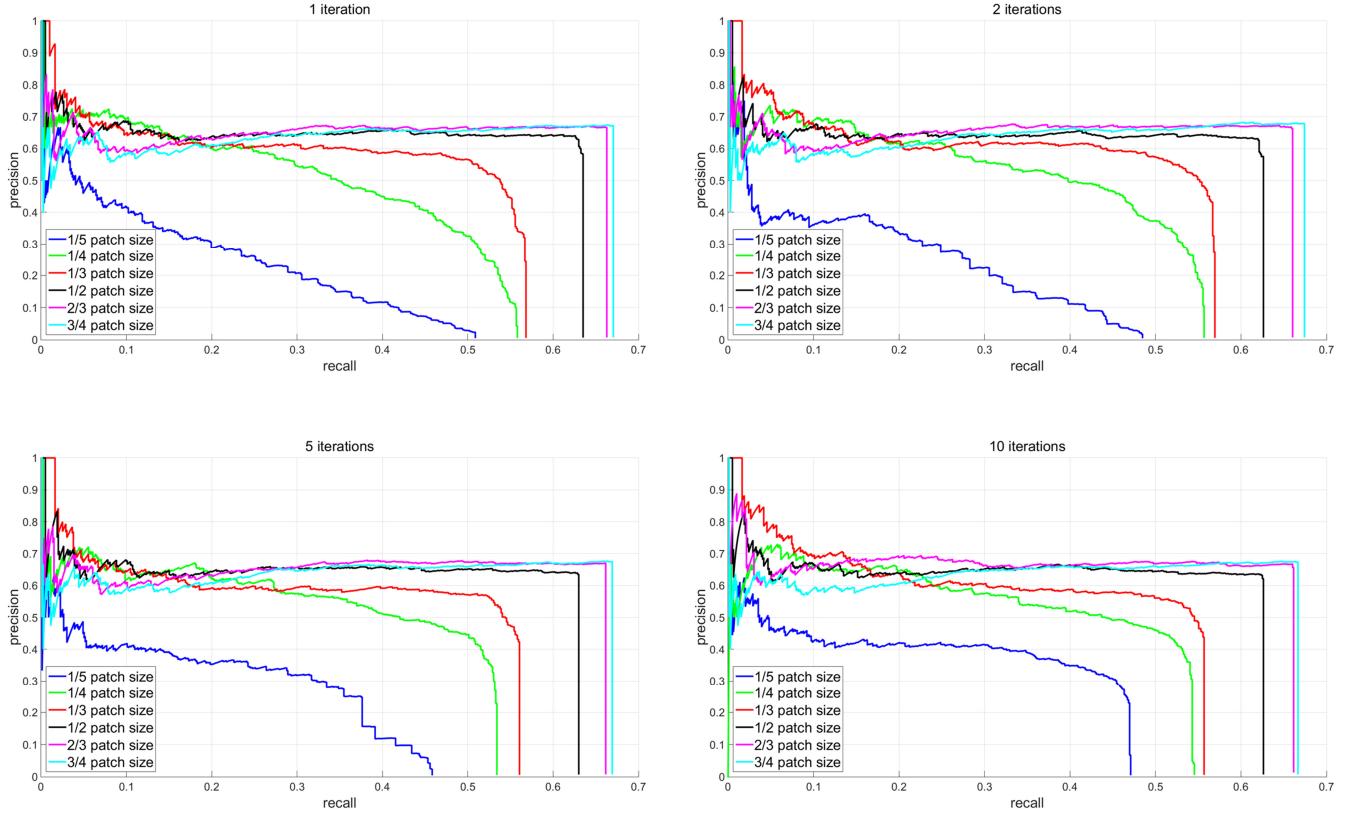
Fig. 8. Precision-Recall curves for different number of iterations and patch sizes.

better in cases we split the forest in the middle, thus, creating two classifiers of equal size.

The quality of the resulted segmentation mask depends on the number of iterations of the co-training process. Since our method makes no use of negative data during training, one straightforward question is whether training with both positive and negative patches would reduce the number of iterations required for extracting the segmentation mask. Towards this end, we used the background dataset of [5] that contains RGB-D images of different cluttered office backgrounds. Patches were sampled from the negative data, pushed through the trees and the classes probability distributions were updated. As expected the overall accuracy was not improved while a slightly faster convergence was noticed only in very few cases. After the first iteration, the extracted consensus pixel set contains both positive and negative data that come directly from the observable unlabelled data in the scene. We believe that this is the reason why training with negative samples offers slightly faster convergence only in very few cases. Fig. 7 illustrates the effect of the iterative part on the performance of our method on our validation dataset. As expected, the iterative method affects only the precision of the system and not the recall, since it does not serve as a false positive rejection framework. Higher number of iterations offer slightly better accuracy with the cost, however, of additional execution time.

In turn, Fig. 8 illustrates the conditional effect of different patch sizes for different number of iterations during inference.

