# RGB-based 3D Hand Pose Estimation
# via Privileged Learning with Depth Images

Shanxin Yuan
Imperial College London
shanxinyuan@gmail.com

Björn Stenger
Rakuten Institute of Technology
bjorn@cantab.net

Tae-Kyun Kim
Imperial College London
tk.kim@imperial.ac.uk

arXiv:1811.07376v1 [cs.CV] 18 Nov 2018

## Abstract

*This paper proposes a method for hand pose estimation from RGB images that uses both external large-scale depth image datasets and paired depth and RGB images as privileged information at training time. We show that providing depth information during training significantly improves performance of pose estimation from RGB images during testing. We explore different ways of using this privileged information: (1) using depth data to initially train a depth-based network, (2) using the features from the depth-based network of the paired depth images to constrain mid-level RGB network weights, and (3) using the foreground mask, obtained from the depth data, to suppress the responses from the background area. By using paired RGB and depth images, we are able to supervise the RGB-based network to learn middle layer features that mimic that of the corresponding depth-based network, which is trained on large-scale, accurately annotated depth data. During testing, when only an RGB image is available, our method produces accurate 3D hand pose predictions. Our method is also tested on 2D hand pose estimation. Experiments on three public datasets show that the method outperforms the state-of-the-art methods for hand pose estimation using RGB image input.*

## 1. Introduction

3D hand pose estimation has been greatly improving in the past few years, especially with the availability of depth cameras. While new methods [18, 45, 8, 40, 34] and datasets [32, 36, 30, 47, 7] have been published, state-of-the-art methods are still lacking in accuracy required for fine manipulations for AR or VR systems. There is a large accuracy gap between pose estimation from RGB and depth image input, which several recent works have aimed to narrow [27, 52, 17, 20]. One of the difficulties has been the lack of large-scale realistic RGB datasets with accurate annotations. Recent papers have addressed this issue by creating synthetic datasets [52], or employing GANs to generate training data [16]. In this paper we propose using depth data as *privileged information* during training. Fully annotated depth datasets [32, 36, 30, 47, 7] are abundant in the literature, but so far no attempt has been made to use this data to support the task of 3D hand pose estimation from RGB images. There are also a few RGB-D datasets proposed recently [52, 48] to tackle the problem of 3D hand pose estimation from RGB images, however all existing methods [52, 16, 48] utilise only RGB images for training. The available depth images, either paired with RGB images [48, 52] or alone in the large-scale *BigHand2.2M* dataset [47] could be used to aid the training.

The use of privileged information in training [39], also called training with hidden information [42], or side information [44], has been shown to improve performance in other domains, such as image classification [4], object detection [12], and action recognition [25]. But the concept of using privileged information to help 3D hand pose estimation from RGB images has not been attempted. To the best of our knowledge, this paper proposes the first solution. Existing methods for 3D hand pose estimation from RGB images pursue two main directions: (1) using only RGB images for 3D hand pose estimation [48, 52, 16], with different CNN models being proposed. Given the limited size of real RGB datasets, a large number of synthetic images [52, 16] are created to help the training, whether they are purely synthetic [52], or using CycleGAN [51] to enforce a certain realism [16]. (2) Using RGB-D images for 3D hand pose tracking [17], where the input is the depth channel in addition to the RGB channels. This works well when the paired RGB and depth images are available at test time. The lack of large-scale annotated training data limits the success of this approach. Our study proposes a new framework for 3D hand pose estimation from RGB images, by using the existing abundant fully annotated depth data in training, as privileged information. This helps improve 3D hand pose estimation using a single RGB image input at test time. Our method transfers supervision from depth images to RGB images. We use two networks, an RGB-based

network and a depth-based network, see Figure 1. We explore different ways to use depth data: (1) initially, we treat a large amount of independent external depth training data as privileged information to train the depth-based network. (2) After the initial training is completed, paired RGB and depth images are used to tune the RGB-based network and the depth-based network. The idea is to let the middle layer activations of the RGB network mimic that of the depth network. (3) We also explore the use of foreground hand masks to suppress background area activations in the middle layers of the RGB network. By doing this, we force the RGB network to extract features only from the foreground area.

Compared to existing methods for 3D hand pose estimation by RGB images, our main contributions are:

- To the best of our knowledge, this paper is the first to introduce the concept of using privileged information (depth images) to help the training of a RGB-based hand pose estimator.

- We propose three ways to use the privileged information: as external training data for a depth-based network, as paired depth data to transfer supervision from the depth-based network to the RGB-based network, as hand masks to suppress the background activations in the RGB-based network.

- Our training strategy can be easily embedded into existing pose estimation methods. We demonstrate this in the experiments of 2D hand pose estimation with an RGB image input by a different CNN model. Results on 2D hand pose estimation, using our training strategy improve over state-of-the-art methods for 2D hand pose estimation with RGB input.

Comprehensive experiments are conducted on three datasets: the Stereo dataset [48], the RHD dataset [52], and the Dexter-Object dataset [29]. The Stereo dataset and RHD dataset are used for evaluating 3D pose estimation from an RGB input. All three datasets are used for evaluating 2D hand pose estimation from a single RGB image.

## 2. Related Work

**3D hand pose estimation.** Hand pose estimation from depth data has made rapid progress in the past years [18, 8, 40, 24, 5], where comprehensive studies [6, 31, 46] have been instrumental in advancing the field. Random forests [32, 33, 41] and CNNs [45, 8, 40, 36] trained on large-scale public depth image datasets [32, 36, 30, 47, 7] have shown good performance. A recent benchmark evaluation [46] showed that modern methods achieve mean 3D joint position errors of less than 10mm. Hand pose estimation from RGB images is significantly more challenging [27, 52, 17, 20].

Due to the difficulty in capturing real RGB datasets with accurate 3D annotations, recent methods employ synthetic CG data [52], or *GANerated* images [16], which are more realistic synthetic images created with a CycleGAN [51]. Mueller *et al.* [16] use an image-to-image translation network to create a large amount of RGB training images and combine a CNN with a kinematic 3D hand model for pose estimation. The method requires a predefined hand model, adapted for each user. Simon *et al.*'s *OpenPose* [27] system generates an annotated RGB dataset using a panoptic studio setup, using multiple views to bootstrap 2D hand pose estimation. Zimmermann and Brox [52] proposed combining hand segmentation and 2D hand pose estimation (using *CPM* [43]), followed by estimating 3D hand pose relative to a canonical pose. Panteleris and Argyros [20] estimate absolute 3D hand pose by first estimating 2D hand pose and then optimizing a 3D hand model with inverse kinematics. Note that there also exists a large body of work on the related task of recovering full 3D human body pose from images. One line of work aims to directly estimate the 3D pose from images [14, 49, 38]. A second approach is to first estimate 2D pose, often in terms of joint locations, and then lift this to 3D pose. 2D key points can be reliably estimated using CNNs and 3D pose is estimated using structured learning or a kinematic model [35, 37, 26, 50].

