Learning to Infer the Depth Map of a Hand from its Color Image

Vassilis C. Nicodemou^{1,2}

Iason Oikonomidis¹

Georgios Tzimiropoulos³

nikodim@ics.forth.gr

oikonom@ics.forth.gr

yorgos.tzimiropoulos@nottingham.ac.uk

Antonis Argyros^{1,2}

argyros@ics.forth.gr

¹Computational Vision and Robotics Laboratory, Institute of Computer Science, FORTH, Greece

²Computer Science Department, University of Crete, Greece

³Computer Vision Laboratory, University of Nottingham, United Kingdom

Abstract

We propose the first approach to the problem of inferring the depth map of a human hand based on a single RGB image. We achieve this with a Convolutional Neural Network (CNN) that employs a stacked hourglass model as its main building block. Intermediate supervision is used in several outputs of the proposed architecture in a staged approach. To aid the process of training and inference, hand segmentation masks are also estimated in such an intermediate supervision step, and used to guide the subsequent depth estimation process. In order to train and evaluate the proposed method we compile and make publicly available HandRGBD, a new dataset of 20,601 views of hands, each consisting of an RGB image and an aligned depth map. Based on HandRGBD, we explore variants of the proposed approach in an ablative study and determine the best performing one. The results of an extensive experimental evaluation demonstrate that hand depth estimation from a single RGB frame can be achieved with an accuracy of 22mm, which is comparable to the accuracy achieved by contemporary low-cost depth cameras. Such a 3D reconstruction of hands based on RGB information is valuable as a final result on its own right, but also as an input to several other hand analysis and perception algorithms that require depth input. Essentially, in such a context, the proposed approach bridges the gap between RGB and RGBD, by making all existing RGBD-based methods applicable to RGB input.

1. Introduction

The task of observing and understanding human activities is of great interest to the field of computer vision. Among other approaches, human activity can be studied by observing and monitoring the state of the human body, ei-





Figure 1: Given an RGB image of a human hand (left) the goal of this work is to produce a depth map of the hand region (right).

ther in 2D or in 3D. Particular emphasis is given to the human hands as the interpretation of their behavior is key to understanding the interaction of humans with their environment. Several efforts have been devoted to this direction and important milestones have been achieved [43, 71]. However, despite the significant progress, a general solution to these problems is still lacking.

This work deals with the problem of estimating the depth map of a hand observed from a regular color camera. Depth information is lost during color image formation and is important for the analysis of hands. Our goal is to develop a method that accepts as input a conventional RGB image of a hand and produces the depth map of the observed hand (see Figure 1). The analysis of the rest of the observed scene (non-hand regions) is out of the scope of this work. Thus, no depth estimation is performed in those regions; the method should only characterize them as background.

Solving the aforementioned problem is interesting and useful, from both a theoretical and a practical point of view. When observing a scene using regular images, it is very appealing to be able to recover the suppressed depth infor-

mation without stereo 3D reconstruction or structure from motion. At the same time, the recovery of this information may have significant impact to the solution of several practical problems. As an example, the hand depth information is useful in the context of interaction applications [15]. On top of this, it can be used to capture and understand hand movement within the 3D space, facilitating tasks such as 3D hand shape and pose estimation, hand-object interaction monitoring, etc, with immediate implications to areas such as human-computer and human-robot interaction, AR/VR, medical rehabilitation, computer games, and more.

Inferring depth information of hands from color information is a problem that presents important challenges. Human hands exhibit large differences in shape and appearance from person to person. On top of this, due to the articulated structure of the hand, there is a very wide range of 3D postures and considerable self-occlusions. Further complications arise when the observed hand interacts with its environment, for example when handling objects. Thus, the recovery of the 3D structure of the hand given color information, only, is a rather demanding task. Depth estimation techniques using regular color input have been proposed for general scenes [61, 26, 10] as well as for specific objects such as human faces [20] and bodies [64, 36]. However, to the best of our knowledge, no existing method has tackled the problem of hand depth estimation.

In this work, we capitalize on recent advances in machine learning and propose a deep neural network architecture to tackle the problem of hand depth estimation of hands from color images. The training of such a method requires aligned RGB and depth information for a large number of hand views. Up to now, there is no publicly available dataset that is suitable for such a training process. This comes to a surprise, but can be explained by the fact that most of the hand-related works dealt with the problem of 3D hand pose estimation and tracking based on depth information. Therefore, annotation involves depth maps and the associated 3D hand poses and does not include color information. We therefore compiled *HandRGBD*¹, a dataset of 20,601 pairs of hand color images aligned with their respective depth maps. We use HandRGBD to evaluate our approach in comparison with variants in the context of an ablation study. We show that hand depth estimation from a single RGB frame can be achieved with an accuracy of 22mm, which is comparable to the accuracy achieved by contemporary low-cost depth cameras. Thus, the proposed method is a significant step towards turning an RGB camera to an RGBD one for hand analysis applications. Moreover, the proposed method makes all depth-based hand analysis methods exploitable on plain RGB input.

