

DenseFusion: 6D Object Pose Estimation by Iterative Dense Fusion

Chen Wang² Danfei Xu¹ Yuke Zhu¹ Roberto Martín-Martín¹
 Cewu Lu² Li Fei-Fei¹ Silvio Savarese¹

¹Department of Computer Science, Stanford University

²Department of Computer Science, Shanghai Jiao Tong University

Abstract

A key technical challenge in performing 6D object pose estimation from RGB-D image is to fully leverage the two complementary data sources. Prior works either extract information from the RGB image and depth separately or use costly post-processing steps, limiting their performances in highly cluttered scenes and real-time applications. In this work, we present DenseFusion, a generic framework for estimating 6D pose of a set of known objects from RGB-D images. DenseFusion is a heterogeneous architecture that processes the two data sources individually and uses a novel dense fusion network to extract pixel-wise dense feature embedding, from which the pose is estimated. Furthermore, we integrate an end-to-end iterative pose refinement procedure that further improves the pose estimation while achieving near real-time inference. Our experiments show that our method outperforms state-of-the-art approaches in two datasets, YCB-Video and LineMOD. We also deploy our proposed method to a real robot to grasp and manipulate objects based on the estimated pose. Our code and video are available at <https://sites.google.com/view/densefusion/>.

1. Introduction

6D object pose estimation is the crux to many important real-world applications, such as robotic grasping and manipulation [7, 34, 43], autonomous navigation [6, 11, 41], and augmented reality [18, 19]. Ideally, a solution should deal with objects of varying shape and texture, show robustness towards heavy occlusion, sensor noise, and changing lighting conditions, while achieving the speed requirement of real-time tasks. The advent of cheap RGB-D sensors has enabled methods that infer poses of low-textured objects even in poorly-lighted environments more accurately than RGB-only methods. Nonetheless, it is difficult for existing methods to satisfy the requirements of accurate pose estimation and fast inference simultaneously.

Classical approaches first extract features from RGB-D

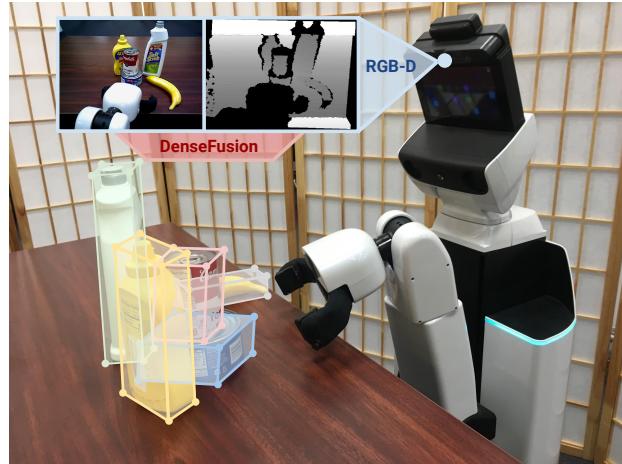


Figure 1. We develop an end-to-end deep network model for 6D pose estimation from RGB-D data, which performs fast and accurate predictions for real-time applications such as robot grasping and manipulation.

data and perform correspondence grouping and hypothesis verification [3, 12, 13, 15, 25, 32, 37]. However, the reliance on handcrafted features and fixed matching procedures have limited their empirical performances in presence of heavy occlusion and lighting variation. Recent success in visual recognition has inspired a family of data-driven methods that use deep networks for pose estimation from RGB-D inputs, such as PoseCNN [40] and MCN [16].

However, these methods require elaborate post-hoc refinement steps to fully utilize the 3D information, such as a highly customized Iterative Closest Point (ICP) [2] procedure in PoseCNN and a multi-view hypothesis verification scheme in MCN. These refinement steps cannot be optimized jointly with the final objective and are prohibitively slow for real-time applications. In the context of autonomous driving, a third family of solutions has been proposed to better exploit the complementary nature of color and depth information from RGB-D data with end-to-end deep models, such as Frustrum PointNet [22] and PointFusion [41]. These models have achieved good performances in driving scenes and the capacity of real-time

inference. However, as we demonstrate empirically, these methods fall short under heavy occlusion, which is common in manipulation domains.

In this work, we propose an end-to-end deep learning approach for estimating 6-DoF poses of known objects from RGB-D inputs. The core of our approach is to embed and fuse RGB values and point clouds at per-pixel level, as opposed to prior work which uses image crops to compute global features [41] or 2D bounding boxes [22]. This per-pixel fusion scheme enables our model to explicitly reason about the local appearance and geometry information, which is essential to handle heavy occlusion. Furthermore, we propose an iterative method which performs pose refinement within the end-to-end learning framework. This greatly enhances model performance while keeping the inference speed real-time.

We evaluate our method in two popular benchmarks for 6D pose estimation, YCB-Video [40] and LineMOD [12]. We show that our method outperforms the state-of-the-art PoseCNN after ICP refinement [40] by 3.5% in pose accuracy while being 200x faster in inference time. In particular, we demonstrate its robustness in highly cluttered scenes thanks to our novel dense fusion method. Last, we also showcase its utility in a real robot task, where the robot estimates the poses of objects and grasp them to clear up a table.

In summary, the contributions of this work are two-fold: First, we present a principled way to combine color and depth information from the RGB-D input. We augment the information of each 3D point with 2D information from an embedding space learned for the task and use this new color-depth space to estimate the 6D pose. Second, we integrate an iterative refinement procedure within the neural network architecture, removing the dependency of previous methods of a post-processing ICP step.

2. Related Work

Pose from RGB images. Classical methods rely on detecting and matching keypoints with known object models [1, 7, 9, 26, 43]. Newer methods address the challenge by learning to predict the 2D keypoints [3, 21, 31, 33, 34] and solve the poses by PnP [10]. Though prevail in speed-demanding tasks, these methods become unreliable given low-texture or low-resolution inputs. Other methods propose to directly estimate objects pose from images using CNN-based architectures [27, 35]. Many such methods focus on orientation estimation: Xiang *et al.* [38, 39] learns a viewpoint-aware pose estimator by clustering 3D features from object models. Mousavian *et al.* [20] predicts 3D object parameters and recovers poses by single-view geometry constraints. Sundermeyer *et al.* [30] implicitly encode orientation in a latent space and in test time find the best match in a codebook as the orientation prediction. However, pose estimation in 3D

remains a challenge for the lack of depth information. Our method leverages both image and 3D data to estimate object poses in 3D in an end-to-end architecture.