**Learning with privileged information and transfer learning.** Privileged information denotes training data that is available only during training but not at test time. The concept to provide teacher-like supervision at training time was introduced by Vapnik and Vashist [39]. The idea has proven useful in other domains [4, 12, 25]. Shi *et al.* [25] treated skeleton data as privileged information in CNN-RNNs for action recognition from depth sequences. Chen *et al.* [4] manually annotated object masks in 10% of the training data and treated these as privileged information for image classification. The idea is related to network compression and mimic learning proposed by Ba and Caruana [1] as well as network distillation by Hinton *et al.* [11], where intermediate layer outputs of one network are approximated by another, possibly smaller, network. These techniques can be used to significantly reduce the number of model parameters without a significant drop in accuracy. In our case, the application target is similar to transfer learning and domain adaptation. Information from one task, prediction from depth images, is shared with another, prediction from RGB images. In transfer learning and domain adaptation information is shared across different data modalities [23, 3, 12]. Chen *et al.* [3] proposed recognition in RGB images by learning from labeled RGB-D data. A common feature representation is learned across two feature modalities. Hoffman *et al.* [12] learned an additional *hallucination* representation, which is informed by the depth data in train-
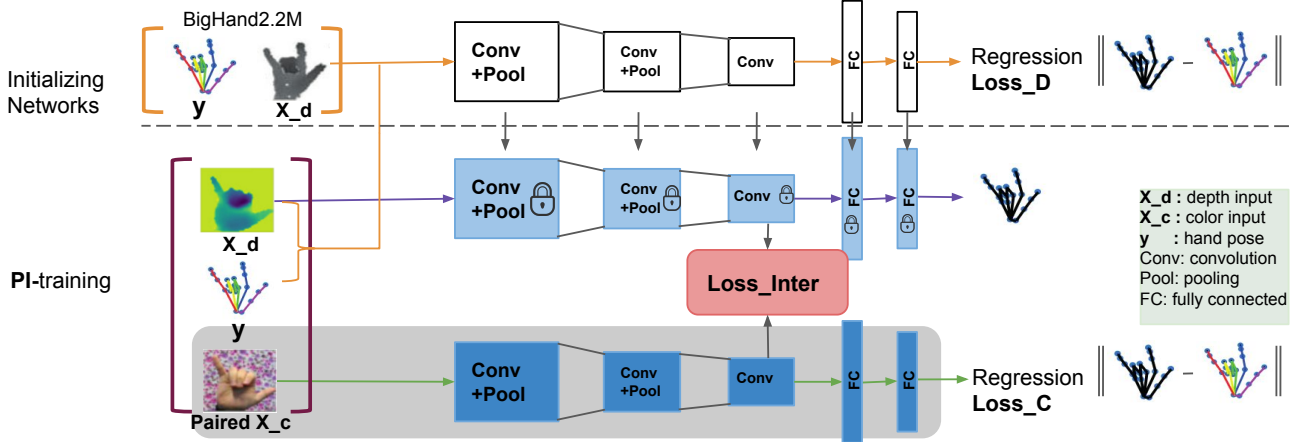
Figure 1. **Proposed framework for 3D hand pose estimation from an RGB image using privileged depth data.** *Training proceeds in two stages, a pre-training stage and privileged information (PI)-training stage. In the first stage, a depth-based network (top) and an RGB-based (bottom) network are trained independently to minimize 3D pose loss Loss_D and Loss_C. In the second stage, we freeze the parameters of the depth-based network and continue training with paired RGB and depth images, by minimizing a joint loss, which includes Loss_C and a mid-level feature regression loss Loss_Inter.*

ing. At testing, it used the softmax to select the final prediction between the predictions from the hallucination representation and the predictions from RGB representation. Luo *et al*. [15] recently proposed graph distillation for action detection with privileged modalities (RGB, depth, skeleton, and flow), where a novel graph distillation layer was used to dynamically learn to distill knowledge from the most effective modality, depending on the type of action. In our case, we use paired depth and RGB images during training. Depth and RGB networks are first trained separately. Subsequently the RGB network are progressively updated, while the depth network parameters remain fixed.

**Learning a latent space representation.** Latent space representation also shows promising for 3D hand pose estimation from RGB images [28, 13]. Spurr *et al*. [28] learned a cross-modal statistical hand model, via learning of a latent space representation that embeds sample points from multiple data sources such as 2D keypoints, images, and 3D hand poses. Multiple encoders were used to project different data modalities into a unified low-dimensional latent space, where a set of decoders reconstruct the hand configuration in either modality. Iqbal *et al*. [13] used latent 2.5D heatmaps, containing the latent 2D heatmaps and latent depth maps, to ensure the scale and translation invariance. Absolute 3D hand poses are reconstructed from the latent 2.5D heatmaps. Cai *et al*. [2] proposed a weakly-supervised method for 3D hand pose estimation from RGB image by introducing an additional depth regularizer module, which rendered a depth image from the estimated 3D hand pose. Training was conducted by minimizing an additional loss term, which is the $L1$ distance between the rendered depth image and the ground truth depth image.

## 3. Methods

We propose a framework to train a hand pose estimation model from RGB images by using depth images as privileged information. The model learns a new RGB representation which is influenced by the paired depth representation through mimicking the mid-level features of a depth network.

As shown in Figure 1, we use depth images in two ways: (1) to train an initial depth-based network with the aim of regressing 3D hand poses. Depth data that is annotated with 3D full hand pose information is abundant in the literature, and we choose the largest real dataset BigHand2.2M [47] to train our depth-based model, see the top row of Figure 1. (2) Paired RGB and depth images are fed into the RGB-based and depth-based network with the parameters of the depth-based network being frozen. The training of the RGB-based network continues with the aim of minimizing a joint loss function. The joint loss function has two parts, the first part being the 3D hand pose regression loss, *Loss_C*, and the second part the mid-level regression loss, *Loss_Inter*.