2. Related Work

Depth from color for general scenes: The general problem of extracting depth information from color images is very interesting [3, 61, 44, 45, 27], and is still a research topic under investigation [6, 37, 29], remaining unsolved in its full generality. Since as early as 1985, methods have been proposed for extracting depth information from color images without prior knowledge of the scene [3, 17, 22, 53, 46, 61, 16]. The exploited cues vary, including local texture, symmetries, edges and their intersections, and fractal dimension. Also, the extracted level of detail ranges from local depth variations [17] to global depth scale of the observed image [61]. The recent success of machine learning, including (most notably) deep neural networks has naturally led to new methods that achieve increasingly better performance and more accurate results [26, 10, 28, 9, 24, 11, 25, 13]. A significant category of methods deals with the problem of estimating the 3D surface of a deformable object [63, 42]. However, the assumed level of deformation in that line of work is much higher than that of human hands. These methods are designed to tackle deformations of paper and cloth, making the deformation model unnecessarily complicated for the case of human hands. Another category of works on monocular scene structure estimation assumes that the input is a sequence of images, coming from a camera that moves [7, 29]. Essentially, this results in stereo pairs of images, enabling the use of the disparity cue. A central assumption in this line of work regards the rigidity of the observed world. Due to the articulated structure and motion of the human hand, this assumption does not hold in our setting. Closer to our work, another category of methods uses priors such as the assumption that the scene contains articulated objects [67, 68]. In these approaches, however, the shape and size variation is larger, and the articulation range of the objects is usually more constrained than the case of the human hand.

Depth from color for human body parts: There is a recent line of work on methods that tackle the problem of depth estimation for specific parts of the human body. The first of these methods [20] estimates the face structure from a single color image. To do so, it uses volumetric information to train a neural network that was based on a stacked hourglass architecture. More works estimating the 3D structure of human faces shortly followed [58, 57]. In parallel to the works on human face, a similar architecture was proposed, targeting the human body [64]. In that work training data are derived from accurate models of pose, shape and appearance of the human body. None of these approaches achieves a direct estimation of the depth model of the observed scene. Instead, intermediate steps with higher-level information are used, such as the estimation of the landmark positions of facial features for the face approaches, or

¹HandRGBD will be made publicly available.

2D and 3D position of the body joints in the second case. Moreover, to the best of our knowledge, currently there is no method to solve the problem of estimating depth information from an RGB image of a hand.

Use of depth for 3D hand pose estimation: 3D hand pose estimation is a long-studied problem [39, 14, 50, 1, 51, 41, 66, 33] that is still of significant interest [71]. Most of the recent works in the area [33, 23, 56, 49, 55, 32, 54, 65, 30, 71] assume the availability of scene depth information, capitalizing on the advent of inexpensive, high-quality depth sensors. Much more recently, a new trend is currently forming [75, 34, 48, 5, 19] that tackles the problem assuming only monocular RGB input. The performance (estimation accuracy and speed) of the older, depth-based approaches is better than that of the more recent RGB-based ones. This is to be expected given the maturity of the older approaches and the richer nature of the depth map as input information.

Our contribution: An important goal of this work is to close the gap between depth-based approaches and the newer trend of works based on RGB input. Until today, in order to extract depth information for hands directly, the only available reliable option is to resort to depth sensors. The proposed method estimates hand depth information of comparable accuracy given only regular color images. This constitutes a significant complexity simplification and cost reduction of the sensing process. At the same time, several robust, depth-based hand perception methods become applicable to regular RGB input. In summary, the major contributions of this work are:

- The first method that estimates depth information from monocular color views of hands. This is achieved with an accuracy that is comparable to the accuracy of a low cost depth sensor.
- The *HandRGBD* dataset (will become publicly available) of 20,601 high resolution RGB hand images that are aligned with their depth maps.

3. Hand Depth Estimation from RGB Input

At the core of the proposed approach, a deep neural network undertakes the task of estimating the geometry of a hand observed in a single RGB image. A stacked hourglass model [31] is inspired from parts of the architecture in [20] and used as the main building block for the proposed approach. The resulting network accepts as input a regular RGB image and outputs the estimated hand depth map. The output of the network is a map of relative depths for all hand pixels of the input image. The absolute depth of the hand is a separate problem [61] that is out of the scope of our work. However, absolute depth estimation can be simplified given a good estimation of the relative depths resulting from the proposed method. Intermediate supervision is used

in several intermediate levels of the proposed architecture in a staged approach. To aid the process of training and inference, hand segmentation masks are also estimated in such an intermediate supervision step, and used as guidance for the subsequent depth estimation process.

3.1. Aligning RGB with Depth

The training data are assumed to be pairs of aligned RGB and depth hand images. The viewpoint of each image pair is assumed to be identical, i.e., each RGB pixel essentially corresponds to the pixel at the same position in the depth map, as if the two streams were captured from the same center of projection. This kind of data can be obtained using common RGBD sensors like Microsoft Kinect2. Most commercially available RGBD cameras have different sensors for each modality, however the viewpoints are very close, and the availability of depth data enables the alignment of the two streams. To achieve this, an accurate intrinsic and extrinsic calibration of the two sensors is required. Given this, a reprojection of the depth map to the RGB image yields the correspondences between the two images.

3.1.1 Ground Truth Annotation

Given the capturing process described in Section 3.1, the training data is already at a usable state and no further manual annotation is required. The only processing that is still required is the segmentation of the RGB and depth channels into foreground (hands) and background (non-hands), and the normalization of the depth range into relative depths. To facilitate this process, it is assumed that the hand is the object closest to the camera. Under this assumption, foreground/background segmentation is easy to perform on the depth map. Towards this end, the minimum value D_{min} in the depth map D(i, j) is estimated, corresponding to the distance of the hands' point that is closest to the camera. The indices i and j run on the horizontal and vertical image dimensions. All pixels with a depth value within a predefined threshold t to this minimum depth D_{min} are considered as the foreground H. The value H(i, j) of the boolean foreground mask H at point (i, j) is defined as:

$$H(i,j) = D(i,j) < (D_{min} + t).$$
 (1)

Working with depth maps in millimeters, it suffices to set t=300mm, a maximum estimation of the possible depth difference within a hand. Since the RGB and depth images correspond pixel-to-pixel, the resulting binary segmentation H is valid for the RGB image, too.