Pose from depth / point cloud. Recent studies have proposed to directly tackle the 3D object detection problem in discretized 3D voxel spaces. For example, Song *et al.* [28, 29] generate 3D bounding box proposals and estimate the poses by featuring the voxelized input with 3D ConvNets. Although the voxel representation effectively encodes geometric information, these methods are often prohibitively expensive: [29] takes nearly 20 seconds for each frame.

More recent 3D deep learning architectures have enabled methods that directly performs 6D pose estimation on 3D point cloud data. As an example, both Frustrum PointNets [22] and VoxelNet [42] use a PointNet-like [23] structure and achieved state-of-the-art performances on the KITTI benchmark [11]. Our method also makes use of similar architecture. However, unlike urban driving applications for which point cloud alone provides enough information, generic object pose estimation tasks such as the YCB-Video dataset [40] demands reasoning over both geometric and appearance information. We address such a challenge by proposing a novel 2D-3D sensor fusion architecture.

Pose from RGB-D data. Classical approaches extract 3D features from the input RGB-D data and perform correspondence grouping and hypothesis verification [3, 12, 13, 15, 25, 32, 37]. However, these features are either hard-coded [12, 13, 25] or learned by optimizing surrogate objectives [3, 32, 37] such as reconstruction [15] instead of the true objective of 6D pose estimation. Newer methods such as PoseCNN [40] directly estimates 6D poses from image data. Li *et al.* [16] further fuses the depth input as an additional channel to a CNN-based architecture. However, these approaches rely on expensive post-processing steps to make full use of 3D input. In comparison, our method fuses 3D data to 2D appearance feature while retaining the geometric structure of the input space, and we show that it outperforms [40] on the YCB-Video dataset [40] without the post-processing step.

Our method is most related to PointFusion [41], in which geometric and appearance information are fused in a heterogeneous architecture. We show that our novel local feature fusion scheme significantly outperforms PointFusion's naive fusion-by-concatenation method. In addition, we use a novel iterative refinement method to further improve the pose estimation.

3. Model

Our goal is to estimate the 6D pose of a set of known objects present in an RGB-D image of a cluttered scene. Without loss of generality, we represent 6D poses as homogeneous transformation matrix, $p \in SE(3)$. In other

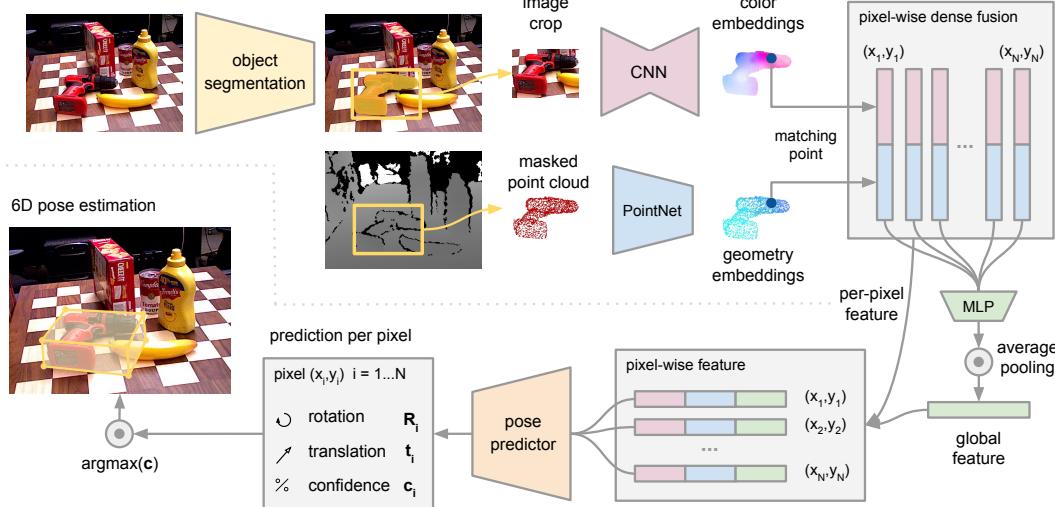


Figure 2. Overview of our 6D pose estimation model. Our model generates object segmentation masks and bounding boxes from RGB images. The RGB colors and point cloud from the depth map are encoded into embeddings and fused at each corresponding pixel. The pose predictor produces a pose estimate for each pixel and the predictions are voted to generate the final 6D pose prediction of the object. (The iterative procedure of our approach is not depicted here for simplicity)

words, a 6D pose is composed by a rotation $R \in SO(3)$ and a translation $t \in \mathbb{R}^3$, $p = [R|t]$. Since we estimate the 6D pose of the objects from camera images, the poses are defined with respect to the camera coordinate frame.

Estimating the pose of a known object in adversarial conditions (e.g., heavy occlusion, poor lighting, ...) is only possible by combining the information contained in the color and depth image channels. However, the two data sources reside in different spaces. Extracting features from heterogeneous data sources and fusing them appropriately is the key technical challenge in this domain.

We address this challenge with (1) a heterogeneous architecture that processes color and depth information differently, retaining the native structure of each data source (Sec. 3.3), and (2) a dense pixel-wise fusion network that performs color-depth fusion by exploiting the intrinsic mapping between the data sources (Sec. 3.4). Finally, the pose estimation is further refined with a differentiable iterative refinement module (Sec. 3.6). In contrast to the expensive post-hoc refinement steps used in [16, 40], our refinement module can be trained jointly with the main architecture and only takes a small fraction of the total inference time.

3.1. Architecture Overview

Fig. 2 illustrates the overall proposed architecture. The architecture contains two main stages. The first stage takes color image as input and performs semantic segmentation for each known object category. Then, for each segmented object, we feed the masked depth pixels (converted to 3D point cloud) as well as an image patch cropped by the bounding box of the mask to the second stage.

The second stage processes the results of the segmentation and estimates the object’s 6D pose. It comprises four components: a) a fully convolutional network that processes the color information and maps each pixel in the image crop to a color feature embedding, b) a PointNet-based [25] network that processes each point in the masked 3D point cloud to a geometric feature embedding, c) a pixel-wise fusion network that combines both embeddings and outputs the estimation of the 6D pose of the object based on an unsupervised confidence scoring, and d) an iterative self-refinement methodology to train the network in a curriculum learning manner and refine the estimation result iteratively. Fig. 2 depicts a), b) and c) and Fig. 3 illustrates d). The details of our architecture are described below.

3.2. Semantic Segmentation

The first step is to segment the objects of interest in the image. Our semantic segmentation network is an encoder-decoder architecture that takes an image as input and generates an $N+1$ -channelled semantic segmentation map. Each channel is a binary mask where active pixels depict objects of each of the N possible known classes. The focus of this work is to develop a pose estimation algorithm. Thus we use an existing segmentation architecture proposed by [40].

