

# Hierarchically Constrained 3D Hand Pose Estimation using Regression Forests from Single Frame Depth Data

Furkan Kirac\*, Yunus Emre Kara, Lale Akarun

*Department of Computer Engineering, Bogazici University, TR-34342 Bebek, Istanbul, Turkey*

---

## Abstract

The emergence of inexpensive 2.5D depth cameras has enabled the extraction of the articulated human body pose. However, human hand skeleton extraction still stays as a challenging problem since the hand contains as many joints as the human body model. The small size of the hand also makes the problem more challenging due to resolution limits of the depth cameras. Moreover, hand poses suffer from self-occlusion which is considerably less likely in a body pose. This paper describes a scheme for extracting the hand skeleton using random regression forests in real-time that is robust to self-occlusion and low resolution of the depth camera. In addition to that, the proposed algorithm can estimate the joint positions even if all of the pixels related to a joint are out of the camera frame. The performance of the new method is compared to the random classification forests based method in the literature. Moreover, the performance of the joint estimation is further improved using a novel hierarchical mode selection algorithm that makes use

---

\*Corresponding author.

*Email addresses:* `kiracmus@boun.edu.tr` (Furkan Kirac),  
`yunus.kara@boun.edu.tr` (Yunus Emre Kara), `akarun@boun.edu.tr` (Lale Akarun)

of constraints imposed by the skeleton geometry. The performance of the proposed algorithm is also tested on a complex dataset where self-occlusion is frequently encountered. The new algorithm which runs in real time using a single depth image is shown to outperform previous methods.

*Keywords:*

hand gesture, articulated hand pose, depth image, kinect, decision tree

---

## 1. Introduction

Interaction of humans and computers has always been an important area of research. Enhancing the communication between humans and computers increases work efficiency and improves the quality of the work being produced. For that reason, new mice and keyboards are still being designed. In spite of the new designs, the use of classical input devices has become a limiting factor when compared to the capabilities of today's computers. As a result, natural interfaces which make use of speech, touch, and gestures are sought. As a natural interface, speech recognition is widely used and has proved its success in recognizing a multitude of languages. After the release of multi-touch enabled devices, touch interfaces have become also mainstream. However, the real breakthrough will be through the use of the human hand as an input device.

Gesture based communication is very intuitive for humans. A simple way to employ gesture recognition is tracking the position of hands and detecting a gesture. Detecting and tracking the centers of naked hands are relatively simpler when compared to detecting the exact configuration of the articulated pose. Unfortunately, most of the hand gestures require the exact detection

19 of the hand pose and tracking the pose variations in a robust fashion. Even  
20 in the presence of powerful CPUs, articulated hand pose extraction is a very  
21 difficult problem. Moreover, the hand is a small limb that produces self-  
22 occluded poses during gesture performance.

23 With the introduction of inexpensive depth cameras, the human computer  
24 interaction field has been revolutionized; and it is now feasible to establish  
25 methods employing computer vision based interaction. Bypassing the illumi-  
26 nation based problems encountered on the images captured by conventional  
27 cameras, the new depth cameras made it possible to establish very fast and  
28 less power hungry methods. Unfortunately, inexpensive depth cameras that  
29 are widely used still work on low resolution.

30 As discussed before, most of the gestures produce depth images where a  
31 large proportion of the hand parts are unseen. The challenge is to extract  
32 the articulated hand pose from self-occluded, noisy, and low resolution hand  
33 depth images, with algorithms that are suitable for real time implementation  
34 on current CPUs/GPUs.

35 There are several surveys on hand pose estimation and gesture recogni-  
36 tion (Erol et al., 2007; Suarez and Murphy, 2012). Erol et al. (2007) review  
37 hand pose estimation methods. They investigate both partial and full pose  
38 estimation methods. They categorize the full pose estimation methods into  
39 single frame and model-based tracking methods. Most of the works in the  
40 literature focus on grayscale or color based methods. These works use ei-  
41 ther single or multiple cameras. Athitsos and Sclaroff (2003) create a large  
42 synthetic hand pose database using an articulated model and estimate 3D  
43 hand pose from a single frame cluttered image by finding the closest match.

de Campos and Murray (2006) recover hand pose from a single frame using an RVM-based learning method. In order to overcome the self-occlusion problem, multiple views are combined. Oikonomidis et al. (2011) use Particle Swarm Optimization for solving the 3D hand pose recovery problem. de La Gorce et al. (2011) work on monocular videos for 3D hand pose estimation. They track hand poses using a generative model-based method. Thippur et al. (2013) use visual shape descriptors for describing the hand shape.