Let us denote with D[H] the depth map D masked with the foreground mask H. D[H] is used to compute the average depth value $\overline{D[H]}$ of hand points. $\overline{D[H]}$ is then subtracted from D and a fixed scaling is applied to the depths,

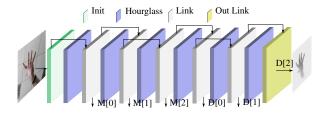


Figure 2: Stacked Hourglass Architecture: The proposed architecture with the intermediate supervision types. The input is preprocessed by some initialization layers ("Init", light green) that include a convolutiona layer and two residual blocks (Figure 3 and compute a feature map to be passed to the first hourglass (see Figure 4) module. Its output is passed to a set of layers that apply some additional convolution layers ("Link", gray) before passing it to the next hourglass module. The Link module also outputs a map to be intermediately guided. Skip connections are used parallel to each hourglass module. The first three outputs of the network target segmentation masks and the remaining three target depth maps with the latter being the final output of the network.

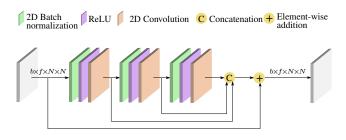


Figure 3: Residual block: The building block of the proposed neural network for hand depth estimation. The input is assumed to be a feature map of spatial dimension $N \times N$. In the figure, the feature count is f. The batch size b is also shown in the tensor dimensions. The output of the block is usually fed to more than one layers, for example to serve a skip connection, as shown here.

bringing the depth values in the range [-1, 1]. Specifically, the relative depth map D_T that will be used for training is

$$D_T(i,j) = c \cdot \left(D(i,j) - \overline{D[H]}\right), \tag{2}$$

where c is a value that scales the depths in the range [-1,1]. For all cases, it suffices to set c as the inverse of the maximum depth difference of an observed hand. When working with depth values in millimeters it suffices to use c=1/200. Finally, the background pixels are set to 1, the largest value in the target range, that is essentially used to denote background areas.

3.2. Stacked Hourglass Architecture

The proposed network is based on the approach of stacked hourglass modules [31, 4]. Additionally, intermediate supervision is applied to the output of each hourglass module, which is a commonly adopted strategy [31]. The architecture of the proposed method is illustrated in Figure 2. In the following description, the intermediate parts of the network will be called stages.

The main building block of the proposed architecture is the hourglass network of [31] built using the residual block of [4]². Figures 3 and 4 illustrate the residual block and the network used graphically. A hourglass module, illustrated in Figure 4 accepts as input a set of feature maps. The residual block of [4] proceeds by applying three successive sets of convolution, batch normalization and ReLU nonlinearity operations, using also skip connections, similar to the DenseNet architecture [18]. This is shown in Figure 3. After these operations, a down-sampling is performed, halving the input dimension. Parallel to this branch with halved spatial dimension, a skip connection runs through another residual block. In total, four repeated residual blocks and resolution halving are applied, and four long-skip connections run in parallel, each at a different spatial resolution. After the last subsampling and application of a residual block, the reverse process is followed, doubling the spatial dimension by upsampling and applying new residual block operations. After each upsampling, the long-skip connection of the appropriate spatial dimension is added to the current feature map. After four upsampling operations in total, the original input spatial and feature dimension is again reached, forming the complete hourglass module.

For the proposed network, we stack 6 such hourglass modules, having therefore in total 6 stages for intermediate supervision. A convolution operation is applied to the input image to compute a feature map of appropriate dimension to be the input of the first hourglass module. The reverse process is followed at the end of the network, and at the end of each hourglass module for intermediate supervision. Specifically, a single 1×1 convolution is applied, yielding a single-channel feature map. Each such output is trained against the foreground mask in the first stages, while the later stages are trained against the depth target. We set equal effort for estimating both the mask and the depth, thus giving 3 stages for the mask estimation and 3 for the depth.

3.3. Loss Function

Each stage has its own target output, therefore each stage has its own loss. The global loss function of the network is the sum of the individual losses. For each stage, regardless of the type of intermediate supervision, its loss is obtained

 $^{^2}$ Implementation available online at https://github.com/ladrianb/face-alignment

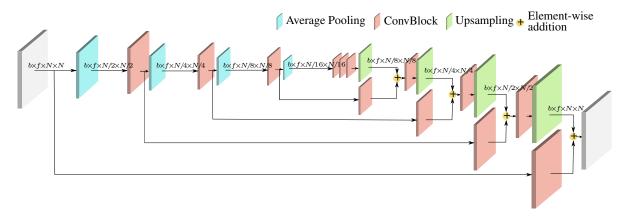


Figure 4: The hourglass building block that is used in the proposed network. Its main building block is the residual block, illustrated in detail in Figure 3. The main idea is to successively lower the spatial input resolution for a total of four input halving steps. After this, the reverse procedure is followed to reach again the input resolution.

by comparing the two images, the predicted and the target image. We define the loss of each stage as the Mean Square Error (MSE) of the target and the output. The final form of the loss function L is:

$$L(m, d, \tilde{m}, \tilde{d}) = \sum_{k=1}^{S_D} \frac{(\tilde{d}_k - d)^2}{||N||} + \sum_{l=1}^{S_M} \frac{(\tilde{m}_l - m)^2}{||N||}, \quad (3)$$

where d and m are the target depth and mask, each having N pixels, S_D and S_M are the total number of depth and mask stages respectively and \tilde{d}_k , \tilde{m}_l are the estimated depths and masks for the kth and lth stage for $k=1\ldots S_D$ and $l=1\ldots S_M$.