At this point we would like to note that, as seen in Table 3, our system is capable of producing accurate 6 DoF pose estimations even without the iterative part. However, depending on the application domain we might switch to different parameters to meet the respective demands. For example, for industrial object manipulation, real-time execution usually constitutes a prerequisite, while robots operating in domestic environments, emphasize more on the accuracy levels. Towards this end, for the *Domestic Environments Dataset* we set the patch size of 2/3, use a forest of 10 trees and set the number of iterations to 10, while for our *Bin-picking Dataset*, we investigate both the aforementioned parameters and the effect of smaller patches (1/2) for forests with only 2 trees without any iteration at all.

Our approach is not guaranteed to converge and it is possible for our method to produce a solution where all the foreground probabilities converge to zero. However, we have not noticed such a case throughout the exhaustive experimental analysis. The extracted consensus pixel set steams from patches that voted for a hypothesis with a score larger than a threshold, which is different for each object and estimated via trial and error. However, false positives play an important role in our method and can generate false segmentation masks, as shown in the last row of Fig. 10.

## 4.2 Comparative study

In the following section we present a comparative analysis of the performance of our method against other state of the art

TABLE 2
Matching score and speed on the dataset of Hinterstoisser *et al.* [16] for LINEMOD [14], the methods of Drost *et al.* [10], Rios Cabrera *et al.* [31], Brachman *et al.* [5] and our approach. In the second column we report which matching score (Eq. 10 or Eq. 11) was used.

| Approach | | LINEMOD [14] | Drost *et al.* [10] | Rios Cabrera *et al.* [31] | Brachmann *et al.* [5] | Our Approach |
|---|---|---|---|---|---|---|
| Sequence (# images) | Metric | Matching Score / Speed (when available) | | | | |
| Ape(1235) | 10 | 95.8% / 127ms | 86.5% / 22.7s | 95.0% / 55.8ms | **95.8% / -** | 95.7% / 1.82s |
| Bench Vise (1214) | 10 | 98.7% / 115ms | 70.7% / 2.94s | 98.9% / 53.3ms | **100% / -** | 99.7% / 2.11s |
| Driller (1187) | 10 | 93.6% / 121ms | 87.3% / 2.65s | 94.3% / 54.6ms | **99.5% / -** | 99.2% / 1.91s |
| Cam (1200) | 10 | 97.5% / 148ms | 78.6% / 2.81s | 98.2% / 58.4ms | **99.6% / -** | **99.6% / 1.98s** |
| Can (1195) | 10 | 95.4% / 122ms | 80.2% / 1.60s | **96.3% / 55.3ms** | 95.9% / - | 96.1% / 2.08s |
| Iron (1151) | 10 | 97.5% / 116ms | 84.9% / 3.18s | 98.4% / 54.3ms | 97.6% / - | **98.5% / 1.97s** |
| Lamp (1226) | 10 | 97.7% / 125ms | 93.9% / 2.29s | 97.9% / 54.8ms | **99.8% / -** | 99.6% / 2.01s |
| Phone (1224) | 10 | 93.3% / 157ms | 80.7% / 4.70s | 95.3% / 58.4ms | **97.6% / -** | 96.7% / 2.05s |
| Cat (1178) | 10 | 99.3% / 111ms | 85.4% / 7.52s | 99.1% / 53.5ms | **100% / -** | 99.8% / 1.97s |
| Hole Punch (1236) | 10 | 95.9% / 110ms | 77.4% / 8.30s | 97.5% / 54.2ms | **99.4% / -** | 99.5% / 1.92s |
| Duck (1253) | 10 | 95.9% / 104ms | 40.0% / 6.97s | 94.2% / 53.6ms | 95.9% / - | **96.1% / 1.74s** |
| Box (1252) | 11 | 99.8% / 101ms | 97.0% / 2.94s | 99.8% / 56.0ms | 98.0% / - | **98.1% / 2.10s** |
| Glue (1219) | 11 | 91.8% / 135ms | 57.2% / 4.03s | 96.3% / 58.5ms | **98.9% / -** | 98.2% / 1.83s |
| **Average** | | 96.3% / 122ms | 78.4% / 5.58s | 97.1% / 55.4ms | **98.3% / -** | 98.2% / 1.96s |

TABLE 3
F1-Scores for LINEMOD [14], the method of Drost *et al.* [10] and our approach for each object class for the dataset of Hinterstoisser *et al.* [16] and our *Domestic Environment Dataset*.