3.3. Dense Feature Extraction

The key technical challenge in this domain is the correct extraction of information from the color and depth channels and their synergistic fusion. Even though color and depth present a similar format in the RGB-D frame, their information resides in different spaces. Therefore, we process

them separately to generate color and geometric features from embedding spaces that retain the intrinsic structure of the data sources.

Dense 3D point cloud feature embedding: Previous approaches have used CNN to process the depth image as an additional image channel [16]. However, such method neglects the intrinsic 3D structure of the depth channel. Instead, we first convert the segmented depth pixels into a 3D point cloud using the known camera intrinsics, and then use a PointNet-like architecture to extract geometric features.

PointNet by Qi *et al.* [23] pioneered the use of a symmetric function (max-pooling) to achieve permutation invariance in processing unordered point sets. The original architecture takes as input a raw point cloud and learns to encode the information about the vicinity of each point and of the point cloud as a whole. The features are shown to be effective in shape classification and segmentation [23] and pose estimation [22, 41]. We propose a geometric embedding network that generates a dense per-point feature by mapping each of the P segmented points to a d_{geo} -dimensional feature space. We implement a variant of PointNet architecture that uses average-pooling as opposed to the commonly used max-pooling as the symmetric reduction function.

Dense color image feature embedding: The goal of the color embedding network is to extract per-pixel features such that we can form dense correspondences between 3D point features and image features. The reason for forming these dense correspondences will be clear in the next section. The image embedding network is a CNN-based encoder-decoder architecture that maps an image of size $H \times W \times 3$ into a $H \times W \times d_{rgb}$ embedding space. Each pixel of the embedding is a d_{rgb} -dimensional vector representing the appearance information of the input image at the corresponding location.

3.4. Pixel-wise Dense Fusion

So far we have obtained dense features from both the image and the 3D point cloud inputs; now we need to fuse the information. A naive approach would be to generate a global feature from the dense color and depth features from the segmented area. However, due to heavy occlusion and segmentation errors, the set of features from previous step may contain features of points/pixels on other objects or parts of the background. Therefore, blindly fusing color and geometric features globally would degrade the performance of the estimation. In the following we describe a novel pixel-wise¹ dense fusion network that effectively combines the extracted features, especially for pose estimation under heavy occlusion and imperfect segmentation.

Pixel-wise dense fusion: The key idea of our dense fusion network is to perform local per-pixel fusion instead

¹Since the mapping between pixels and 3D points is unique, we will use interchangeably *pixel-fusion* and *point-fusion*.

of global fusion so that we can make predictions based on *each* fused feature. In this way, we can potentially select the predictions based on the visible part of the object and minimize the effects of occlusion and segmentation noise. Concretely, our dense fusion procedure first associates the geometric feature of each point to its corresponding image feature pixel based on a projection onto the image plane using the known camera intrinsic parameters. The obtained pairs of features are then concatenated and fed to another network to generate a fixed-size global feature vector using a symmetric reduction function. While we refrained from using a single global feature for the estimation, here we enrich each dense pixel-feature with the global densely-fused feature to provide a global context.

We feed each of the resulting per-pixel features into a final network that predicts the object’s 6D pose. In other words, we will train this network to predict one pose from each densely-fused feature. The result is a set of P predicted poses, one per feature. This defines our first learning objective, as we will see in Sec. 3.5. We will now explain our approach to learn to choose the best prediction in a self-supervised manner, inspired by the work by Xu *et al.* [41]. **Per-pixel self-supervised confidence:** We would like to train our pose estimation network to decide which pose estimation is likely to be the best hypothesis based on the specific context. To do so, we modify the network to output a confidence score c_i for each prediction in addition to the pose estimation predictions. We will have to reflect this second learning objective in the overall learning objective, as we will see at the end of the next section.

3.5. 6D Object Pose Estimation

Having defined the overall network structure, we now take a closer look at the learning objective. We define the pose estimation loss as the distance between the points sampled on the objects model in ground truth pose and corresponding points on the same model transformed by the predicted pose. Specifically, the loss to minimize for the prediction per dense-pixel is defined as

$$L_i^p = \frac{1}{M} \sum_j \|(Rx_j + t) - (\hat{R}_i x_j + \hat{t}_i)\| \quad (1)$$

where x_j denotes the j^{th} point of the M randomly selected 3D points from the object’s 3D model, $p = [R|t]$ is the ground truth pose, and $\hat{p}_i = [\hat{R}_i|\hat{t}_i]$ is the predicted pose generated from the fused embedding of the i^{th} dense-pixel.

The above loss function is only well-defined for asymmetric objects, where the object shape and/or texture determines a unique canonical frame. Symmetric objects have more than one and possibly an infinite number of canonical frames, which leads to ambiguous learning objectives. Therefore, for symmetric objects, we instead minimize the

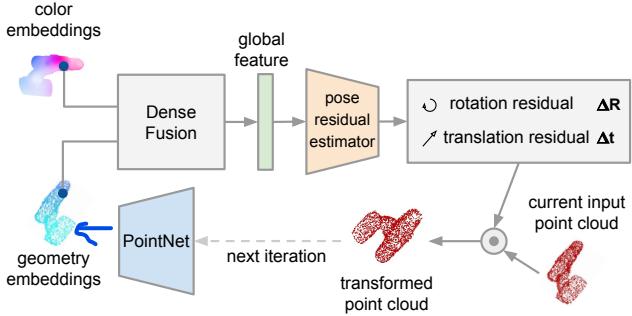


Figure 3. Iterative Pose Refinement. We introduce a network module that refines the pose estimation in an iterative procedure.

distance between each point on the estimated model orientation and the closest point on the ground truth model. The loss function becomes:

$$L_i^p = \frac{1}{M} \sum_j \min_{0 < k < M} \|(Rx_j + t) - (\hat{R}_i x_k + \hat{t}_i)\| \quad (2)$$

Optimizing over all predicted per dense-pixel poses would be to minimize the sum of the per dense-pixels losses: $L = \frac{1}{N} \sum_i L_i^p$. However, as explained before, we would like our network to learn to balance the confidence among the per dense-pixel predictions. To do that we weight the per dense-pixel loss with the dense-pixel confidence, and add a second confidence regularization term:

$$L = \frac{1}{N} \sum_i (L_i^p c_i - w \log(c_i)), \quad (3)$$

where N is the number of randomly sampled dense-pixel features from the P elements of the segment and w is a balancing hyperparameter. Intuitively, low confidence will result in low pose estimation loss but would incur high penalty from the second term, and vice versa. We use the pose estimation that has the highest confidence as the final output.