3.4. Data Augmentation

A commonly employed strategy during training is data augmentation which aims at enriching the diversity of the training set and at increasing the generalization capability of the trained network. In our case, the input is regular RGB images, and common augmentation practices apply. Specifically, we apply (a) random horizontal flip (so that we don't have to capture both hands from a subject) (b) random rotation, (c) random crop and (d) random color jittering (to capture the widest possible range of skin tones and different illumination cases). The geometric transformations are applied to both the RGB and the depth maps, ensuring pixel-to-pixel correspondence. The color transformations (e.g., color jittering) is only applied to the RGB channel. Finally, all data are resized to fit the network's input and output dimensions.

4. Experimental Evaluation

Due to the unavailability of an appropriate, published dataset, the quantitative evaluation of the proposed hand depth estimation method was performed on the basis of *HandRGBD*, a dataset that is introduced in this paper. A first category of experiments assessed the adopted design choices in an ablation study. We also assessed the potential of the proposed method to support methods that perform depth-based estimation of the 3D hand pose.

4.1. Hand-related Datasets

Works related to hand appearance and shape modeling as well as pose estimation, require datasets appropriately annotated with ground truth for the purposes of objective, quantitative comparison of competitive approaches, and also - whenever applicable - for training. Therefore, numerous datasets have been proposed so far in the relevant literature (see Table 1). Input modalities such as monocular RGB, stereo, multiview, and depth are covered. Also, scenarios including egocentric viewpoint, hand-object interaction, and hand-hand interaction are available.

The training and evaluation tasks of the problem we are addressing in this work call for a dataset that includes aligned RGB and depth observations of hands. The RGB input should be unaltered, since the goal is to apply our method to regular color input. Some datasets [60, 59] warp the RGB image to the depth map, introducing big black holes in the images that defeat this goal. Another dataset [21] segments the hand in the image and masks the background with a black color. Given that one of our goals is also to learn this segmentation, the dataset becomes unusable. Two additional requirements are the presence of multiple actors and close-up views of the depicted hand(s), so that details on the variation of hand shapes across humans and under articulation are adequately captured.

Table 1 presents a list of the most relevant datasets to our work. Columns of the matrix list some of the requirements listed above, specifically the availability of RGB and depth data, and of their alignment. Evidently, only the datasets

Table 1: Datasets on human hands. For the purposes of this work, aligned pairs of RGB and depth data are required.

Dataset	RGB	Depth	Alignment
Gomez [12]	√	-	-
Simon [47]	√	-	-
Bambach [2]	✓	-	-
Dreuw [8]	✓	-	-
Yuen [73]	√	-	-
Tang [55]	-	\checkmark	-
Sun [52]	-	\checkmark	-
Yuan [72]	-	\checkmark	-
Xu [70]	-	\checkmark	-
Tompson [60]	Warped	✓	-
Tkatch [59]	Warped	√	-
Zhang [74]	√	√	-
Rogez [40]	✓	✓	-
Sridhar [49]	✓	✓	-
Kanhangad [21]	No BG	√	√
Zimmermann [75]	√	√	√
Tzionas [62]	√	√	√

by Zimmerman and Brox [75], called "Rendered Handpose Dataset" (RHD) and Tzionas et al. [62], called "Hands in Action" have aligned RGB and depth data. Unfortunately, the RHD dataset [75] is synthetic, and, although it has a large variation on hand sizes, shapes and appearances, it is of rather low resolution (320×240) and contains distant views of a hand. The Hands in Action dataset [62] contains real world data, and the depth is captured by a structured light sensor. The actor diversity is small and the view is not closeup, in images of resolution 640×480 . Overall, this dataset comes closest to fulfilling our requirements, however it is still unsuitable due to the somewhat small resolution, the low actor diversity, and (less importantly) the use of a structured light sensor.

4.2. The HandRGBD Dataset

Despite the existence of several hand datasets, it turns out that none of them covers the requirements of this work. Consequently, we resorted to creating *HandRGBD*, our own dataset of aligned RGB and depth hand images.

As the capturing device, we employed a Kinect V2 [38] sensor because of its high quality color camera, and the Time of Flight depth sensor. Among the available options, this sensor provided the best combination of image and depth resolution and quality. The native SDK does not provide an alignment of the depth data to the RGB image, only the opposite, resulting in black holes in the RGB image. Therefore, we used the library libfreenect2 [69]³ that sup-

ports this functionality, simultaneously scaling and aligning the depth information on the color image.

The captured dataset contains 20, 601 images along with their respective depth maps. The depicted hands are in closeup view, in distances ranging from 40cm to 100cm from the sensor. Some of the captured images contain two hands that interact (strongly, in some cases). 17 subjects, 13 male and 4 female, contributed to the dataset. The subjects were instructed to keep their hand(s) roughly in the center of the camera field of view, but some images were also captured with hands close to the image edges. The subjects were also instructed to perform free hand gestures and articulations, exploring as much as possible the hand articulation space. Special care was taken to capture the hands in front of different background scenes, facilitating the generalization of foreground/background segmentation. Also, some of the images contain two hands, that are both annotated as foreground areas.

4.3. Training Details

We implemented the proposed approach using the Py-Torch framework [35]. The Adam optimizer was used to train it for 100 epochs, with a learning rate value of 10^{-3} , weight decay of 10^{-5} and a learning rate scheduler with $\gamma=0.5$ applied every 30 epochs. For training, we employed an Nvidia GTX 1080 Ti GPU. On that machine, each epoch took about 825 seconds. For all the experiments, the input size to the network was a 256×256 RGB image, and the output a 64×64 depth map.

We split *HandRGBD* into training and test sets to train the proposed method. The training set is composed of 19, 104 samples, while the test set contains 1, 497 samples from sequences that are not included in the training set.

