# Single-Stage 6D Object Pose Estimation

Yinlin Hu,    Pascal Fua,    Wei Wang,    Mathieu Salzmann

CVLab, EPFL, Switzerland

{firstname.lastname}@epfl.ch

## Abstract

*Most recent 6D pose estimation frameworks first rely on a deep network to establish correspondences between 3D object keypoints and 2D image locations and then use a variant of a RANSAC-based Perspective-n-Point (PnP) algorithm. This two-stage process, however, is suboptimal: First, it is not end-to-end trainable. Second, training the deep network relies on a surrogate loss that does not directly reflect the final 6D pose estimation task.*

*In this work, we introduce a deep architecture that directly regresses 6D poses from correspondences. It takes as input a group of candidate correspondences for each 3D keypoint and accounts for the fact that the order of the correspondences within each group is irrelevant, while the order of the groups, that is, of the 3D keypoints, is fixed. Our architecture is generic and can thus be exploited in conjunction with existing correspondence-extraction networks so as to yield single-stage 6D pose estimation frameworks. Our experiments demonstrate that these single-stage frameworks consistently outperform their two-stage counterparts in terms of both accuracy and speed.*

## 1. Introduction

Detecting 3D objects in images and computing their 6D pose must be addressed in a wide range of applications [11, 31, 48, 29], ranging from robotics to augmented reality. State-of-the-art approaches [39, 41, 32, 16, 13, 36, 54, 34, 24] follow a two-stage paradigm: First use a deep network to establish correspondences between 3D object points and their 2D image projections, then use a RANSAC-based Perspective-n-Point (PnP) algorithm to compute the 6 pose parameters [9, 20, 40, 47, 21, 18, 7, 46].

While effective, this paradigm suffers from several weaknesses. First, the loss function used to train the deep network does not reflect the true goal of pose estimation, but encodes a surrogate task, such as minimizing the 2D errors of the detected image projections. The relationship between such errors and the pose accuracy, however, is not one-to-
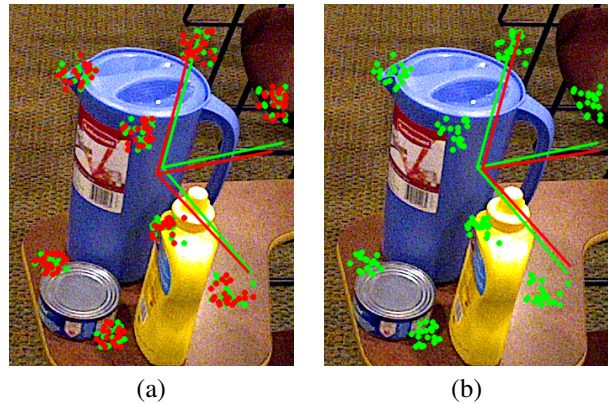
Figure 1: **Motivation.** Consider the modern 6D pose estimation algorithm of [13] that uses a deep network to predict several 2D correspondences for each of the eight 3D corners of the pitcher's bounding box. **(a)** Because it minimizes the average 2D error of these correspondences, two instances of such a framework could produce correspondences that differ but have the same *average* accuracy, such as the green and the red ones. As evidenced by the projected green and red reference frames, applying a RANSAC-based PnP algorithm to these two sets of correspondences can yield substantially different poses. **(b)** Even when using only the set of green correspondences, simply changing their order causes a RANSAC-based PnP algorithm to return different solutions.

one. As shown in Fig. 1 (a) for the state-of-the-art framework of [13], two sets of correspondences with the same *average* 2D error can result in different pose estimates. Second, the two-stage process is not end-to-end trainable. Finally, the iterative RANSAC is time-consuming when there are many correspondences that need to be handled.

In principle, an end-to-end framework could be designed by exploiting a deep version of RANSAC [1, 2], followed by another network performing pose estimation from correspondences [5]. However, the time-consuming character of RANSAC in the presence of many outliers, and the poor repeatability of its solution, arising from the fact that, as shown in Fig. 1 (b), the order of the correspondences affects the resulting pose, do not make it a good candidate for inclusion into an end-to-end trainable network. Further-

more, the approach of [5] relies on using a Direct Linear Transform (DLT) [9] to compute the pose, which is known to be imprecise and would exacerbate in the poor repeatability problem.

As a result, there are still no end-to-end frameworks that can handle jointly keypoint localization and 6D pose estimation. In this paper, we overcome this by introducing a simple but effective network that directly regresses the 6D pose from groups of 3D-to-2D correspondences associated to each 3D object keypoint. Its architecture explicitly encodes that the order of the correspondences in each group is irrelevant, while exploiting the fact that the order of the groups is fixed and corresponds to that of the 3D keypoints.

We then demonstrate the generality of this network by combining it with two state-of-the-art correspondence-extraction frameworks [13, 36]. This yields end-to-end trainable 6D pose estimation frameworks that are both accurate and repeatable. We show that these single-stage frameworks systematically outperform the original two-stage ones [13, 36], in terms of both accuracy and runtime.

## 2. Related Work

Detecting keypoints in the input image followed by running a RANSAC-based PnP algorithm on the established 3D-to-2D correspondences is a classical way for solving the 6D object pose estimation problem. Over the years, many methods have been proposed to improve 3D-to-2D matching [28, 42, 43, 44, 35, 33], relying on diverse techniques, such as template-matching [10, 11], edge-matching [22, 27], and 3D model-based matching [14, 26, 12]. However, these traditional methods still often fail in the presence of severe occlusions and cluttered background.

As in many other areas, the modern take on 6D object pose estimation from an RGB image involves deep neural networks. The simplest approach is to directly regress from the image to the pose parameters [17, 50]. However, this tends to be less accurate than first establishing 3D-to-2D correspondences [39, 41, 32, 16, 13, 36, 54, 34, 24] and then running a RANSAC-based Perspective-n-Point (PnP) algorithm [9] to estimate the object position and orientation given the camera intrinsic parameters. What these methods all have in common is that the correspondences are established independently from each other and consistency is only imposed after the fact by the RANSAC PnP algorithm, which is not part of the deep network. As shown in [53], albeit in a different context, this fails to exploit the fact that all correspondences are constrained by the camera pose and are therefore *not* independent from each other.

Our goal in this paper is to turn the two-stage process described above into a single-stage one by implementing the RANSAC-based PnP part of the process as a deep network that can be combined with the one that establishes the correspondences. This is not a trivial problem because the

standard approach to PnP involves performing a Singular Value Decomposition (SVD), which can be embedded in a deep network but often results in numerical instabilities. In [5], this was addressed by avoiding the explicit use of SVD and instead treating PnP as a least-square fitting problem via the Direct Linear Transform (DLT) approach [9]. This, however, does not guarantee that the result describes a true rotation and further post processing is still needed.