Recent advances have been made on the depth camera front. Usually, time-of-flight depth cameras are used to acquire depth (range) images (Malassiotis and Srinivas, 2008; Doliotis and Athitsos, 2012; Gallo and Placitelli, 2012; Billiet et al., 2013). Some works make use of disparities to get depth information, and some works use color and depth data together for estimating hand poses (Yao and Fu, 2012; Liang et al., 2012). Liu and Fujimura (2004) use time-of-flight cameras to acquire depth images for hand gesture recognition. Their hand detection method is based on measuring the shape similarity by thresholding the depth data and using Chamfer distance. They recognize gestures using shape, location, trajectory, orientation and speed features of the hand. Mo and Neumann (2006) work on low-resolution depth images acquired from a laser-based camera. Their algorithm requires manual initialization and uses basic sets of finger poses for interpolating a hand pose. Malassiotis and Srinivas (2008) use depth images generated from synthetic 3D hand models. Suryanarayan et al. (2010) use depth for dynamically recognizing scale and rotation invariant poses. Using a volumetric shape descriptor formed by augmenting a 2D image, they classify six signature hand poses. Lee et al. (2012) model hand shape variations using manifolds. Shot-

ton et al. (2011, 2013) use classification forests for human pose estimation using depth data. Keskin et al. (2011) adapted this method to hand pose estimation.

Recent advances have shown that using random classification or regression forests on depth images is a suitable choice for hand pose skeleton extraction, since the recognition phase is very fast and requires minimal algorithmic complexity.

In this work, we investigate the use of random decision forests for the hand pose estimation problem. We adapt the regression forests method to the hand pose estimation task. This method was initially proposed for the human pose estimation area by Girshick et al. (2011). We compare the results of the newly proposed method to the results of Keskin et al. (2011). In addition, we introduce a hierarchical mode selection method that makes use of constraints imposed by the hand skeleton geometry. Regression forests of Girshick et al. (2011) extract multi-modal joint position distributions per joint. They only consider the global mode. However, disregarded local modes of the joint distribution provide invaluable information in the case of self-occlusions and missing data. Knowing all joint configurations are sampled from the same skeletal constraints, provides a strong prior knowledge about the hierarchy of the modes over different joints. We investigate possible skeleton configurations that fit not only on the global modes but also on the local modes. The probable configurations are then filtered out using distance constraints based on a priori positions from the hand skeleton model. At the end, the best skeletal configuration is selected to be the hand pose using dynamic programming.

94 In Section 2, we give the details of random decision forests for classifica-  
 95 tion and regression along with the features we used within these methods.  
 96 Moreover, we introduce our novel hierarchical mode selection method which  
 97 provides a significant improvement to the joint estimation step. Section 3  
 98 includes the details of our synthetic data generation method and the datasets  
 99 along with our cross validated parameter selection and the results on the test  
 100 sets. We conclude with the discussion of the results in Section 4.

## 101 **2. Methodology**

102 This section gives the details of the hand pose estimation methods. The  
 103 outline of the methods discussed are shown in Figure 1. Section 2.1 includes  
 104 description of the random decision forests approach. Sections 2.2 and 2.3  
 105 explain the use of random decision forests for hand pose estimation based on  
 106 classification (RDF-C) and regression (RDF-R), respectively. In Section 2.4,  
 107 we introduce a new hierarchical mode selection algorithm which we call RDF-  
 108 R+.

### 109 *2.1. Randomized Decision Trees*

110 A decision tree is used for inferring a set of posterior probabilities for the  
 111 input. It consists of internal nodes and leaf nodes; where the internal nodes  
 112 propagate the data to one of its children. In the binary case, decisions to  
 113 split the data are simply yes/no decisions. Leaf nodes do not make a decision  
 114 but give statistical information about the nature of the data. The type of  
 115 statistical information depends on the application.

116 In randomized decision trees, the decisions on internal nodes are made  
 117 by selecting a random subset of the features. The aim is to reach leaf nodes