As already mentioned, data augmentation was used in order to increase the generalization of the network. Specifically, each training sample was randomly flipped horizontally with probability 0.5. Also, a random rotation in the range of $[-90^{\circ}, 90^{\circ}]$ was applied. For the random cropping, a bounding box of size 0.8 of the original size was selected. Finally, a random intensity value in the range of [-20, 20] for each color channel was added for color jittering.

4.4. Evaluation Metrics

Assessing depth estimation accuracy: For each hand pixel we consider the absolute difference between ground truth and estimated depth. The first error metric E (in mm) is the average of all these differences for all actual hand pixels and all frames of a test set. A second error metric considers the percentage F(e) of hand pixels in the test set for which the absolute difference between ground truth and estimated depth is less that a threshold e.

Assessing hand/background segmentation: The proposed method also produces a segmentation of the hand re-

 $^{^3}$ Source Code available online at https://github.com/OpenKinect/libfreenect2

Table 2: Ablative study for the proposed hand depth estimation.

Variants of the architecture	Error E (mm)	IoU
0 Mask Stages, 1 Depth Stage	39.75	0.62
1 Mask Stage, 2 Depth Stages	33.16	0.65
1 Mask Stage, 3 Depth Stages	29.04	0.70
1 Mask Stage, 4 Depth Stages	28.05	0.73
1 Mask Stage, 5 Depth Stages	28.83	0.72
2 Mask Stages, 4 Depth Stages	25.00	0.73
3 Mask Stages, 3 Depth Stages	24.64	0.81
4 Mask Stages, 2 Depth Stages	34.42	0.68

gions from the background. To assess this, we compute the *IoU* (Intesection over Union) criterion for this classification.

4.5. Ablative Study

We evaluate different architectural choices (Section 3.2) based on a subset of *HandRGBD*. Specifically, variants of the proposed method were trained on 4,500 images and tested on 500 separate images of the dataset.

An important hyper-parameter of the proposed network is the number of intermediate supervision stages that target the mask segmentation. Experimenting with different training strategies for the proposed network, it became apparent that the hand segmentation mask is an important cue for the task at hand. In a preliminary experiment, the ground truth segmentation mask was provided as a fourth channel concatenated along with the RGB image to the network. This experiment lowered significantly the depth estimation error, indicating that the segmentation mask is indeed useful. It is therefore important to use this cue as an intermediate supervision target, since it aids the task of the network.

In a network with a fixed number of hourglass modules, some of the first hourglass module outputs target segmentation masks and the rest target depths. We performed an experiment to determine the optimal number of stages for each of the two tasks, experimenting also with the total number of hourglass modules. The results of this experiment are presented in Table 2. The best results are highlighted with bold font. From this experiment we can conclude that, in fact, the segmentation cue is equally important to the depth map itself. The best performing network with six stages was trained with the first three of them targeted as segmentation mask and the rest targeting depth.

4.6. Hand Depth Estimation Accuracy

We explored the performance of the best performing variant (line 7 in Table 2) when trained in a larger subset of *HandRGBD*, compared to the experiments in Section 4.5. Specifically, we used our training set of 19, 104 images of

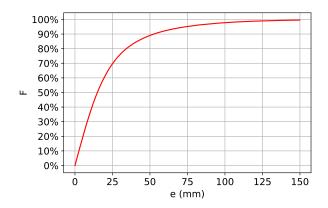


Figure 5: The error metric F(e) for the depth accuracy estimation experiment on the HandRGBD test set (see text for details).

the dataset for training and the rest 1,497 images for testing. In this experiment, the depth estimation error E was equal to E=22.88mm. Figure 5 shows the metric F(e). For this experiment, the estimated IoU was equal to 0.84.

4.7. Supporting 3D Hand Pose Estimation

We assessed the quality of the depth estimated by the proposed method by evaluating the extend at which it can support depth-based hand pose estimation. To do so, we employed the test set part of *HandRGBD* on which we applied a depth-based 3D hand pose estimation method in two different experimental conditions: C1, on the actual depth information of the testset as this was measured by the Kinect2 sensor and C2, the depth that has been estimated by our method. From the available 3D hand pose estimation methods, we chose to employ the tracking approach of Oikonomidis et al. [33]. We selected this method because it depends explicitly on the quality of the employed depth map on which it fits a synthetic hand model. This is contrasted to more recent approaches like the one in [32], where hand pose estimation relies on a learned, indirect function of the hand's depth map.

By comparing the performance of [33] under C1 and C2, we can assess the potential of the proposed method to provide depth maps that are usable by higher level hand perception methods. Ideally, this comparison can be performed by quantifying the 3D hand pose estimation error in C1 and C2 based on 3D hand pose ground truth. However, due to the lack of such ground truth, we follow a different strategy. Specifically, we measured the average distance of the corresponding hand joints as those were estimated by [33] in conditions C1 and C2. This distance was estimated as 24.71mm. Thus, it turns out that the 3D pose discrepancy between C1 and C2 is very similar to the error in depth es-

timation of our approach on this testset. This is very important, as it provides quantitative evidence that an improvement in RGB-based depth estimation will translate directly to an improvement in 3D hand pose estimation.

5. Qualitative Results

Figure 6 shows representative depth estimation results on three sequences of the test set of *HandRGBD*. For each sequence, we show the input RGB image, the ground truth depth map, the estimated one, and their color-coded difference. It can be verified that the depth maps estimated by our method are very close to the ones measured by the depth sensor.

6. Discussion

We presented the first method that has been specifically designed to estimate the depth map of a human hand based on a single RGB frame. The proposed method consists of a specially designed convolutional neural network that has been trained and evaluated on HandRGBD, a new dataset of aligned RGB and depth images. Extensive experiments evaluated design choices behind the proposed method, verified its depth estimation accuracy and provided evidence on the potential of the method to support, providing input, existing depth-based hand pose estimation methods. The obtained results demonstrate that for the specific context of hands observation, the proposed method constitutes an important step towards turning a conventional RGB camera to an RGBD one. Future plans include the consideration of the absolute depth estimation problem and the investigation of the suitability of the proposed approach for other RGBbased depth estimation tasks.

