

Recent Advances in 3D Object and Hand Pose Estimation

Vincent Lepetit

vincent.lepetit@enpc.fr

LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France

Abstract

3D object and hand pose estimation have huge potentials for Augmented Reality, to enable tangible interfaces, natural interfaces, and blurring the boundaries between the real and virtual worlds. In this chapter, we present the recent developments for 3D object and hand pose estimation using cameras, and discuss their abilities and limitations and the possible future development of the field.

1 3D Object and Hand Pose Estimation for Augmented Reality

An Augmented Reality system should not be limited to the visualization of virtual elements integrated to the real world, it should also *perceive* the user and the real world surrounding the user. As even early Augmented Reality applications demonstrated, this is required to provide rich user experience, and unlock new possibilities. For example, the magic book application developed by Billinghurst and Kazo [Billinghurst *et al.*, 2001] and illustrated by Figure 1, shows the importance of being able to manipulate real objects in Augmented Reality. It featured a real book, from which virtual objects would pop up, and these virtual objects could be manipulated by the user using a small real “shovel”.

This early system relied on visual markers on the real book and the shovel. These markers would become the core of the popular ARToolkit, and were essential to robustly estimate the locations and orientations of objects in the 3D space, a geometric information required for the proper integration of the virtual objects with the real book, and their manipulation by the shovel.

Visual markers, however, require a modification of the real world, which is not always possible nor desirable in practice. Similarly, magnetic sensors have also been used to perceive the spatial positions of real objects, in order to integrate them to the augmented world. However, they share the same drawbacks as visual markers.

Being able to perceive the user’s hands is also very important, as the hands can act as an interface between the user’s intention and the virtual world. Haptic gloves or various types of joysticks have been used to capture the hands’ locations and motions. In particular, input pads are still popular with current augmented and virtual reality platforms, as they ensure robustness and accuracy for the perception of the user’s actions. Entirely getting rid of such hardware is extremely desirable, as it makes Augmented Reality user interfaces intuitive and natural.

Researchers have thus been developing Computer Vision approaches to perceiving the real world using simple cameras and without having to engineer the real objects with markers or sensors, or the user’s hands with gloves or handles. Such perception problem also appears in fields such as robotics, and the scientific literature on this topic is extremely rich. The goal of this chapter is to introduce the



Figure 1: Augmented Reality (left) and external (right) views of the early Magic book developed by Kazo and Billinghurst. This visionary experiment shows the importance of real object manipulation in AR applications.

reader to 3D object and hand perception based on cameras for Augmented Reality applications, and to the scientific literature on the topic, with a focus on the most recent techniques.

In the following, we will first detail the problem of 3D pose estimation, explain why it is difficult, and review early approaches to motivate the use of Machine Learning techniques. After a brief introduction to Deep Learning, we will discuss the literature on 3D pose estimation for objects, hands, and hands manipulating objects through representative works.

2 Formalization

In the rest of this chapter, we will define the 3D pose of a rigid object as its 3D location and orientation in some coordinate system, as shown in Figure 2. In the case of Augmented Reality, this coordinate system should be directly related to the headset, or to the tablet or smartphone, depending on the visualization system, *and* to the object. This is thus different from Simultaneous Localization and Mapping (SLAM), where the 3D pose of the camera can be estimated in an arbitrary coordinate system as long as it is consistent over time.

In general, this 3D pose has thus 6 degrees of freedom: 3 for the 3D translation, and 3 for the 3D orientation, which can be represented for example with a 3D rotation matrix or a quaternion. Because of these 6 degrees of freedom, the 3D pose is also called '6D pose' in the literature, the terms '3D pose' and '6D pose' thus refer to the same notion.

Articulated objects, such as scissors, or even more deformable objects such as clothes or sheets of papers have also been considered in the literature [Salzmann and Fua, 2010; Pumarola *et al.*, 2018]. Their positions and shapes in space have many more degrees of freedom, and specific representations of these poses have been developed. Estimating these values from images robustly remains very challenging.

The 3D pose of the hand is also very complex, since we would like to also consider the positions of the individual fingers. This is by contrast with gesture recognition for example, which aims at assigning a distinctive label such as 'open' or 'close' to a certain hand pose. Instead, we would like to estimate continuous values in the 3D space. One option among others to represent the 3D pose of a hand is to consider the 3D locations of each joint in some coordinate system. Given a hand model with bone

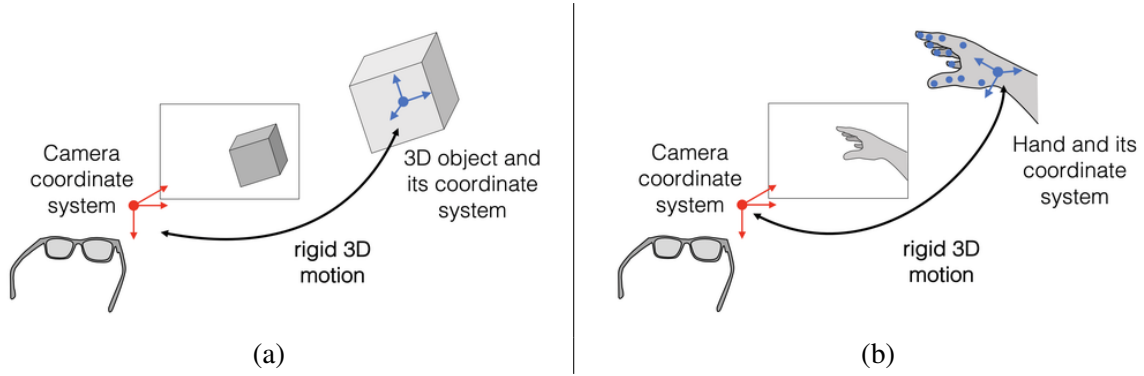


Figure 2: (a) The 3D pose of a rigid object can be defined as a 3D rigid motion between a coordinate system attached to the object and another coordinate system, for example one attached to an Augmented Reality headset. (b) One way to define the 3D pose of a hand can be defined as a rigid motion between a coordinate system attached to one of the joints (for example the wrist) and another coordinate system, plus the 3D locations of the joints in the first coordinate system. 3D rigid motions can be defined as the composition of a 3D rotation and a 3D translation.

length and rotation angles for all DoF of the joints, forward kinematics [Gustus *et al.*, 2012] can be used to calculate the 3D joint locations. Reciprocally, inverse kinematics [Tolani *et al.*, 2000] can be applied to obtain joint angles and bone length from the 3D joint locations. In addition, one may also want to estimate the shapes of the user’s hands, such as their sizes or the thickness of the fingers.

To be useful for Augmented Reality, 3D pose estimation has to run continuously, in real-time. When using computer vision, it means that it should be done from a single image, or stereo images captured at the same time, or RGB-D images that provide both color and depth data, maybe also relying on the previous images to guide or stabilize the 3D poses estimations. Cameras with large fields of view help, as they limit the risk of the target objects to leave the field of view compared to more standard cameras.

Stereo camera rigs, made of two or more cameras, provide additional spatial information that can be exploited for estimating the 3D pose, but make the device more complex and more expensive. RGB-D cameras are currently a good trade-off, as they also provide, in addition to a color image, 3D information in the form of a depth map, *i.e.* a depth value for each pixel, or at least most of the pixels of the color image. “Structured light” RGB-D cameras measure depth by projecting a known pattern in infrared light and capturing the projection with an infrared camera. Depth can then be estimated from the deformation of the pattern. “Time-of-flight” cameras are based on pulsed light sources and measure the time a light pulse takes to travel from the emitter to the scene and come back after reflection.

With structured-light technology, depth data can be missing at some image locations, especially along the silhouettes of objects, or on dark or specular regions. This technology also struggles to work outdoor, because of the ambient infrared sunlight. Time-of-flight cameras is also disturbed outdoor, as the high intensity of sunlight causes a quick saturation of the sensor pixels. Multiple reflections produced by concave shapes can also affect the time-of-flight.

But despite these drawbacks, RGB-D cameras remain however very useful in practice when they can be used, and algorithms using depth maps perform much better than algorithms relying only on color information—even though the gap is decreasing in recent research.

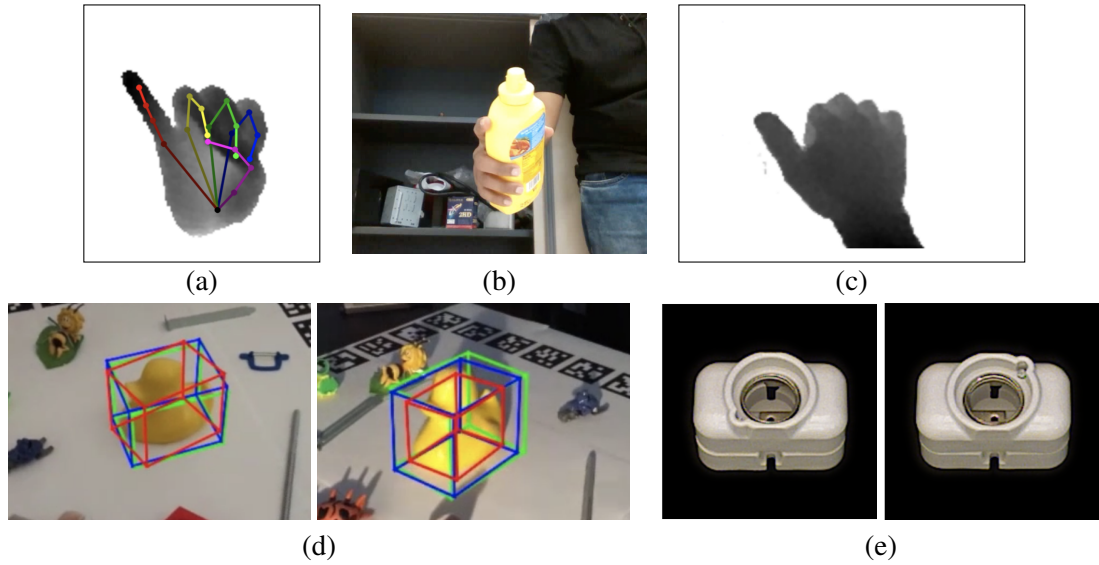


Figure 3: Some challenges when estimating the 3D poses of objects and hands. (a) Degrees of freedom. The human hand is highly articulated, and many parameters have to be estimated to correctly represent its pose. (b) Occlusions. In this case, the hand partially occludes the object, and the object partially occludes the hand. (c) Self-occlusions can also occur in the case of hands, especially in egocentric views. (d) Illuminations. Like the rubber duck in these examples, the appearance of objects can vary dramatically with illumination. (e) Ambiguities. Many manufactured objects have symmetric or almost symmetric shapes, which may make pose estimation ambiguous (object and images from the T-Less dataset [Hodan *et al.*, 2016]).

3 Challenges of 3D Pose Estimation using Computer Vision

There are many challenges that need to be tackled in practice when estimating the 3D poses of objects and hands from images. We list below some of these challenges, illustrated in Figure 3.

High degrees of freedom. As mentioned above, the pose of a rigid object can be represented with 6 scalar values, which is already a large number of degrees of freedom to estimate. The 3D pose of a human hand lives in an even much higher dimensional space, as it is often represented with about 32 degrees of freedom. The risk of making an error when estimating the pose increases with this number of degrees of freedom.

Occlusions. Occlusions of the target objects or hands, even partial, often disturb pose estimation algorithms. It can be difficult to identify which parts are visible and which parts are occluded, and the presence of occlusions may make pose estimation completely fail, or be very inaccurate. Occlusions often happen in practice. For example, when a hand manipulates an object, both the hand and the object are usually partially occluded.

Cluttered background. Objects in the background can act as distractors, especially if they look like target objects, or have similar parts.

Changing illumination conditions. In practice, illumination cannot be controlled, and will change the appearance of the objects, not only because the lighting is different but also because shadows can be cast on them. This requires robust algorithms. Moreover, cameras may struggle to keep a good balance between bright glares and dark areas, and using high dynamic range cameras becomes appealing under such conditions. Moreover, sunlight may make depth cameras fail as discussed above.

Material and Textures. Early approaches to pose estimation relied on the presence of texture or pattern on the objects' surfaces, because stable features can be detected and matched relatively easily and efficiently on patterns. However, in practice, many objects lack textures, making such approach fail. Non-Lambertian surfaces such as metal and glass make the appearance of objects change with the objects' poses because of specularities and reflections appearing on the objects' surfaces. Transparent objects are also of course problematic, as they do not appear clearly in images.

Ambiguities. Many manufactured objects are symmetrical or almost symmetrical, or have repetitive patterns. This generates possible ambiguities when estimating the poses of these objects that need to be handled explicitly.

3D Model Requirement. Most of the existing algorithms for pose estimation assume the knowledge of some 3D models of the target objects. Such 3D models provide very useful geometric constraints, which can be exploited to estimate the objects' 3D poses. However, building such 3D models takes time and expertise, and in practice, a 3D model is not readily available. Hand pose estimation suffers much less from this problem, since the 3D geometry of hands remains very similar from one person to another. Some recent works have also considered the estimation of the objects' geometry together with their poses [Grabner *et al.*, 2018]. These works are still limited to a few object categories, such as chairs or cars, but are very interesting for future developments.

4 Early Approaches to 3D Pose Estimation and Their Limits

3D pose estimation from images has a long history in computer vision. Early methods were based on simple image features, such as edges or corners. Most importantly, they strongly relied on some prior knowledge about the 3D pose, to guide the pose estimation in the high-dimensional space of possible poses. This prior may come from the poses estimated for the previous images, or could be provided by a human operator.

3D tracking methods. For example, pioneer works such as [Harris and Stennett, 1990] and [Lowe, 1991] describe the object of interest as a set of 3D geometric primitives such as lines and conics, which were matched with contours in the image to find a 3D pose estimate by solving an optimization problem to find the 3D pose that reprojects the 3D geometric primitives to the matched image contours. The entire process is computationally light, and careful implementations were able to achieve high frame rates with computers that would appear primitive to us.

Unfortunately, edge-based approaches are quite unreliable in practice. Matching the reprojection of the 3D primitives with image contours is difficult to achieve: 3D primitives do not necessarily appear as strong image contours, except for carefully chosen objects. As a result, these primitives are likely

to be matched with the wrong image contours, especially in case of cluttered background. Incorrect matches will also occur in case of partial occlusions. Introducing robust estimators into the optimization problem [Drummond and Cipolla, 2002] helps, but it appears that it is impossible in general to be robust to such mismatches in edge-based pose estimation: Most of the time, a match between two contours provides only limited constraints on the pose parameters, as the two contours can “slide” along each other and still be a good match. Contours thus do not provide reliable constraints for pose estimation.

Relying on feature points [Harris and Stephens, 1988] instead of contours provides more constraints as point-to-point matching does not have the “sliding ambiguity” of contour matching. Feature point matching has been used for 3D object tracking with some success, for example in [Vacchetti *et al.*, 2004]. However, such approach assumes the presence of feature points that can be detected on the object’s surface, which is not true for all the objects.

The importance of detection methods. The two approaches described above assume that prior knowledge on the object pose is available to guide the matching process between contours or points. Typically, the object pose estimated at time t is exploited to estimate the pose at time $t + 1$. In practice, this makes such approaches fragile, because the poses at time t and $t + 1$ can be very different if the object moves fast, or because the pose estimated at time t can be wrong.

Being able to estimate the 3D pose of an object without relying too much on prior knowledge is therefore very important in practice. As we will see, this does not mean that methods based on strong prior knowledge on the pose are not useful. In fact, they tend to be much faster and/or more accurate than “detection methods”, which are more robust. A natural solution is thus to combine both, and this is still true in the modern era of object and pose estimation.