3.6. Iterative Refinement

The iterative closest point algorithm (ICP) [2] is a powerful refinement approach used by many 6D pose estimation methods [14, 30, 40]. However, the best-performing ICP implementations are often not efficient enough for real-time applications. Here we propose a neural network-based iterative refinement module that can improve the final pose estimation result in a fast and robust manner.

The goal is to enable the network to correct its own pose estimation error in an iterative manner. The challenge here is training the network to refine the previous prediction as opposed to making new predictions. To do so, we must include the prediction made in a previous iteration as part of the input to the next iteration. Our key idea is to consider the previously predicted pose as an estimate of canonical frame of the target object and transform the input point cloud into

this estimated canonical frame. This way, the transformed point cloud implicitly encodes the estimated pose. We then feed the transformed point cloud back into the network and predict a residual pose based on the previously estimated pose. This procedure can be applied iteratively and generate potentially finer pose estimation each iteration.

The procedure is illustrated in Fig. 3. Concretely, we train a dedicated pose residual estimator network to perform the refinement given the initial pose estimation from the main network. At each iteration, we reuse the image feature embedding from the main network and perform dense fusion with the geometric features computed for the new transformed point cloud. The pose residual estimator uses as input a global feature from the set of fused pixel features. After K iterations, we obtain the final pose estimation as the concatenation of the per-iteration estimations:

$$\hat{p} = [R_K | t_K] \cdot [R_{K-1} | t_{K-1}] \cdots [R_0 | t_0] \quad (4)$$

The pose residual estimator can be trained jointly with the main network. However, the pose estimation at the beginning of the training is too noisy for it to learn anything meaningful. Thus in practice, the joint training starts after the main network has converged.

4. Experiments

In the experimental section, we would like to answer the following questions: (1) How does the dense fusion network compare to naive global fusion-by-concatenation? (2) Is the dense fusion and prediction scheme robust to heavy occlusion and segmentation errors? (3) Does the iterative refinement module improve the final pose estimation? (4) Is our method robust and efficient enough for downstream tasks such as robotic grasping?

To answer the first three questions, we evaluate our method on two challenging 6D object pose estimation datasets: YCB-Video Dataset [40] and LineMOD [12]. The YCB-Video Dataset features objects of varying shapes and texture levels under different occlusion conditions. Hence it's an ideal testbed for our occlusion-resilient multi-modal fusion method. The LineMOD dataset is a widely-used dataset that allows us to compare with a broader range of existing methods. We compare our method with state-of-the-art methods [14, 30] as well as model variants. To answer the last question, we deploy our model to a real robot platform and evaluate the performance of a robot grasping task that uses the predictions from our model.

4.1. Datasets

YCB-Video Dataset. The YCB-Video Dataset Xiang *et al.* [40] features 21 YCB objects Calli *et al.* [5] of varying shape and texture. The dataset contains 92 RGB-D videos, where each video shows a subset of the 21 objects in different indoor scenes. The videos are annotated with 6D poses

Table 1. Quantitative evaluation of 6D pose (ADD-S[40]) on YCB-Video Dataset. Objects with bold name are symmetric.

	PointFusion [41]		PoseCNN+ICP [40]		Ours (single)		Ours (per-pixel)		Ours (iterative)	
	AUC	<2cm	AUC	<2cm	AUC	<2cm	AUC	<2cm	AUC	<2cm
002_master_chef_can	90.9	99.8	95.8	100.0	93.9	100.0	95.2	100.0	96.4	100.0
003_cracker_box	80.5	62.6	92.7	91.6	90.8	98.4	92.5	99.3	95.5	99.5
004_sugar_box	90.4	95.4	98.2	100.0	94.4	99.2	95.1	100.0	97.5	100.0
005_tomato_soup_can	91.9	96.9	94.5	96.9	92.9	96.7	93.7	96.9	94.6	96.9
006_mustard_bottle	88.5	84.0	98.6	100.0	91.2	97.8	95.9	100.0	97.2	100.0
007_tuna_fish_can	93.8	99.8	97.1	100.0	94.9	100.0	94.9	100.0	96.6	100.0
008_pudding_box	87.5	96.7	97.9	100.0	88.3	97.2	94.7	100.0	96.5	100.0
009_gelatin_box	95.0	100.0	98.8	100.0	95.4	100.0	95.8	100.0	98.1	100.0
010_potted_meat_can	86.4	88.5	92.7	93.6	87.3	91.4	90.1	93.1	91.3	93.1
011_banana	84.7	70.5	97.1	99.7	84.6	62.0	91.5	93.9	96.6	100.0
019_pitcher_base	85.5	79.8	97.8	100.0	86.9	80.9	94.6	100.0	97.1	100.0
021_bleach_cleanser	81.0	65.0	96.9	99.4	91.6	98.2	94.3	99.8	95.8	100.0
024_bowl	75.7	24.1	81.0	54.9	83.4	55.4	86.6	69.5	88.2	98.8
025_mug	94.2	99.8	95.0	99.8	90.3	94.7	95.5	100.0	97.1	100.0
035_power_drill	71.5	22.8	98.2	99.6	83.1	64.2	92.4	97.1	96.0	98.7
036_wood_block	68.1	18.2	87.6	80.2	81.7	76.0	85.5	93.4	89.7	94.6
037_scissors	76.7	35.9	91.7	95.6	83.6	75.1	96.4	100.0	95.2	100.0
040_large_marker	87.9	80.4	97.2	99.7	91.2	88.6	94.7	99.2	97.5	100.0
051_large_clamp	65.9	50.0	75.2	74.9	70.5	77.1	71.6	78.5	72.9	79.2
052_extra_large_clamp	60.4	20.1	64.4	48.8	66.4	50.2	69.0	69.5	69.8	76.3
061_foam_brick	91.8	100.0	97.2	100.0	92.1	100.0	92.4	100.0	92.5	100.0
MEAN	83.9	74.1	93.0	93.2	88.2	87.9	91.2	95.3	93.1	96.8

and segmentation masks. We follow prior work [40] and split the dataset into 80 videos for training and 2,949 key frames chosen from the rest 12 videos for testing and include the same 80,000 synthetic images released by [40] in our training set. In our experiments, we compare with the result of [40] after depth refinement(ICP) and learning-based depth method [41].

LineMOD Dataset. The LineMOD dataset Hinterstoisser *et al.* [12] consists of 13 low-textured objects in 13 videos. It is widely adopted by both classical methods [4, 8, 36] and recent learning-based approaches [17, 30, 33]. We use the same training and testing set as prior learning-based works [17, 24, 33] without additional synthetic data and compare with the best ICP-refined results of the state-of-the-art algorithms.