By contrast, the backpropagation-friendly eigendecomposition method of [49] performs explicit SVD, and could in principle used to perform PnP. Doing so, however, would fail to account for the RANSAC part of the algorithm to select the correct correspondences. While RANSAC can be implemented via a deep network [1, 2], its poor repeatability, evidenced in Fig. 1(b), makes it ill-suited to train an end-to-end 6D pose estimation network. In short, no one yet has proposed a satisfying solution to designing a single-stage 6D pose estimation network, which is the problem we address here.

Our architecture is inspired by PointNet [37, 38]. However, PointNet was designed to deliver invariance to rigid transformations, which is the opposite of what we need. Furthermore, we introduce a grouped feature aggregation scheme to effectively hande correspondence clusters in 6D object pose estimation.

## 3. Approach

Given an RGB image captured by a calibrated camera, our goal is to simultaneously detect objects and estimate their 6D pose. We assume them to be rigid and their 3D model to be available. In this section, we first formalize the 6D pose estimation problem assuming that sets of 2D correspondence are given *a priori* for each 3D keypoint on the target object and propose a network architecture that yields 6D poses from such inputs. This network is depicted by Fig. 3. We then discuss how to obtain a single-stage 6D pose estimation framework when these correspondences are the output of another network.

### 3.1. 6D Pose from Correspondence Clusters

Let us assume that we are given the $3 \times 3$ camera intrinsic parameter matrix $\mathbf{K}$ and $m$ potential 2D correspondences $\mathbf{u}_{ik}$ for each one of $n$ 3D object keypoints $\mathbf{p}_i$, with $1 \leq i \leq n$ and $1 \leq k \leq m$. The $\mathbf{p}_i$ is expressed in a coordinate system linked to the object, as shown in Fig. 2(a). For each *valid* 3D to 2D correspondence $\mathbf{p}_i \leftrightarrow \mathbf{u}_{ik}$, we have

$$\lambda_{ik} \begin{bmatrix} \mathbf{u}_{ik} \\ 1 \end{bmatrix} = \mathbf{K}(\mathbf{R}\mathbf{p}_i + \mathbf{t}), \quad (1)$$

where $\lambda_i$ is a scale factor, and $\mathbf{R}$ and $\mathbf{t}$ are the rotation matrix and translation vector that define the camera pose. Because $\mathbf{R}$ is a rotation, it only has three degrees of freedom and $\mathbf{t}$ likewise, for a total of 6.

Note that the 3D-to-2D correspondences above are not restricted to 3D point to 2D point correspondences. In particular, as shown in Fig. 2(b), our formalism can handle 3D point to 2D *vector* correspondences, which have been shown to be better-suited to use in conjunction with a deep network [36]. In that case, the 2D locations can be infered as the crosspoint of two 2D vectors, and Eq. 1 still holds on crosspoints. Our approach as discussed below also still applies, and we therefore do not explicitly distinguish between these two types of 3D-to-2D correspondences unless necessary.

Classical PnP methods [21, 7, 46] try to recover $\mathbf{R}$ and $\mathbf{t}$ given several correspondences, which typically involves using RANSAC to find the valid ones. In the process, an SVD has to be performed on the many randomly chosen subsets of correspondences that must be tried before one containing only valid correspondences is found. In this work, we propose to replace this cumbersome process by a non-linear regression implemented by an appropriately designed deep network $g$ with parameters $\Theta$. In other words, we have

$$(\mathbf{R}, \mathbf{t}) = g(\{(\mathbf{p}_i \leftrightarrow \mathbf{u}_{ik})\}_{1 \leq i \leq n, 1 \leq k \leq m}; \Theta) . \quad (2)$$

We now turn to the actual implementation of $g_\theta$. In the remainder of this section, we first discuss the properties of the set of 3D to 2D correspondences $\mathcal{C}_2^3 = \{(\mathbf{p}_i \leftrightarrow \mathbf{u}_{ik})\}_{1 \leq i \leq n, 1 \leq k \leq m}$ that the network takes as input and then the architecture we designed to account for them.

### 3.1.1 Properties of the Correspondence Set

We will refer to all the 2D points associated to a specific 3D point as a *cluster* because, assuming that the algorithm used to find them is a good one, they tend to cluster around the true location of the 3D point's projection, as can be seen in Fig. 1. Our implementation choice were driven by the following considerations:

**Cluster ordering.** The order of the correspondences within a cluster is irrelevant and should not affect the result. However, the order of the clusters corresponds to the order of the 3D points, which is given and fixed.

**Interaction within a cluster and across clusters.** Although the points in the same cluster correspond to the same 3D point, the 2D location estimate for each point should be expected to be noisy. Thus the model needs to capture the noise distribution within each cluster. More importantly, one single cluster can tell us nothing about the pose, and the final pose can only be inferred by capturing the global structure for multiple clusters.

**Rigid transformations matter.** When processing 3D point clouds with a deep network, one usually wants the result to be invariant to rigid transformations. By contrast, here, we want our 2D points to represent projections of 3D
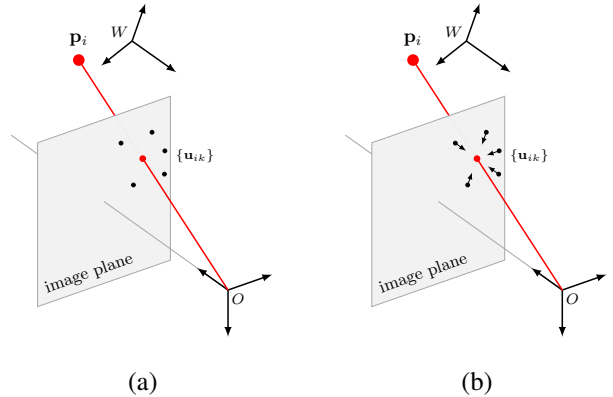


Figure 2: **3D to 2D correspondences.** **(a)** Given $m$ potential 2D correspondences $\mathbf{u}_{ik}$ for each one of $n$ 3D object keypoints $\mathbf{p}_i$, $\{(\mathbf{p}_i \leftrightarrow \mathbf{u}_{ik})\}_{1 \leq i \leq n, 1 \leq k \leq m}$, the pose can be computed based on these 3D-to-2D correspondences. Here, we only show the correspondence cluster for $\mathbf{p}_i$. The camera and object coordinate systems are denoted by $O$ and $W$ respectively. **(b)** The pose can also be obtained from point-to-vector correspondences, in which case a 3D-to-2D correspondence is defined between a 3D point and a 2D vector. Our method can handle both cases.

points, and the features that we extract from them should depend on their absolute positions, which are critical to pose estimation.