| Approach | LINEMOD [14] | Drost *et al.* [10] | SI LINEMOD | Without Iterations | With Iterations |
|---|---|---|---|---|---|
| Dataset of Hinterstoisser *et al.* [16] | | | | | |
| Sequence (# images) | F1-Score | | | | |
| Ape(1235) | 0.533 | 0.628 | 0.631 | 0.799 | **0.855** |
| Bench Vise (1214) | 0.846 | 0.237 | 0.869 | 0.941 | **0.961** |
| Driller (1187) | 0.691 | 0.597 | 0.744 | 0.899 | **0.905** |
| Cam (1200) | 0.640 | 0.513 | 0.711 | 0.636 | **0.718** |
| Can (1195) | 0.512 | 0.510 | 0.550 | 0.708 | **0.709** |
| Iron (1151) | 0.683 | 0.405 | 0.749 | 0.705 | **0.735** |
| Lamp (1226) | 0.675 | 0.776 | 0.790 | 0.911 | **0.921** |
| Phone (1224) | 0.563 | 0.471 | 0.655 | 0.660 | **0.728** |
| Cat (1178) | 0.656 | 0.566 | 0.773 | 0.884 | **0.888** |
| Hole Punch (1236) | 0.516 | 0.500 | 0.601 | 0.819 | **0.875** |
| Duck (1253) | 0.580 | 0.313 | 0.659 | 0.888 | **0.907** |
| Box (1252) | 0.860 | 0.826 | **0.933** | 0.736 | 0.740 |
| Glue (1219) | 0.438 | 0.382 | 0.462 | 0.643 | **0.678** |
| **Average** | 0.630 | 0.517 | 0.702 | 0.788 | **0.817** |
| Domestic Environment Dataset | | | | | |
| Sequence (# images) | F1-Score | | | | |
| Coffee Cup (708) | 0.819 | 0.867 | 0.831 | 0.821 | **0.877** |
| Shampoo (1058) | 0.625 | 0.651 | 0.649 | 0.712 | **0.759** |
| Joystick (1032) | 0.454 | 0.277 | 0.491 | 0.511 | **0.534** |
| Camera (708) | 0.422 | 0.407 | **0.498** | 0.291 | 0.372 |
| Juice Carton (859) | 0.494 | 0.604 | 0.506 | 0.812 | **0.870** |
| Milk (860) | 0.176 | 0.259 | 0.228 | 0.315 | **0.385** |
| **Average** | 0.498 | 0.511 | 0.533 | 0.582 | **0.633** |

TABLE 1
Impact of different modalities on the performance of our method based on a smaller validation subset.

| Object | RGB | Depth | RGB-D |
|---|---|---|---|
| Coffee Cup | 0.748 | **0.838** | 0.828 |
| Camera | 0.309 | 0.324 | **0.369** |

techniques. We perform a self comparisons to investigate the role of our iterative method, while we compare against several recently published methods in three different datasets.

### 4.2.1 Self Comparisons

We perform two self comparisons on the dataset of Hinterstoisser *et al.* [16] and our *Domestic Environment Dataset*. Firstly we compare the results of our method with and without the iterative process. As can be seen in Table 3 for the dataset of Hinterstoisser *et al.* [16], our approach with the iterative process improves the F1-Score by 2.8% on average and up to 6.4% on some objects. The biggest gains are seen in objects that have large amounts of indistinct regions (image patches) for which background clutter can easily be confused. For example, the biggest improvements are seen in the Glue, Holepuncher and Phone objects that contain large
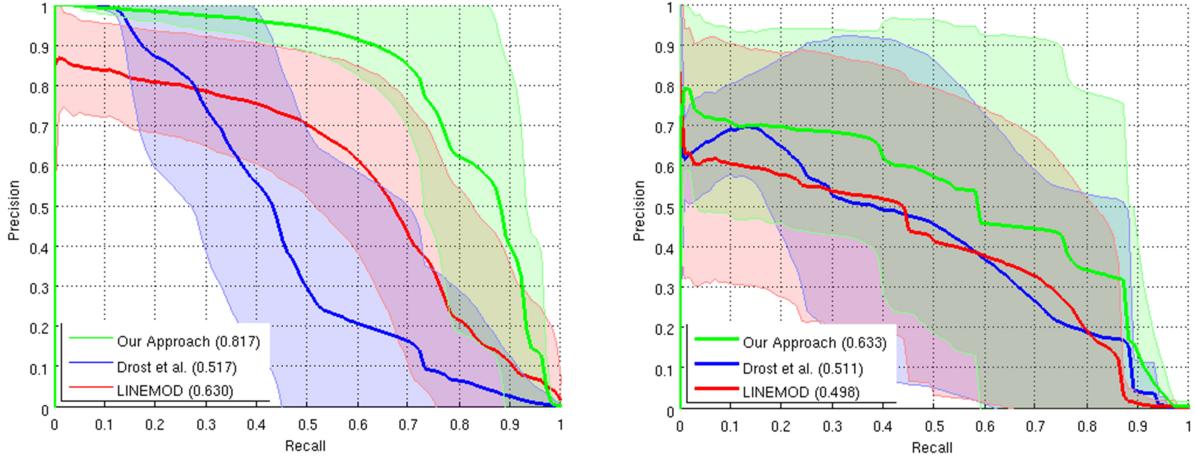
Fig. 9. Average Precision-Recall curve over all objects in the dataset of LINEMOD [16] (left) and our *Domestic Environment Dataset* (right). The shaded region represents one standard deviation above and below the precision value at a given recall value.

planar regions. Furthermore, in Table 3 we compare the results of using holistic LINEMOD templates, to scale-invariant (SI) LINEMOD templates. As the scale-invariant version is trained using only one scale, the performance is increased 6-fold (623 templates opposed to 3738). Furthermore, the performance is also increased by 6.9% on average, this is due to the fact that templates are able to matched at scales not seen in the template learning stage of the original LINEMOD. The impact of our iterative method is more visible in testing scenes that contain small objects with significant amount of foreground occlusions, e.g. Camera object in *Domestic Environment Dataset* or Coffee Cup object in *Bin-Picking Dataset*.