References

- [1] V. Athitsos and S. Sclaroff. Estimating 3D hand pose from a cluttered image. 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition 2003 Proceedings, 2:II–432–9, 2003.
- [2] S. Bambach, S. Lee, D. J. Crandall, and C. Yu. Lending a hand: Detecting hands and recognizing activities in complex egocentric interactions. *Proceedings of the IEEE International Conference on Computer Vision*, 2015 Inter:1949– 1957, 2015.
- [3] H. G. Barrow and J. M. Tenenbaum. Interpreting Line Drawings as Three-Dimensional Surfaces. *Artificial Intelligence*, 17(3):75–116, 1981.
- [4] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Computer Vision (ICCV)*, 2017 IEEE International Conference on, pages 1021–1030. IEEE, 2017.
- [5] Y. Cai, L. Ge, J. Cai, and J. Yuan. Weakly-supervised 3D Hand Pose Estimation from Monocular RGB Images. In *European Conference on Computer Vision*, pages 1–17, 2018.

- [6] Z. Chen, V. Badrinarayanan, G. Drozdov, and A. Rabinovich. Estimating Depth from RGB and Sparse Sensing. arXiv preprint arXiv:1804.02771, pages 1–20, 2018.
- [7] F. Chhaya, D. Reddy, S. Upadhyay, V. Chari, M. Z. Zia, and K. M. Krishna. Monocular reconstruction of vehicles: Combining SLAM with shape priors. *Proceedings - IEEE International Conference on Robotics and Automation*, 2016– June(April):5758–5765, 2016.
- [8] P. Dreuw, T. Deselaers, D. Keysers, and H. Ney. Modeling image variability in appearance-based gesture recognition. Proc. of the ECCV 2006 3rd Workshop on Statistical Methods in Multi-Image and Video Processing (SMVP), 12 May, Graz, Austria, pages 7–18, 2006.
- [9] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. *Proceedings of the IEEE International Conference on Computer Vision*, 2015 Inter:2650–2658, 2015.
- [10] D. Eigen, C. Puhrsch, and R. Fergus. Depth Map Prediction from a Single Image using a Multi-Scale Deep Network. NIPS, pages 1–9, 2014.
- [11] R. Garg, B. G. Vijay Kumar, G. Carneiro, and I. Reid. Unsupervised CNN for single view depth estimation: Geometry to the rescue. Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 9912 LNCS:740–756, 2016.
- [12] F. Gomez-donoso, S. Orts-escolano, and M. Cazorla. Largescale Multiview 3D Hand Pose Dataset. pages 1–23.
- [13] L. He, G. Wang, and Z. Hu. Learning Depth from Single Images with Deep Neural Network Embedding Focal Length. *IEEE Transactions on Image Processing*, (April), 2018.
- [14] T. Heap and D. Hogg. Towards 3D Hand Tracking using a Deformable Model. In *Ieee*, volume 9, pages 140–145, 1996.
- [15] O. Hilliges, D. Kim, S. Izadi, M. Weiss, and A. Wilson. HoloDesk: Direct 3D Interactions with a Situated See-Through Display. *Proceedings of the 2012 ACM annual con*ference on Human Factors in Computing Systems - CHI '12, page 2421, 2012.
- [16] D. Hoiem, A. A. Efros, and M. Herbert. Geometric Context from a Single Image. In *ICCV*, 2005.
- [17] B. K. P. Horn and M. J. Brooks. The Variational Approach to Shape from Shading. *Computer Vision, Graphics, and Image Processing*, 208:174–208, 1986.
- [18] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. Proceedings 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, 2017-Janua:2261–2269, 2017.
- [19] 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.
- [20] A. S. Jackson, A. Bulat, V. Argyriou, and G. Tzimiropoulos. Large Pose 3D Face Reconstruction from a Single Image via Direct Volumetric CNN Regression. *Proceedings of the IEEE International Conference on Computer Vision*, 2017-Octob:1031–1039, 2017.

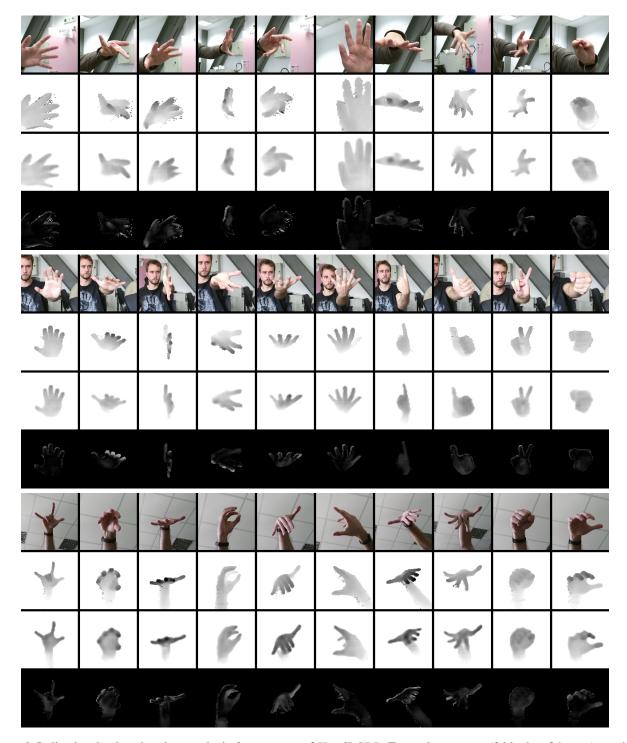


Figure 6: Indicative depth estimation results in 3 sequences of *HandRGBD*. For each sequence (3 blocks of 4 rows) we show the RGB input (1st row), ground truth depth (2nd row), estimated depth (3rd row) and difference between ground truth and estimated depth (4th row).