### 3.1.2 Network Architecture

We construct a simple network architecture, depicted by Fig. 3, that utilizes the properties discussed above to predict the pose from correspondence clusters. It comprises three main modules: A local feature extraction module with shared network parameters, a feature aggregation module within individual clusters, and a global inference module made of simple fully-connected layers.

**Local feature extraction.** We use an MLP with three layers to extract local features for each correspondence, with weights shared across the correspondences and across the clusters.

**Grouped feature aggregation.** As the order of the clusters is given but the points within each cluster are orderless, to extract the representation for each cluster, we design a grouped feature aggregation method that insensitive to the correspondence order. In theory, we could have used an architecture similar to that of PointNet [37, 38]. However, PointNet is designed to deliver invariance to rigid transformations, which is the opposite of what we need. Instead, given $n$ clusters, each containing $m$ 2D points $\{\mathbf{u}_{ik}\}, 1 \leq i \leq n, 1 \leq k \leq m$, we define a set function $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}^{nD}$ that maps correspondences $\{\mathbf{u}_{ik}\}_{1 \leq k \leq m}$ to the $nD$-dimensional vector

$$CAT\Big(\underset{k}{MAX}(\{\mathbf{f}_{1k}\}), \underset{k}{MAX}(\{\mathbf{f}_{2k}\}), .., \underset{k}{MAX}(\{\mathbf{f}_{nk}\})\Big), \quad (3)$$
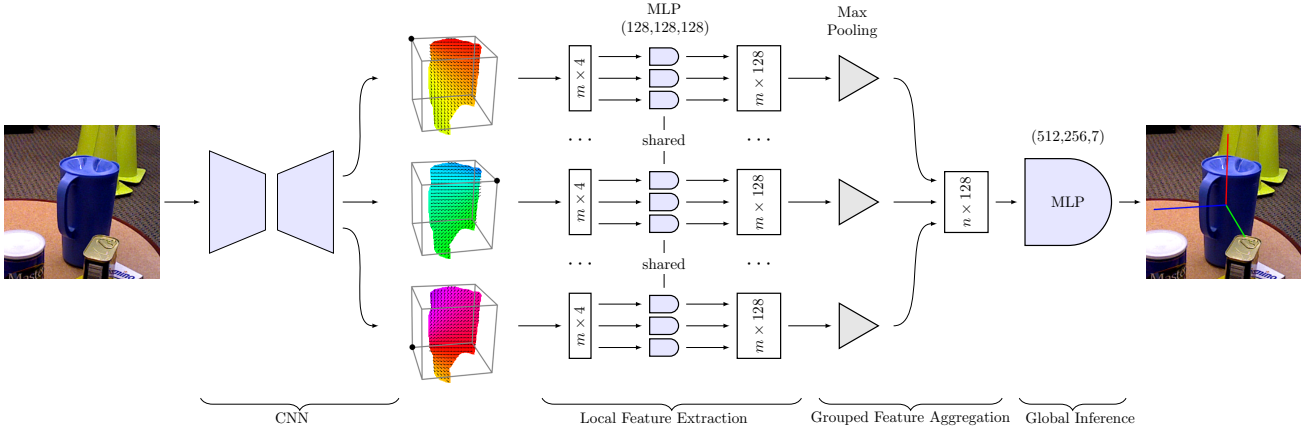
Figure 3: **Overall architecture for single-stage 6D object pose estimation.** After establishing 3D-to-2D correspondences by some segmentation-driven CNN for 6D pose [13, 36], we use three main modules to infer the pose from these correspondence clusters directly: a local feature extraction module with shared network parameters, a feature aggregation module operating within the different clusters, and a global inference module consisting of simple fully-connected layers to estimate the final pose as a quaternion and a translation. The color in the CNN outputs indicates the direction of the 2D offset from the grid cell center to the corresponding projected 3D bounding box corner.

where $\mathbf{f}_{ik}$ is the $D$-dimensional feature representation of $\mathbf{u}_{ik}$ obtained via the above-mentioned fully-connected layers, $MAX()$ is the max-pooling operation and $CAT()$ is the concatenation operation. In our experiments, we found that neither instance normalization [45, 52] nor batch normalization [15] improved the performance here. Therefore, we do not use these operations in our network $g_\theta$.

In principle, one could use a single max-pooling operation, without accounting for the order of the groups, just as PointNet [37] does to achieve permutation invariance for all points. In our case, however, this would mean ignoring the property that the order of the groups is fixed. By contrast, Eq. 3 is invariant to any permutation within a cluster but still accounts for the pre-defined cluster order. We demonstrate the benefits of this approach in the results section.

**Global inference.** We then pass the $nD$-dimensional vector aggregating the group features through another MLP which outputs the 6D pose. To this end, we use three fully-connected layers and encode the final pose as a quaternion and a translation.

### 3.2. Single-Stage 6D Object Pose Estimation

The deep network described above gives us a differentiable way to predict the 6D pose from correspondence clusters for a given object. Given the input image, we therefore still need to detect each object and establish the 3D to 2D correspondences. To do so, we use another deep regressor $f$ with parameters $\Phi$, which, for one object, lets us write

$$[\mathbf{u}_{i1}, \ldots, \mathbf{u}_{im}] = f(\mathbf{p}_i, \mathbf{I}; \Phi), \quad 1 \le i \le n \quad (4)$$

where $\mathbf{I}$ is the input RGB image. To implement $f$, we use the recent encoder-decoder architecture of either [13]

or [36].

In practice, the $\{\mathbf{p}_i\}$ are often taken to be the eight corners of the 3D bounding box of the object's 3D model [39, 32, 13], which leads to different 3D points $\{\mathbf{p}_i\}$, for different object types. In our experiments, we have observed that using the same $\{\mathbf{p}_i\}$ for every object has little impact on the accuracy of $f_\phi$ and makes the subsequent training of $g_\theta$ much easier. We therefore use a single cube for all dataset objects, defined as the largest cube contained by a sphere whose radius is the average of that of the bounding spheres of all object 3D models. This means that the 3D keypoint coordinates are implicitly given by the order of the clusters and do not need to be explicitly specified as network inputs. We therefore a use of 4D representation for each input correspondence, which does *not* include the 3D coordinates. Instead, because the network of [13] operates on an image grid, when we use it to find the correspondences, we take the input to be the $x$ and $y$ coordinates of the center of the grid cell in which the 2D projections are and the $dx$ and $dy$ offsets from that center. In other words, the image coordinates of a 2D correspondence are $x + dx$ and $y + dy$. We tried using these directly as input but we found out experimentally that giving the network what amounts to a first order expansion works better. When using the network of [36] instead of that of [13] to find the correspondences, we use the same input format but normalize the $dx$ and $dy$ so that they represent an orientation.