### 3.1. Architecture

Figure 1 shows our training architecture. There are two base models, each for one input channel. We use deep convolutional neural networks (CNNs), which have been widely used in hand pose estimation and have proven useful in transferring information from one network to another [11]. Prior work [17] has been shown useful in combining RGB and depth images as a four-dimensional RGB-D input to a single CNN model to estimate 3D hand pose. In our architecture, we share information in the middle lay-

ers of our two CNN models, one is a depth-based network and the other one is an RGB-based network. Each CNN model takes an input (a depth image or an RGB image) and produces a 3D hand pose estimation result.

For clarity, we denote the depth-based network *Depth_Net*, the RGB-based network *RGB_Net* when this is trained before privileged information is used. When privileged information is introduced in the training, we denote the RGB-based network *RGB_PI_Net*. In summary, *RGB_Net* and *RGB_PI_Net* are the same CNN model trained before and after the paired RGB and images are used to train the RGB channel.

We aim at sharing information between the middle layers of our two CNN models, and in particular using *Depth_Net* to inform *RGB_PI_Net* in the training time when paired RGB and depth images are available. To let the *Depth_Net* channel share information with *RGB_PI_Net*, we introduce an intermediate regression loss between the paired layers in the two models. This intermediate regression loss is inspired by prior works [12, 11], where similar techniques are used for model distillation [11], supervision transfer from well labeled RGB images to depth images with limited annotation [10], and hallucination of different modalities [12]. We therefore introduce an intermediate loss, which helps *RGB_PI_Net* to extract middle level features that mimic the responses of the corresponding layer of the *Depth_Net* using the paired depth image.

The intermediate loss (or *Loss_Inter* as shown in Figure 1) is defined as:

$$Loss\_Inter(k) = \|A_k^{Depth} - A_k^{RGB}\|_2^2, \quad (1)$$

where $A_k^{Depth}$ and $A_k^{RGB}$ are the $k_{th}$ layer activations for Depth Network and RGB Network, respectively. During testing, where only an RGB image is available, we feed the RGB image into *RGB_PI_Net* to estimate the 3D hand pose.

### 3.2. Training with privileged information

This section explains the details of training the proposed architecture. We choose a base CNN for *Depth_Net* and *RGB_PI_Net* for 3D hand pose estimation. For the base model, we build on Convolutional Pose Machine (CPM)'s [43] feature extraction layers with two fully connected layers to regress a 63 dimensional 3D hand pose with 21 joints.

In this initial stage, we call this external depth images as privileged information. Our *Depth_Net* is initially independently trained on BigHand2.2M [47] dataset, which has 2.2 million fully annotated (21 joints) depth images. After training, the model is further trained on the depth images of a smaller dataset (*e.g.*, Stereo [48] and RHD [52] datasets) that has fully annotated paired RGB and depth images. The *RGB_Net* is initially trained on the RGB images from the same dataset.

| Dataset | No. Training | No. Test | No. Joints | Annotation | Type |
|---|---|---|---|---|---|
| Stereo [48] | 15,000 | 3,000 | 21 | 2D, 3D | real |
| RHD [52] | 41,258 | 2,728 | 21 | 2D, 3D | synthetic |
| Dexter-Object [29] | - | 3,111 | 5 (tips) | 2D, 3D | real |

Table 1. **Public datasets used in our experiments.**

When the initial training is completed for both CNN models, we freeze the parameters of the *Depth_Net* and start training *RGB_PI_Net* with privileged information. In this stage, our privileged information is the paired depth images, and comes into use in the form of the middle layer activations of the *Depth_Net*. During the privileged training stage, we want the *RGB_PI_Net*'s middle level layer's activations to match the activations of the corresponding layers of the *Depth_Net*. We have two losses to optimize: (1) *Loss_Inter* (Eqn. 1) is used to match the middle layer activations of the two CNN models. (2) *Loss_C* (see Figure 1) is the L2 loss between the ground truth and the estimated 3D hand pose. Here we use a joint loss:

$$Loss\_Joint(k) = Loss\_Inter(k) + \lambda \cdot Loss\_C, \quad (2)$$

where $\lambda$ is used to balance the two losses, a larger value of $\lambda$ means less supervision is required from the privileged information, a smaller value means that the model depends more on the supervision. We set $\lambda$ to 100 for all experiments.

### 3.3. Foreground mask as privileged information

In addition to the supervision from depth images, we also explore the idea of extracting hand masks from depth images and embedding the hand masks into CNN layers of *RGB_PI_Net* to suppress the background features. As shown in Figure 2, we treat the hand mask $M_h$ as privileged information. At test time, when the hand mask is not available, the CNN model is viewed as a standard CNN with convolutional layers, pooling layers and full-connected layers, where the *Loss_Mask* is not used. In the training stage, the foreground hand mask is introduced in the last convolutional layer, as shown in Figure 2. Pixels of the mask $M_h$ are zero on the hand region, and one otherwise. We suppress background features by minimizing the regression loss *Loss_Mask*:

$$Loss\_Mask = \|A_k^{RGB} \odot M_h\|_2^2, \quad (3)$$

where $\odot$ denotes element-wise multiplication.

By minimizing the regression loss, where the response on the hand is multiplied by zero and the response outside the hand is multiplied by one, the response from outside the hand area is suppressed, focusing the response on the hand region.

## 4. Experiments

We carried out experiments on both 3D and 2D hand pose estimation from RGB images. Our experiments are
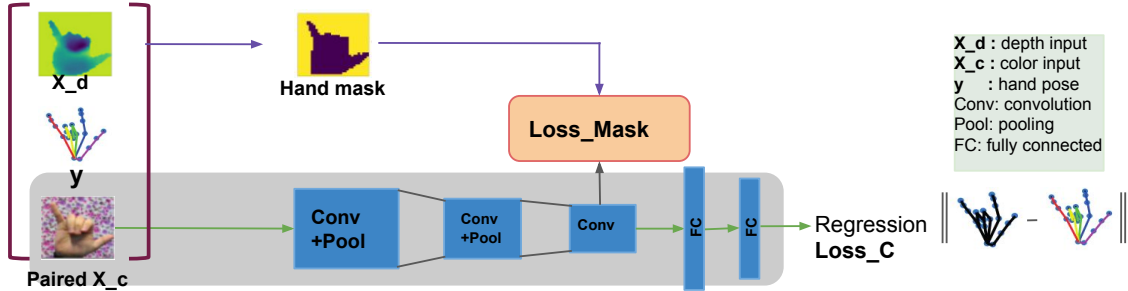
Figure 2. **Treating hand mask as privileged information**. Hand mask are used as privileged information to suppress the responses from the background area in the middle layers.