Additionally, we investigate the impact of different modalities on the performance of our system. We used a much smaller validation subset comprised of images of two of the objects contained in both our *Domestic Environment Dataset* and *Bin-Picking Dataset*. Table 1 summarizes the calculated F1-Scores, while one interesting finding is that geometrical cues seem to be more important from appearance ones when trying to detect the object Coffee Cup. In turn, although the detection levels are relatively low, the Camera object favors the usage of both RGB and D modalities.

### 4.2.2  $1^{st}$ *Dataset: LINEMOD*

The performance of our method was evaluated also on the famous dataset of of Hinterstoisser *et al.* [16] against the state of the art works of LINEMOD [14], Drost *et al.* [10], Rios Cabrera *et al.* [31] and Brachmann *et al.* [5]. This dataset contains only one object instance per testing image in a cluttered background, however, without occlusions. Table 2 summarizes the results of the conducted comparison with respect to the performance levels and the execution times. In the second column of Table 2 we report which matching score (Eq. 10 or Eq. 11) was used per object, while for our method we used the same $k_m$ threshold values with the ones

presented in [16]. The performance of our method is almost identical to the one of Brachmann *et al.* [5] on a dataset that could be apprehended as obsolete, since the challenges it offers, i.e. background clutter with no foreground occlusions, are experimentally proven to be easily addressed. As far as execution times are concerned, the scalable method of Rios Cabrera *et al.* [31] provided the lowest runtimes followed by LINEMOD [14] and our method, respectively. Compared to the aforementioned works, our method provides higher performance levels with the additional computational cost. Unfortunately, we could not identify reported runtimes for the method of Brachmann *et al.* [5]. We would like to note that since the evaluation metric for symmetric objects is not very representative, one could transform every hypothesis using the rotation matrices $R_x$ that convert the object to its symmetric shapes and evaluate the detection performance using the non-symmetric measure of Eq. 10.

*Implementation details*: We have implemented our own version of LINEMOD that produced the same exactly results presented in [16]. For the method of Drost *et al.* [10], we use a binary version kindly provided by the author and set the parameters as described in [10]. Source code of the method of Brachmann *et al.* [5] is publicly available. As far as the method of Rios Cabrera *et al.* [31] is concerned, we show in Table 2 the respective performance levels as they were presented in the original paper. Regarding runtimes, reported numbers are taken from [31].

### 4.2.3  $2^{nd}$ *Dataset: Domestic Environments*

Our *Domestic Environment Dataset* consists of 6 objects placed on top of a table simulating, thus, a kitchen table in a house. We provide a dense 3D reconstruction of each object obtained via a commercially available 3D scanning tool [35]. Additionally, we provide for each object, similarly to [16], an individual testing sequence containing over 700 images

TABLE 4
Comparison of two versions of our *Latent-Class Hough Forests* and the method of Brachmann *et al.* [5].

**Domestic Environment Dataset**

| Object | [5] | Latent-Class Hough Forests 5 trees 1/2 patch | Latent-Class Hough Forests 10 trees 2/3 patch |
|---|---|---|---|
| Coffee Cup | 91.2% | 92.4% | **94.6%** |
| Shampoo | 82.4% | 86.9% | **88.2%** |
| Joystick | **75.9%** | 71.3% | 74.1% |
| Camera | 69.1% | 73.8% | **78.4%** |
| Juice Carton | 89.7% | 91.3% | **93.5%** |
| Milk | 47.6% | 50.1% | **51.6%** |
| Average | 75.9% | 77.6% | **80.0%** |

annotated with ground truth position and 3D pose. Testing sequences were obtained by a freely moving handheld RGB-D camera and ground truth was calculated using marker boards and verified manually. The testing images were sampled to produce sequences that are uniformly distributed in the pose space by $[0° − 360°]$, $[−80° − 80°]$ and $[−70° − 70°]$ in the yaw, roll and pitch angles, respectively. Unlike the dataset of [16], our testing sequences contain *multiple object instances* and *foreground occlusions* along with near and far range 2D and 3D clutter, making it more challenging for the task of 3D object detection and pose estimation. Some example frames from this dataset can be seen in Fig 10.

In Fig. 9 we show the average precision-recall curves across all objects and in Table 3 we show the F1-Score per object for each dataset. All methods provided evidence of worse performance levels on the new dataset, which is to be suspected due to the introduction of occlusions as well as multiple object instances. As can be seen we outperform both state of the arts in both datasets. The method of Drost *et al.* [10] has considerably lower precision values due to the fact that it does not take object boundaries into consideration, thus large planar regions of the target object can have a large surface overlap in the background clutter causing many false positives in addition to the true positives. Conversely, our method maintains high levels of precision at high recall which is due to the inference process simplifying the Hough space.

We have also compared our method with the one of Brachmann *et al.* [5] that is designed to estimate only one object per image. For a fair comparison, we compare the scores of the top hypothesis produced by the respective methods per image. Table 4 shows the detection results for the method of Brachmann *et al.* [5] and two different versions of our *Latent-Class Hough Forests*. In Fig. 10 we present some qualitative results on both our *Domestic Environment Dataset* and the one of LINEMOD [16]. A video demonstrating the efficiency of our method in this dataset is also available[1].