One early popular method for object pose estimation without pose prior was also based on feature points, often referred to as keypoints in this context. This requires the ability to match keypoints between an input image, and a reference image of the target object, which is captured offline and in which the 3D pose and the 3D shape of the object is known. By using geometry constraints, it is then possible to estimate the object’s 3D pose. However, wide baseline point matching is much more difficult than short baseline matching used by tracking methods. SIFT keypoints and descriptors [Lowe, 2001; Lowe, 2004] were a breakthrough that made many computer vision applications possible, including 3D pose estimation. They were followed by faster methods, including SURF [Bay *et al.*, 2008] and ORB [Rublee *et al.*, 2011].

As for tracking methods based on feature points, the limitation of keypoint-based detection and pose estimation is that it is limited to objects exhibiting enough keypoints, which is not the case in general. This approach was still very successful for Augmented Reality in magazines for example, where the “object” is an image printed on paper—so it has a simple planar geometry—and selected to guarantee that the approach will work well by exhibiting enough keypoints [Kim *et al.*, 2010].

To be able to detect and estimate the 3D pose of objects with almost no keypoints, sometimes referred to as “texture-less” objects, some works attempted to use “templates”: The templates of [Hinterstoisser *et al.*, 2012b] aim at capturing the possible appearances of the object by discretizing the 3D pose space, and representing the object’s appearance for each discretized pose by a template covering the full object, in a way that is robust to lighting variations. Then, by scanning the input images looking for templates, the target object can be detected in 2D and its pose estimated based on the template that matches best the object appearance. However, such approach requires the creation of many templates, and is poorly robust to occlusions.

Conclusion. We focused in this section on object pose estimation rather than hand pose estimation, however the conclusion would be the same. Early approaches were based on handcrafted methods to extract features from images, with the goal of estimating the 3D pose from these features. This is however very challenging to do, and almost doomed to fail in the general case. Since then, Machine Learning-based methods have been shown to be more adapted, even if they also come with their drawbacks, and will be discussed in the rest of this chapter.

5 Machine Learning and Deep Learning

Fundamentally, 3D pose estimation of objects or hands can be seen as a mapping from an image to a representation of the pose. The input space of this mapping is thus the space of possible images, which is an incredibly large space: For example, a RGB VGA image is made of almost 1 million values of pixel values. Not many fields deal with 1 million dimension data! The output space is much smaller, since it is made of 6 values for the pose of a rigid object, or a few tens for the pose of a hand. The natures of the input and output spaces of the mapping sought in pose estimation are therefore very different, which makes this mapping very complex. From this point of view, we can understand that it is very difficult to hope for a pure “algorithmic” approach to code this mapping.

This is why Machine Learning techniques, which use data to improve algorithms, have become successful for pose estimation problems, and computer vision problems in general. Because they are data-driven, they can find automatically an appropriate mapping, by contrast with previous approaches that required hardcoding mostly based on intuition, which can be correct or wrong.

Many Machine Learning methods exist, and Random Forests, also called Random Decision Forests or Randomized Trees, were an early popular method in the context of 3D pose estimation [Lepetit *et al.*, 2005; Shotton *et al.*, 2013; Brachmann *et al.*, 2014]. Random Forests can be efficiently applied to image patches and discrete (multi-class) and continuous (regression) problems, which makes them flexible and suitable to fast, possibly real-time, applications.

Deep Learning. For multiple reasons, Deep Learning, a Machine Learning technique, took over the other methods almost entirely during the last years in many scientific fields, including 3D pose estimation. It is very flexible, and in fact, it has been known for a long time that any continuous mapping can be approximated by two-layer networks, as finely as wanted [Hornik *et al.*, 1989; Pinkus, 1999]. In practice, networks with more than two layers tend to generalize better, and to need dramatically less parameters than two-layer networks [Eldan and Shamir, 2016], which makes them a tool of choice for computer vision problems.

Many resources can now be easily found to learn the fundamentals of Deep Learning, for example [Goodfellow *et al.*, 2016]. To stay brief, we can say here that a Deep Network can be defined as a composition of functions (“layers”). These functions may depend on parameters, which need to be estimated for the network to perform well. Almost any function can be used as layer, as long as it is useful to solve the problem at hand, and if it is *differentiable*. This differentiable property is indeed required to find good values for the parameters by solving an optimization problem.

Deep Network Training. For example, if we want to make a network F predict the pose of a hand visible in an image, one way to find good values $\hat{\Theta}$ for the network parameters is to solve the following

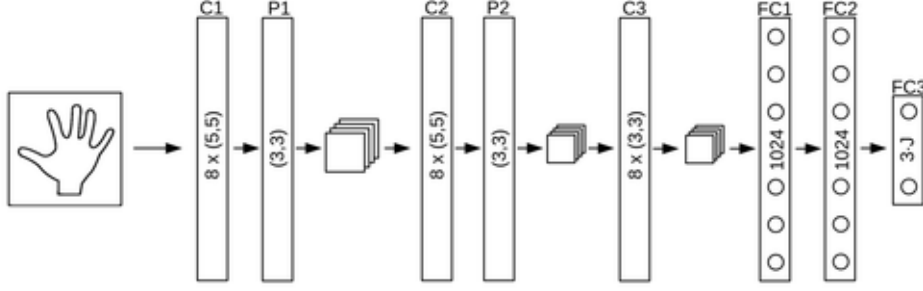


Figure 4: A basic Deep Network architecture applied to hand pose prediction. C1, C2, and C3 are convolutional layers, P1 and P2 are pooling layers, and FC1, FC2, and FC3 are fully connected layers. The numbers in each bar indicates either: The number and size of the linear filters in case of the convolutional layers, the size of the pooling region for the pooling layers, and the size of the output vector in case of the fully connected layers.

optimization problem:

$$\hat{\Theta} = \arg \min_{\Theta} \mathcal{L}(\Theta) \text{ with} \quad (1)$$

$$\mathcal{L}(\Theta) = \frac{1}{N} \sum_{i=1}^N \|F(I_i; \Theta) - \mathbf{e}_i\|^2,$$

where $\{(I_1, \mathbf{e}_1), \dots, (I_N, \mathbf{e}_N)\}$ is a *training set* containing pairs made of images I_i and the corresponding poses \mathbf{e}_i , where the \mathbf{e}_i are vectors that contain, for example, the 3D locations of the joints of the hand as was done in [Oberweger *et al.*, 2017]. Function \mathcal{L} is called a loss function. Optimizing the network parameters Θ is called training. Many optimization algorithms and tricks have been proposed to solve problems like Eq. (1), and many software libraries exist to make the implementation of the network creation and its training an easy task.

Because optimization algorithms for network training are based on gradient descent, any differentiable loss function can be used in principle, which makes Deep Learning a very flexible approach, as it can be adapted easily to the problem at hand.

Supervised Training and the Requirement for Annotated Training Sets. Training a network by using a loss function like the one in Eq. 1 is called supervised training, because we assume the availability of an annotated dataset of images. Supervised training tends to perform very well, and it is used very often in practical applications.

This is however probably the main drawback of Deep Learning-based 3D pose estimation. While early methods relied only on a 3D model, and for some of them only on a small number of images of the object, modern approaches based on Deep Learning require a large dataset of images of the target objects, annotated with the ground truth 3D poses

Such datasets are already available (see Section 6), but they are useful mostly for evaluation and comparison purposes. Applying current methods to new objects requires creating a dataset for these new objects, and this is a cumbersome task.

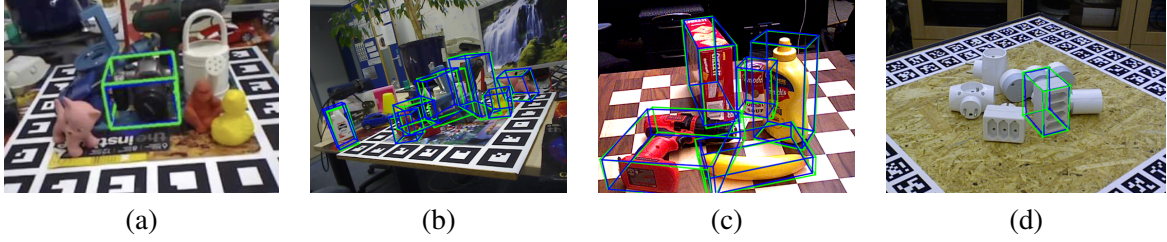


Figure 5: Images from popular datasets for 3D object pose estimation. (a) LineMOD; (b) Occluded LineMOD; (c) YCB-Video; (d) T-Less.

It is also possible to use synthetic data for training, by generating images using Computer Graphics techniques. This is used in many works, often with special care to take into account the differences between real and synthetic images, as we will see below.

6 Datasets

Datasets have become an important aspect of 3D pose estimation, for training, evaluating, and comparing methods. We describe some for object, hand, and hand+object pose estimation below.

6.1 Datasets for Object Pose Estimation

LineMOD Dataset. The LineMOD dataset [Hinterstoisser *et al.*, 2012b] predates most machine learning approaches and as such, it is not divided into a training and a test sets. It is made of 15 small objects, such as a camera, a lamp, and a cup. For each object, it offers a set of 1200 RGB-D images of the object surrounded by clutter. The other objects are often visible in the clutter, but only the 3D pose of the target object is provided for each set. The 3D models of the objects are also provided.

Occluded LineMOD Dataset. The Occluded LineMOD dataset was created by the authors of [Brachmann *et al.*, 2014] from LineMOD by annotating the 3D poses of the objects belonging to the dataset but originally not annotated because they were considered as part of the clutter. This results into a sequence of 1215 frames, each frame labeled with the 3D poses of eight objects in total, as well as the objects’ masks. The objects show severe occlusions, which makes pose estimation challenging.

YCB-Video Dataset. This dataset [Xiang *et al.*, 2018a] consists of 92 video sequences, where 12 sequences are used for testing and the remaining 80 sequences for training. In addition, the dataset contains 80k synthetically rendered images, which can be used for training as well. There are 21 “daily life” objects in the dataset, from cereal boxes to scissors or plates. These objects were selected from the YCB dataset [Calli *et al.*, 2017] and are available for purchase. The dataset is captured with two different RGB-D sensors. The test images are challenging due to the presence of significant image noise, different illumination levels, and large occlusions. Each image is annotated with the 3D object poses, as well as the objects’ masks.

T-Less Dataset. The T-Less dataset [Hodan *et al.*, 2016] is made from 30 “industry-relevant” objects. These objects have no discriminative color nor texture. They present different types of symmetries and similarities between them, making pose estimation often almost ambiguous. The images were captured using three synchronized sensors: two RGB-D cameras, one structured-light based and one time-of-flight based, and one high-resolution RGB camera. The test images (10K from each sensor) are from 20 scenes with increasing difficulties, with partial occlusions and contacts between objects. This dataset remains extremely challenging.

6.2 Datasets for Hand Pose Estimation

Early datasets. The NYU dataset [Tompson *et al.*, 2014] contains over 72k training and 8k test RGB-D images data, captured from three different viewpoints using a structured-light camera. The images were annotated with 3D joint locations with a semi-automated approach, by using a standard 3D hand tracking algorithm reinitialized manually in case of failure. The ICVL dataset [Tang *et al.*, 2014] contains over 180k training depth frames showing various hand poses, and two test sequences with each approximately 700 frames, all captured with a time-of-flight camera. The depth images have a high quality with hardly any missing depth values and sharp outlines with little noise. Unfortunately, the hand pose variability of this dataset is limited compared to other datasets, and annotations are rather inaccurate [Supancic *et al.*, 2015]. The MSRA dataset [Sun *et al.*, 2015] contains about 76k depth frames, captured using a time-of-flight camera from nine different subjects.

BigHands Dataset. The BigHands dataset [Yuan *et al.*, 2017] contains an impressive 2.2 millions RGB-D images captured with a structured-light camera. The dataset was automatically annotated by using six 6D electromagnetic sensors and inverse kinematics to provide 21 3D joint locations per frame. A significant part is made of egocentric views. The depth images have a high quality with hardly any missing depth values, and sharp outlines with little noise. The labels are sometimes inaccurate because of the annotation process, but the dataset has a large hand pose variability, from ten users.

CMU Panoptic Hand Dataset. The CMU Panoptic hand dataset [Simon *et al.*, 2017] is made of 15k real and synthetic RGB images, from a third-person point of view. The real images were recorded in the CMU’s Panoptic studio and annotated by the method proposed in the paper, based on multiple views. The annotations are only in 2D but can still be useful for multi-view pose estimation.

6.3 Datasets for Object and Hand Pose Estimation

GANerated Hand Dataset. The GANerated hand dataset [Mueller *et al.*, 2018b] is a large dataset made of 330K synthetic images of hands, sometimes holding an object, in front of a random background. The images are annotated with the 3D poses of the hand. The images were made more realistic by extending CycleGAN [Zhu *et al.*, 2017].

First-Person Hand Action dataset. The First-Person Hand Action dataset [Garcia-Hernando *et al.*, 2018] provides a dataset of hand and object interactions with 3D annotations for both hand joints and object pose. They used a motion capture system made of magnetic sensors attached to the user’s hand and to the object in order to obtain hand 3D pose annotations in RGB-D video sequences. Unfortunately,

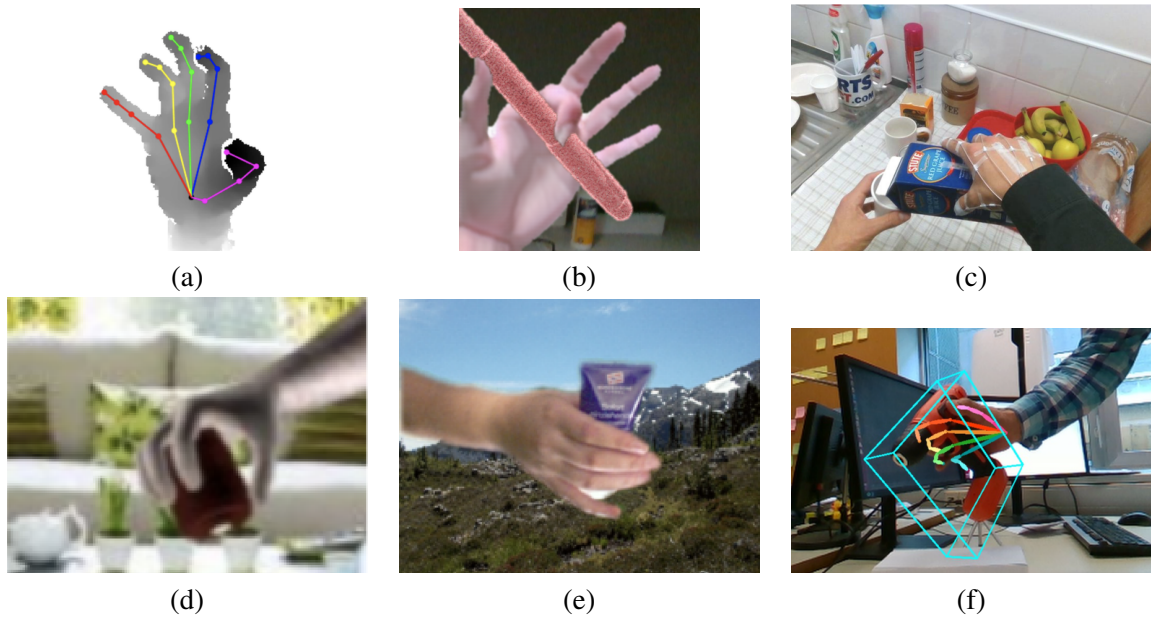


Figure 6: Images from popular datasets for 3D hand (a) and hand+object (b-f) pose estimation. (a) BigHand; (b) GANerated hand dataset; (c) First-Person Hand Action dataset; (d) Obman dataset; (e) FreiHAND dataset; (f) HO-3D dataset.

this changes the appearance of the hand in the color images as the sensors and the tape attaching them are visible, but the dataset proposes a large number of frames under various conditions (more than 100K egocentric views of 6 subjects doing 45 different types of daily-life activities).