conducted on three public RGB-D datasets: the RHD dataset [52], the Stereo dataset [48], and the Dexter-Object dataset [29], as shown in Table 1. The RHD dataset is created synthetically and contains 41,238 training and 2,728 test images, with a resolution of $320 \times 320$. Each pair of RGB and depth images contains 3D annotations for 21 hand joints, and intrinsic camera parameters. The RHD dataset is built from 20 different subjects performing 39 actions. The training set has 16 subjects performing 31 actions, while the test set has 4 subjects performing 8 actions. The dataset contains diverse backgrounds sampled from 1,231 Flickr images. The Stereo [48] dataset is a real RGB-D dataset, which has 18,000 pairs of RGB and depth images with a resolution of $640 \times 480$ pixels. Each pair is fully annotated with 21 joints. The dataset contains six different backgrounds with respect to different difficulties (*e.g.* textured/textureless, dynamic/static, near/far, hightlights/no-highlights). For each background, there are two sequences, each containing 1,500 image pairs. The dataset is manually annotated. In our experiments, we follow the evaluation protocol of [52], *i.e.*, we train on 10 sequences (15,000 images) and test on the remaining 2 sequences (3,000 images). The Dexter-Object [29] dataset contains 3,111 images of two subjects performing manipulations with a cuboid. The dataset provides RGB and depth images, but only fingertips are annotated. The RGB images have a resolution of $640 \times 320$ pixels. Due to the incomplete hand annotation, we use this dataset for cross-dataset generalization.

During testing on a GTX 1080 Ti, the network forward steps take 6ms for 3D pose estimation and 8ms for the 2D case. The image cropping and normalization is the same as in [52]. To crop the hand region, we use ground truth annotations to obtain an axis aligned crop, resized to $256 \times 256$ pixels by bilinear interpolation. Examples are shown in the first row of Figure 4. For 3D hand pose estimation, we use the root joint's world coordinates and the hand's scale to normalize the results.

### 4.1. 3D hand pose estimation from RGB

In this section, we investigate the usefulness of depth images to improve the performance of 3D hand pose estimation from an RGB image. Our base CNN model is built upon the feature extraction layers of Convolutional Pose Machine (CPM) [43] with two fully connected layers. The final output is a 63 dimensioal vector denoting the 21 joint 3D locations. Specifically, our base CNN model contains 14 convolutional layers, 4 pooling layers, and 2 fully-connected layers. At training stage, we have access to paired RGB and depth images. Initially the *Depth_Net* is trained on *BigHand2.2M* [47]. We continue to train the *Depth_Net* using the depth images from the small dataset, *e.g.*, Stereo dataset or RHD dataset. We train the *RGB_Net* with the RGB images from the small dataset. When the initial training is completed, we start PI-training with the paired RGB and depth images. We freeze the weights of the *Depth_Net* and add the intermediate regression loss *Loss_Inter* among the mid-level features of *Depth_Net* and *RGB_PI_Net*, then we continue the training of *RGB_PI_Net* by minimizing the joint loss $Loss\_Joint$. We apply the intermediate loss to the last convolutional layers of both branches, where the parameter $k$ is set to 18 in Equation 1 and Equation 2.

**Effect of PI-Learning:** We conduct experiments with the two baseline CNNs and the CNN after PI training, see the accuracy curves in Figure 3 (top-left plot). Our networks only estimate relative 3D pose from a cropped RGB image patch containing the hand, to yield 3D hand pose in world coordinates, we follow a similar procedure of [52], *i.e.*, by adding the absolute position of the root joint to our estimated results. For comparison we choose the Percentage of Correct Keypoints (PCK) over a varying threshold. Training with depth data significantly improves the performance of the RGB-based network, narrowing the gap to the depth-based network.

**Comparison with the state of the art:** We compare our results with state-of-the-art methods, including PSO [19], ICPPSO [22], Zhang *et al*. [48], Z&B [52],
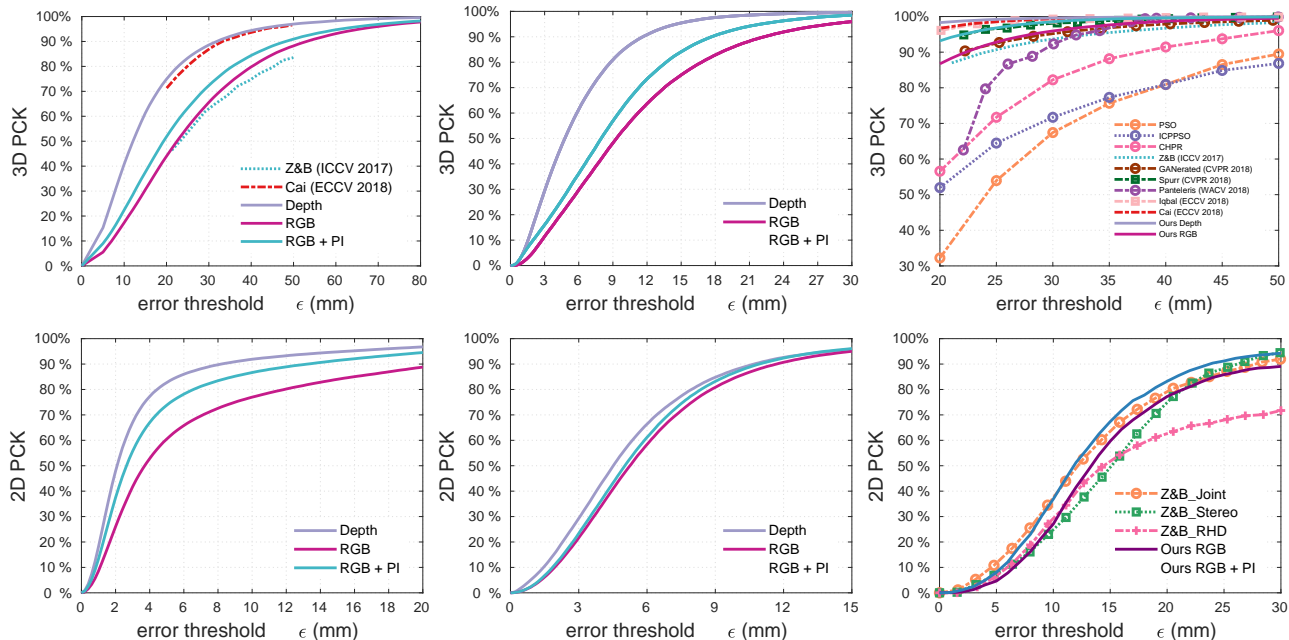
Figure 3. **Results one Stereo and RHD dataset for 3D hand pose and 2D pose accuracy**.*Top row shows the comparisons of 3D hand pose accuracy, bottom row shows the comparisons of 2D hand pose accuracy. Top-left is self-comparison on RHD dataset, top-middle is self-comparison on Stereo dataset, top-right is comparison with state-of-the-art on Stereo dataset. Bottom-left is self comparison on RHD dataset, bottom-middle is self comparison on Stereo dataset, bottom-right is comparison with state-of-the-art on the Dexter-Object dataset.*