### 4.2.4 $3^{rd}$ Dataset: Bin-picking

One of the most widely encountered application in industry is the one of robotic manipulation and, specifically, the one of manipulating similar objects placed in bins (e.g. grasping bolts and screws from a box). It is apparent that, this particular

---

1. https://www.youtube.com/watch?v=idY3Q7wg5rk

working scenario is very challenging for any vision algorithm since several cascading issues arise (e.g. severe occlusions, foreground clutter etc.). Annotating objects stacked in bins can not be done via the methodology we used for our *Domestic Environments Dataset*, where we placed markers beneath each testing object in order to acquire the necessary ground truth measurements. Since there is no simple way of obtaining the latter, Liu *et al.* [24] devised a statistical evaluation of the pose estimation consistency across multiple viewpoints of the camera. They utilized a robot arm to register the different viewpoints, while the final metric is a histogram of the deviations from the median pose estimation.

Unilke [24], we put manual labor to provide, to the best of our knowledge, the first fully annotated *Bin-picking Dataset*. The building of the dataset was divided into two phases, firstly, registering the viewpoints and secondly, annotating the objets. A freely moving handheld RGB-D camera was used to capture the testing frames, which, in turn, are registered by manually selecting key-point correspondences across all sequences. Afterwards, we annotate each object in the testing scenes by manually projecting the 3D mesh of the targets onto the respective point clouds of the testing sequences. The dataset comprises of multiple instances of two objects (Juice Carton and Coffee Cup) in more than 150 testing images.

In Table 5 we show the efficiency of our method in this challenging dataset in the form of accurate detection and F1-Scores. Similar to our previous comparison with the method of Brachmann *et al.* [5], we evaluate the scores of the top hypothesis produced by the respective methods per image. F1-Scores are calculated only for our method which is designed to work with multi-instance objects. In Fig. 11 we present several qualitative results on the *Bin-picking Dataset*. Our system is capable of achieving near real-time execution without any GPU processing. A video demonstrating the efficiency of our method in the *Bin-picking Dataset* is also available[2].

TABLE 5
Percentages of accurate detections and F1-Scores for two versions of our *Latent-Class Hough Forests* and the method of Brachmann *et al.* [5].

**Bin-Picking Dataset**

| Object | [5] | Latent-Class Hough Forests 5 trees 1/2 patch | Latent-Class Hough Forests 10 trees 2/3 patch |
|---|---|---|---|
| Coffee Cup | 89.4% | 90.1% (0.521) | **91.2%** (0.542) |
| Juice Carton | 87.6% | 89.6% (0.484) | **90.4%** (0.492) |
| Average | 88.5% | 89.8% (0.502) | **90.8%** (0.517) |

## 5 CONCLUSION

In this paper we have introduced a novel framework for accurate 3D detection and pose estimation of multiple object instances in cluttered and occluded scenes. We have demonstrated that these challenges can be efficiently met via the adoption of a state of the art template matching feature into a patch-based regression forest. During training we employ a one-class learning scheme, i.e. training with positive samples

---

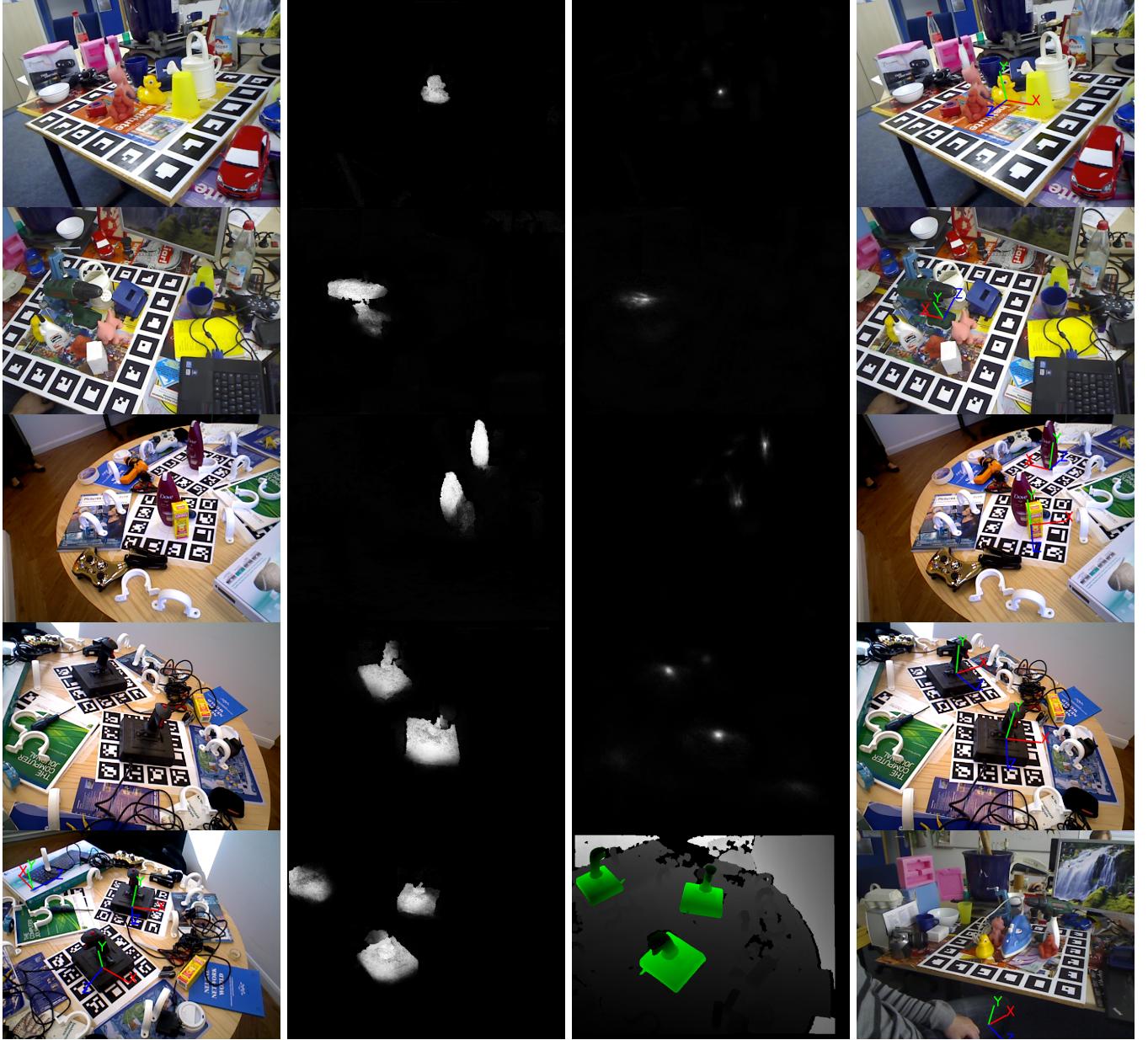2. https://www.youtube.com/watch?v=dh2VtnnsGuY