Our complete model can therefore be written as

$$(\mathbf{R}, \mathbf{t}) = g\big(f(\mathbf{p}_1, \mathbf{I}; \Phi), \cdots, f(\mathbf{p}_8, \mathbf{I}; \Phi); \Theta\big) . \quad (5)$$

To train it, we minimize the loss function

$$\mathcal{L} = \mathcal{L}_s + \mathcal{L}_k + \mathcal{L}_p , \quad (6)$$

Figure 4: **Synthetic data.** We create synthetic data by randomly changing the pose of a unit sphere in 3D space relative to the camera. We capture 20K images for training and 2K for testing.

which combines segmentation term $\mathcal{L}_s$ aiming to assign each grid cell to an object class of to the background, a keypoint regression term $\mathcal{L}_k$, and a pose estimation term $\mathcal{L}_p$. We take $\mathcal{L}_s$ to be the Focal Loss of [25], and $\mathcal{L}_k$ to be the regression term of either [13] or [36] depending on which of the two architectures we use. As in [50, 23], we take $\mathcal{L}_p$ to be the 3D space reconstruction error, that is

$$\mathcal{L}_p = \frac{1}{n} \sum_{i=1}^{n} \|(\hat{\mathbf{R}}\mathbf{p}_i + \hat{\mathbf{t}}) - (\mathbf{R}\mathbf{p}_i + \mathbf{t})\| , \qquad (7)$$

where $\hat{\mathbf{R}}$ and $\hat{\mathbf{t}}$ are the estimated rotation matrix and translation vector, $\mathbf{R}$ and $\mathbf{t}$ are the ground-truth ones. The rotations are estimated from the estimated and ground-truth quaternions, which can be done in a differentiable manner [55]. We also normalize the translations to make sure the regression targets all have a comparable range.

Our architecture simultaneously outputs a segmentation mask and potential 2D locations for a set of predefined 3D keypoints. More specifically, for a dataset with $S$ object classes and an input image $\mathbf{I}$ of size $h \times w \times 3$, it outputs a 3D tensor of size $H \times W \times C$. The dimensions $H$ and $W$ are proportional to the input resolution and $C = (S+1)+2*n$ with $(S + 1)$ channels for segmentation, including one for the background class, and $2 * n$ for the 2D locations (or 2D direction vectors) corresponding to the $n$ 3D points $\mathbf{p}_i$. To obtain correspondence clusters for a given object, we randomly sample $m = 200$ grid cells on the output feature tensor that fall under the segmentation mask of a particular class label.

## 4. Experiments

We compare our single-stage approach to more traditional but state-of-the-art two-stage frameworks [13, 36], first on synthetic data and then on real data from the challenging Occluded-LINEMOD [19] and YCB-Video [50] datasets. Our source code is publicly available at https://github.com/cvlab-epfl/single-stage-pose.

### 4.1. Synthetic Data

As in [21, 7], we create synthetic 3D-to-2D correspondences using a virtual calibrated camera, with image size $640 \times 480$, focal length 800, and principal point at the image center. We take our target object to be a unit 3D sphere,
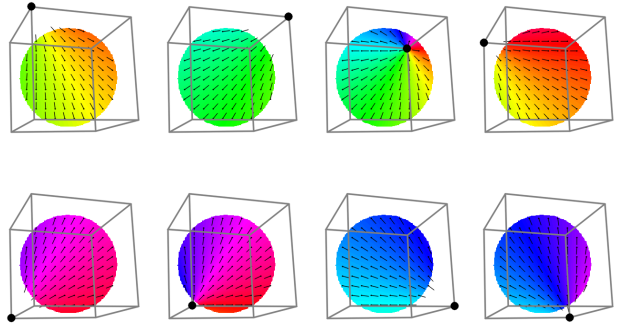


Figure 5: **Generating correspondences.** We project each corner of the sphere's 3D bounding box in the image and, for each grid cell within the object mask, create a correspondence by recording the center $x, y$ of the grid cell and the offset $dx, dy$ to the projected corner.

which we randomly rotate and whose center we randomly translate within the interval $[-2, 2] \times [-2, 2] \times [4, 8]$ expressed in the camera coordinate system, as shown in Fig. 4.

Recall from Section 3.2, that $g_\theta$, the network that regresses poses from the correspondence clusters, expects 4D inputs in the form $[x, y, dx, dy]$, where $x, y$ represent the center of an image grid location and $dx, dy$ a shift from that center. Here, each one should represent a potential image correspondence for a specific corner of the sphere's bounding box for a particular object. Given the segmentation mask of a particular object obtained by projecting the object's 3D model in the image, we create correspondences in the following manner. We project each corner of the sphere's 3D bounding box in the image and, for each grid cell in the segmentation mask, record the cell center $x, y$ and the displacement $dx, dy$ to the projected corner. We then take the resulting correspondences from 200 randomly sampled grid cells within the mask. We add Gaussian noise to their $dx, dy$ values as well as create outliers by setting some percentage of the $dx, dy$ to values uniformly sampled in the image. Fig. 5 demonstrates this procedure.

We trained $g_\theta$ for 300 epoch on 20K synthetic training images with batch size 32, and a learning rate of 1e-3 using the Adam optimizer. During training we randomly add 2D noise with variance $\sigma$ in the range of [0, 15] and create from 0% to 30% of outliers. To test the accuracy obtained with different noise levels and outlier rates, we use 2K synthetic test images and report the mean pose accuracy in terms of the ratio of the 3D space reconstruction error of Eq. 7 to the diameter of the target object.