ObMan Dataset. Very recently, [Hasson *et al.*, 2019a] introduced ObMan, a large dataset of images of hands grasping objects. The images are synthetic, but the grasps are generated using an algorithm from robotics and the grasps still look realistic. The dataset provides the 3D poses and shapes of the hand as well as the object shapes.

FreiHand Dataset. [Zimmermann *et al.*, 2019] proposed a multi-view RGB dataset, FreiHAND, which includes hand-object interactions and provides the 3D poses and shapes of the hand. It relies on a green-screen background environment so that it is easy to change the background for training purposes.

HO-3D Dataset. [Hampali *et al.*, 2020] proposed a method to automatically annotate video sequences captured with one or more RGB-D cameras with the object and hand poses and shapes. This results in a dataset made of 75,000 real RGB-D images, from 10 different objects and 10 different users. The objects come from the YCB dataset (see Section 6.1). The backgrounds are complex, and the mutual occlusions are often large, which makes the pose estimation realistic but very challenging.

6.4 Metrics

Metrics are important to evaluate and compare methods. Many metrics exist, and we describe here only the main ones for object pose estimation. Discussions on metrics for 3D object pose estimation can be found in [Hodan *et al.*, 2016] and [Brégier *et al.*, 2018].

ADD, ADI, ADD-S, and the 6D Pose metrics. The ADD metric [Hinterstoisser *et al.*, 2012a] calculates the average distance in 3D between the model points, after applying the ground truth pose and the predicted pose. This can be formalized as:

$$\text{ADD} = \frac{1}{|\mathcal{V}|} \sum_{\mathbf{M} \in \mathcal{V}} \|\text{Tr}(\mathbf{M}; \hat{\mathbf{p}}) - \text{Tr}(\mathbf{M}; \bar{\mathbf{p}})\|_2, \quad (2)$$

where \mathcal{V} is the set of the object’s vertices, $\hat{\mathbf{p}}$ the estimated pose and $\bar{\mathbf{p}}$ the ground truth pose, and $\text{Tr}(\mathbf{M}; \mathbf{p})$ the rigid transformation in \mathbf{p} applied to 3D point \mathbf{M} . In the 6D Pose metric, a pose is considered when the ADD metric is less than 10% of the object’s diameter.

For the objects with ambiguous poses due to symmetries, [Hinterstoisser *et al.*, 2012a] replaces the ADD metric by the ADI metric, also referred to as the ADD-S metric in [Xiang *et al.*, 2018b], computed as follows:

$$\text{ADD-S} = \frac{1}{|\mathcal{V}|} \sum_{\mathbf{M}_1 \in \mathcal{V}} \min_{\mathbf{M}_2 \in \mathcal{V}} \|\text{Tr}(\mathbf{M}_1; \hat{\mathbf{p}}) - \text{Tr}(\mathbf{M}_2; \bar{\mathbf{p}})\|_2, \quad (3)$$

which averages the distances from points after applying the predicted pose to the *closest* points under the ground truth pose. The advantage of this metric is that it is indeed equal to zero when the pose is retrieved up to a symmetry, even if it does not exploit the symmetries of the object.

7 Modern Approaches to 3D Object Pose Estimation

Over the past years, many authors realise that Deep Learning is a powerful tool for 3D object pose estimation from images. We discuss here the development of Deep Learning applied to 3D object pose estimation over time. This development was and is still extremely fast, with improving accuracy, robustness, and computation times. We present this development through several milestone methods, but much more methods could also be included here.

7.1 BB8

One of the first Deep Learning methods for 3D object pose estimation is probably BB8 [Rad and Lepetit, 2017]. As it is the first method we describe, we will present it in some details.

This method proceeds in three steps. It first detect the target objects in 2D using coarse object segmentation. It then applies a Deep Network on each image window centered on detected objects. Instead of predicting the 3D pose of the detected objects in the form of a 3D translation and a 3D rotation, it predict the 2D projections of the corners of the object’s bounding box, and compute the 3D pose from these 2D-3D correspondences with a PnP algorithm [Gao *et al.*, 2003]—hence the name for the method, from the 8 corners of the bounding box, as illustrating in Figure 7. Compared to the direct prediction of the pose, this avoids the need for a meta-parameter to balance the translation and rotation

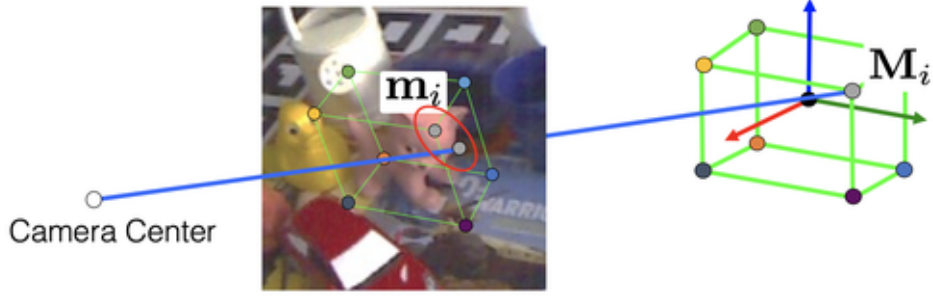


Figure 7: Some 3D object pose estimation methods predict the pose by first predicting the 2D reprojections \mathbf{m}_i of some 3D points \mathbf{M}_i , and then computing the 3D rotation and translation from the 2D-3D correspondences between the \mathbf{m}_i and \mathbf{M}_i points using a PnP algorithm.

terms. It also tends to make network optimization easier. This representation was used in some later works.

Since it is the first Deep Learning method for pose estimation we describe, we will detail the loss function (see Section 5) used to train the network. This loss function \mathcal{L} is a least-squares error between the predicted 2D points and the expected ones, and can be written as:

$$\mathcal{L}(\Theta) = \frac{1}{8} \sum_{(W, \mathbf{p}) \in \mathcal{T}} \sum_i \|\text{Proj}_{\mathbf{p}}(\mathbf{M}_i) - F(W; \Theta)_i\|^2, \quad (4)$$

where F denotes the trained network and Θ its parameters. \mathcal{T} is a training set made of image windows W containing a target object under a pose \mathbf{p} . The \mathbf{M}_i are the 3D coordinates of the corners of the bounding box of for this object, in the object coordinate system. $\text{Proj}_{\mathbf{e}, \mathbf{t}}(\mathbf{M})$ projects the 3D point \mathbf{M} on the image from the pose defined by \mathbf{e} and \mathbf{t} . $F(W; \Theta)_i$ returns the two components of the output of F corresponding to the predicted 2D coordinates of the i -th corner.

The problem of symmetrical and “almost symmetrical” objects. Predicting the 3D pose of objects a standard least-squares problem, using a standard representation of the pose or point reprojections as in BB8, yields to large errors on symmetrical objects, such as many objects in the T-Less dataset (Figure 5(d)). This dataset is made of manufactured objects that are not only similar to each other, but also have one axis of rotational symmetry. Some objects are not perfectly symmetrical but only because of small details, like a screw.

The approach described above fails on these objects because it tries to learn a mapping from the image space to the pose space. Since two images of a symmetrical object under two different poses look identical, the image-pose correspondence is in fact a one-to-many relationship.

For objects that are perfectly symmetrical, [Rad and Lepetit, 2017] proposed a solution that we will not describe here, as simpler solutions have been proposed since. Objects that are “almost symmetrical” can also disturb pose prediction, as the mapping from the image to the 3D pose is difficult to learn even though this is a one-to-one mapping. Most recent methods ignore the problem and consider that these objects are actually perfectly symmetrical and consider a pose recovered up to the symmetries as

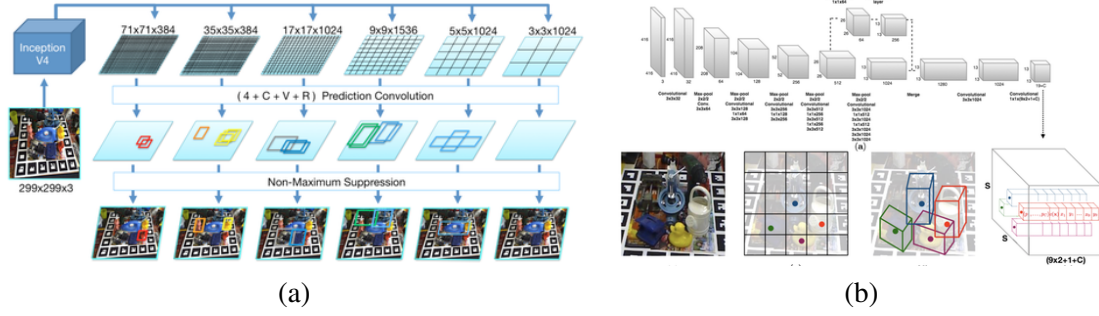


Figure 9: (a) The architecture of SSD, extended by SSD-6D to predict 3D poses. For each bounding box corresponding to a detected object, a discretized 3D pose is also predicted by the network. Image from [Kehl *et al.*, 2017]. (b) The architecture of YOLO-6D. Image from [Tekin *et al.*, 2018].

the original training images, which can be done using the ground truth poses and the objects’ 3D models. The background is replaced by a patch extracted from a randomly picked image from the ImageNet dataset [Russakovsky *et al.*, 2015]. Note that this procedure removes context by randomly replacing the surroundings of the objects. Some information is thus lost, as context could be useful for pose estimation.

7.2 SSD-6D

BB8 relies on two separate networks to first detect the target objects in 2D and then predict their 3D poses, plus a third one if refinement is performed. Instead, as shown in Figure 9(a), SSD-6D [Kehl *et al.*, 2017] extends a deep architecture (the SSD architecture [Liu *et al.*, 2016]) developed for 2D object detection to 3D pose estimation (referred in the paper as 6D estimation). SSD had already been extended to pose estimation in [Poirson *et al.*, 2016], but SS6-6D performs full 3D pose prediction.

This prediction is done by first discretizing the pose space. Each discretized pose is considered as a class, to turn the pose prediction into a classification problem rather than a regression one, as this was performing better in the authors’ experience. This discretization was done on the 3D pose decomposition into direction of view over a half-sphere and in-plane rotation. The 3D translation can be computed from the 2D bounding box. To deal with symmetrical objects, views on the half-sphere corresponding to identical appearance for the object are merged into the same class.

A single network is therefore trained to perform both 2D object detection and 3D pose prediction, using a single loss function made of a weighted sum of different terms. The weights are hyperparameters and have to be tuned, which can be difficult. However, having a single network is an elegant solution, and more importantly for practical applications, this allows to save computation times: Image feature extraction, the slowest part of the network, is performed only once even if it is used to predict the 2D bounding boxes and the 3D pose estimation.

A refinement procedure is then run on each object detection to improve the pose estimate predicted by the classifier. SSD-6D relies on a method inspired by an early approach to 3D object tracking based on edges [Drummond and Cipolla, 2002]. Data augmentation was done by rendering synthetic views of the objects using their 3D models over images from COCO [Lin *et al.*, 2014], which helps but only to some extent as the differences between the real and synthetic images (the “domain gap”) remain high.

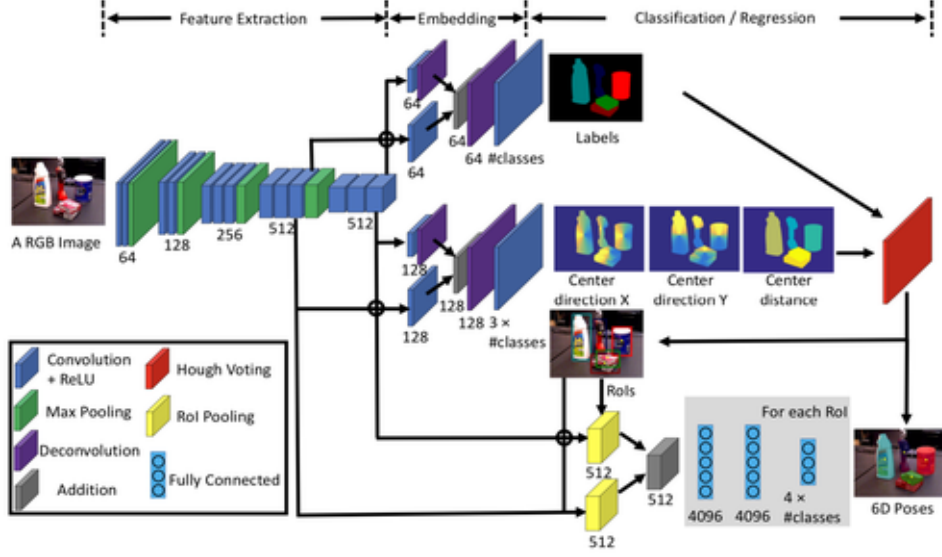


Figure 10: The architecture of PoseCNN. Image from [Xiang *et al.*, 2018b].

7.3 YOLO-6D

The method proposed in [Tekin *et al.*, 2018], sometimes referred as YOLO-6D, is relatively similar to SSD-6D, but makes different choices that makes it faster and more accurate. As shown in Figure 9(b), it relies on the YOLO architecture [Redmon *et al.*, 2016; Redmon and Farhadi, 2017] for 2D object detection, and predicts the 3D object poses in a form similar to the one of BB8. The authors report that training the network using this pose representation was much simpler than when using quaternions to represent the 3D rotation. As YOLO is much faster and as they do not discretize the pose space, [Tekin *et al.*, 2018] also reports much better performance than SSD-6D in terms of both computation times and accuracy.

7.4 PoseCNN

PoseCNN is a method proposed in [Xiang *et al.*, 2018b]. It relies on a relatively complex architecture, based on the idea of decoupling the 3D pose estimation task into different sub-tasks. As shown in Figure 10, this architecture predicts for each pixel in the input image 1) the object label, 2) a unit vector towards the 2D object center, and 3) the 3D distance between the object center and the camera center, plus 4) 2D bounding boxes for the objects and 5) a 3D rotation for each bounding box in the form of a quaternion. From 1), 2), and 3), it is possible to compute the 3D translation vector for each visible object, which is combined with the 3D rotation to obtain the full 3D pose for each visible object.

Maybe the most interesting contribution of PoseCNN is the loss function. It uses the ADD metric as the loss (see Section 6.4), and even more interestingly, the paper shows that the ADD-S metric, used to evaluate the pose accuracy for symmetrical objects (see Section 6.1), can be used as a loss function to deal with symmetrical objects. The ADI metric makes use of the min operator, which makes it maybe an unconventional loss function, but it is still differentiable and can be used to train a network. This results in an elegant and efficient way of handling object symmetries.

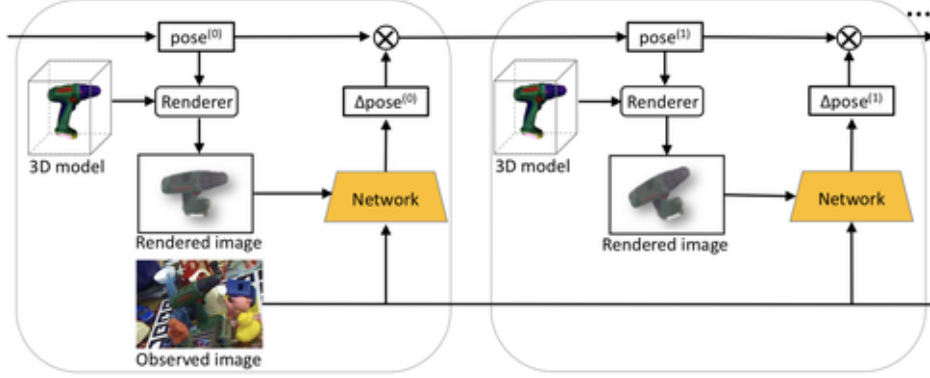


Figure 11: The DeepIM [Xiang *et al.*, 2018b] refinement architecture. It is similar to the one of BB8 (Figure 8) but predicts a pose update in a carefully chosen coordinate system for better generalization.