Fig. 10. Rows 1-4 show, from left to right, the original RGB image, the final segmentation mask, the final Hough vote map and the augmented 3D axis of the estimated result. The final row shows some incorrect results, from left to right: one false positive leading to a false segmentation mask and wrong 3D rendering and finally a false negative.

only rather than involving negative examples. In turn, during inference, we engage the proposed *Latent-Class Hough Forests* that iteratively produce a more accurate estimation of the clutter / occluder distribution by considering class distribution as latent variables. As a result, apart from accurate detection results we can, further, obtain an highly representative occlusion-aware masks facilitating further tasks such as scene layout understanding, occlusion aware ICP or online domain adaption to name a few. Our method is evaluated using both the public dataset of Hinterstoisser *et al.* [16] and our new challenging ones containing foreground occlusion (severe in cases of the *Bin-Picking Dataset*) and multiple object instances. Experimental evaluation provides evidence of our

novel *Latent-Class Hough Forest* outperforming all baselines highlighting the potential benefits of part-based strategies to address the issues of such a challenging problem.

## REFERENCES

[1] A. Aldoma, F. Tombari, L. Di Stefano, and M. Vincze, "A global hypotheses verification method for 3d object recognition," in *ECCV*, 2012.

[2] P. Bariya and K. Nishino, "Scale-hierarchical 3d object recognition in cluttered scenes," in *CVPR*, 2010.

[3] C. M. Bishop, "Novelty detection and neural network validation," in *VISP*, 1994.