GANerated [16], Cai *et al.* [2], Spurr *et al.* [28], Iqbal *et al.* [13], Panteleris *et al.* [21], see Figure 3 (top-right plot). Our method out-performs all existing state-of-the-art methods. We outperform (Z&B) [52] and [16]. While both [52] and [16] used extra training data, [52] used both Stereo (real) and RHD (synthetic) data to train their network. [16] used synthetic (GANerated) data to train their network. The proposed method uses less RGB training data and achieved the best performance. we significantly outperformed both methods with our privileged training strategy.

**Feature activation maps:** To give more intuitions on the effectiveness of training using additional privileged information, we visualize the activations of the mid-level feature for the three networks. Feeding an RGB image into each network, we aggregate all the mid-level feature maps into feature map by taking the maximum across all feature maps (similar to the maxout operation [9]). As shown in Figure 4, training with privileged information helps to select more representative features, where the visualized activations are close to the foreground (the hand).

**Loss function evolution:** We keep a record of the loss during our training on the Stereo dataset, see Figure 5. The loss for 3D hand pose (left plot) of the RGB network on the test data converges at iteration 15,000, we continue training for another 5,000 iterations. From iteration 20,000, we fix the depth network parameters and connect mid-level features between the RGB and depth networks, and continue

training by minimizing the joint loss (right plot) using RGB-D image pairs. The intermediate loss (middle plot) is used to suppress the difference between the mid-level feature between the RGB and depth networks. Loss for 3D hand pose of the RGB network, and the joint loss stop decreasing at around iteration 30,000.

## 4.2. 2D hand pose estimation from RGB

In this section, we choose the base CNN model as CPM [43], which has shown great performance for 2D human pose estimation [43], and 2D hand pose estimation [52]. Results are reported in Table 2, where 'EPE' stands for the 'average end point error' in pixels, where an end point is a hand joint. Qualitative examples are shown in Figure 7 and Figure 8. In this part of experiments, we treat the hand mask as privileged data, the CNN base model is CPM [27]. The baseline is obtained by the normal training procedure, *i.e.*, feeding the pre-processed hand image into CPM and obtaining the 2D hand pose by finding out the maximum location in each of the 21 heatmaps. For training with privileged information, we randomly select a certain proportion of RGB training data and use the hand masks, which are obtained by thresholding the depth images, in the *Loss_Mask* to suppress the background responses. As shown in Table 2, where 0.2 and 0.8 denotes the percentage of images when the *Loss_Mask* is used during the training for 2D hand pose estimation.
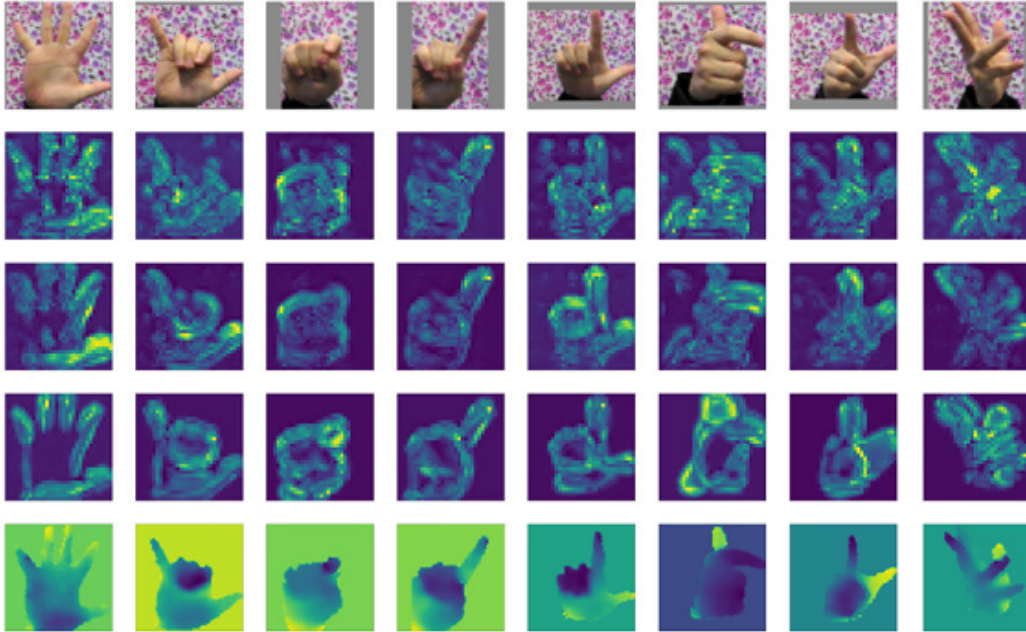
Figure 4. **Feature activation maps**. *(top row) input images, (row 2) activations of the RGB network trained on RGB only, (row 3) activations of the RGB network trained with additional depth data, (row 4) activations of the depth network, and (row 5) depth images. During training, depth data helps the RGB network focus on the region of interest, reducing the influence of background regions.*
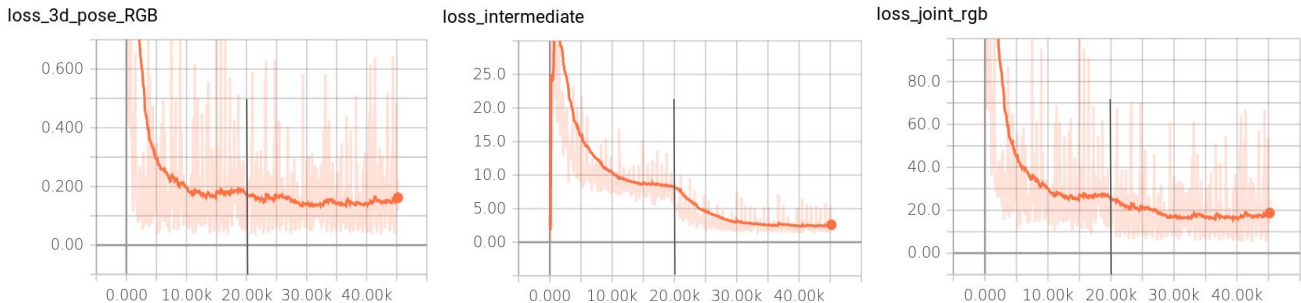