4.2. Metrics

We use two metrics to report on the YCB-Video Dataset. The average closest point distance (ADD-S) [40] is an ambiguity-invariant pose error metric which takes care of both symmetric and non-symmetric objects into an overall evaluation. Given the estimated pose $[\hat{R}|\hat{t}]$ and ground truth pose $[R|t]$, ADD-S calculates the mean distance from each 3D model point transformed by $[\hat{R}|\hat{t}]$ to its closest neighbour on the target model transformed by $[R|t]$. We report the area under the ADD-S curve (AUC) following PoseCNN [40]. We follow prior work and set the maximum threshold of AUC to be 0.1m. We also report the percentage of ADD-S smaller than 2cm (<2cm), which measures

the predictions under the minimum tolerance for robot manipulation (2cm for most of the robot grippers).

For the LineMOD dataset, we use the Average Distance of Model Points (ADD) [13] for non-symmetric objects and ADD-S for the two symmetric objects (*eggbox* and *glue*) following prior works [13, 30, 33].

4.3. Implementation Details

The image embedding network consists of a Resnet-18 encoder followed by 4 up-sampling layers as the decoder. The PointNet architecture is an MLP followed by an average-pooling reduction function. Both color and geometric dense feature embedding are of dimension 128. We choose $w = 0.01$ for Eq. 3 by empirical evaluation. The iterative pose refinement module consists of a 4 fully connected layers that directly output the pose residual from the global dense feature. We use the 2 refinement iterations for all experiments.

4.4. Architectures

We compare four model variants that showcase the effectiveness of our design choices.

- *PointFusion* [41] uses a CNN to extract a fixed-size feature vector and fuse by directly concatenating the image feature with the geometry feature. The rest of the network is similar to our architecture. The comparison to this baseline demonstrates the effectiveness of our dense fusion network.

- *Ours (single)* uses our dense fusion network, but instead

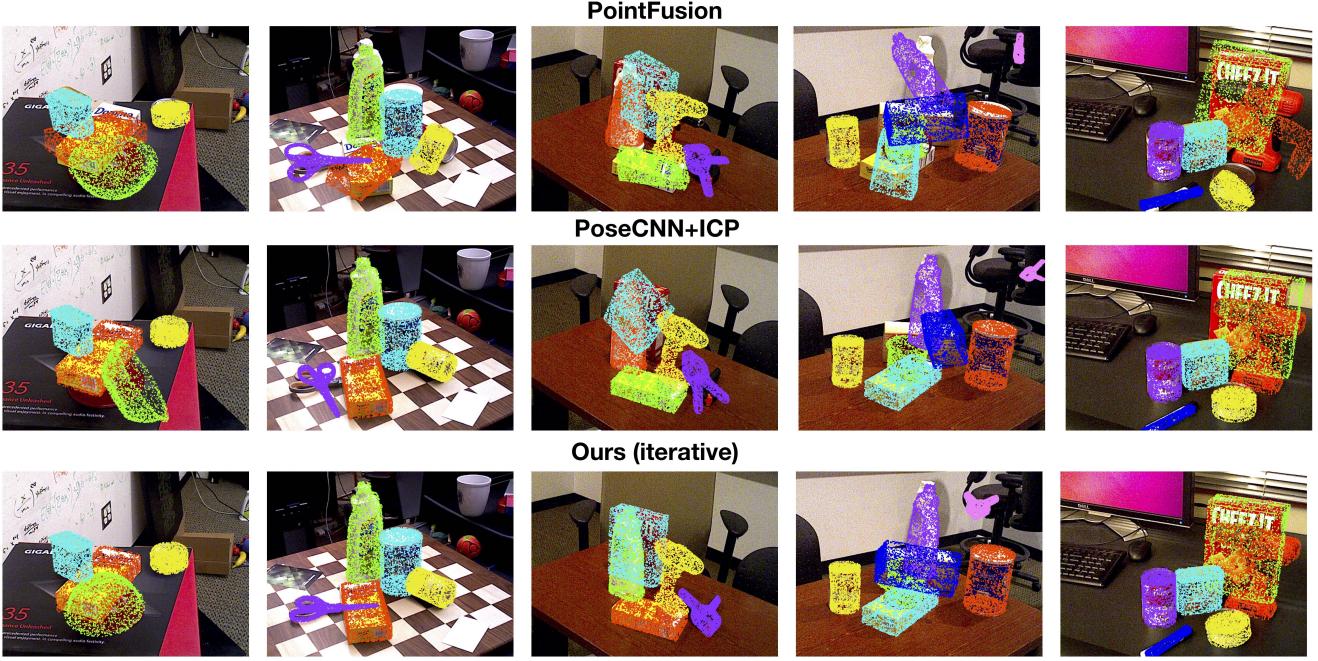


Figure 4. **Qualitative results on the YCB-Video Dataset.** All three methods shown here are tested with the same segmentation masks as in PoseCNN. Each object point cloud in different color are transformed with the predicted pose and then projected to the 2D image frame. The first two rows are former RGB-D methods and the last row is our approach with dense fusion and iterative refinement (2 iterations).

of performing per-point prediction, it only outputs a single prediction using the global feature vector.

- *Ours (per-pixel)* performs per-pixel prediction based on each densely fused feature.
- *Ours (iterative)* is our complete model that uses the iterative refinement (Sec. 3.6) on top of *Ours (per-pixel)*.

4.5. Evaluation on YCB-Video Dataset

Table 1 shows the evaluation results for all the 21 objects in the YCB-Video Dataset. We report the ADD-S AUC ($<0.1\text{m}$) and the ADD-S $<2\text{cm}$ metrics on PoseCNN [40] and our four model variants. To ensure a fair comparison, all methods use the same segmentation masks as in PoseCNN [40]. Among our model variants, *Ours (Iterative)* achieves the best performance. Our method is able to outperform PoseCNN + ICP[40] even without iterative refinement. In particular, *Ours (Iterative)* outperforms PoseCNN + ICP by 3.5% on the ADD-S $<2\text{cm}$ metric.

Effect of dense fusion Both of our dense fusion baselines (*Ours (single)* and *Ours (per-pixel)*) outperform *PointFusion* by a large margin, which shows that dense fusion has a clear advantage over the global fusion-by-concatenation method used in *PointFusion*.

Effect of iterative refinement Table 1 shows that our iterative refinement improves the overall pose estimation performance. In particular, it significantly improves the performances for texture-less symmetric object, e.g., bowl (29%), banana (6%), and extra_large_clamp (6%)

which suffer from orientation ambiguity.