**Comparing with RANSAC PnP.** Combining a PnP algorithm with RANSAC is the most widespread approach to handling noisy correspondences [39, 41, 13, 54]. Fig. 6 shows that RANSAC-based EPnP [21] and RANSAC-based P3P [8] yield similar performance. While they are more accurate than our learning-based method when there is very
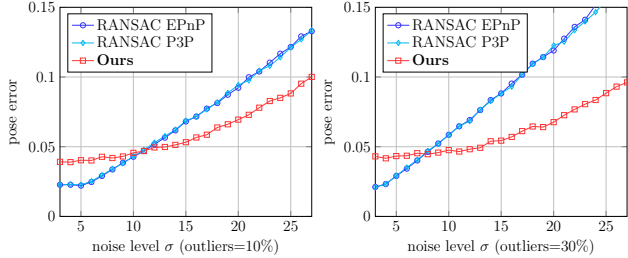
**Figure 6: Comparison with RANSAC PnP.** We compare our network with two classical RANSAC-based PnP methods, EPnP [21] and P3P [8]. The two RANSAC-based methods have very similar performance. More importantly, our method is much more accurate and robust when the noise increases. The pose error is reported as the ratio of the 3D space reconstruction error to the diameter of the target object.
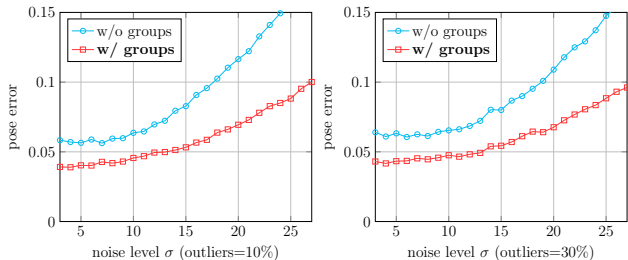


**Figure 7: Importance of correspondence clustering.** We compare our network with one having a single max-pooling operation, thus not accounting for the order of the clusters. Ignoring this property clearly degrades the performance.

little noise, our method quickly becomes much more accurate when the noise level increases.

**Importance of correspondence clustering.** To showcase the importance to structure our network in the way we did, we implemented a simplified version that uses a single max-pooling operation to achieve permutation invariance for all correspondences, without accounting for the order of the clusters that matches that of the keypoints. To make this work, we had to incorporate explicitly the 3D keypoint coordinates associated to each correspondence as input to the network. As shown in Fig. 7, not modeling the fixed order of the keypoints yields a significant decreases in accuracy.

**Comparing with PVNet's voting-based PnP.** In the above experiments, the 2D correspondences were expressed in terms of 2D locations of image points. Since one of the best current techniques [36] uses directions instead and infers poses from those using a voting-based PnP scheme, we feed the same 3D point to 2D vector correspondences to our own network. In this setting, as shown in Fig. 8, the pose is more sensitive to the correspondence noise. However, as in the previous case, while voting-based PnP yields more accurate results when there is little noise, our method is much more robust and accurate when the noise level increases.
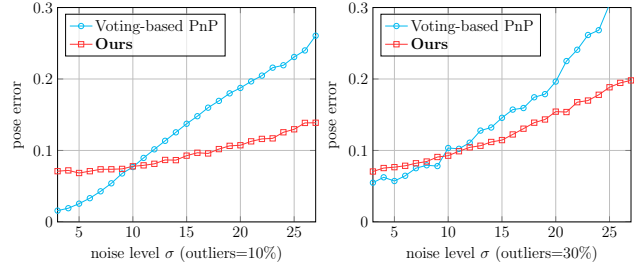


**Figure 8: Comparison with PVNet's voting-based PnP [36].** When using 3D point to 2D vector correspondences, we compare our network with the voting-based PnP used by PVNet. Our method is much more robust to noise than voting-based PnP.

## 4.2. Real Data

We evaluate our method on real data from two challenging datasets, Occluded-LINEMOD [19] and YCB-Video [50].

**Occluded-LINEMOD** consists of 8 objects and is a subset of the older LINEMOD dataset [11]. Unlike LINEMOD in which only one object per image is annotated, Occluded-LINEMOD features multiple annotated objects. This makes it more meaningful for evaluating methods that perform both instance detection and pose estimation. In addition to the cluttered backgrounds, textureless objects, and changing lighting conditions of LINEMOD, Occluded-LINEMOD also has severe occlusions between multiple object instances. As there are only 1214 testing images and no explicit training data in Occluded-LINEMOD, we train our network based on the LINEMOD training data.

**YCB-Video** is more recent and even more challenging. It features 21 objects taken from the YCB dataset [4, 3] and comprises about 130K real images from 92 video sequences. It offers all the challenges of Occluded-LINEMOD plus more diverse object sizes, including several tiny textures-less objects.

**Data preparation.** For Occluded-LINEMOD, as in [41, 13, 36], we first use the Cut-and-Paste synthetic technique [6] to generate 20K images from LINEMOD data and random background data [51], with 4 to 10 different instances for each image. Then, we generate 10K rendering images for each object type from the textured 3D mesh, as in [36]. The pose range during the rendering procedure is the same as in LINEMOD except for one thing: To handle pose ambiguities when encountering symmetry objects [30], we restrict the pose range to a subrange according to the symmetry type of the object during training to avoid confusing the network [39]. In the end, our training data consists of 20K synthetic images with multiple instances and 10K rendered images with only one instance for each object, a total of $(20 + 10 \times 8)$K images.

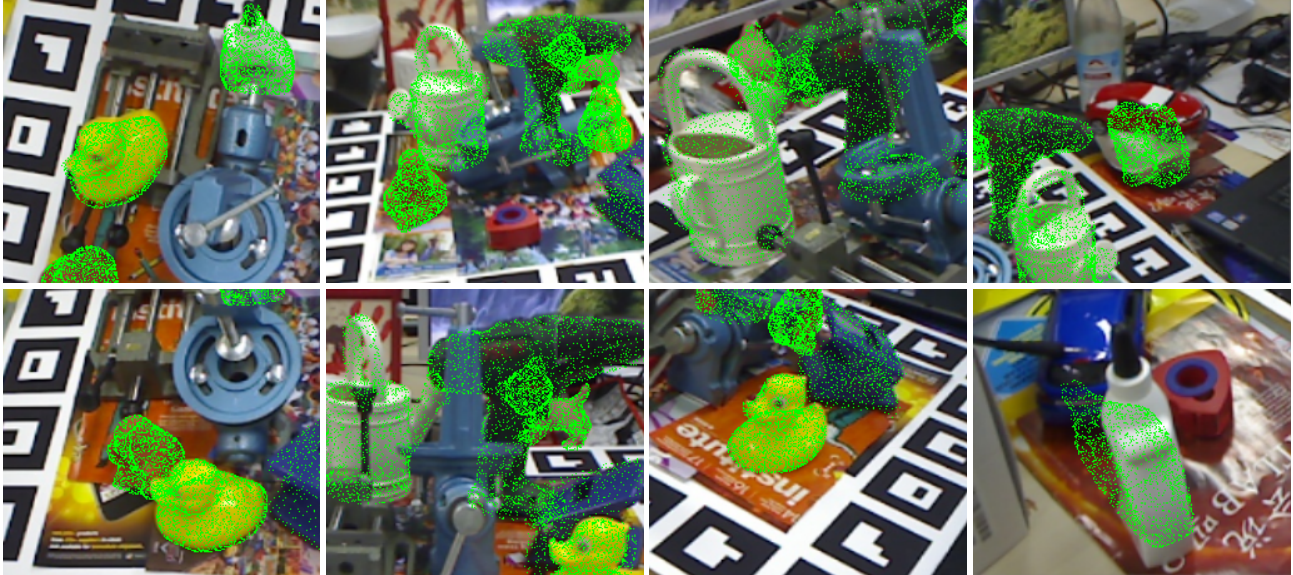For YCB-Video, we follow a similar procedure. We ren-