[4] U. Bonde, V. Badrinarayanan, and R. Cipolla, "Robust instance recognition in presence of occlusion and clutter," in *ECCV*, 2014.
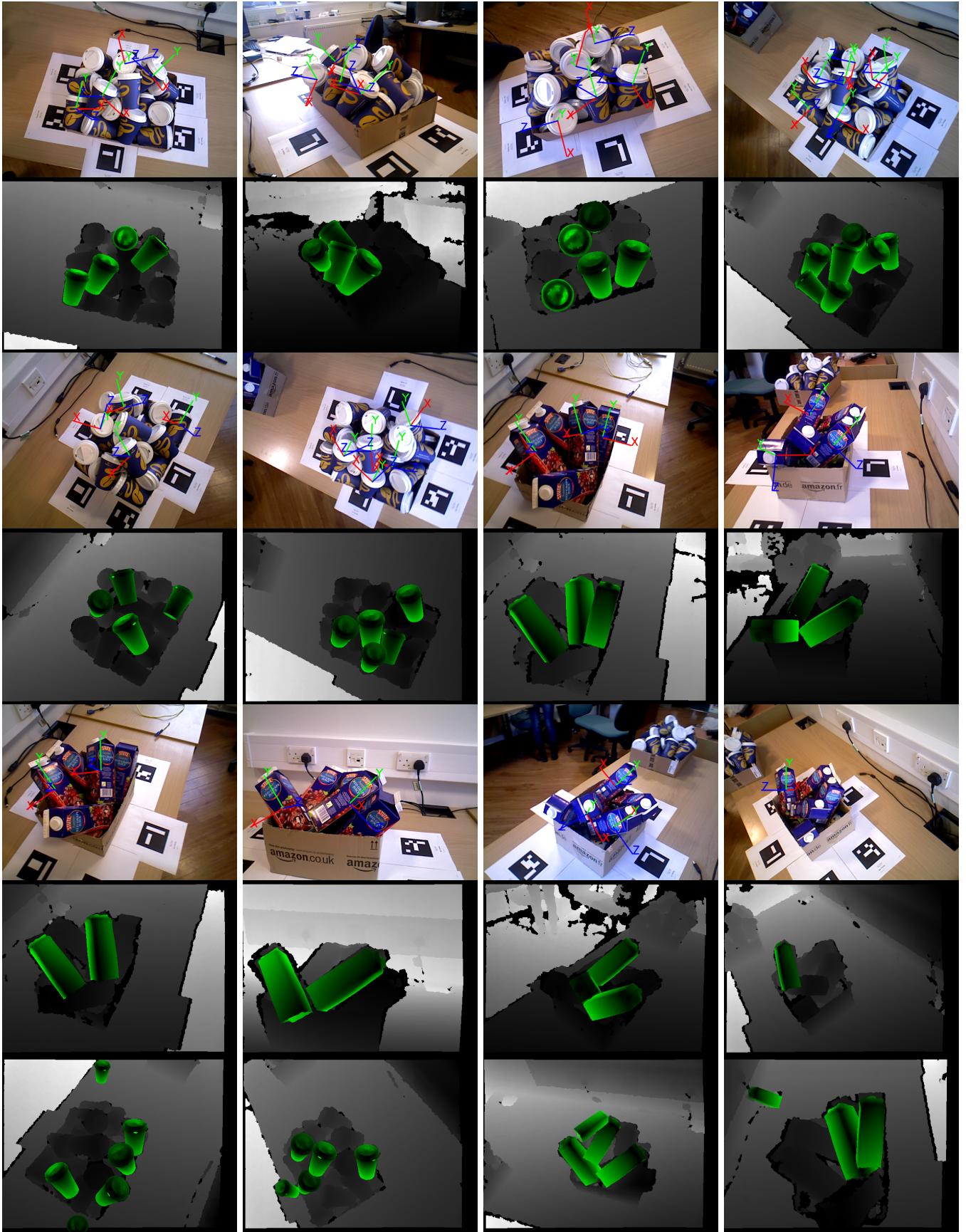
Fig. 11. Rows 1, 3 and 5 show the augmented 3D axis of the estimated result for our *Bin-picking Dataset*. Rows 2, 4 and 6 illustrate the overlaid mesh of the registered objects for the scene right above. The final row shows cases of false positives.

[5] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, "Learning 6d object pose estimation using 3d object coordinates," in *ECCV*, 2014.

[6] L. Breiman, "Random forests," *Machine learning*, 2001.

[7] A. G. Buch, Y. Yang, N. Kruger, and H. G. Petersen, "In search of inliers: 3d correspondence by local and global voting," in *CVPR*, 2014.

[8] A. Crivellaro, M. Rad, Y. Verdie, K. Moo Yi, P. Fua, and V. Lepetit, "A novel representation of parts for accurate 3d object detection and tracking in monocular images," in *ICCV*, 2015.

[9] A. Doumanoglou, R. Kouskouridas, S. Malassiotis, and T.-K. Kim, "6d object detection and next-best-view prediction in the crowd," *arXiv preprint arXiv:1512.07506*, 2015.

[10] B. Drost, M. Ulrich, N. Navab, and S. Ilic, "Model globally, match locally: Efficient and robust 3d object recognition," in *CVPR*, 2010.

[11] G. Fanelli, J. Gall, and L. Van Gool, "Real time head pose estimation with random regression forests," in *CVPR*, 2011.

[12] J. Gall, A. Yao, N. Razavi, L. Van Gool, and V. Lempitsky, "Hough forests for object detection, tracking, and action recognition," *PAMI*, 2011.

[13] R. Girshick, J. Shotton, P. Kohli, A. Criminisi, and A. Fitzgibbon, "Efficient regression of general-activity human poses from depth images," in *ICCV*, 2011.

[14] S. Hinterstoisser, S. Holzer, C. Cagniart, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of textureless objects in heavily cluttered scenes," in *ICCV*, 2011.

[15] S. Hinterstoisser, V. Lepetit, S. Ilic, P. Fua, and N. Navab, "Dominant orientation templates for real-time detection of texture-less objects," in *CVPR*, 2010.

[16] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes," in *ACCV*, 2012.

[17] E. Hsiao and M. Hebert, "Occlusion reasoning for object detection under arbitrary viewpoint," in *CVPR*, 2012.

[18] I. Kokkinos and P. Maragos, "Synergy between object recognition and image segmentation using the expectation-maximization algorithm," *PAMI*, 2009.

[19] R. Kouskouridas, K. Charalampous, and A. Gasteratos, "Sparse pose manifolds," *Autonomous Robots*, 2014.

[20] R. Kouskouridas, A. Gasteratos, and C. Emmanouilidis, "Efficient representation and feature extraction for neural network-based 3d object pose estimation," *Neurocomputing*, 2013.

[21] A. Krull, F. Michel, E. Brachmann, S. Gumhold, S. Ihrke, and C. Rother, "6-dof model based tracking via object coordinate regression," in *ACCV*, 2014.