Figure 5. **Loss function evolution on the stereo dataset [48].** *The loss for 3D hand pose (left plot) of the RGB network on the test data converges at iteration 15,000, we continue training for another 5,000 iterations. From iteration 20,000, we fix the depth network parameters and connect mid-level features between the RGB and depth networks, and continue training by minimizing the joint loss (right plot) using RGB-D image pairs. The intermediate loss (middle plot) is used to suppress the difference between the mid-level feature between the RGB and depth networks. Loss for 3D hand pose of the RGB network, and the joint loss stop decreasing at around iteration 30k.*

**Performance on hand-object interaction dataset:** In Figure 3 (bottom-right plot), we show a comparison in terms of 2D PCK (in pixels) on the Dexter-Object [29] dataset. *Z&B_Joint* denotes the method of Z&B [52] trained on both RHD and Stereo datasets, which is better than *Z&B_Stereo* (trained on Stereo) and *Z&B_RHD* (trained on RHD). Our approach outperformed *Z&B_Joint* even though we used less RGB training data.

## 5. Conclusions

In this paper, we proposed a framework for 3D hand pose estimation from RGB images, with the training stage aided

| Method | Testing | Training | EPE median | EPE mean |
|---|---|---|---|---|
| Z&B [52] | RHD | RHD+Stereo | 5.001 | 9.135 |
| Baseline RGB | RHD | RHD | 3.708 | 7.841 |
| Baseline Depth | RHD | RHD | 2.087 | 3.902 |
| RGB + PI training | RHD | RHD | 2.642 | 5.223 |
| Z&B [52] | Stereo | RHD+Stereo | 5.522 | 5.013 |
| Baseline RGB | Stereo | Stereo | 5.250 | 6.533 |
| Baseline Depth | Stereo | Stereo | 4.775 | 5.883 |
| RGB + PI training (0.2) | Stereo | Stereo | 5.068 | 6.280 |
| RGB + PI training (0.8) | Stereo | Stereo | 4.515 | 5.801 |
| Z&B [52] | Dexter-Object | RHD+Stereo | 13.684 | 25.160 |
| Baseline RGB | Dexter-Object | RHD | 13.360 | 18.278 |
| RGB + PI training | Dexter-Object | RHD | 11.809 | 14.593 |

Table 2. **2D Hand Pose Accuracy**. *Results when training on the RHD and Stereo datasets.*
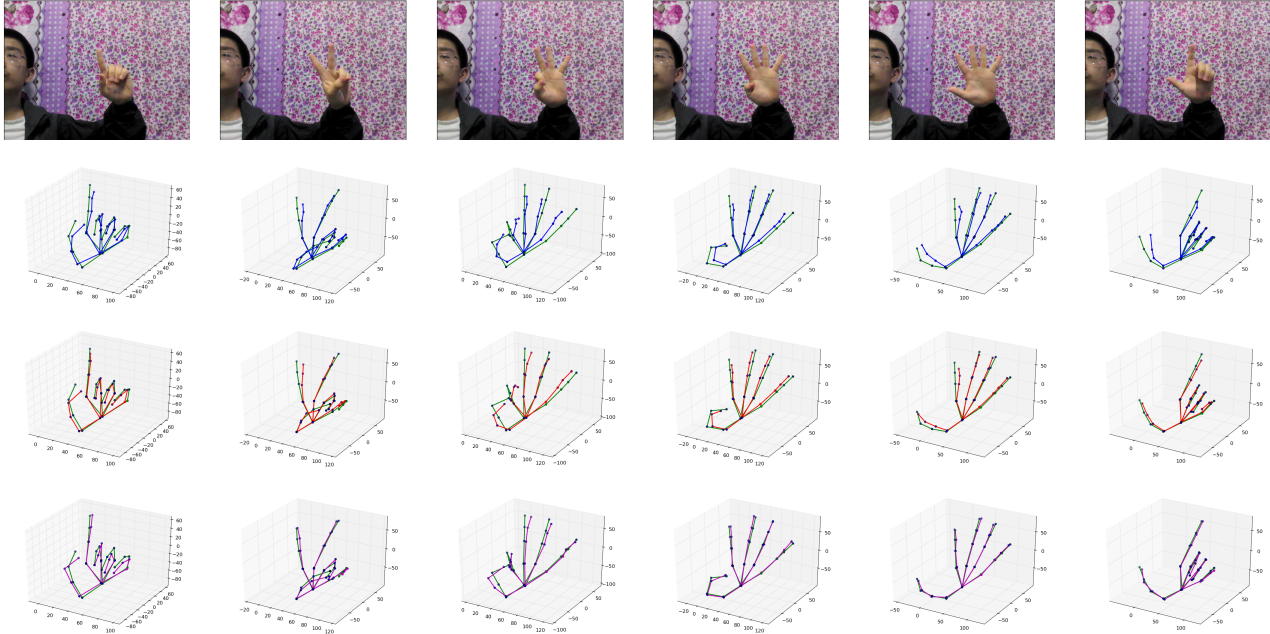
Figure 6. **Qualitative 3D pose estimation results**. *Comparing the outputs of the RGB network (blue, second row), the RGB network with PI training (red, third row), and the Depth network (magenta, bottom row) with the ground truth 3D pose (green) on the Stereo dataset. Top row are the original images.*
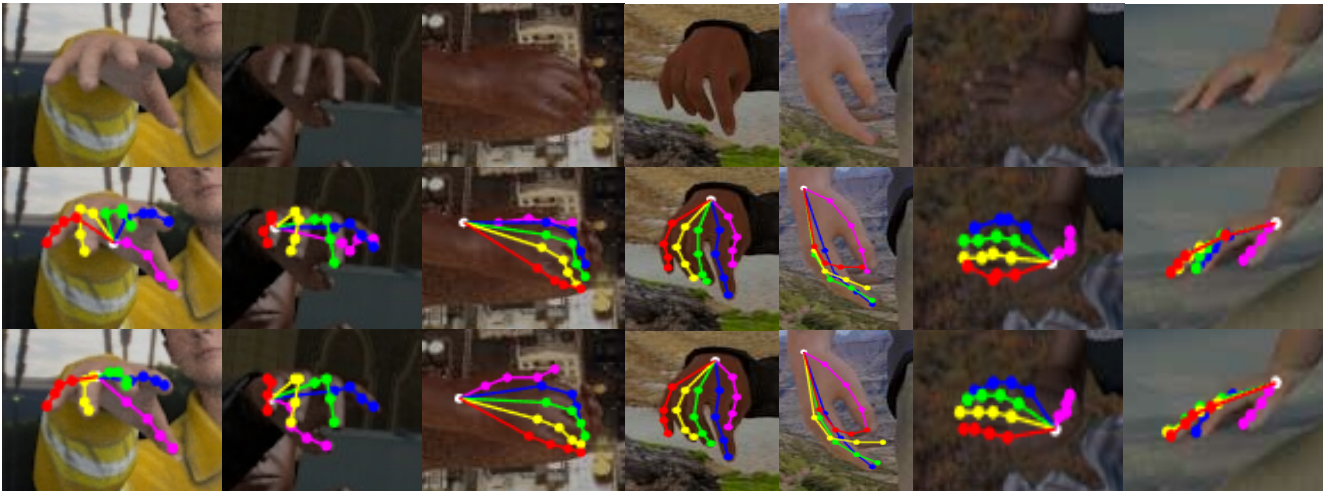