Figure 9: **Qualitative results on Occluded-LINEMOD.** Our method yields accurate results even in the presence of large occlusions, as shown in the first three columns. The last column shows two failure cases, where the target egg box is occluded too much and the target glue exhibits subtle symmetry ambiguities, making it not easy for the correspondence-extraction network [36] to establish stable correspondences. Here, the pose is visualized as the reprojection of the 3D mesh for each object.

| | [13] | [13] + **Ours** | [36] | [36] + **Ours** |
|---|---|---|---|---|
| Ape | 12.1 | **14.8** | 15.8 | **19.2** |
| Can | 39.9 | **45.5** | 63.3 | **65.1** |
| Cat | 8.2 | **12.1** | 16.7 | **18.9** |
| Driller | 45.2 | **54.6** | 65.7 | **69.0** |
| Duck | 17.2 | **18.3** | 25.2 | **25.3** |
| Eggbox* | 22.1 | **30.2** | 50.2 | **52.0** |
| Glue* | 35.8 | **45.8** | 49.6 | **51.4** |
| Holepun. | 36.0 | **37.4** | 39.7 | **45.6** |
| Average | 27.0 | **32.3** | 40.8 | **43.3** |

Table 1: **Evaluation with different correspondence-extraction networks on Occluded-LINEMOD.** We evaluate two state-of-the-art correspondence-extraction networks: SegDriven [13] and PVNet [36], by replacing their original RANSAC-based post processing with our small network. Our method consistently outperforms the original versions in both cases. Here, we report the ADD-0.1d.

der 10K images for each of the 21 objects using the 3D mesh models that are provided and according to the pose statistic of the dataset. However, we do not use the Cut-and-Paste technique to generate images with multiple instances because in the original YCB-Video images are already annotated with multiple objects and we use that directly.

**Training Procedure.** For both datasets, we use an input image resolution of $640 \times 480$ for both training and testing, as in [36]. We use Adam to optimize with the initial learning rate set to 1e-3 and divided by 10 after processing 50%, 75%, and 90% of the total number of data samples. We set the batch size to be 8 and rely on the usual data augmentation techniques, that is as random luminance, Gaussian noise, translation, scaling, and also occlusions [56]. We train the network on 5M training samples through online data augmentation.

**Metrics.** We quantify the pose error in both 3D and 2D as in [50, 13]. In 3D, it use the average distance between the 3D model points transformed using the predicted pose and those obtained with the ground-truth one, and we refer to it as ADD [50]. In 2D, we use the usual 2D reprojection error of the 3D model points, and we refer it as REP [13]. We measure the pose accuracy in terms the percentage of recovered poses that are correct. In the tables below, we report ADD-0.1d and REP-5px, for which the predicted pose are considered to be correct if ADD is smaller than 10% of the model diameter and REP is below 5 pixel, respectively. For each metric, we use the symmetric version for symmetric objects, which we denote by a * superscript.

### 4.2.1 Occluded-LINEMOD Results

As discussed before, to demonstrate that our method is generic, we test it in conjunction with two correspondence-extraction networks SegDriven [13] and PVNet [36]. Table 1 shows that, by replacing the original RANSAC-based post processing by our network to turn the approach into a single-stage one we improve performance in both cases.

|  | ADD-0.1d | | | | REP-5px | | | |
|---|---|---|---|---|---|---|---|---|
|  | PoseCNN | SegDriven | PVNet | **Ours** | PoseCNN | SegDriven | PVNet | **Ours** |
| Ape | 9.6 | 12.1 | 15.8 | **19.2** | 34.6 | 59.1 | 69.1 | **70.3** |
| Can | 45.2 | 39.9 | 63.3 | **65.1** | 15.1 | 59.8 | **86.1** | 85.2 |
| Cat | 0.9 | 8.2 | 16.7 | **18.9** | 10.4 | 46.9 | 65.1 | **67.2** |
| Driller | 41.4 | 45.2 | 65.7 | **69.0** | 7.4 | 59.0 | **73.1** | 71.8 |
| Duck | 19.6 | 17.2 | 25.2 | **25.3** | 31.8 | 42.6 | 61.4 | **63.6** |
| Eggbox* | 22.0 | 22.1 | 50.2 | **52.0** | 1.9 | 11.9 | 8.4 | **12.7** |
| Glue* | 38.5 | 35.8 | 49.6 | **51.4** | 13.8 | 16.5 | 55.4 | **56.5** |
| Holepun. | 22.1 | 36.0 | 39.7 | **45.6** | 23.1 | 63.6 | 69.8 | **71.0** |
| Average | 24.9 | 27.0 | 40.8 | **43.3** | 17.2 | 44.9 | 61.1 | **62.3** |

Table 2: **Comparison with the state of the art on Occluded-LINEMOD.** We compare our results with those of PoseCNN [50], Seg-Driven [13], and PVNet [36] in terms of both ADD-0.1d and REP-5px. Our method outperforms the state of the art, especially in ADD-0.1d.

|  | correspondence extraction | fusion | total time | FPS |
|---|---|---|---|---|
| PoseCNN | - | - | >250 | <4 |
| SegDriven | 30 | 20 | 50 | 20 |
| PVNet | **14** | 26 | 40 | 25 |
| **Ours** | **14** | **8** | **22** | **45** |

Table 3: **Comparing speed.** We compare the running times (in milliseconds) of PoseCNN [50], SegDriven [13], PVNet [36] and our method on a modern GPU (GTX1080 Ti). Except for PoseCNN, these methods first extract correspondences and then fuse them. With the same correspondence-extraction backbone as in PVNet, our method runs about 2 times faster, thanks to our network that prevents the need for RANSAC-based fusion.

|  | ADD-0.1d | REP-5px |
|---|---|---|
| PoseCNN | 21.3 | 3.7 |
| SegDriven | 39.0 | 30.8 |
| PVNet | - | 47.4 |
| **Ours** | **53.9** | **48.7** |

Table 4: **Comparison with the state of the art on YCB-Video.** We compare our results with those of PoseCNN [50], Seg-Driven [13], and PVNet [36] in terms of ADD-0.1d and REP-5px. We denote by "-" the result missing from the original PVNet paper.