[22] B. Leibe, A. Leonardis, and B. Schiele, "Combined object categorization and segmentation with an implicit shape model," in *ECCV*, 2004.

[23] J. J. Lim, A. Khosla, and A. Torralba, "Fpm: Fine pose parts-based model with 3d cad models," in *ECCV*, 2014.

[24] M.-Y. Liu and et al., "Fast object localization and pose estimation in heavy clutter for robotic bin picking," *IJRR*, 2012.

[25] T. Malisiewicz and et al., "Ensemble of exemplar-svms for object detection and beyond," in *ICCV*, 2011.

[26] M. Mathias and et al., "Handling occlusions with franken-classifiers," in *ICCV*, 2013.

[27] M. M. Moya and D. R. Hush, "Network constraints and multi-objective optimization for one-class classification," *Neural Networks*, 1996.

[28] M. Moya, M. Koch, and L. Hostetler, "One-class classifier networks for target recognition applications," Sandia National Labs., Albuquerque, NM (United States), Tech. Rep., 1993.

[29] L. Parra and et al., "Statistical independence and novelty detection with information preserving nonlinear maps," *Neural Computation*, 1996.

[30] B. Pepikj and et al, "Occlusion patterns for object class detection," in *CVPR*, 2013.

[31] R. Rios-Cabrera and T. Tuytelaars, "Discriminatively trained templates for 3d object detection: A real time scalable approach," in *ICCV*, 2013.

[32] M. Ristin, M. Guillaumin, J. Gall, and L. Van Gool, "Incremental learning of ncm forests for large-scale image classification," in *CVPR*, 2014.

[33] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural computation*, vol. 13, no. 7, 2001.

[34] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore, "Real-time human pose recognition in parts from single depth images," *ACM*, 2013.

[35] Skanect, "skanect.manctl.com/," 2014.

[36] S. Song and J. Xiao, "Sliding shapes for 3d object detection in depth images," in *ECCV*, 2014.

[37] D. Tang, Y. Liu, and T.-K. Kim, "Fast pedestrian detection by cascaded random forest with dominant orientation templates." in *BMVC*, 2012.

[38] D. Tang, T.-H. Yu, and T.-K. Kim, "Real-time articulated hand pose estimation using semi-supervised transductive regression forests," in *ICCV*, 2013.

[39] D. M. Tax, *One-class classification*. TU Delft, Delft University of Technology, 2001.

[40] A. Tejani, D. Tang, R. Kouskouridas, and T.-K. Kim, "Latent-class hough forests for 3d object detection and pose estimation," in *ECCV*, 2014.

[41] T. Wang and et al., "Learning structured hough voting for joint object detection and occlusion reasoning," in *CVPR*, 2013.

[42] P. Wohlhart and V. Lepetit, "Learning descriptors for object recognition and 3d pose estimation," in *CVPR*, 2015.

[43] Z. Wu, S. Song, A. Khosla, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *CVPR*, 2015.

**Rigas Kouskouridas** received his Ph.D. degree from the Department of Production and Management Engineering of the Democritus University of Thrace in 2013. He is currently a Postdoctoral research associate in the Imperial Computer Vision and Learning Lab of the Electrical and Electronic Engineering Department of Imperial College London. His areas of interest include computer vision, machine learning and robotics.

**Alykhan Tejani** received his MEng degree in Computing in 2010 and his MPhil whilst at the Computer Vision and Learning Lab in 2014, both at Imperial College London. He currently works at Blippar developing their visual search and augmented reality platforms. His current research interests include large-scale object recognition and 3D object pose estimation.

**Andreas Doumanoglou** received his MEng in Electrical and Electronic Engineering from the Aristotle University of Thessaloniki in 2009. Currently, he is Ph.D candidate in the Imperial Computer Vision and Learning Lab of the Electrical and Electronic Engineering Department of Imperial College London. He works in collaboration with CERTH-ITI and his areas of interest include robot vision and machine learning.

**Danhang Tang** is a PhD candidate in the Imperial Computer Vision and Learning Lab of the Electrical and Electronic Engineering Department of Imperial College London. Prior to his PhD, he received a 1st honor MSc degree from University College London and a BSc degree from Sun Yat-sen University. From 2007 to 2009, he worked as system architect for Evryx Technologies Ltd., in support for SnapNow, one of the first image recognition apps in the world.

**Tae-Kyun (T-K) Kim** is an Assistant Professor and leader of Computer Vision and Learning Lab at Imperial College London, UK, since Nov 2010. He received the B.Sc. and M.Sc. degrees from Korea Advanced Institute of Science and Technology in 1998 and 2000, respectively, and worked at Samsung Advanced Institute of Technology in 2000-2004. He obtained his PhD from Univ. of Cambridge in 2008 and Junior Research Fellowship (governing body) of Sidney Sussex College, Univ. of Cambridge for 2007-2010. His research interests primarily lie in decision forests (tree-structure classifiers) and linear methods for: articulated hand pose estimation, face analysis and recognition by image sets and videos, 6D object pose estimation, active robot vision, activity recognition and object detection/tracking. He has co-authored over 40 academic papers in top-tier conferences and journals in the field, his co-authored algorithm for face image retrieval is an international standard of MPEG-7 ISO/IEC.