[21] V. Kanhangad, A. Kumar, and D. Zhang. Contactless and pose invariant biometric identification using hand surface. *IEEE Transactions on Image Processing*, 20(5):1415–1424,

2011.

[22] J. M. Keller, R. M. Crownover, and R. Y. U. Chen. Characteristics of Natural Scenes Related to the Fractal Dimension.

- *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (5):621–627, 1987.
- [23] C. Keskin, F. Kirac, Y. E. Kara, and L. Akarun. 3D hand pose estimation and classification using depth sensors. In Signal Processing and Communications Applications Conference (SIU), 2012 20th, pages 1–4, 2012.
- [24] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. *Proceedings - 2016 4th International Conference on 3D Vision, 3DV 2016*, pages 239–248, 2016.
- [25] J. Li, R. Klein, and A. Yao. A Two-Streamed Network for Estimating Fine-Scaled Depth Maps from Single RGB Images. In Proceedings of the IEEE International Conference on Computer Vision, 2017.
- [26] Y.-H. Lin, W.-H. Cheng, H. Miao, T.-H. Ku, and Y.-H. Hsieh. Single image depth estimation from image descriptors. In 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 809–812, 2012.
- [27] B. Liu, S. Gould, and D. Koller. Single image depth estimation from predicted semantic labels. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1253–1260, 2010.
- [28] F. Liu, C. Shen, and G. Lin. Deep Convolutional Neural Fields for Depth Estimation from a Single Image. *Computer Vision and Pattern Recognition*, pages 1–13, 2015.
- [29] R. Mahjourian, M. Wicke, and A. Angelova. Unsupervised Learning of Depth and Ego-Motion from Monocular Video Using 3D Geometric Constraints. In CVPR 2018, pages 5667–5675, 2018.
- [30] 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. CVPR 2018, pages 29–31, 2018.
- [31] A. Newell, K. Yang, and J. Deng. Stacked Hourglass Networks for Human Pose Estimation. In *European Conference* on Computer Vision, 2016.
- [32] M. Oberweger, P. Wohlhart, and V. Lepetit. Training a feed-back loop for hand pose estimation. *Proceedings of the IEEE International Conference on Computer Vision*, 2015 Inter:3316–3324, 2015.
- [33] I. Oikonomidis, N. Kyriazis, and A. Argyros. Efficient model-based 3D tracking of hand articulations using Kinect. In *British Machine Vision Conference*, pages 101.1–101.11, Dundee, UK, 2011.
- [34] P. Panteleris, I. Oikonomidis, and A. Argyros. Using a single RGB frame for real time 3D hand pose estimation in the wild. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2018.
- [35] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. De-Vito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- [36] G. Pavlakos, L. Zhu, X. Zhou, and K. Daniilidis. Learning to Estimate 3D Human Pose and Shape from a Single Color Image. In *CVPR 2018*, 2018.
- [37] 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. *Cvpr*, 2018.

- [38] M. C. Redmond. Kinect for XBox One.
- [39] J. Reng and T. Kanade. Visual tracking of high dof articulated structures: an application to human hand tracking. Proceedings of the Third European Conference on Computer Vision (ECCV '94), 801(May):35–46, 1994.
- [40] G. Rogez, M. Khademi, J. S. Supancic, J. Montiel, and D. Ramanan. 3D Hand Pose Detection in Egocentric RGB-D Images. In European Conference on Computer Vision Workshop, 2014.
- [41] J. Romero, H. Kjellström, and D. Kragic. Monocular real-time 3D articulated hand pose estimation. 9th IEEE-RAS International Conference on Humanoid Robots, HU-MANOIDS09, pages 87–92, dec 2009.
- [42] M. Salzmann and P. Fua. Linear Local Models for Monocular Reconstruction of Deformable Surfaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):931–944, 2011.
- [43] N. Sarafianos, B. Boteanu, B. Ionescu, and I. A. Kakadiaris. 3D Human pose estimation: A review of the literature and analysis of covariates. *Computer Vision and Image Under*standing, 152:1–20, 2016.
- [44] A. Saxena, S. H. Chung, and A. Y. Ng. Learning Depth from Single Monocular Images. In Advances in neural information processing systems, 2006.
- [45] A. Saxena, M. Sun, and A. Y. Ng. Make3D: Depth Perception from a Single Still Image. *Aaai*, pages 1571–1576, 2008.
- [46] I. Shimshoni, Y. Moses, and M. Lindenbaum. Shape Reconstruction of 3D Bilaterally Symmetric Surfaces. *Interna*tional Journal of Computer Vision, 15, 2000.
- [47] T. Simon, H. Joo, I. Matthews, and Y. Sheikh. Hand keypoint detection in single images using multiview bootstrapping. Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, 2017-Janua:4645– 4653, 2017.
- [48] A. Spurr, J. Song, S. Park, O. Hilliges, and E. Zurich. Cross-modal Deep Variational Hand Pose Estimation. In CVPR 2018, 2018.
- [49] S. Sridhar, A. Oulasvirta, and C. Theobalt. Interactive markerless articulated hand motion tracking using RGB and depth data. In *Proceedings of the IEEE International Conference* on Computer Vision, pages 2456–2463, 2013.
- [50] B. Stenger, P. Mendonca, and R. Cipolla. Model-based 3D tracking of an articulated hand. *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, 2:II–310–II–315, 2001.
- [51] E. Sudderth, M. Mandel, W. Freeman, and A. Willsky. Visual Hand Tracking Using Nonparametric Belief Propagation. In Computer Vision and Pattern Recognition Workshop, page 189, 2004.
- [52] X. Sun, Y. Wei, S. Liang, X. Tang, and J. Sun. Cascaded hand pose regression. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June, pages 824–832, 2015.
- [53] B. J. Super and A. Bovik. Shape from Texture Using Local Spectral Moments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (May), 1995.