In Table 2, we shown that our single-stage network out-perform the state-of-the-art methods, PoseCNN [50], Seg-Driven [13] and PVNet [36]. Fig. 9 provides qualitative results. In Table 3, we report runtimes for an input image containing about 4-5 objects. Our method is also faster than the others because it does away for the iterative RANSAC procedure.

### 4.2.2 YCB-Video Results

Table 4 summarizes the results comparing against PoseCNN [50], SegDriven [13], and PVNet [36]. It shows that our method consistently also outperforms the others on this dataset. Furthermore, note that it runs nearly 10 times faster than PoseCNN and also nearly 2 times faster than SegDriven and PVNet.

### 4.3. Limitations

While our method is accurate and fast when used in conjunction with state-of-the-art correspondence-extraction networks [13, 36], the network that estimates the poses from the correspondences is still not as accurate as traditional geometry-based PnP algorithms when very precise correspondences can be obtained by other means, as shown in Fig. 6. Furthermore, it does not address the generic PnP problem because we only trained it for fixed sets of 3D coordinates. Addressing this will be the focus of our future work.

### 5. Conclusion

We have introduced a single-stage approach approach to 6D detection and pose estimation. Its key ingredient is a small network that takes candidate 3D-to-2D correspondences and returns a 6D pose. When combined with state-of-the-art approaches to establish the correspondences, it boosts performance by allowing end-to-end training and eliminating the cumbersome RANSAC style procedure that they normally require.

Future work will focus on making the pose estimation network more accurate and more generic so that it can be used in a broader context.

# References

[1] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. DSAC – Differentiable RANSAC for Camera Localization. *Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2

[2] Eric Brachmann and Carsten Rother. Neural-Guided RANSAC: Learning Where to Sample Model Hypotheses. In *International Conference on Computer Vision*, 2019. 1, 2

[3] Berk Calli, Arjun Singh, James Bruce, Aaron Walsman, Kurt Konolige, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M Dollar. Yale-Cmu-Berkeley Dataset for Robotic Manipulation Research. In *International Journal of Robotics Research*, 2017. 6

[4] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron Dollar. The YCB Object and Model Set: Towards Common Benchmarks for Manipulation Research. In *International Conference on Advanced Robotics*, 2015. 6

[5] Zheng Dang, Kwang Moo Yi, Yinlin Hu, Fei Wang, Pascal Fua, and Mathieu Salzmann. Eigendecomposition-Free Training of Deep Networks with Zero Eigenvalue-Based Losses. In *European Conference on Computer Vision*, 2018. 1, 2

[6] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, Paste and Learn: Surprisingly Easy Synthesis for Instance Detection. In *International Conference on Computer Vision*, 2017. 6

[7] Luis Ferraz, Xavier Binefa, and Francesc Moreno-Noguer. Very Fast Solution to the PnP Problem with Algebraic Outlier Rejection. In *Conference on Computer Vision and Pattern Recognition*, pages 501–508, 2014. 1, 3, 5

[8] Xiao-Shan Gao, Xiao-Rong Hou, Jianliang Tang, and Hang-Fei Cheng. Complete Solution Classification for the Perspective-Three-Point Problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(8):930–943, 2003. 5, 6

[9] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2000. 1, 2

[10] Stefan Hinterstoißer, Cedric Cagniart, Slobodan Ilic, Peter F. Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient Response Maps for Real-Time Detection of Textureless Objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34, May 2012. 2

[11] Stefan Hinterstoißer, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary R. Bradski, Kurt Konolige, and Nassir Navab. Model Based Training, Detection and Pose Estimation of Texture-Less 3D Objects in Heavily Cluttered Scenes. In *Asian Conference on Computer Vision*, 2012. 1, 2, 6

[12] Edward Hsiao, Sudipta N. Sinha, Krishnan Ramnath, Simon Baker, C. Lawrence Zitnick, and Richard Szeliski. Car Make and Model Recognition Using 3D Curve Alignment. In *IEEE Winter Conference on Applications of Computer Vision*, 2014. 2

[13] Yinlin Hu, Joachim Hugonot, Pascal Fua, and Mathieu Salzmann. Segmentation-Driven 6D Object Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 4, 5, 6, 7, 8

[14] Daniel P. Huttenlocher, Gregory A. Klanderman, and William Rucklidge. Comparing Images Using the Hausdorff Distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 850–863, 1993. 2

[15] Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning*, 2015. 4

[16] Omid Hosseini Jafari, Siva Karthik Mustikovela, Karl Pertsch, Eric Brachmann, and Carsten Rother. Ipose: Instance-Aware 6D Pose Estimation of Partly Occluded Objects. In *Asian Conference on Computer Vision*, 2018. 1, 2

[17] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. SSD-6D: Making Rgb-Based 3D Detection and 6D Pose Estimation Great Again. In *International Conference on Computer Vision*, 2017. 2

[18] Laurent Kneip, Hongdong Li, and Yongduek Seo. UPnP: An Optimal O(n) Solution to the Absolute Pose Problem with Universal Applicability. In *European Conference on Computer Vision*, 2014. 1

[19] Alexander Krull, Eric Brachmann, Frank Michel, Michael Ying Yang, Stefan Gumhold, and Carsten Rother. Learning Analysis-By-Synthesis for 6D Pose Estimation in RGB-D Images. In *International Conference on Computer Vision*, 2015. 5, 6

[20] Vincent Lepetit and Pascal Fua. *Monocular Model-Based 3D Tracking of Rigid Objects: A Survey*. Now Publishers, September 2005. 1

[21] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. EPnP: An Accurate O(n) Solution to the PnP Problem. *International Journal of Computer Vision*, 2009. 1, 3, 5, 6

[22] Dengwang Li, Hongjun Wang, Yong Yin, and Xiuying Wang. Deformable Registration Using Edge-preserving Scale Space for Adaptive Image-guided Radiation Therapy. In *Journal of Applied Clinical Medical Physics*, 2011. 2