Robustness towards occlusion The main advantage of our dense fusion method is its robustness towards occlusions. To quantify the effect of occlusion on final performance, we calculate the visible surface ratio of each object instance (further detail available in supplementary material). Then we calculate how the accuracy (ADD-S $<2\text{cm}$ percentage) changes with extent of occlusion. As shown in Fig. 5, the performances of *PointFusion* and PoseCNN+ICP degrade significantly as the occlusion increases. In contrast, none of our methods experiences notable performance drop. In particular, the performance of both *Ours (per-pixel)* and *Ours (iterative)* only decrease by 2% overall.

Time efficiency We compare the time efficiency of our model with PoseCNN+ICP in Table 3. We can see that our method is two order of magnitude faster than PoseCNN+ICP. In particular, PoseCNN+ICP spends most of time on the post processing ICP. In contrast, all of our computation component, namely segmentation (Seg), pose estimation (PE), and iterative refinement (Refine), are equally efficient, and the overall runtime is fast enough for real-time application (16 FPS, about 5 objects in each frame).

Qualitative evaluation Fig. 4 visualizes some sample predictions made by PoseCNN+ICP, PointFusion, and our iterative refinement model. As we can see, PoseCNN+ICP and PointFusion fail to estimate the correct pose of the bowl in the leftmost column and the cracker box in the middle col-

Table 2. Quantitative evaluation of 6D pose (ADD[13]) on the LineMOD dataset. Objects with bold name are symmetric.

	RGB		RGB-D				
	BB8 [24] w ref.	PoseCNN +DeepIM [17, 40]	Implicit [30]+ICP	SSD-6D [14]+ICP	PointFusion [41]	Ours (per-pixel)	Ours (iterative)
ape	40.4	77.0	20.6	65	70.4	79.5	92.3
bench vi.	91.8	97.5	64.3	80	80.7	84.2	93.2
camera	55.7	93.5	63.2	78	60.8	76.5	94.4
can	64.1	96.5	76.1	86	61.1	86.6	93.1
cat	62.6	82.1	72.0	70	79.1	88.8	96.5
driller	74.4	95.0	41.6	73	47.3	77.7	87.0
duck	44.3	77.7	32.4	66	63.0	76.3	92.3
eggbox	57.8	97.1	98.6	100	99.9	99.9	99.8
glue	41.2	99.4	96.4	100	99.3	99.4	100.0
hole p.	67.2	52.8	49.9	49	71.8	79.0	92.1
iron	84.7	98.3	63.1	78	83.2	92.1	97.0
lamp	76.5	97.5	91.7	73	62.3	92.3	95.3
phone	54.0	87.7	71.0	79	78.8	88.0	92.8
MEAN	62.7	88.6	64.7	79	73.7	86.2	94.3

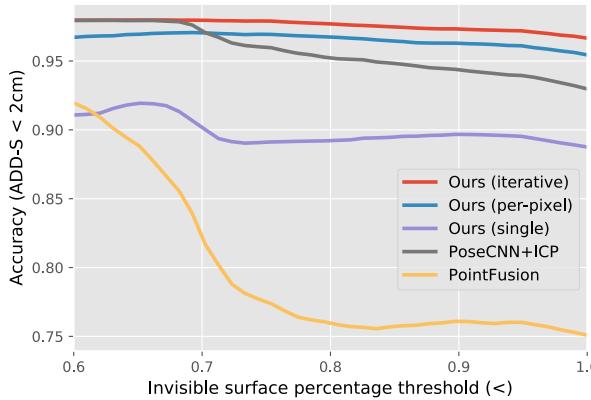


Figure 5. **Model performance under increasing levels of occlusion.** Here the levels of occlusion is estimated by calculating the invisible surface percentage of each object in the image frame. Our methods work more robustly under heavy occlusion compared to baseline methods.

Table 3. Runtime breakdown (second per frame on YCB-Video Dataset). Our method is approximately 200x faster than PoseCNN+ICP. Seg means Segmentation, and PE means Pose Estimation.

PoseCNN+ICP [40]				Ours			
Seg	PE	ICP	ALL	Seg	PE	Refine	ALL
0.03	0.17	10.4	10.6	0.03	0.02	0.01	0.06

umn due to heavy occlusion, whereas our method remains robust. Another challenging case is the clamp in the middle row due to poor segmentation (not shown in the figure). Our approach localizes the clamp from only the visible part of the object and effectively reduces the dependency on accurate segmentation result.

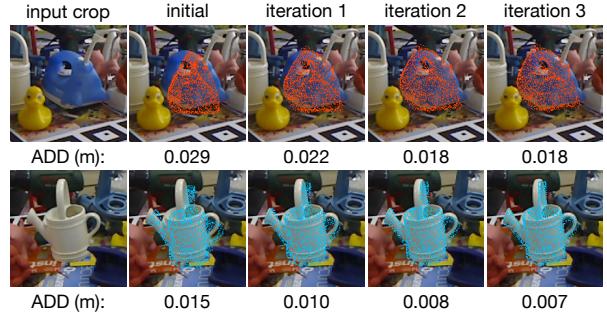


Figure 6. **Iterative refinement performance on LineMOD dataset** We visualize how our iterative refinement procedure corrects initially sub-optimal pose estimation.

4.6. Evaluation on LineMOD Dataset

Table 2 compares our method with previous RGB methods with depth refinement(ICP) (results from [30, 33]) on the ADD metric [13]. Even without the iterative refinement step, our method can outperform 7% over the state-of-the-art depth refinement method. After processing the iterative refinement approach, the final result has another 8% improvement, which proves that our learning-based depth method is superior to the sophisticated application of ICP in both accuracy and efficiency. We visualize the estimated 6D pose after each refinement iteration in Fig.6, where our pose estimation improves by an average of 0.8 cm (ADD) after 2 refinement iterations. The results of some other color-only methods are also listed in Table 2 for reference.

4.7. Robotic Grasping Experiment

In our last experiment, we evaluate whether the poses estimated by our approach are accurate enough to enable robot grasping and manipulation. As shown in Fig. 1, we place 5 YCB objects on a table and command the robot to grasp them using the estimated pose. We follow a similar

procedure to Tremblay *et al.* [34]: we place the five objects in four different random locations on the table, at three random orientations, including configurations with partial occlusions. Since the order of picking the objects is not optimized, we do not allow configurations where objects lay on top of each other. The robot attempts 12 grasps on each object, 60 attempts in total. The robot uses the estimated object orientation to compute an alignment of the gripper’s fingers to the object narrower dimension.