Another contribution of [Xiang *et al.*, 2018b] is the introduction of the YCB-Video dataset (see Section 6.1).

7.5 DeepIM

DeepIM [Li *et al.*, 2018] proposes a refinement step that resembles the one in BB8 (see the discussion on refinement steps in Section 7.1), but with several main differences. As BB8, it trains a network to compare the input image with a rendering of an object that has already been detected, and for which a pose estimate is already available. The network outputs an update for the 3D object pose that will improve the pose estimate, and this process can be iterated until convergence. The first difference with BB8 is that the rendered image has a high-resolution and the input image is centered on the object and upsampled to the same resolution. This allows a higher precision when predicting the pose update. The loss function also includes terms to make the network output the flow and the object mask, to introduce regularization.

The second difference is more fundamental, as it introduces a coordinate system well suited to define the rotation update that must be predicted by the network. The authors first remark that predicting the rotation in the camera coordinate system because this rotation would also *translate* the object. They thus set the center of rotation to the object center. For the axes of the coordinate system, the authors remark that using those of the object’s 3D model is not a good option, as they are arbitrary and this would force the network to learn them for each object. They therefore propose to use the axes of the camera coordinate system, which makes the network generalize much better. The translation update is predicted as a 2D translation on the image plane, plus a delta along the z axis of the camera in a log-scale. To learn this update, DeepIM uses the same loss function as [Xiang *et al.*, 2018b] (ADD, or ADD-S for symmetrical objects). This approach is even shown to generalize to unseen objects, but this is admittedly demonstrated in the paper only on very simple object renderings.

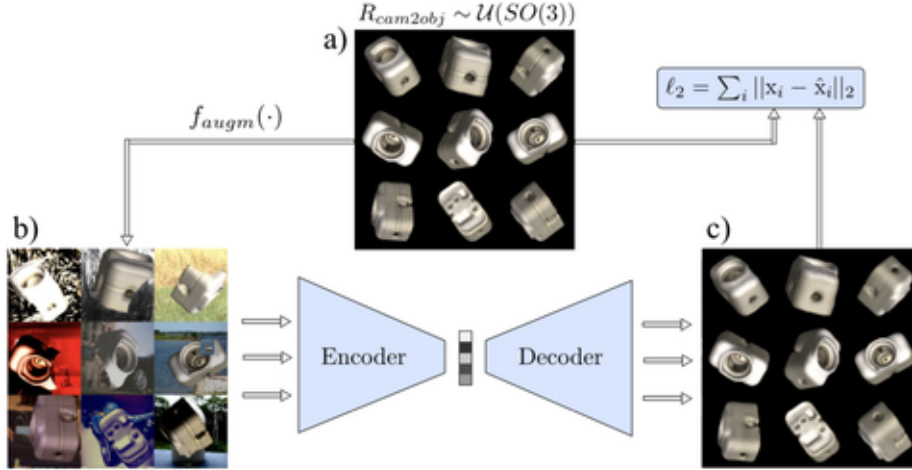


Figure 12: Autoencoders trained to learn an image embedding. This embedding is robust to nuisances such as background and illumination changes, and the possible rotations for the object in the input image can be retrieved based on this embedding. Image from [Sundermeyer *et al.*, 2019].

7.6 Augmented Autoencoders

The most interesting aspect of the Augmented Autoencoders method proposed in [Sundermeyer *et al.*, 2019] is the way it deals with symmetrical objects. It proposes to first learn an embedding that can be computed from an image of the object. This embedding should be robust to imaging artefacts (illumination, background, etc.) and depends only on the object’s appearance: The embeddings for two images of the object under ambiguous rotations should thus be the same. At run-time, the embedding for an input image can then be mapped to all the possible rotations. This is done efficiently by creating a codebook offline, by sampling views around the target objects, and associating the embeddings of these views to the corresponding rotations. The translation can be recovered from the object bounding box’ 2D location and scale.

To learn to compute the embeddings, [Sundermeyer *et al.*, 2019] relies on an autoencoder architecture, shown in Figure 12. An Encoder network predicts an embedding from an image of an object with different image nuisances; a Decoder network takes the predicted embedding as input and should generate the image of the object under the same rotation but without the nuisances. The encoder and decoder are trained together on synthetic images of the objects, to which nuisances are added. The loss function encourages the composition of the two networks to output the original synthetic images, despite using the images with nuisances as input. Another motivation for this autoencoder architecture is to learn to be robust to nuisances, even though a more standard approach outputting the 3D pose would probably do the same when trained on the same images.

7.7 Robustness to Partial Occlusions: [Oberweger *et al.*, 2018], [Hu *et al.*, 2019], PVNet

[Oberweger *et al.*, 2018], [Hu *et al.*, 2019], and PVNet [Peng *et al.*, 2019] developed almost in parallel similar methods that provide accurate 3D poses even under large partial occlusions. Indeed, [Oberweger *et al.*, 2018] shows that Deep Networks can be robust to partial occlusions when predicting a pose from

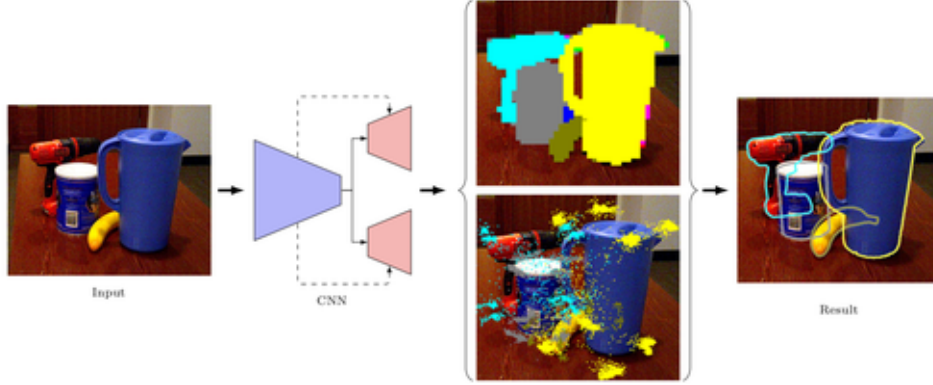


Figure 13: Being robust to large occlusions. [Hu *et al.*, 2019] makes predictions for the 2D reprojections of the objects’ 3D bounding boxes from many image locations. Image locations that are not occluded will result in good predictions, the predictions made from occluded image locations can be filtered out using a robust estimation of the pose. Image from [Hu *et al.*, 2019].

a image containing the target object when trained on examples with occlusions, but only to some extent. To become more robust and more accurate under large occlusions, [Oberweger *et al.*, 2018; Hu *et al.*, 2019; Peng *et al.*, 2019] proposed to predict the 3D pose in the form of the 2D reprojections of some 3D points as in BB8, but by combining multiple predictions, where each prediction is performed from different local image information. As shown in Figure 13, the key idea is that local image information that is not disturbed by occlusions will result into good predictions; local image information disturbed by occlusions will predict erroneous reprojections, but these can be filtered out with a robust estimation.

[Oberweger *et al.*, 2018] predicts the 2D reprojections in the form of heatmaps, to handle the ambiguities of mapping between image locations and the reprojections, as many image locations can look the same. However, this makes predictions relatively slow. Instead, [Hu *et al.*, 2019] predicts a single 2D displacement between the image location and each 2D reprojection. These results in many possible reprojections, some noisy, but the correct 3D pose can be retrieved using RANSAC. PVNet [Peng *et al.*, 2019] chose to predict the directions towards the reprojections, rather than a full 2D displacement, but relies on a similar RANSAC procedure to estimate the final 3D pose. [Hu *et al.*, 2019] and [Peng *et al.*, 2019] also predict the masks of the target objects, in order to consider the predictions from the image locations that lie on the objects, as they are the only informative ones.

7.8 DPOD and Pix2Pose

DPOD [Zakharov *et al.*, 2019] and Pix2Pose [Park *et al.*, 2019b] are also two methods that have been presented at the same conference and present similarities. They learn to predict for each pixel of an input image centered on a target object its 3D coordinates in the object coordinate system, see Figure 14. More exactly, DPOD predicts the pixels’ texture coordinates but this is fundamentally the same thing as they both provide 2D-3D correspondences between image locations and their 3D object coordinates. Such representation is often called ‘object coordinates’ and more recently ‘location field’ and has already been used for 3D pose estimation in [Brachmann *et al.*, 2014; Brachmann *et al.*, 2016] and before that in [Taylor *et al.*, 2012] for human pose estimation. From these 2D-3D correspondences, it

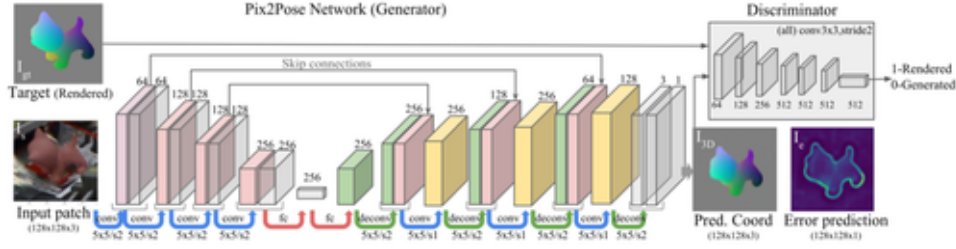


Figure 14: Pix2Pose [Park *et al.*, 2019b] predicts the object coordinates of the pixels lying on the object surface, and computes the 3D pose from these correspondences. It relies on a GAN [Goodfellow *et al.*, 2014] to perform this prediction robustly even under occlusion. Image from [Park *et al.*, 2019b].

is then possible to estimate the object pose using RANSAC and PnP. [Park *et al.*, 2019b] relies on a GAN [Goodfellow *et al.*, 2014] to learn to perform this prediction robustly even under occlusion. Note that DPOD has been demonstrated to run in real-time on a tablet.

Pix2Pose [Park *et al.*, 2019b] also uses a loss function that deals with symmetrical objects and can be written:

$$\mathcal{L} = \min_{\mathbf{R} \in \text{Sym}} \frac{1}{|\mathcal{V}|} \sum_{\mathbf{M} \in \mathcal{V}} \|\text{Tr}(\mathbf{M}; \hat{\mathbf{p}}) - \text{Tr}(\mathbf{M}; \mathbf{R} \cdot \bar{\mathbf{p}})\|_2, \quad (5)$$

where $\hat{\mathbf{p}}$ and $\bar{\mathbf{p}}$ denote the predicted and ground truth poses respectively, Sym is a set of rotations corresponding to the object symmetries, and $\mathbf{R} \cdot \mathbf{p}$ denotes the composition of such a rotation and a 3D pose. This deals with symmetries in a much more satisfying way than the ADD-S loss (Eq. (3)).

7.9 Discussion

As can be seen, the last recent years have seen rich developments in 3D object pose estimation, and methods have become even more robust, more accurate, and faster. Many methods rely on 2D segmentation to detect the objects, which seems pretty robust, but assumes that only several instances of the same object do not overlap in the image. Most methods also rely on a refinement step, which may slow things down, but relax the need for having a single strong detection stage.

There are still several caveats though. First, it remains difficult to compare methods based on Deep Learning, even though the use of public benchmarks helped to promote fair comparisons. Quantitative results depend not only on the method itself, but also on how much effort the authors put into training their methods and augmenting training data. Also, the focus on the benchmarks may make the method overfit to their data, and it is not clear how well current methods generalize to the real world. Also, these methods rely on large numbers of registered training images and/or on textured models of the objects, which can be cumbersome to acquire.

8 3D Pose Estimation for Object Category

So far, we discussed methods that estimate the 3D pose of specific objects, which are known in advance. As we already mentioned, this requires the cumbersome step of capturing a training set for each object.

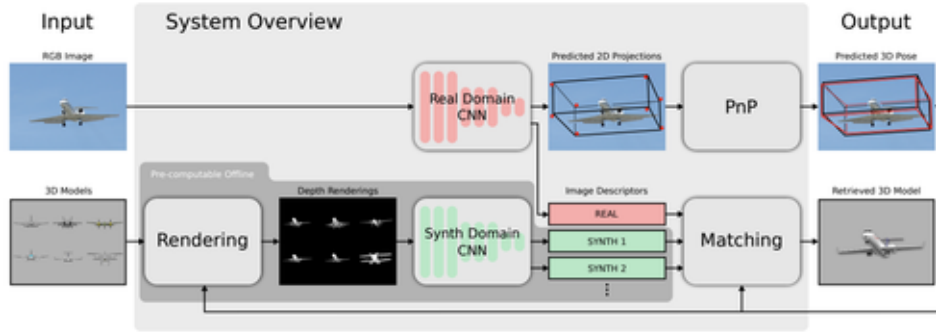


Figure 15: Estimating the 3D pose of an unknown object from a known category. [Grabner *et al.*, 2018] predicts the 2D reprojections of the corners of the 3D bounding box *and* the size of the bounding box. For this information, it is possible to compute the pose of the object using a PnP algorithm. Image from [Grabner *et al.*, 2018].

One possible direction to avoid this step is to consider a “category”-level approach, for objects that belong to a clear category, such as ‘car’ or ‘chair’. The annotation burden moves then to images of objects from the target categories, but we can then estimate the 3D pose of new, unseen objects from these categories.

Some early category-level methods only estimate 3 degrees-of-freedom (2 for the image location and 1 for the rotation over the ground plane) of the object pose using regression, classification or hybrid variants of the two. For example, [Xiang *et al.*, 2016] directly regresses azimuth, elevation and in-plane rotation using a Convolutional Neural Network (CNN), while [Tulsiani *et al.*, 2015; Tulsiani and Malik, 2015] perform viewpoint classification by discretizing the range of each angle into a number of disjoint bins and predicting the most likely bin using a CNN.

To estimate a full 3D pose, many methods rely on ‘semantic keypoints’, which can be detected in the images, and correspond to 3D points on the object. For example, to estimate the 3D pose of a car, one may one to consider the corners of the roof and the lights as semantic keypoints. [Pepik *et al.*, 2015] recovers the pose from keypoint predictions and CAD models using a PnP algorithm, and [Pavlakos *et al.*, 2017] predicts semantic keypoints and trains a deformable shape model which takes keypoint uncertainties into account. Relying on such keypoints, however, is even more depending in terms of annotations as keypoints have to be careful chosen and manually located in many images, and as they do not generalize across categories.

In fact, if one is not careful, the concept of 3D pose for object categories is ill-defined, as different objects from the same category are likely to have different sizes. To properly define this pose, [Grabner *et al.*, 2018] considers that the 3D pose of an object from a category is defined as the 3D pose of its 3D bounding box. It then proposes a method illustrated in Figure 15 that extends BB8 [Rad and Lepetit, 2017]. BB8 estimates a 3D pose by predicting the 2D reprojections of the corners of its 3D bounding box. [Grabner *et al.*, 2018] predicts similar 2D reprojections, plus the size of the 3D bounding box in the form of 3 values (length, height, width). Using these 3 values, it is possible to compute the 3D coordinates of the corners in a coordinate system related to the object. From these 3D coordinates, and from the predicted 2D reprojections, it is possible to compute the 3D pose of the 3D bounding box.