- [54] A. Tagliasacchi, M. Schröder, A. Tkach, S. Bouaziz, M. Botsch, and M. Pauly. Robust Articulated-ICP for Real-Time Hand Tracking. In *Computer Graphics Forum*, 2015.
- [55] D. Tang, H. Jin Chang, A. Tejani, and T.-K. Kim. Latent regression forest: Structured estimation of 3d articulated hand posture. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3786–3793, 2014.
- [56] D. Tang, T.-H. Yu, and T.-K. Kim. Real-Time Articulated Hand Pose Estimation Using Semi-supervised Transductive Regression Forests. In 2013 IEEE International Conference on Computer Vision, pages 3224–3231, 2013.
- [57] A. Tewari, M. Zollhöfer, P. Garrido, F. Bernard, H. Kim, P. Pérez, and C. Theobalt. Self-supervised Multi-level Face Model Learning for Monocular Reconstruction at over 250 Hz. In CVPR 2018, 2018.
- [58] A. Tewari, M. Zollhofer, H. Kim, P. Garrido, F. Bernard, P. Perez, and C. Theobalt. MoFA: Model-based Deep Convolutional Face Autoencoder for Unsupervised Monocular Reconstruction. In *ICCV*, pages 1274–1283, 2017.
- [59] A. Tkach, M. Pauly, and A. Tagliasacchi. Sphere-meshes for real-time hand modeling and tracking. ACM Transactions on Graphics, 35(6):1–11, 2016.
- [60] J. Tompson, M. Stein, Y. Lecun, and K. Perlin. Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks. ACM Transactions on Graphics (SIG-GRAPH 2014), 33(5):1–10, 2014.
- [61] A. Torralba and A. Oliva. Depth estimation from image structure. *IEEE Transactions on Pattern Analysis and Ma*chine Intelligence, 24(9):1226–1238, 2002.
- [62] D. Tzionas, L. Ballan, A. Srikantha, P. Aponte, M. Pollefeys, and J. Gall. Capturing Hands in Action Using Discriminative Salient Points and Physics Simulation. *International Journal* of Computer Vision, 118(2):172–193, 2016.
- [63] A. Varol, M. Salzmann, E. Tola, and P. Fua. Template-Free Monocular Reconstruction of Deformable Surfaces. Number ICCV, pages 1811–1818, 2009.
- [64] G. Varol, D. Ceylan, B. Russell, J. Yang, E. Yumer, I. Laptev, and C. Schmid. BodyNet: Volumetric Inference of 3D Human Body Shapes. 2018.
- [65] C. Wan, A. Yao, and L. Van Gool. Direction matters: hand pose estimation from local surface normals. In *European Conference on Computer Vision*, pages 554–569. Springer, 2016.
- [66] R. Y. Wang and J. Popović. Real-time hand-tracking with a color glove. ACM Transactions on Graphics, 28(3):1, jul 2009.
- [67] J. Wu, T. Xue, J. J. Lim, Y. Tian, J. B. Tenenbaum, A. Torralba, and W. T. Freeman. Single Image 3D Interpreter Network. In *ECCV* 2016, volume 1, pages 1–16, 2016.
- [68] J. Wu, T. Xue, J. J. Lim, Y. Tian, J. B. Tenenbaum, A. Torralba, and W. T. Freeman. 3D Interpreter Networks for Viewer-Centered Wireframe Modeling. *International Journal of Computer Vision*, 126(9):1009–1026, 2018.
- [69] L. Xiang, F. Echtler, C. Kerl, T. Wiedemeyer, Lars, hanyazou, R. Gordon, F. Facioni, laborer2008, R. Wareham, M. Goldhoorn, alberth, gaborpapp, S. Fuchs, jmtatsch,

- J. Blake, Federico, H. Jungkurth, Y. Mingze, vinouz, D. Coleman, B. Burns, R. Rawat, S. Mokhov, P. Reynolds, P. Viau, M. Fraissinet-Tachet, Ludique, J. Billingham, and Alistair. libfreenect2: Release 0.2, Apr. 2016.
- [70] C. Xu and L. Cheng. Efficient Hand Pose Estimation from a Single Depth Image. In 2013 IEEE International Conference on Computer Vision, pages 3456–3462, 2013.
- [71] S. Yuan, G. Garcia-Hernando, B. Stenger, G. Moon, J. Y. Chang, K. M. Lee, P. Molchanov, J. Kautz, S. Honari, L. Ge, J. Yuan, X. Chen, G. Wang, F. Yang, K. Akiyama, Y. Wu, Q. Wan, M. Madadi, S. Escalera, S. Li, D. Lee, I. Oikonomidis, A. Argyros, and T.-K. Kim. Depth-Based 3D Hand Pose Estimation: From Current Achievements to Future Goals. In CVPR 2018, 2018.
- [72] S. Yuan and T.-k. Kim. BigHand2 . 2M Benchmark: Hand Pose Dataset and State of the Art Analysis. In *Computer Vision and Pattern Recognition*, pages 15–20, 2017.
- [73] K. Yuen. VIVA Hand Tracking Challenge, 2015.
- [74] 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.
- [75] C. Zimmermann and T. Brox. Learning to Estimate 3D Hand Pose from Single RGB Images. *arXiv preprint arXiv:1705.01389*, 2017.