The robot succeeds on 73% of the grasps using our proposed approach to estimate the pose of the objects. The most difficult object to grasp is the banana (7 out of 12 successful attempts). One possible reason is that our banana model is not exactly the same as in the dataset – ours is plain yellow. This characteristic hinders the estimation, especially of the orientation, and leads to some failed grasp attempts along the longer axis of the object. In spite of this less accurate case, our results indicate that our approach is robust enough to be deployed in real-world robotic tasks without explicit domain adaptation, even with a different RGB-D sensor and in a different background than the ones in the training data.

5. Conclusion

We presented a novel approach to estimating 6D poses of known objects from RGB-D images. Our approach fuses a dense representation of features that include color and depth information based on the confidence of their predictions. With this dense fusion approach, our method outperforms previous approaches in several datasets, and is significantly more robust against occlusions. Additionally, we demonstrated that a robot can use our proposed approach to grasp and manipulate objects.

Acknowledgement

This work has been partially supported by JD.com American Technologies Corporation (“JD”) under the SAIL-JD AI Research Initiative and by an ONR MURI award (1186514-1-TBCJE). This article solely reflects the opinions and conclusions of its authors and not JD or any entity associated with JD.com.

References

- [1] M. Aubry, D. Maturana, A. A. Efros, B. C. Russell, and J. Sivic, “Seeing 3d chairs: Exemplar part-based 2d-3d alignment using a large dataset of cad models,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 3762–3769.
- [2] P. J. Besl and N. D. McKay, “A method for registration of 3-d shapes,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 14, pp. 239–256, 1992.
- [3] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, “Learning 6d object pose estimation using 3d object coordinates,” in *European conference on computer vision*, Springer, 2014, pp. 536–551.
- [4] A. G. Buch, L. Kiforenko, and D. Kraft, “Rotational subgroup voting and pose clustering for robust 3d object recognition,” in *Computer Vision (ICCV), 2017 IEEE International Conference on*, IEEE, 2017, pp. 4137–4145.
- [5] B. Calli, A. Singh, A. Walsman, S. S. Srinivasa, P. Abbeel, and A. M. Dollar, “The ycb object and model set: Towards common benchmarks for manipulation research,” in *2015 International Conference on Advanced Robotics (ICAR)*, pp. 510–517, 2015.
- [6] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, “Multi-view 3d object detection network for autonomous driving,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [7] A. Collet, M. Martinez, and S. S. Srinivasa, “The moped framework: Object recognition and pose estimation for manipulation,” *The International Journal of Robotics Research*, vol. 30, no. 10, pp. 1284–1306, 2011.
- [8] B. Drost, M. Ulrich, N. Navab, and S. Ilic, “Model globally, match locally: Efficient and robust 3d object recognition,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, Ieee, 2010, pp. 998–1005.
- [9] V. Ferrari, T. Tuytelaars, and L. Van Gool, “Simultaneous object recognition and segmentation from single or multiple model views,” *International Journal of Computer Vision*, vol. 67, no. 2, pp. 159–188, 2006.
- [10] M. A. Fischler and R. C. Bolles, “Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography,” *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [11] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2012, pp. 3354–3361.
- [12] S. Hinterstoisser, S. Holzer, C. Cagniart, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, “Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes,” *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 858–865, 2011.
- [13] 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 *Asian conference on computer vision*, Springer, 2012, pp. 548–562.
- [14] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, “Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 22–29.

- [15] W. Kehl, F. Milletari, F. Tombari, S. Ilic, and N. Navab, “Deep learning of local rgb-d patches for 3d object detection and 6d pose estimation,” in *European Conference on Computer Vision*, Springer, 2016, pp. 205–220.
- [16] C. Li, J. Bai, and G. D. Hager, “A unified framework for multi-view multi-class object pose estimation,” *ArXiv preprint arXiv:1803.08103*, 2018.
- [17] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, “Deepim: Deep iterative matching for 6d pose estimation,” *ArXiv preprint arXiv:1804.00175*, 2018.
- [18] E. Marchand, H. Uchiyama, and F. Spindler, “Pose estimation for augmented reality: A hands-on survey,” *IEEE transactions on visualization and computer graphics*, vol. 22, no. 12, pp. 2633–2651, 2016.
- [19] E. Marder-Eppstein, “Project tango,” in *ACM SIGGRAPH 2016 Real-Time Live!*, ser. SIGGRAPH ’16, Anaheim, California: ACM, 2016, 40:25–40:25.
- [20] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, “3d bounding box estimation using deep learning and geometry,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [21] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, “6-dof object pose from semantic keypoints,” *ArXiv preprint arXiv:1703.04670*, 2017.
- [22] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, “Frustum pointnets for 3d object detection from rgb-d data,” *ArXiv preprint arXiv:1711.08488*, 2017.
- [23] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” *ArXiv preprint arXiv:1612.00593*, 2016.
- [24] M. Rad and V. Lepetit, “Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth.”
- [25] R. Rios-Cabrera and T. Tuytelaars, “Discriminatively trained templates for 3d object detection: A real time scalable approach,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 2048–2055.
- [26] F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce, “3d object modeling and recognition using local affine-invariant image descriptors and multi-view spatial constraints,” *International Journal of Computer Vision*, vol. 66, no. 3, pp. 231–259, 2006.
- [27] M. Schwarz, H. Schulz, and S. Behnke, “Rgb-d object recognition and pose estimation based on pre-trained convolutional neural network features,” in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*, IEEE, 2015, pp. 1329–1335.
- [28] S. Song and J. Xiao, “Sliding shapes for 3d object detection in depth images,” in *European conference on computer vision*, Springer, 2014, pp. 634–651.
- [29] ———, “Deep sliding shapes for amodal 3d object detection in rgb-d images,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 808–816.
- [30] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and R. Triebel, “Implicit 3d orientation learning for 6d object detection from rgb images,” in *European Conference on Computer Vision*, Springer, 2018, pp. 712–729.
- [31] S. Suwajanakorn, N. Snavely, J. Tompson, and M. Norouzi, “Discovery of latent 3d keypoints via end-to-end geometric reasoning,” *ArXiv preprint arXiv:1807.03146*, 2018.
- [32] A. Tejani, D. Tang, R. Kouskouridas, and T.-K. Kim, “Latent-class hough forests for 3d object detection and pose estimation,” in *Proceedings of the European Conference on Computer Vision*, Springer, 2014, pp. 462–477.
- [33] B. Tekin, S. N. Sinha, and P. Fua, “Real-Time Seamless Single Shot 6D Object Pose Prediction,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [34] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, “Deep object pose estimation for semantic robotic grasping of household objects,” *ArXiv preprint arXiv:1809.10790*, 2018.
- [35] S. Tulsiani and J. Malik, “Viewpoints and keypoints,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1510–1519.
- [36] J. Vidal, C.-Y. Lin, and R. Martí, “6d pose estimation using an improved method based on point pair features,” in *2018 4th International Conference on Control, Automation and Robotics (ICCAR)*, IEEE, 2018, pp. 405–409.
- [37] P. Wohlhart and V. Lepetit, “Learning descriptors for object recognition and 3d pose estimation,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3109–3118.
- [38] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, “Data-driven 3d voxel patterns for object category recognition,” in *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1903–1911.
- [39] ———, “Subcategory-aware convolutional neural networks for object proposals and detection,” in *Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on*, IEEE, 2017, pp. 924–933.
- [40] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, “Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes,” *ArXiv preprint arXiv:1711.00199*, 2017.
- [41] D. Xu, D. Anguelov, and A. Jain, “Pointfusion: Deep sensor fusion for 3d bounding box estimation,” *ArXiv preprint arXiv:1711.10871*, 2017.
- [42] Y. Zhou and O. Tuzel, “Voxelenet: End-to-end learning for point cloud based 3d object detection,” *ArXiv preprint arXiv:1711.06396*, 2017.
- [43] M. Zhu, K. G. Derpanis, Y. Yang, S. Brahmbhatt, M. Zhang, C. Phillips, M. Lecce, and K. Daniilidis, “Single image 3d object detection and pose estimation for grasping,” in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, IEEE, 2014, pp. 3936–3943.