Figure 7. **Qualitative 2D pose estimation results**. *Comparing the outputs of (middle) the RGB network and (bottom) the RGB network with PI training on the RHD dataset. Top row are the original images.*

with privileged information, *i.e.* depth data. To the best of our knowledge, our method is the first to introduce the concept of using privileged information (depth images) to support training a RGB-based 3D hand pose estimator. We proposed three ways to use the privileged information: as external training data for a depth-based network branch, as paired depth data to transfer supervision from the depth-based network to the RGB-based network, and as a hand mask to suppress background activations in the RGB-based

network. Our training strategy can be easily embedded into existing pose estimation methods. As an illustration, we estimate 2D hand pose from an RGB image using a different CNN model. Results on 2D hand pose estimation, using our training strategy, are improved over state-of-the-art methods for 2D hand pose estimation from RGB input. During testing, when only RGB images are available, our model significantly outperforms the same model trained only using RGB images. This training strategy can be incorpo-
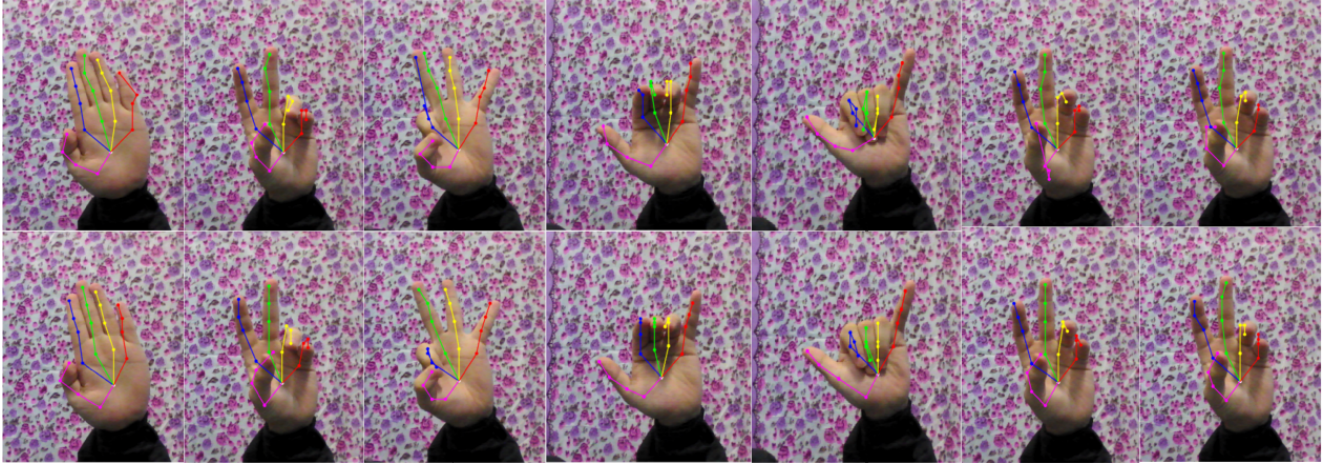
Figure 8. **Qualitative 2D pose estimation results on Stereo dataset**. *Comparing the outputs of (top) the RGB network and (bottom) the RGB network with PI training.*

rated into existing models to boost the performance of hand pose estimation from an RGB image. One limitation of our method is the difficulty of handling occlusion by objects, which can be addressed by systematically adding synthetic objects in the depth data (privileged information).

## References

[1] J. Ba and R. Caruana. Do deep nets really need to be deep? In *NIPS*, 2014.

[2] Y. Cai, L. Ge, J. Cai, and J. Yuan. Weakly-supervised 3d hand pose estimation from monocular rgb images. In *ECCV, 2018*.

[3] L. Chen, W. Li, and D. Xu. Recognizing rgb images by learning from rgb-d data. In *CVPR 2014*.

[4] Y. Chen, X. Jin, J. Feng, and S. Yan. Training group orthogonal neural networks with privileged information. *arXiv:1701.06772*, 2017.

[5] C. Choi, A. Sinha, J. Hee Choi, S. Jang, and K. Ramani. A collaborative filtering approach to real-time hand pose estimation. In *ICCV*, 2015.

[6] A. Erol, G. Bebis, M. Nicolescu, R. D. Boyle, and X. Twombly. Vision-based hand pose estimation: A review. *CVIU*, 108(1):52–73, 2007.

[7] 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. *CVPR*, 2018.

[8] L. Ge, H. Liang, J. Yuan, and D. Thalmann. 3D convolutional neural networks for efficient and robust hand pose estimation from single depth images. CVPR, 2017.

[9] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. *ICML*, 2013.

[10] S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *CVPR 216*.

[11] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *arXiv:1503.02531*, 2015.

[12] J. Hoffman, S. Gupta, and T. Darrell. Learning with side information through modality hallucination. In *CVPR*, 2016.

[13] U. Iqbal, P. Molchanov, T. Breuel, J. Gall, and J. Kautz. Hand pose estimation via latent 2.5 d heatmap regression. In *ECCV, 2018*.