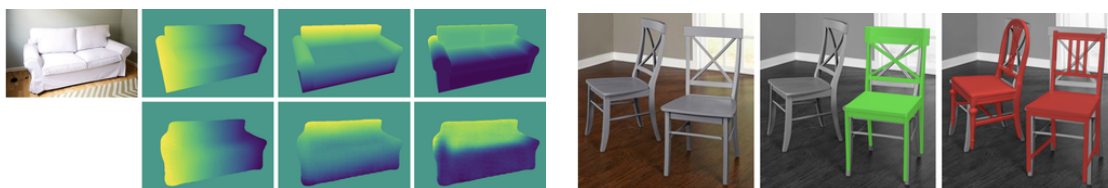


Figure 16: Recovering a 3D model for an unknown object from a known category using an image. (a) [Grabner *et al.*, 2019] generates renderings of the object coordinates (called location fields in the paper) from 3D models (top), and predicts similar object coordinates from the input images (bottom). Images of objects and 3D models are matched based on embeddings computed from the object coordinates. (b) An example of recovered 3D poses and 3D models for two chairs. The 3D models capture the general shapes of the real chairs, but do not fit perfectly.

3D model retrieval. Since objects from the same category have various 3D models, it may also be interesting to recover a 3D model for these objects, in addition to their 3D poses. Different approaches are possible. One is to recover an existing 3D model from a database that fits well the object. Large datasets of light 3D models exist for many categories, which makes this approach attractive [Chang *et al.*, 2015]. To make this approach efficient, it is best to rely on metric learning, so that an embedding computed for the object from the input image can be matched against the embeddings for the 3D models [Aubry and Russell, 2015; Izadinia *et al.*, 2017; Grabner *et al.*, 2018]. If the similarity between these embeddings can be estimated based on their Euclidean distance or their dot product, then it is possible to rely on efficient techniques to perform this match.

The challenge is that the object images and the 3D models have very different natures, while this approach needs to compute embeddings that are comparable from these two sources. The solution is then to replace the 3D models by image renderings; the embeddings for the image renderings can then be compared more easily with the embeddings for the input images. Remains the domain gap between the real input images and the synthetic image renderings. In [Grabner *et al.*, 2019], as shown in Figure 16, this is solved by predicting the object coordinates for the input image using a deep network, and by rendering the object coordinates for the synthetic images, rather than a regular rendering. This brings the representations for the input images and the 3D models even closer. From these representations, it is then easier to compute suitable embeddings.

Another approach is to predict the object 3D geometry directly from images. This approach learns a mapping from the appearance of an object to its 3D geometry. This is more appealing than recovering a 3D model from a database that has no guarantee to fit perfectly the object, since this latter method can potentially adapt better to the object’s geometry [Shin *et al.*, 2018; Groueix *et al.*, 2018; Park *et al.*, 2019a]. This is however much more challenging, and current methods are still often limited to clean images with blank background, and sometimes even synthetic renderings, however this is a very active area, and more progress can be expected in the near future.

9 3D Hand Pose Estimation from Depth Maps

While related to 3D object pose estimation, hand pose estimation has its own specificities. On one hand (no pun intended), it is more challenging, if only because more degrees-of-freedom must be estimated.

On the other hand, while considering new 3D objects still requires acquiring new data and/or specific 3D models, a single model can generalize to unseen hands, despite differences in size, shape, and skin color, as these differences remain small, and considerably smaller than differences between two arbitrary objects. In practice, this means that a model does not need new data nor retraining to adapt to new users, which is very convenient.

However, the motivation for 3D hand pose in AR applications is mostly for making possible natural interaction with virtual objects. As we will briefly discuss in Section 10.6, convincing interaction requires very high accuracy, which current developments are aiming to achieve.

We will start with 3D hand pose estimation from a depth map. As for 3D object pose estimation, the literature is extremely large, and we will focus here as well on a few representative methods. This is the direction currently taken in the industry, because it can provide a robust and accurate estimate. It is currently being deployed on commercial solutions, such as HoloLens 2 and Oculus. We will then turn to hand pose estimation from color images, which is more challenging, but has also the potential to avoid the drawbacks of using depth cameras. We will turn to pose estimation from color images in the next section.

9.1 DeepPrior++

Estimating the 3D pose of a hand from a depth map is in fact now relatively easy and robust, when using a well engineered solution based on Deep Learning. For example, [Oberweger *et al.*, 2018] discusses the implementation of a method called DeepPrior++, which is simple in principle and performs well.

DeepPrior++ first selects a bounding box on the hand, and then predicts the 3D joint locations from the depth data within this bounding box. An accurate localization of the bounding appears to improve the final accuracy, and the method from the center of mass of the depth data within a threshold, and then uses a refinement step that can be iterated. This refinement step is performed by a regression network predicting, from a 3D bounding box centered on the center of mass, the 3D location of one of the joints used as a referential. Then, the 3D joint locations can be predicted using a second Network that is not very different from the one presented in Figure 4, except that it relies on ResNet block [He *et al.*, 2016] for more accurate results. The input to this network is therefore the input depth map cropped to contain only the hand, and its output is made of the 3D coordinates of all the hand joints.

Simple data augmentation appears to improve accuracy, even when starting from an already large dataset. [Oberweger *et al.*, 2018] applies random in-plane rotations, scales, 3D offsets to the input data.

9.2 V2V-PoseNet

V2V-PoseNet [Moon *et al.*, 2018] is a method that has been identified by [Yuan *et al.*, 2017] as one of the best performing methods at the time, based on the results on a public benchmark. [Moon *et al.*, 2018] argues that using a cropped depth map as input to the network makes training difficult, as the same hand pose appears (slightly) differently in a depth map depending on where it reprojects.

To avoid distortion, V2V-PoseNet converts the depth map information in the hand bounding box into a voxel-based representation (it uses a refinement step similar to the one in [Oberweger *et al.*, 2018], so accuracy of the hand detection seems important). The conversion is done by voxelizing the 3D bounding box for the hand, and each voxel that contains at least one pixel from the depth map is set to occupied. The voxel-based representation therefore corresponds to a discretization of the hand surface.

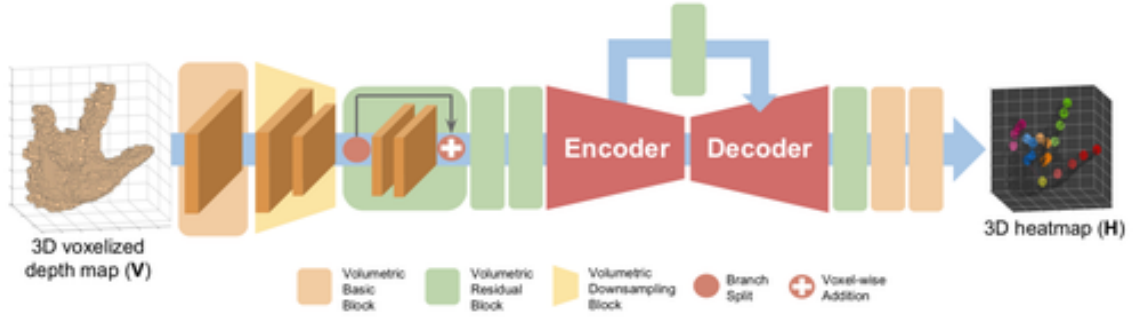


Figure 17: Architecture of V2V-PoseNet. The network takes a voxel-based representation of the depth map centered on the hand, and predicts the 3D joint locations in the form of heatmaps.

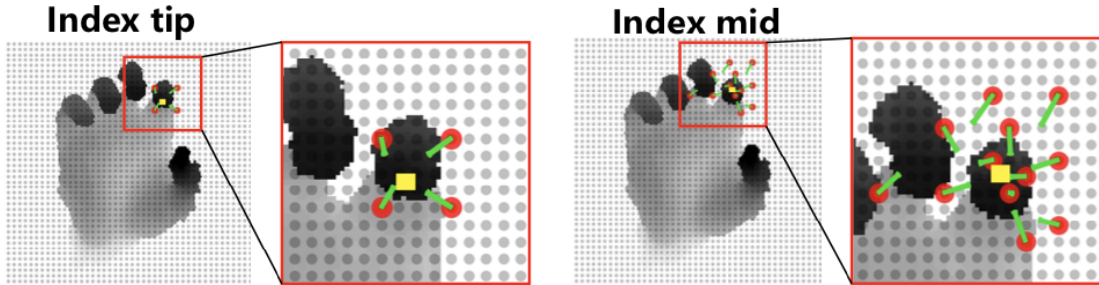


Figure 18: The A2J method predicts, for regularly sampled “anchor points” (gray dots), and for each joint, a 2D displacement towards the 2D joint location (green line segments), its depth, and the confidence for these predictions. Red dots correspond to anchor points with high confidences. The 3D joint locations can be estimated from these as weighted sums (yellow squares). Images from [Xiong *et al.*, 2019].

Then, a network is trained to take this voxel-based representation as input, and predicts the 3D joint locations in the form of 3D heatmaps. A 3D heatmap is a 3D array containing confidence values. There is one 3D heatmap for each joint, and the joint location is taken as the 3D location with the higher confidence.

9.3 A2J

The A2J method [Xiong *et al.*, 2019] performs similarly or better than V2V-PoseNet, but faster. The prediction of the 3D joint locations is performed via what are called in the paper “anchor points”. Anchor points are 2D locations regularly sampled over the bounding box of the hand, obtained as in [Moon *et al.*, 2018] and [Oberweger *et al.*, 2018]. As shown in Figure 18, a network is trained to predict, for each anchor point, and for each joint, 1) a 2D displacement from the anchor point to the joint, 2) the depth of the joint, and 3) a weight reflecting the confidence for the prediction. The prediction of 2D displacements makes this method close to [Hu *et al.*, 2019] and PVNet [Peng *et al.*, 2019], which were developed for 3D object pose estimation under occlusions.

From the output of the network, it is possible to estimate the 3D joint locations as weighted sums of the 2D locations and depths, where the weights of the sums are the ones predicted by the network. Special care is taken to make sure the weights sum to 1. The paper argues that this is similar to predictions by an “ensemble”, which are known to generalize well [Breiman, 1996]. The weights play a critical role here: Without them, the predicted locations would be a linear combination of the network output, which could be done by a standard layer without having to introduce anchor points. The product between weights and 2D displacements and depths make the predicted joint locations a more complex, non-linear transformation of the network output.

9.4 Discussion

As discussed in detail in [Armagan *et al.*, 2020], the coverage of the space of possible hand poses is critical to achieve good performance. This should not be surprising: Deep Networks essentially learn a mapping between the input (here the depth map) and the output (here the hand pose) from the training set. They can interpolate very well between samples, but poorly extrapolate: If the input contains a pose too different from any sample in the training set, the output is likely to be wrong. When using as input depth maps, which do not suffer from light changes or cluttered background like color images, having a good training set is the main important issue.

That was in fact already discussed in [Oberweger *et al.*, 2018], and even before that in [Supancic *et al.*, 2015]. Combining the DeepPrior++ that relies on a simple mapping with a simple domain adaptation method [Rad *et al.*, 2018] in order to use more synthetic data appears to perform very well, and still outperforms more complex methods on the NYU dataset, for example. While it is not entirely clear, the direction taken by HoloLens 2 and Oculus is to train a simple model on a large set of synthetic training data, and refine the pose estimate using an ICP algorithm [Sharp *et al.*, 2015] aligning a 3D model of the hand to the depth data.

10 3D Hand Pose Estimation from an RGB Image

We now turn to modern approaches to 3D hand pose estimation from an RGB image, which is significantly more challenging than from a depth map, but which is a field that has been impressively fast progress over the past few years.

10.1 [Zimmermann and Brox, 2017]

[Zimmermann and Brox, 2017] was one of the first methods predicting the 3D hand pose from a single color image, without using depth information. An overview of the method is shown in Figure 19. The hand is first segmented in 2D, as in recent approaches for 3D object pose estimation, to obtain a 2D bounding box centered on it. From this window, a 2D heatmap is predicted for each joint, using Convolutional Pose Machines [Wei *et al.*, 2016]. This provides information on the 2D joint locations. To obtain the final 3D pose, the heatmaps are lifted in 3D by using 2 networks: One network predicts the 3D joint locations in a canonical system attached to the hand; the other network predicts the 3D pose (3D rotation and translation) of the hand in the camera coordinate system. From these 2 set of information, it is possible to compute the 3D joint coordinates in the camera coordinate system. This

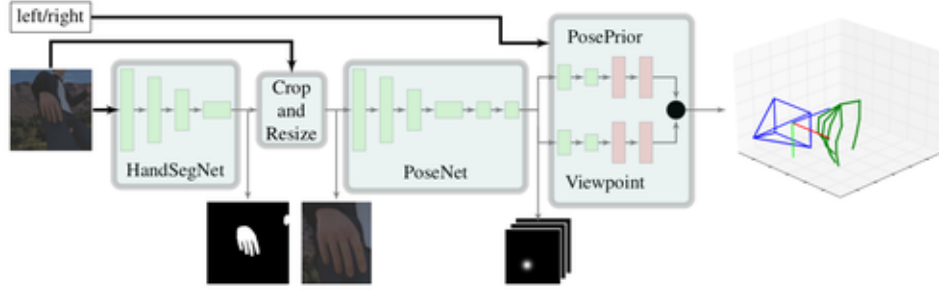


Figure 19: Overview of [Zimmermann and Brox, 2017]. The hand is first segmented in 2D to obtain a 2D bounding box. From the image information within this bounding box, a set of heatmaps for the hand joints is predicted and lifted to 3D. Image from [Zimmermann and Brox, 2017].

is shown to work slightly better than a single network directly predicting the 3D joint locations in the camera coordinate system.

To train the segmentation, heatmap prediction, and 3D pose predictions, [Zimmermann and Brox, 2017] created a dataset of synthetic images of hands, as the 3D pose can be known exactly. To do so, they used freely available 3D models of humans with corresponding animations from Mixamo, a company specialized in character animation, and Blender to render the images over random background. Lighting, viewpoint, and skin properties were also randomized. The dataset is publicly available online.

10.2 [Iqbal *et al.*, 2018]

Like [Zimmermann and Brox, 2017], [Iqbal *et al.*, 2018] relies on 2D heatmaps from a color image to lift them to 3D, but with several differences that improve the pose accuracy. They oppose two approaches to 3D pose prediction: The first approach relies on heatmaps; it is generally accurate but keeping this accuracy with 3D heatmaps is still intractable, as a finely sampled 3D volume can be very large. The second approach is to predict the 3D joints in a “holistic” way, meaning that the values of the 3D joint locations are directly predicted from the input image, as was done in [Toshev and Szegedy, 2014; Sun *et al.*, 2017] for human body pose prediction, and [Oberweger *et al.*, 2017] for hand pose prediction from depth data. Unfortunately, “holistic” approaches tend to be less accurate.

To keep the accuracy of heatmaps, and to predict the joint depths with a tractable method, [Iqbal *et al.*, 2018] also predicts, in addition to the 2D heatmaps, a depth map for each joint, where the depth values are predicted depths for the corresponding joint relative to a reference joint. This is shown in Figure 20. The final 3D pose is computed using geometric constraints from this information. Instead of learning to predict the heatmaps in a supervised way, [Iqbal *et al.*, 2018] learns to predict “latent” heatmaps. This means that the deep model is pretrained to predict heatmaps and depth maps, but the constraints on the predicted heatmaps and depth maps are removed from the full loss, and the model is trained “end-to-end”, with a loss involving only the final predicted 3D pose. This lets the optimization find heatmaps and depth maps that are easier to predict, while being more useful to predict the correct pose.