6. Supplementary Materials

6.1. Invisible surface percentage calculation

The invisible surface percentage is a measurement that quantifies how occluded an object is given the camera viewpoint. The measurement is used in Sec.4.5 of the main manuscript. Following are the details of how to compute the invisible surface percentage.

First, we transform the ground truth model of an object to its target pose. Then, the 3D points on the surface of the model are sampled and projected back to a 2D image plane as depth pixels according to the camera intrinsic parameters. The projected depth pixels should be close to the depth measured by a depth sensor if there is *no occlusion*. In other words, if the distance between the measured depth of a pixel and the model-projected depth is larger than a margin, we consider the pixel as being occluded and thus invisible. Concretely, suppose a projected depth pixel p has depth value $d(p)$, and the measured depth of p is $\hat{d}(p)$. p is considered *invisible* if $|d(p) - \hat{d}(p)| > h$. The margin h is set to be 20mm in the experiment. The invisible surface percentage is thus the percentage of the points that are *invisible* out of all sampled points on the object model surface. Since around half of the points on an object model are always invisible due to self-occlusion, Fig.5 in the main manuscript shows results starting from 60 invisible surface percentage.

6.2. Details of the robotic grasping experiment

The robot used in the experiment is a Toyota HSR (Human Support Robot). The robot is equipped with an Asus Xtion RGB-D sensor, a holonomic mobile base, and a two-finger gripper. We deployed our pose estimation model trained on YCB-Video dataset without finetuning. Note that our camera (Asus Xtion) is different from the one used to capture the YCB-Video dataset (Kinect-v2). Our experiment shows that our model is able to tolerate the difference in camera and perform accurate pose estimation. The evaluation includes five YCB objects: 005_tomato_soup_can, 006_mustard_bottle, 007_tuna_fish_can, 011_banana, and 021_bleach_cleanser.

6.3. Additional iterative refinement examples

See Fig. 7.

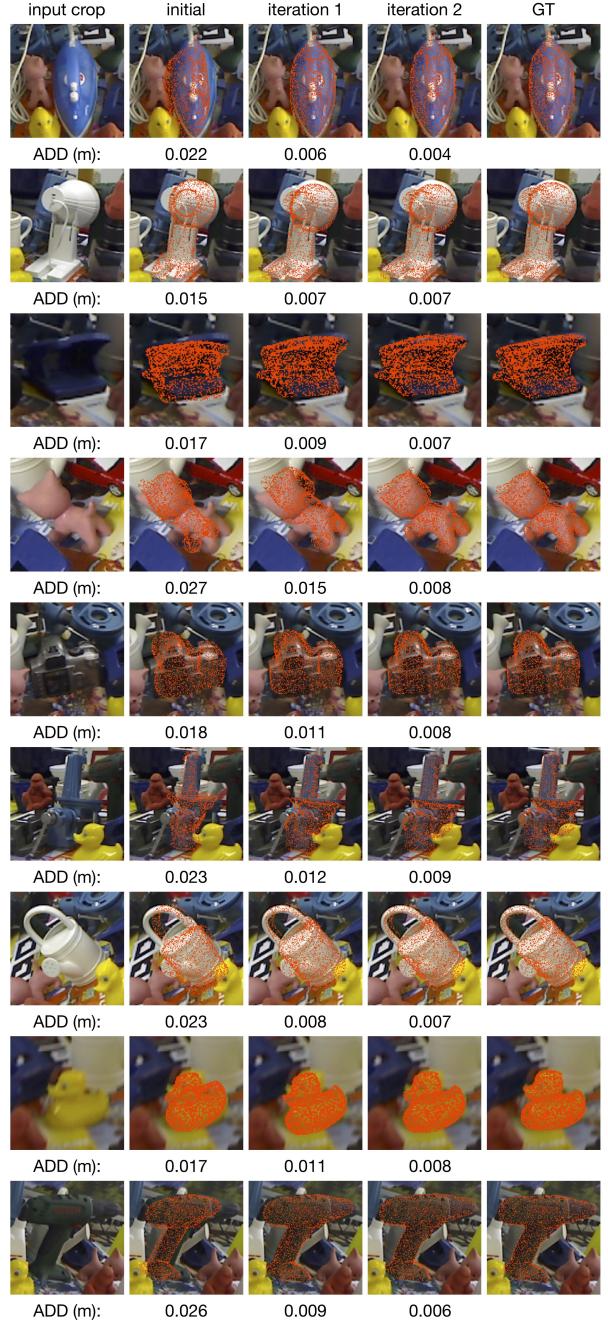


Figure 7. Iterative refinement performance on LineMOD dataset The initial estimation is outputted by Ours (per-pixel). We first transform the object model with the estimated pose and ground truth pose into the 3D space. The ADD distance is the average distance between each corresponding point pair on the two transformed model point clouds. Here we show our iterative refinement performance in more situations includes blurring and low light conditions, where we can see clear improvement on accuracy by using our neural network based iterative refinement method.