[23] Yi Li, Gu Wang, Xiangyang Ji, Yu Xiang, and Dieter Fox. DeepIM: Deep Iterative Matching for 6D Poseestimation. In *European Conference on Computer Vision*, 2018. 5

[24] Zhigang Li, Gu Wang, and Xiangyang Ji. CDPN: Coordinates-Based Disentangled Pose Network for Real-Time RGB-Based 6-DoF Object Pose Estimation. In *International Conference on Computer Vision*, 2019. 1, 2

[25] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollr. Focal Loss for Dense Object Detection. In *International Conference on Computer Vision*, 2017. 5

[26] Ming-Yu Liu, Oncel Tuzel, Ashok Veeraraghavan, and Rama Chellappa. Fast Directional Chamfer Matching. In *Conference on Computer Vision and Pattern Recognition*, 2010. 2

[27] David G. Lowe. Fitting Parameterized Three-Dimensional Models to Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(5):441–450, June 1991. 2

[28] David G. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, 20(2):91–110, November 2004. 2

[29] Subhransu Maji and Jitendra Malik. Object Detection using a Max-margin Hough Transform. In *Conference on Computer Vision and Pattern Recognition*, 2009. 1

[30] Fabian Manhardt, Diego Martin Arroyo, Christian Rupprecht, Benjamin Busam, Tolga Birdal, Nassir Navab, and Federico Tombari. Explaining the Ambiguity of Object Detection and 6D Pose From Visual Data. In *International Conference on Computer Vision*, 2019. 6

[31] Frank Michel, Alexander Kirillov, Eric Brachmann, Alexander Krull, Stefan Gumhold, Bogdan Savchynskyy, and Carsten Rother. Global Hypothesis Generation for 6D Object Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2017. 1

[32] Markus Oberweger, Mahdi Rad, and Vincent Lepetit. Making Deep Heatmaps Robust to Partial Occlusions for 3D Object Pose Estimation. In *European Conference on Computer Vision*, 2018. 1, 2, 4

[33] Yuki Ono, Eduard Trulls, Pascal Fua, and Kwang Moo Yi. LF-Net: Learning Local Features from Images. In *Advances in Neural Information Processing Systems*, 2018. 2

[34] Kiru Park, Timothy Patten, and Markus Vincze. Pix2Pose: Pixel-Wise Coordinate Regression of Objects for 6D Pose Estimation. In *International Conference on Computer Vision*, 2019. 1, 2

[35] Georgios Pavlakos, Xiaowei Zhou, Aaron Chan, Konstantinos G. Derpanis, and Kostas Daniilidis. 6-DoF Object Pose from Semantic Keypoints. In *International Conference on Robotics and Automation*, 2017. 2

[36] Sida Peng, Yuan Liu, Qixing Huang, Hujun Bao, and Xiaowei Zhou. PVNet: Pixel-Wise Voting Network for 6DoF Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 3, 4, 5, 6, 7, 8

[37] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Conference on Computer Vision and Pattern Recognition*, 2017. 2, 3, 4

[38] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In *Advances in Neural Information Processing Systems*, 2017. 2, 3

[39] Mahdi Rad and Vincent Lepetit. Bb8: A Scalable, Accurate, Robust to Partial Occlusion Method for Predicting the 3D Poses of Challenging Objects Without Using Depth. In *International Conference on Computer Vision*, 2017. 1, 2, 4, 5, 6

[40] Fred Rothganger, Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. 3D Object Modeling and Recognition Using Local Affine-Invariant Image Descriptors and Multi-View Spatial Constraints. *International Journal of Computer Vision*, 66(3), 2006. 1

[41] Bugra Tekin, Sudipta N. Sinha, and Pascal Fua. Real-Time Seamless Single Shot 6D Object Pose Prediction. In *Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2, 5, 6

[42] Engin Tola, Vincent Lepetit, and Pascal Fua. DAISY: An Efficient Dense Descriptor Applied to Wide Baseline Stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(5):815–830, 2010. 2

[43] Tomasz Trzcinski, C. Mario Christoudias, Vincent Lepetit, and Pascal Fua. Learning Image Descriptors with the Boosting-Trick. In *Advances in Neural Information Processing Systems*, December 2012. 2

[44] Shubham Tulsiani and Jitendra Malik. Viewpoints and Keypoints. In *Conference on Computer Vision and Pattern Recognition*, 2015. 2

[45] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance Normalization: The Missing Ingredient for Fast Stylization. In *arXiv Preprint*, 2016. 4

[46] Steffen Urban, Jens Leitloff, and Stefan Hinz. MLPnP-A Real-Time Maximum Likelihood Solution to the Perspective-N-Point Problem. *arXiv Preprint*, 2016. 1, 3

[47] Daniel Wagner, Gerhard Reitmayr, Alessandro Mulloni, Tom Drummond, and Dieter Schmalstieg. Pose Tracking from Natural Features on Mobile Phones. In *International Symposium on Mixed and Augmented Reality*, September 2008. 1

[48] Chen Wang, Danfei Xu, Yuke Zhu, Roberto Martn-Martn, Cewu Lu, Li Fei-Fei, and Silvio Savarese. DenseFusion: 6D Object Pose Estimation by Iterative Dense Fusion. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1

[49] Wei Wang, Zheng Dang, Yinlin Hu, Pascal Fua, and Mathieu Salzmann. Backpropagation-Friendly Eigendecomposition. In *Advances in Neural Information Processing Systems*, 2019. 2

[50] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes. In *Robotics: Science and Systems Conference*, 2018. 2, 5, 6, 7, 8

[51] Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *Conference on Computer Vision and Pattern Recognition*, 2010. 6

[52] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to Find Good Correspondences. In *Conference on Computer Vision and Pattern Recognition*, 2018. 4

[53] Kwang Moo Yi, Yannick Verdie, Pascal Fua, and Vincent Lepetit. Learning to Assign Orientations to Feature Points. In *Conference on Computer Vision and Pattern Recognition*, 2016. 2

[54] Sergey Zakharov, Ivan Shugurov, and Slobodan Ilic. DPOD: 6D Pose Object Detector and Refiner. In *International Conference on Computer Vision*, 2019. 1, 2, 5

[55] Yinqiang Zheng, Yubin Kuang, Shigeki Sugimoto, Kalle strm, and Masatoshi Okutomi. Revisiting the PnP Problem: A Fast, General and Optimal Solution. In *International Conference on Computer Vision*, 2013. 5

[56] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random Erasing Data Augmentation. In *arXiv Preprint*, 2017. 7