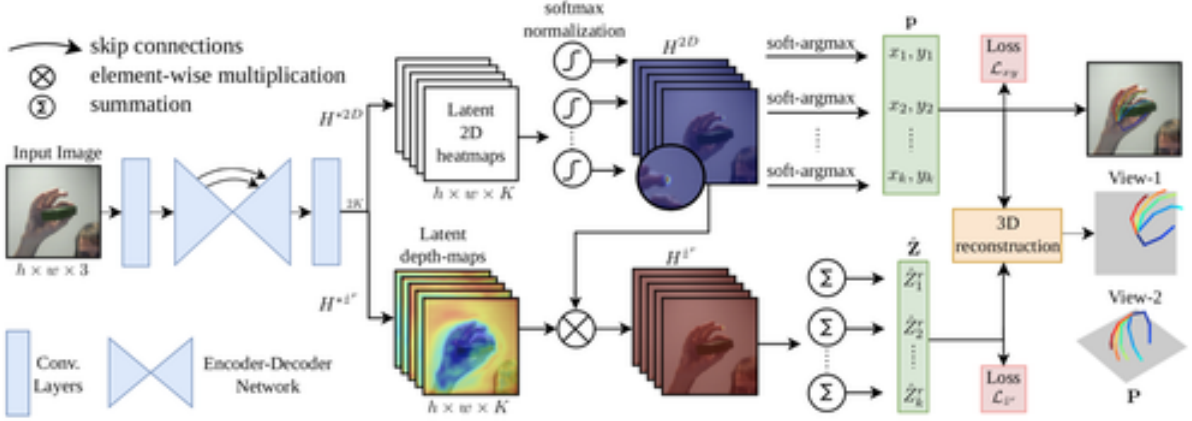


Figure 20: Architecture of [Iqbal *et al.*, 2018]. The method predicts a heatmap and a depth map for each joint, and estimates the 3D pose using geometric constraints. The model is trained end-to-end, which lets the model find optimal heatmaps and depth maps. Image from [Iqbal *et al.*, 2018].

10.3 GANerated Hands

As already mentioned above, training data and its diversity is of crucial importance in Deep Learning. For 3D hand pose estimation from color images, this means that training data need to span the hand pose space, the hand shape space, possible lighting, possible background, etc. Moreover, this data has to be annotated in 3D, which is very challenging. Generating synthetic images is thus attractive, but because of the “domain gap” between real and synthetic images, this may impact the accuracy of a model trained on synthetic images and applied to real images.

To bridge the domain gap between synthetic and real images, [Mueller *et al.*, 2018a] extends CycleGAN [Zhu *et al.*, 2017], which is itself based on Generative Adversarial Networks (GANs). In original GANs [Goodfellow *et al.*, 2014], a network is trained to generate images from random vectors. In order to ensure the generated images are “realistic”, a second network is trained jointly with the first one to distinguish the generated images from real ones. When this network fails and cannot classify the images generated by the first network as generated, it means these generated images are similar to real ones. CycleGAN [Zhu *et al.*, 2017] builds on this idea to change an image from a domain into an image with the *same content* but similar to the images from another domain. An image generator is thus train to take an image from some domain as input and to generate as output another image. A second network ensures that the output looks similar to the target domain. To make sure the content of the image is preserved, a second pair of networks is jointly trained with the first one: This pair is trained jointly with the first one, to re-generate the original input image from the image generated by the first generator.

[Mueller *et al.*, 2018a] extends CycleGAN by adding a constraint between the original image (here a synthetic image) and the image returned by the first generator. The synthetic image is a hand rendered over a uniform background, and a segmentation network is trained jointly to segment the generated image, so that the predicted segmentation matches the segmentation of the hand in the synthetic image. This helps preserving, at least to some extent, the 3D pose of the hand in the transformed image. Figure 21 shows two synthetic images transformed by CycleGAN, with and without this constraint. Another advantage is that the plain background can be changed by a random background to augment the

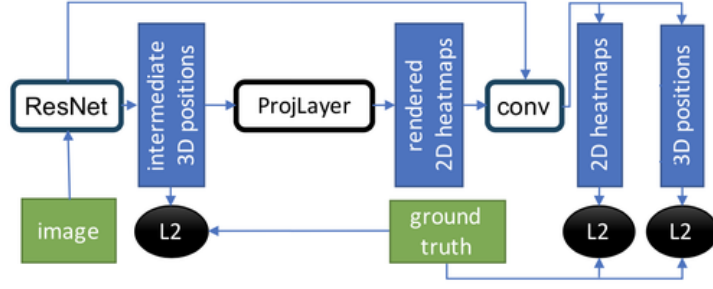


Figure 21: The architecture used in GANerated [Mueller *et al.*, 2018a] predicts both 2D heatmaps and 3D joints locations, and estimates the 3D hand pose with a (non-learned) optimization. Image from [Mueller *et al.*, 2018a].

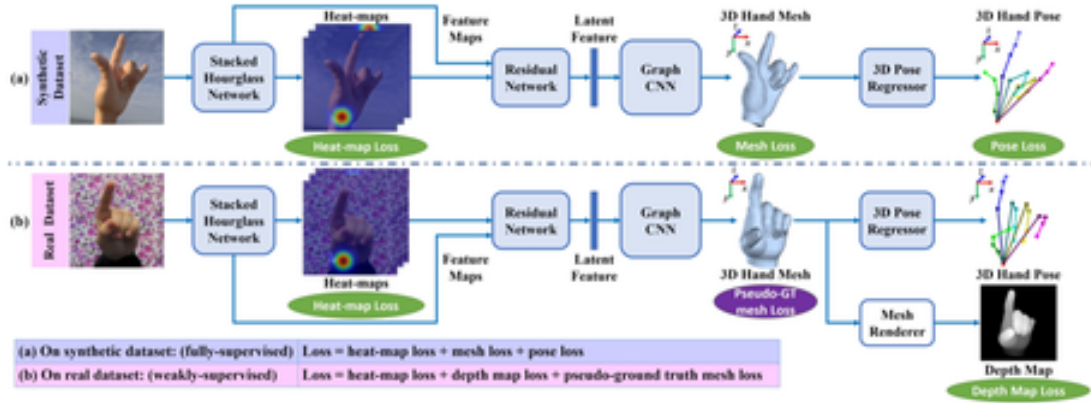


Figure 22: The architecture used in [Ge *et al.*, 2019] predicts 2D heatmaps for the joints and a 3D mesh for the hand. Image from [Ge *et al.*, 2019].

dataset of transformed images.

Besides this data generation, the 3D hand pose prediction in [Mueller *et al.*, 2018a] is also interesting. A network is trained to predict *both* 2D heatmaps for the 2D joint locations, and the 3D joint locations. Given the output of this network over a video stream, an optimization is performed to fit a kinematic model to these 2D and 3D predictions under joint angle constraints and temporal smoothness constraints. The output of the network is therefore not used directly as the 3D hand pose, but used as an observation in this final optimization.

10.4 3D Hand Shape and Pose Estimation [Ge *et al.*, 2019; Boukhayma *et al.*, 2019]

[Ge *et al.*, 2019] and [Boukhayma *et al.*, 2019] were published at the same conference (CVPR'19) and both propose a method for predicting the hand shape in addition to the hand pose.

[Ge *et al.*, 2019] first generates a training set of hands under various poses, where each image is annotated with the pose and shape parameters, with special care to make the rendered images as realistic as possible. Using this dataset, they train a network (shown in Figure 22) to predict a “latent feature

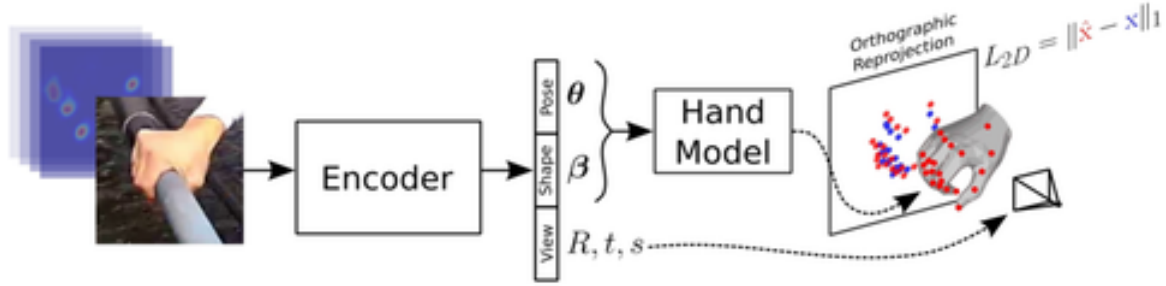


Figure 23: [Boukhayma *et al.*, 2019] learns to predict the pose and shape parameters of the hand MANO model [Romero *et al.*, 2017] without 3D annotations on the pose and shape, by minimizing on the reprojection error of the joints. Image from [Boukhayma *et al.*, 2019].

vector”, which is fed to a Graph CNN [Defferrard *et al.*, 2016]. The motivation for using a Graph CNN is to predict a 3D mesh for the hand, which can be represented as a graph. The loss function thus a loss that compares, for synthetic training images, the mesh predicted by the GraphCNN and the ground truth mesh, and the predicted 3D pose and the ground truth pose. For real training images, which do not have ground truth mesh available, they use the pose loss term, but also a term that compares the predicted 3D mesh with the depth map for the input RGB image when available, and another term that compares the same predicted 3D mesh with a mesh predicted from the ground truth 2D heatmaps.

The method of [Boukhayma *et al.*, 2019] is simpler. It relies on the MANO model [Romero *et al.*, 2017], which is a popular deformable and parameterized model of the human hand, created from many captures of real hands. The MANO model provides a differentiable function that generates a 3D model of a hand, given pose parameters and shape parameters. It does not require a dataset annotated with the shape parameters, but instead it can re-uses existing datasets, annotated only with the 3D hand pose. To do so, they train a network (shown in Figure 23) to predict the hand shape and pose parameters, and the loss function is the 2D distances between the reprojections of the 3D joints, as computed from the predicted shape and pose parameters and the annotated 2D joints.

10.5 Implementation in MediaPipe

We can also mention a Google Project [Bazarevsky and Zhang, 2019] for real-time 3D hand pose implementation on mobile devices. They first detect the palm, rather than the full hand: This allows them to only consider squared bounding boxes, and to avoid confusion between multiple hand detections. The 3D hand pose is directly predicted as the 3D joint locations, using a model trained on real and synthetic images.

Unfortunately, the work is not presented in a formal research publication, and it is difficult to compare their results with the ones obtained by other methods. However, the implementation is publicly available, and appears to work well. This may show again that good engineering can make simple approaches perform well.

10.6 Manipulating Virtual Objects

One of the main motivations to track the 3D pose of a hand for an AR system is to make possible natural interaction with virtual objects. It is therefore important to be able to exploit the estimated 3D hand pose to compute how the virtual objects should move when they are in interaction with the user's hand.

This is not a trivial problem. In the real world, our ability to interact with objects is due to the presence of friction between the surfaces of the objects and our hands, as friction is a force that resists motion. However, friction is very challenging to simulate correctly. Also, nothing prevents the real hand to penetrate the virtual objects, and make realism fail.

Some approaches rely on heuristics to perform a set of object interactions. It may be reasonable to focus on object grasping for AR systems, as it is the most interesting interaction. Some methods are data-driven and synthesize prerecorded, real hand data to identify the most similar one during runtime from a predefined database [Li *et al.*, 2007; Miller *et al.*, 2003; Rijpkema and Girard, 1991].

A friction model from physics, the *Coulomb-Contensou* model was used in [Talvas *et al.*, 2015] to reach more accurate results, but the computational complexity becomes very high, and not necessarily compatible with real-time constraints. [Hoell *et al.*, 2018] relies on the simpler Coulomb model that appears to be a good trade-off between realism and tractability. The force applied by the user are taken proportional to how much their real hand penetrates the virtual objects. Computing accurately the forces requires very accurate 3D hand poses, otherwise they can become very unstable.

Very recently, videos demonstrating the real-time interaction system of Oculus using depth-based hand tracking appeared. No technical details are available yet, but the system seems to allow for realistic grasping-based interactions.

11 3D Object+Hand Pose Estimation

We finally turn to the problem of estimating a 3D pose for both a hand and an object, when the hand directly manipulates the object. Even if the existing methods are still preliminary, the potential application to AR is very appealing, as this could offer tangible interfaces by manipulating objects, or even the possibility to bring real objects into the virtual world.

This problem still remains very challenging, as the close contact between the hand and the object results in mutual occlusions that can be large. Dealing with egocentric views, which are more relevant for AR applications, is often even more challenging. Fortunately, there are also physical constraints between the hand and the object, such as impossibility of penetration but also natural grasping poses, that may help solving this problem.

Pioneered approaches joint hand+object pose estimation [Oikonomidis *et al.*, 2011; Wang *et al.*, 2011; Ballan *et al.*, 2012] typically relied on frame-by-frame tracking, and in the case of [Kyriazis and Argyros, 2013; Tzionas *et al.*, 2016], also on a physics simulator to exploit the constraints between hand and object. However, frame-to-frame tracking, when used alone, may require careful initialization from the user, and may drift over time. The ability to estimate the 3D poses from a single image without prior from previous frames or from the user makes tracking more robust. Given the difficulty of the task, this has been possible only with the use of Deep Learning.

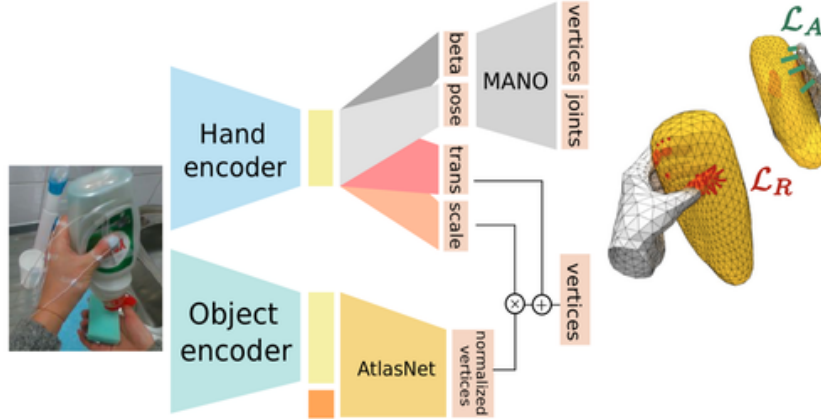


Figure 24: [Hasson *et al.*, 2019b] learns to directly predict pose and shape parameters for the hand, and shape for the object, using a term that encourages the hand and the object to be in contact without interpenetration. Image from [Hasson *et al.*, 2019b].

11.1 ObMan and HOPS-Net

ObMan [Hasson *et al.*, 2019b] and HOPS-Net [Kokic *et al.*, 2019] almost simultaneously proposed methods with some similarities for estimating the 3D pose *and shape* of both a hand and of the object it manipulates from a single RGB image. Creating realistic training data is the first challenge, and they both created a dataset of synthetic images of hands and objects using the *GraspIt!* simulation environment [Miller and Allen, 2004] for generating realistic grasps given 3D models for the hand and for the object.

[Hasson *et al.*, 2019b] train their model on their synthetic dataset before fine-tuning it on [Garcia-Hernando *et al.*, 2018], which provide annotated real images. To make the synthetic images more realistic, [Kokic *et al.*, 2019] use “Augmented CycleGAN” [Almahairi *et al.*, 2018]: The motivation to use Augmented CycleGAN rather than CycleGAN is that the former can deal with many-to-many mappings, while the latter considers only one-to-one mappings, and that a synthetic image can correspond to multiple real images.