[14] S. Li and A. B. Chan. 3D human pose estimation from monocular images with deep convolutional neural network. In *ACCV*, 2014.

[15] Z. Luo, L. Jiang, J.-T. Hsieh, J. C. Niebles, and L. Fei-Fei. Graph distillation for action detection with privileged information. In *ECCV, 2018*.

[16] 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 *CVPR, 2018*.

[17] F. Mueller, D. Mehta, O. Sotnychenko, S. Sridhar, D. Casas, and C. Theobalt. Real-time hand tracking under occlusion from an egocentric RGB-D sensor. In *ICCV*, 2017.

[18] M. Oberweger, P. Wohlhart, and V. Lepetit. Training a feedback loop for hand pose estimation. In *ICCV*, 2015.

[19] I. Oikonomidis, N. Kyriazis, and A. A. Argyros. Efficient model-based 3D tracking of hand articulations using kinect. In *BMVC*, 2011.

[20] P. Panteleris and A. A. Argyros. Back to RGB: 3D tracking of hands and hand-object interactions based on short-baseline stereo. In *ICCV Workshop*, 2017.

[21] P. Panteleris, I. Oikonomidis, and A. A. Argyros. Using a single rgb frame for real time 3d hand pose estimation in the wild. In *WACV 2018*.

[22] C. Qian, X. Sun, Y. Wei, X. Tang, and J. Sun. Realtime and robust hand tracking from depth. In *CVPR, 2014*.

[23] M. Rad, M. Oberweger, and V. Lepetit. Domain transfer for 3d pose estimation from color images without manual annotations. In *ACCV 2018*.

[24] T. Sharp, C. Keskin, D. Robertson, J. Taylor, J. Shotton, D. Kim, C. Rhemann, I. Leichter, A. Vinnikov, Y. Wei, D. Freedman, P. Kohli, E. Krupka, A. Fitzgibbon, and

S. Izadi. Accurate, robust, and flexible real-time hand tracking. In *CHI*, 2015.

[25] Z. Shi and T.-K. Kim. Learning and refining of privileged information-based RNNs for action recognition from depth sequences. In *CVPR*, 2017.

[26] E. Simo-Serra, A. Ramisa, G. Alenyà, C. Torras, and F. Moreno-Noguer. Single image 3D human pose estimation from noisy observations. In *CVPR*, 2012.

[27] T. Simon, H. Joo, I. Matthews, and Y. Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *CVPR*, 2017.

[28] A. Spurr, J. Song, S. Park, and O. Hilliges. Cross-modal deep variational hand pose estimation. In *CVPR, 2018*.

[29] S. Sridhar, F. Mueller, M. Zollhöfer, D. Casas, A. Oulasvirta, and C. Theobalt. Real-time joint tracking of a hand manipulating an object from rgb-d input. In *ECCV*, pages 294–310. Springer, 2016.

[30] X. Sun, Y. Wei, S. Liang, X. Tang, and J. Sun. Cascaded hand pose regression. In *CVPR*, 2015.

[31] J. S. Supancic, G. Rogez, Y. Yang, J. Shotton, and D. Ramanan. Depth-based hand pose estimation: methods, data, and challenges. *ICCV*, 2015.

[32] D. Tang, H. J. Chang, A. Tejani, and T.-K. Kim. Latent regression forest: Structured estimation of 3D articulated hand posture. In *CVPR*, 2014.

[33] D. Tang, J. Taylor, P. Kohli, C. Keskin, T.-K. Kim, and J. Shotton. Opening the black box: Hierarchical sampling optimization for estimating human hand pose. In *ICCV*, 2015.

[34] D. Tang, Q. Ye, S. Yuan, J. Taylor, P. Kohli, C. Keskin, T.-K. Kim, and J. Shotton. Opening the black box: Hierarchical sampling optimization for hand pose estimation. *TPAMI*, 2018.

[35] D. Tome, C. Russell, and L. Agapito. Lifting from the deep: Convolutional 3D pose estimation from a single image. *CVPR*, 2017.

[36] J. Tompson, M. Stein, Y. Lecun, and K. Perlin. Real-time continuous pose recovery of human hands using convolutional networks. In *TOG*, 2014.

[37] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS*, 2014.

[38] A. Toshev and C. Szegedy. DeepPose: Human pose estimation via deep neural networks. In *CVPR*, 2014.

[39] V. Vapnik and A. Vashist. A new learning paradigm: Learning using privileged information. *Neural networks*, 22(5-6):544–557, 2009.

[40] C. Wan, T. Probst, L. Van Gool, and A. Yao. Crossing nets: Combining GANs and VAEs with a shared latent space for hand pose estimation. In *CVPR*, 2017.

[41] C. Wan, A. Yao, and L. Van Gool. Hand pose estimation from local surface normals. In *ECCV*, 2016.

[42] Z. Wang and Q. Ji. Classifier learning with hidden information. In *CVPR*, 2015.

[43] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *CVPR*, 2016.

[44] M. Xu, R. Jin, and Z.-H. Zhou. Speedup matrix completion with side information: Application to multi-label learning. In *NIPS*, 2013.

[45] Q. Ye, S. Yuan, and T.-K. Kim. Spatial attention deep net with partial PSO for hierarchical hybrid hand pose estimation. In *ECCV*, 2016.

[46] S. Yuan, G. Garcia-Hernando, B. Stenger, G. Moon, J. Y. Chang, K. M. Lee, P. Molchanov, J. Kautz, S. Honari, L. Ge, et al. 3D hand pose estimation: From current achievements to future goals. In *CVPR*, 2018.

[47] S. Yuan, Q. Ye, B. Stenger, S. Jain, and T.-K. Kim. Big-Hand2.2M benchmark: Hand pose dataset and state of the art analysis. *CVPR*, 2017.

[48] J. Zhang, J. Jiao, M. Chen, L. Qu, X. Xu, and Q. Yang. 3D hand pose tracking and estimation using stereo matching. *arXiv:1610.07214*, 2016.

[49] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei. Deep kinematic pose regression. In *ECCV*, 2016.

[50] X. Zhou, M. Zhu, S. Leonardos, K. G. Derpanis, and K. Daniilidis. Sparseness meets deepness: 3D human pose estimation from monocular video. In *CVPR*, 2016.

[51] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017.

[52] C. Zimmermann and T. Brox. Learning to estimate 3D hand pose from single RGB images. In *ICCV*, 2017.