As shown in Figure 24, [Hasson *et al.*, 2019b] does not go through an explicit 2D detection of the object or the hand, by contrast with the other methods we already discussed. Instead, the method directly predicts from the input image the pose and shape parameters for the hand (using the MANO model), and the shape for the object. To predict the object shape, they use AtlasNet [Groueix *et al.*, 2018], a method that can predict a set of 3D points from a color image. This shape prediction does not require to know the object in advance, and predicts the shape in the camera coordinates system, there is therefore no notion of object pose. The loss function to train the architecture includes a term that prevents intersection between the hand and the object, and a term that penalizes cases in which the hand is close to the surface of the object without being in contact.

[Kokic *et al.*, 2019] focuses on objects from specific categories (bottle, mug, knife, and bowl), and uses Mask R-CNN [He *et al.*, 2017] to detect and segment the object: The object does not have to be known as long as it belongs to a known category. The method also relies on [Zimmermann and Brox, 2017] to predict the 3D hand pose. A network is trained to predict, from the bounding box predicted by

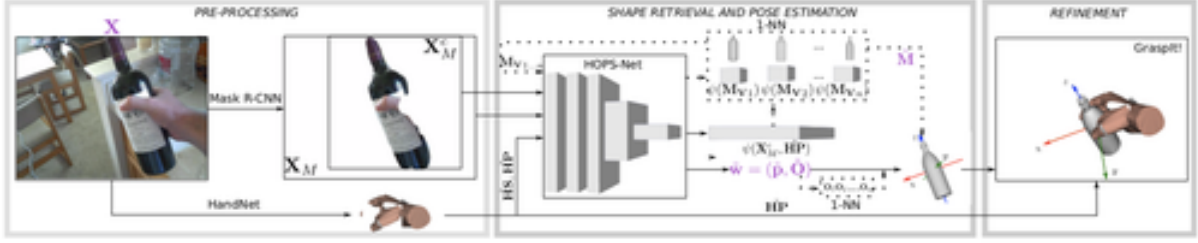


Figure 25: [Kokic *et al.*, 2019] first predicts the hand pose and detects the object in 2D, then predicts the object’s pose for its 2D bounding pose and the hand pose. It also finds a 3D model for the object from its appearance, and uses it to refine the object pose so that the object is in contact with the hand. Image from [Kokic *et al.*, 2019].

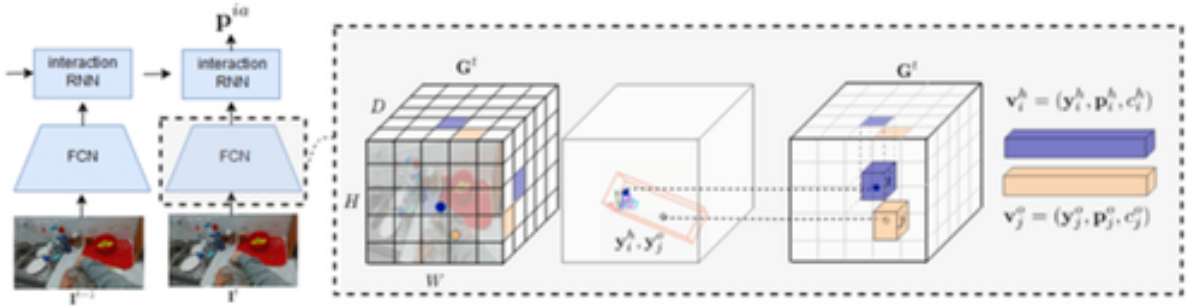


Figure 26: The H+O method [Tekin *et al.*, 2019a] predicts the hand and object poses both as a set of 3D locations, to get a unified representation and simplify the loss function. These 3D locations are predicted *via* a 3D grid. Moreover, the predicted 3D locations are provided to a recurrent neural network, to propagate this information over time and recognize the action performed by the hand. Image from [Tekin *et al.*, 2019a].

Mask-RCNN and the hand pose (to provide additional constraints), the 3D object pose and a description vector for the shape. The 3D pose representation depends on the category to handle symmetries in the object shape, and the description vector is matched against a database of 3D models to retrieve a 3D model with a similar shape (see [Grabner *et al.*, 2019] discussed in Section 8). A final refinement step optimizes the object pose so that it is in contact with the hand.

11.2 H+O [Tekin *et al.*, 2019a]

H+O [Tekin *et al.*, 2019a] is a method that can predict the 3D poses of a hand and an object from an image, but it also learns a temporal model to recognize the performed actions (for example, pouring, opening, or cleaning) and to introduce temporal constraints to improve the pose estimation. Moreover, the architecture remains relatively simple, as shown in Figure 26. It is trained on the First-Person Hand Action dataset (see Section 6.3).

The 3D pose of the object is predicted as the 3D locations of the corners of the 3D bounding box— from this information, it is possible to compute a 3D rotation and translation, using [Gower, 1975] for

example. Since the 3D hand pose can also be predicted as a set of 3D locations, this makes the 2 output of the network consistent, and avoids a weight to tune between the loss term for the hand and the loss term for the object. More exactly, the 3D locations are not predicted directly. The 3D space is split into cells, and for each cell and each 3D location, it is predicted if the 3D location is within this cell, plus an offset between a reference corner of the cell and the 3D location—this offset ensures that the 3D location can be predicted accurately, even with the space discretization into cells. A confidence is also predicted for each 3D location prediction. The 3D hand and object poses are computed by taking the 3D points with the highest confidences.

To enforce both temporal constraints and recognize the actions, the method relies on a Recurrent Neural Network, and more exactly on an Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber, 1997]. Such network can propagate information from the predicted 3D locations over time, and it is used here to predict the action and the nature of the object, in addition to be a way to learn constraints between the hand and the object.

11.3 HOnnotate [Hampali *et al.*, 2020]

[Hampali *et al.*, 2020] proposed a method to automatically annotate real images of hands grasping objects with their 3D poses (see Figure 6(f)), which works with a single RGB-D camera, but can exploit more cameras if available for better robustness and accuracy. The main idea is to optimize jointly all the 3D poses of the hand and the object over the sequence to exploit temporal consistency. The method is automated, which means that it can be used easily to labeled new sequences.

The authors use the resulting dataset called HO-3D to learn to predict from a single RGB image the 3D pose of a hand manipulating an object. This is done by training a Deep Network to predict the 2D joint locations of the hand along with the joint direction vectors and lift them to 3D by fitting a MANO model to these predictions. This reaches good accuracy despite occlusions by the object, even when the object was not seen during training.

12 The Future of 3D Object and Hand Pose Estimation

This chapter aimed at demonstrating and explaining the impressive development of 3D pose estimation in computer vision since the early pioneer works, and in the context of potential applications to Augmented Reality. Methods are becoming more robust and accurate, while benefiting on fast implementations on GPUs. We focused on some popular works, but they are only entries to many other works we could not present here, and that the reader can explore on their own.

However, as we pointed out at the end of Section 7, it is still too early to draw conclusions on what is the best methodology for 3D pose estimation. Quantitative results not only reflect the contributions of the methods, and also depend on how much effort the authors put in their implementation. As a result, there are sometimes contractory conclusions between papers. For example, is it important to first detect the object or the hand in 2D first, or not? Are 2D heatmaps important for accuracy, or not? What is the best pose representation for a 3D pose: A 3D rotation and translation, or a set of 3D points? Also, performance is not the only aspect here: For example, using a set of 3D points for both the object and hand poses as in [Tekin *et al.*, 2019b] relaxes for tuning the weights in the loss function, which is also interesting in practice.

The focus on some public benchmarks may also bias the conclusions on the performance of the methods in the real world: How well do they perform under poor lighting? How well do they perform on different cameras without retraining?

Another aspect that needs to be solved is the dependence on training sets, which are complex to create, and to a training stage each time a new object should be considered. The problem does not really occur for hands, as a large enough training set for hands would ensure generalization to unseen hands. For objects, variability in terms of shape and appearance is of course much more important. Some methods have shown some generalization power within known categories [Grabner *et al.*, 2019; Kokic *et al.*, 2019]. [Xiang *et al.*, 2018b] has shown some generalization to new objects, but only on simple synthetic images and for known 3D models. Being able to consider instantly (without retraining nor capturing a 3D model) any object in an AR system is a dream that will probably become accessible in the next years.

References

- [Almahairi *et al.*, 2018] A. Almahairi, S. Rajeswar, A. Sordoni, P. Bachman, and A. Courville. Augmented CycleGAN: Learning Many-To-Many Mappings from Unpaired Data. In *arXiv Preprint*, 2018.
- [Armagan *et al.*, 2020] A. Armagan, G. Garcia-Hernando, S. Baek, S. Hampali, M. Rad, Z. Zhang, S. Xie, M. Chen, B. Zhang, F. Xiong, Y. Xiao, Z. Cao, J. Yuan, P. Ren, W. Huang, H. Sun, M. Hruz, J. Kanis, Z. Krnoul, Q. Wan, S. Li, L. Yang, D. Lee, A. Yao, W. Zhou, S. Mei, Y. Liu, A. Spurr, U. Iqbal, P. Molchanov, P. Weinzaepfel, R. Brégier, G. Rogez, V. Lepetit, and T.-K. Kim. Measuring Generalisation to Unseen Viewpoints, Articulations, Shapes and Objects for 3D Hand Pose Estimation Under Hand-Object Interaction. In *arXiv Preprint*, 2020.
- [Aubry and Russell, 2015] M. Aubry and B. Russell. Understanding Deep Features with Computer-Generated Imagery. In *Conference on Computer Vision and Pattern Recognition*, pages 2875–2883, 2015.
- [Ballan *et al.*, 2012] L. Ballan, A. Taneja, J. Gall, L. Van Gool, and M. Pollefeys. Motion Capture of Hands in Action Using Discriminative Salient Points. In *European Conference on Computer Vision*, pages 640–653, October 2012.
- [Bay *et al.*, 2008] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool. SURF: Speeded Up Robust Features. *Computer Vision and Image Understanding*, 10(3):346–359, 2008.
- [Bazarevsky and Zhang, 2019] V. Bazarevsky and F. Zhang. Google Project on Hand Pose Prediction, 2019.
- [Billinghurst *et al.*, 2001] M. Billinghurst, H. Kato, and I. Poupyrev. The Magicbook: A Transitional AR Interface. *Computers & Graphics*, 25:745–753, 10 2001.
- [Boukhayma *et al.*, 2019] A. Boukhayma, R. de Bem, and P. H.S. Torr. 3D Hand Shape and Pose from Images in the Wild. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Brachmann *et al.*, 2014] 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*, 2014.

- [Brachmann *et al.*, 2016] E. Brachmann, F. Michel, A. Krull, M. M. Yang, S. Gumhold, and C. Rother. Uncertainty-Driven 6D Pose Estimation of Objects and Scenes from a Single RGB Image. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [Brégier *et al.*, 2018] R. Brégier, F. Devernay, L. Leyrit, and J.L. Crowley. Defining the Pose of Any 3D Rigid Object and an Associated Distance. *International Journal of Computer Vision*, 2018.
- [Breiman, 1996] L. Breiman. Bagging Predictors. *Machine Learning*, 24(2):123–140, 1996.
- [Calli *et al.*, 2017] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dollar. Yale-Cmu-Berkeley Dataset for Robotic Manipulation Research. In *International Journal of Robotics Research*, 2017.
- [Chang *et al.*, 2015] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. In *arXiv Preprint*, 2015.
- [Defferrard *et al.*, 2016] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *Advances in Neural Information Processing Systems*, pages 3844–3852, 2016.
- [Drummond and Cipolla, 2002] T. Drummond and R. Cipolla. Real-Time Visual Tracking of Complex Structures. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(7):932–946, July 2002.
- [Eldan and Shamir, 2016] R. Eldan and O. Shamir. The Power of Depth for Feedforward Neural Networks. In *Conference on Learning Theory*, 2016.
- [Gao *et al.*, 2003] X. Gao, X. Hou, J. Tang, and H. Cheng. Complete Solution Classification for the Perspective-Three-Point Problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(8):930–943, 2003.
- [Garcia-Hernando *et al.*, 2018] G. Garcia-Hernando, S. Yuan, S. Baek, and T.-K. Kim. First-Person Hand Action Benchmark with RGB-D Videos and 3D Hand Pose Annotations. In *Conference on Computer Vision and Pattern Recognition*, pages 409–419, 2018.
- [Ge *et al.*, 2019] L. Ge, Z. Ren, Y. Li, Z. Xue, Y. Wang, J. Cai, and J. Yuan. 3D Hand Shape and Pose Estimation from a Single RGB Image. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Goodfellow *et al.*, 2014] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and A. Y. Bengio. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, 2014.
- [Goodfellow *et al.*, 2016] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.
- [Gower, 1975] J. C. Gower. Generalized Procrustes Analysis. *Psychometrika*, 40(1):33–51, 1975.
- [Grabner *et al.*, 2018] A. Grabner, P. M. Roth, and V. Lepetit. 3D Pose Estimation and 3D Model Retrieval for Objects in the Wild. In *Conference on Computer Vision and Pattern Recognition*, 2018.

- [Grabner *et al.*, 2019] A. Grabner, P. M. Roth, and V. Lepetit. Location Field Descriptors: Single Image 3D Model Retrieval in the Wild. In *International Conference on 3D Vision*, 2019.
- [Groueix *et al.*, 2018] T. Groueix, M. Fisher, V. Kim, B. Russell, and M. Aubry. Atlasnet: A Papier-Mâché Approach to Learning 3D Surface Generation. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Gustus *et al.*, 2012] A. Gustus, G. Stillfried, J. Visser, H. Jörntell, and P. van der Smagt. Human Hand Modelling: Kinematics, Dynamics, Applications. *Biological Cybernetics*, 106(11):741–755, 2012.
- [Hampali *et al.*, 2020] S. Hampali, M. Rad, M. Oberweger, and A. V. Lepetit. HOnnotate: A Method for 3D Annotation of Hand and Object Poses. In *Conference on Computer Vision and Pattern Recognition*, 2020.
- [Harris and Stennett, 1990] C. Harris and C. Stennett. RAPiD-A Video Rate Object Tracker. In *British Machine Vision Conference*, 1990.
- [Harris and Stephens, 1988] C.G. Harris and M.J. Stephens. A Combined Corner and Edge Detector. In *Fourth Alvey Vision Conference*, 1988.
- [Hasson *et al.*, 2019a] Y. Hasson, G. Varol, D. Tzionas, I. Kalevatykh, M. J. Black, I. Laptev, and C. Schmid. Learning Joint Reconstruction of Hands and Manipulated Objects. In *Conference on Computer Vision and Pattern Recognition*, pages 11807–11816, 2019.
- [Hasson *et al.*, 2019b] Y. Hasson, G. Varol, D. Tzionas, I. Kalevatykh, M. J. Black, I. Laptev, and C. Schmid. Learning Joint Reconstruction of Hands and Manipulated Objects. In *Conference on Computer Vision and Pattern Recognition*, pages 11807–11816, 2019.
- [He *et al.*, 2016] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [He *et al.*, 2017] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *International Conference on Computer Vision*, 2017.
- [Hinterstoisser *et al.*, 2012a] S. Hinterstoisser, C. Cagniart, S. Ilic, P. F. Sturm, N. Navab, P. Fua, and V. Lepetit. Gradient Response Maps for Real-Time Detection of Textureless Objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34, May 2012.
- [Hinterstoisser *et al.*, 2012b] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. R. 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*, 2012.
- [Hochreiter and Schmidhuber, 1997] S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [Hodan *et al.*, 2016] T. Hodan, P. Haluza, S. Obdrzalek, J. Matas, M. Lourakis, and X. Zabulis. On Evaluation of 6D Object Pose Estimation. In *European Conference on Computer Vision*, 2016.
- [Hoell *et al.*, 2018] M. Hoell, M. Oberweger, C. Arth, and A. V. Lepetit. Efficient Physics-Based Implementation for Realistic Hand-Object Interaction in Virtual Reality. In *IEEE Conference on Virtual Reality*, 2018.

- [Hornik *et al.*, 1989] K. Hornik, M. Stinchcombe, and H. White. Multilayer Feedforward Networks Are Universal Approximators. *Neural Networks*, 1989.
- [Hu *et al.*, 2019] Y. Hu, J. Hugonot, P. Fua, and M. Salzmann. Segmentation-Driven 6D Object Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Iqbal *et al.*, 2018] U. Iqbal, P. Molchanov, T. Breuel, J. Gall, and J. Kautz. Hand Pose Estimation via Latent 2.5D Heatmap Regression. In *European Conference on Computer Vision*, 2018.
- [Izadinia *et al.*, 2017] H. Izadinia, Q. Shan, and S. Seitz. IM2CAD. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [Kehl *et al.*, 2017] 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 *International Conference on Computer Vision*, 2017.
- [Kim *et al.*, 2010] K. Kim, M. Grundmann, A. Shamir, I. Matthews, J. Hodgins, and I. Essa. Motion Fields to Predict Play Evolution in Dynamic Sport Scenes. In *Conference on Computer Vision and Pattern Recognition*, 2010.
- [Kokic *et al.*, 2019] M. Kokic, D. Kragic, and J. Bohg. Learning to Estimate Pose and Shape of Hand-Held Objects from RGB Images. In *International Conference on Intelligent Robots and Systems*, pages 3980–3987, 11 2019.
- [Kyriazis and Argyros, 2013] N. Kyriazis and A. A. Argyros. Physically Plausible 3D Scene Tracking: The Single Actor Hypothesis. In *Conference on Computer Vision and Pattern Recognition*, 2013.
- [Lepetit *et al.*, 2005] V. Lepetit, P. Laguerre, and P. Fua. Randomized Trees for Real-Time Keypoint Recognition. In *Conference on Computer Vision and Pattern Recognition*, June 2005.
- [Li *et al.*, 2007] Y. Li, J. L. Fu, and N. S. Pollard. Data-Driven Grasp Synthesis Using Shape Matching and Task-Based Pruning. *IEEE Transactions on Visualization and Computer Graphics*, 13(4):732–747, July 2007.
- [Li *et al.*, 2018] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox. DeepIM: Deep Iterative Matching for 6D Pose Estimation. In *European Conference on Computer Vision*, 2018.
- [Lin *et al.*, 2014] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C.L. Zitnick. Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision*, pages 740–755, 2014.
- [Liu *et al.*, 2016] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S.E. Reed, C.-Y. Fu, and A.C. Berg. SSD: Single Shot Multibox Detector. In *European Conference on Computer Vision*, 2016.
- [Lowe, 1991] D. G. Lowe. Fitting Parameterized Three-Dimensional Models to Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(5):441–450, June 1991.
- [Lowe, 2001] D.G. Lowe. Local Feature View Clustering for 3D Object Recognition. In *Conference on Computer Vision and Pattern Recognition*, pages 682–688, 2001.
- [Lowe, 2004] D.G. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, 20(2), 2004.

- [Miller and Allen, 2004] A.T. Miller and P.K. Allen. Graspit! a Versatile Simulator for Robotic Grasping. *IEEE Robotics & Automation Magazine*, 2004.
- [Miller *et al.*, 2003] A. T. Miller, S. Knoop, H. I. Christensen, and P. K. Allen. Automatic Grasp Planning Using Shape Primitives. In *Proc. ICRA*, pages 1824–1829, 2003.
- [Moon *et al.*, 2018] G. Moon, J. Y. Chang, and K. M. Lee. V2v-Posenet: Voxel-To-Voxel Prediction Network for Accurate 3D Hand and Human Pose Estimation from a Single Depth Map. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Mueller *et al.*, 2018a] F. Mueller, F. Bernard, O. Sotnychenko, D. Mehta, S. Sridhar, D. Casas, and A. C. Theobalt. GANerated Hands for Real-Time 3D Hand Tracking from Monocular RGB. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Mueller *et al.*, 2018b] F. Mueller, F. Bernard, O. Sotnychenko, D. Mehta, S. Sridhar, D. Casas, and C. Theobalt. GANerated Hands for Real-Time 3D Hand Tracking from Monocular RGB. In *Conference on Computer Vision and Pattern Recognition*, pages 49–59, 2018.
- [Oberweger *et al.*, 2015] M. Oberweger, P. Wohlhart, and V. Lepetit. Training a Feedback Loop for Hand Pose Estimation. In *International Conference on Computer Vision*, 2015.
- [Oberweger *et al.*, 2017] M. Oberweger, P. Wohlhart, and V. Lepeti. DeepPrior++: Improving Fast and Accurate 3D Hand Pose Estimation. In *International Conference on Computer Vision*, page 2, 2017.
- [Oberweger *et al.*, 2018] M. Oberweger, M. Rad, and V. Lepetit. Making Deep Heatmaps Robust to Partial Occlusions for 3D Object Pose Estimation. In *European Conference on Computer Vision*, 2018.
- [Oikonomidis *et al.*, 2011] I. Oikonomidis, N. Kyriazis, and A. A. Argyros. Full DoF Tracking of a Hand Interacting with an Object by Modeling Occlusions and Physical Constraints. In *International Conference on Computer Vision*, pages 2088–2095, 2011.
- [Park *et al.*, 2019a] J. J. Park, P. Florence, J. Straub, R. Newcombe, and A. S. Lovegrove. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Park *et al.*, 2019b] K. Park, T. Patten, and M. Vincze. Pix2pose: Pixel-Wise Coordinate Regression of Objects for 6D Pose Estimation. In *International Conference on Computer Vision*, 2019.
- [Pavlakos *et al.*, 2017] G. Pavlakos, X. Zhou, A. Chan, K. Derpanis, and K. Daniilidis. 6-DoF Object Pose from Semantic Keypoints. In *International Conference on Robotics and Automation*, pages 2011–2018, 2017.
- [Peng *et al.*, 2019] S. Peng, Y. Liu, Q. Huang, H. Bao, and X. Zhou. Pvnnet: Pixel-Wise Voting Network for 6DoF Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Pepik *et al.*, 2015] B. Pepik, M. Stark, P. Gehler, T. Ritschel, and B. Schiele. 3D Object Class Detection in the Wild. In *Conference on Computer Vision and Pattern Recognition*, 2015.
- [Pinkus, 1999] A. Pinkus. Approximation Theory of the MLP Model in Neural Networks. *Acta Numerica*, 1999.

- [Poirson *et al.*, 2016] P. Poirson, P. Ammirato, C.-Y. Fu, W. Liu, J. Kosecka, and A. C. Berg. Fast Single Shot Detection and Pose Estimation. In *International Conference on 3D Vision*, 2016.
- [Pumarola *et al.*, 2018] A. Pumarola, A. Agudo, L. Porzi, A. Sanfeliu, V. Lepetit, and F. Moreno-Noguer. Geometry-Aware Network for Non-Rigid Shape Prediction from a Single View. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Rad and Lepetit, 2017] 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. In *International Conference on Computer Vision*, 2017.
- [Rad *et al.*, 2018] M. Rad, M. Oberweger, and V. Lepetit. Feature Mapping for Learning Fast and Accurate 3D Pose Inference from Synthetic Images. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Redmon and Farhadi, 2017] J. Redmon and A. Farhadi. Yolo9000: Better, Faster, Stronger. In *arXiv Preprint*, 2017.
- [Redmon *et al.*, 2016] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You Only Look Once: Unified, Real-Time Object Detection. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [Rijpkema and Girard, 1991] H. Rijpkema and M. Girard. Computer Animation of Knowledge-Based Human Grasping. *SIGGRAPH Computer Graphics*, 25(4):339–348, July 1991.
- [Romero *et al.*, 2017] J. Romero, D. Tzionas, and M. J. Black. Embodied Hands: Modeling and Capturing Hands and Bodies Together. *ACM Transactions on Graphics (TOG)*, 2017.
- [Rublee *et al.*, 2011] E. Rublee, V. Rabaud, K. Konolidge, and G. Bradski. ORB: An Efficient Alternative to SIFT or SURF. In *International Conference on Computer Vision*, 2011.
- [Russakovsky *et al.*, 2015] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A.C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [Salzmann and Fua, 2010] M. Salzmann and P. Fua. *Deformable Surface 3D Reconstruction from Monocular Images*. Morgan-Claypool, 2010.
- [Sharp *et al.*, 2015] T. Sharp, C. Keskin, D. P. Robertson, J. Taylor, J. Shotton, D. Kim, C. Rhemann, I. Leichter, A. Vinnikov, Y. Wei, D. Freedman, P. Kohli, E. Krupka, A. W. Fitzgibbon, and A. S. Izadi. Accurate, Robust, and Flexible Real-Time Hand Tracking. In *ACM Conference on Human Factors in Computing Systems*, 2015.
- [Shin *et al.*, 2018] D. Shin, C. C. Fowlkes, and D. Hoiem. Pixels, Voxels, and Views: A Study of Shape Representations for Single View 3D Object Shape Prediction. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Shotton *et al.*, 2013] J. Shotton, R. Girshick, A. Fitzgibbon, T. Sharp, M. Cook, M. Finocchio, R. Moore, P. Kohli, A. Criminisi, A. Kipman, and A. Blake. Efficient Human Pose Estimation from Single Depth Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2821–2840, 2013.

- [Simon *et al.*, 2017] T. Simon, H. Joo, I. A. Matthews, and A. Y. Sheikh. Hand Keypoint Detection in Single Images Using Multiview Bootstrapping. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [Sun *et al.*, 2015] X. Sun, Y. Wei, S. Liang, X. Tang, and J. Sun. Cascaded Hand Pose Regression. In *Conference on Computer Vision and Pattern Recognition*, 2015.
- [Sun *et al.*, 2017] X. Sun, J. Shang, S. Liang, and Y. Wei. Compositional Human Pose Regression. In *International Conference on Computer Vision*, 2017.
- [Sundermeyer *et al.*, 2019] M. Sundermeyer, Z.-C. Marton, M. Durner, and R. Triebel. Augmented Autoencoders: Implicit 3D Orientation Learning for 6D Object Detection. *International Journal of Computer Vision*, 2019.
- [Supancic *et al.*, 2015] J. S. Supancic, G. Rogez, Y. Yang, J. Shotton, and D. Ramanan. Depth-Based Hand Pose Estimation: Data, Methods, and Challenges. In *International Conference on Computer Vision*, 2015.
- [Talvas *et al.*, 2015] A. Talvas, M. Marchal, C. Duriez, and M. A. Otaduy. Aggregate Constraints for Virtual Manipulation with Soft Fingers. *IEEE Transactions on Visualization and Computer Graphics*, 21(4):452–461, April 2015.
- [Tang *et al.*, 2014] D. Tang, H. J. Chang, A. Tejani, , and T.-K. Kim. Latent Regression Forest: Structured Estimation of 3D Articulated Hand Posture. In *Conference on Computer Vision and Pattern Recognition*, 2014.
- [Taylor *et al.*, 2012] J. Taylor, J. Shotton, T. Sharp, and A. Fitzgibbon. The Vitruvian Manifold: Inferring Dense Correspondences for One-Shot Human Pose Estimation. In *Conference on Computer Vision and Pattern Recognition*, pages 103–110, 2012.
- [Tekin *et al.*, 2018] B. Tekin, S. N. Sinha, and P. Fua. Real-Time Seamless Single Shot 6D Object Pose Prediction. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- [Tekin *et al.*, 2019a] B. Tekin, F. Bogo, and M. Pollefeys. H+o: Unified Egocentric Recognition of 3D Hand-Object Poses and Interactions. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Tekin *et al.*, 2019b] B. Tekin, F. Bogo, and M. Pollefeys. H+o: Unified Egocentric Recognition of 3D Hand-Object Poses and Interactions. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [Tolani *et al.*, 2000] D. Tolani, A. Goswami, and N. I. Badler. Real-Time Inverse Kinematics Techniques for Anthropomorphic Limbs. *Graphical Models*, 62(5):353–388, 2000.
- [Tompson *et al.*, 2014] J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation. In *Advances in Neural Information Processing Systems*, 2014.
- [Toshev and Szegedy, 2014] A. Toshev and C. Szegedy. DeepPose: Human Pose Estimation via Deep Neural Networks. In *Conference on Computer Vision and Pattern Recognition*, 2014.

- [Tulsiani and Malik, 2015] S. Tulsiani and J. Malik. Viewpoints and Keypoints. In *Conference on Computer Vision and Pattern Recognition*, pages 1510–1519, 2015.
- [Tulsiani *et al.*, 2015] S. Tulsiani, J. Carreira, and J. Malik. Pose Induction for Novel Object Categories. In *International Conference on Computer Vision*, pages 64–72, 2015.
- [Tzionas *et al.*, 2016] D. Tzionas, L. Ballan, A. Srikantha, P. Aponte, M. Pollefeys, and J. Gall. Capturing Hands in Action Using Discriminative Salient Points and Physics Simulation. In *International Journal of Computer Vision*, 2016.
- [Vacchetti *et al.*, 2004] L. Vacchetti, V. Lepetit, and P. Fua. Stable Real-Time 3D Tracking Using On-line and Offline Information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(10), October 2004.
- [Wang *et al.*, 2011] R. Wang, S. Paris, and J. Popović. 6D Hands: Markerless Hand-Tracking for Computer Aided Design. In *ACM Symposium on User Interface Software and Technology*, pages 549–558, 2011.
- [Wei *et al.*, 2016] S. E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional Pose Machines. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [Xiang *et al.*, 2016] Y. Xiang, W. Kim, W. Chen, J. Ji, C. Choy, H. Su, R. Mottaghi, L. Guibas, and S. Savarese. ObjectNet3D: A Large Scale Database for 3D Object Recognition. In *European Conference on Computer Vision*, pages 160–176, 2016.
- [Xiang *et al.*, 2018a] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox. Posecnn: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes. In *Robotics: Science and Systems Conference*, 2018.
- [Xiang *et al.*, 2018b] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox. PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes. *Robotics: Science and Systems Conference*, 2018.
- [Xiong *et al.*, 2019] F. Xiong, B. Zhanga, Y. Xiao, Z. Cao, T. Yu, J. T. Zhou, and J. Yuan. A2J: Anchor-To-Joint Regression Network for 3D Articulated Pose Estimation from a Single Depth Image. In *International Conference on Computer Vision*, 2019.
- [Yuan *et al.*, 2017] S. Yuan, Q. Ye, B. Stenger, S. Jain, and T.-K. Kim. BigHand2.2M Benchmark: Hand Pose Dataset and State of the Art Analysis. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [Zakharov *et al.*, 2019] S. Zakharov, I. Shugurov, and S. Ilic. DPOD: 6D Pose Object Detector and Refiner. In *International Conference on Computer Vision*, 2019.
- [Zhu *et al.*, 2017] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired Image-To-Image Translation Using Cycle-Consistent Adversarial Networks. In *International Conference on Computer Vision*, 2017.
- [Zimmermann and Brox, 2017] C. Zimmermann and T. Brox. Learning to Estimate 3D Hand Pose from Single RGB Images. In *International Conference on Computer Vision*, 2017.

[Zimmermann *et al.*, 2019] C. Zimmermann, D. Ceylan, J. Yang, B. C. Russell, M. J. Argus, and T. Brox. FreiHAND: A Dataset for Markerless Capture of Hand Pose and Shape from Single RGB Images. In *International Conference on Computer Vision*, pages 813–822, 2019.