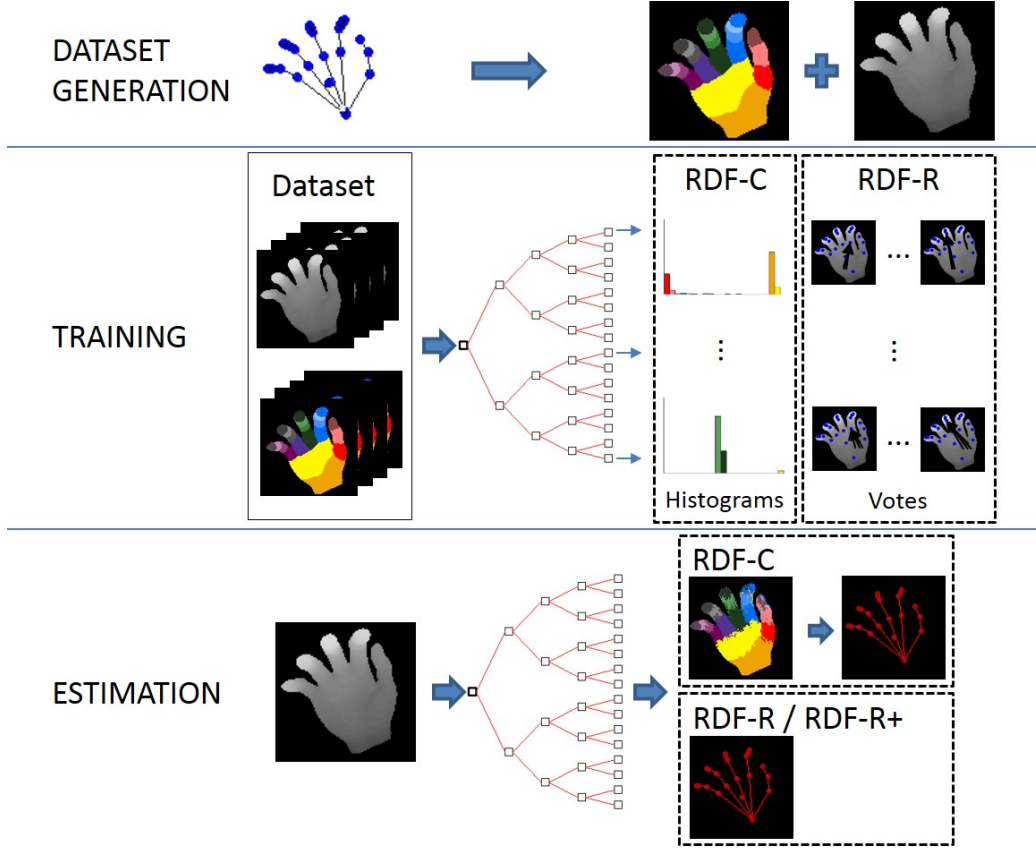


Figure 1: Overview of the methods. The only difference of RDF-C and RDF-R training phases is at the leaves of the trees. RDF-C stores histograms whereas RDF-R stores relative votes. RDF-C classifies each pixel and estimates the joint positions using the estimated pixel labels. However, RDF-R estimates joint positions directly by using the relative votes stored at the leaves.

118 that are as pure as possible. A pure node consists of samples of only one  
 119 class. Thus, the features are selected to yield maximum information gain, in  
 120 other words, minimum entropy. The decision rule is usually of the form:

$$f_n(v) < \tau_n \quad (1)$$

121 where  $f_n(v)$  is a split function,  $v$  is the feature vector and  $\tau_n$  is a threshold,  
 122 at split node  $n$ . A split function is a real valued function on a subset of  
 123 features.

124 During training, the split functions and thresholds of nodes are chosen to  
 125 satisfy the minimum entropy condition. On the leaf nodes, statistics are col-  
 126 lected using the data associated with that node. In the case of classification,  
 127 this is usually a histogram of the class labels of the leaf node data.

128 Decision Forests are ensembles of decision trees. Each tree can be trained  
 129 on the same or slightly different datasets. During testing, the given sample  
 130 is processed in each tree separately. The statistics on the reached leaves are  
 131 combined for a common response. In classification problems, this is usually  
 132 done by accumulating normalized histograms on the leaves.

## 133 2.2. *Hand Pose Estimation using Randomized Decision Forest for Classifi-* 134 *cation (RDF-C)*

135 Shotton et al. (2011) used RDF-C for human body pose estimation. Ke-  
 136 skin et al. (2011) adapted that method to the hand pose estimation problem.  
 137 The aim of the method is to find the pixels closest to each joint. Centroids  
 138 of those pixels would give the joint position. However, since some classifica-  
 139 tion errors are anticipated, it would be better to find the mode of the pixel



positions instead of the mean. For this purpose, first the training data pixels are labeled to define the area around the joints. A decision forest is trained using this data. On the leaves, histograms are calculated using the classes of the pixels.

#### 2.2.1. Pixel Classification using RDF-C

At each node of a randomized decision tree, a random subset of features must be selected and a decision must be made using this subset. Keskin et al. (2011) used the same feature family for hand pose estimation as Shotton et al. (2011). The training data consists of large number of pixels of different depth images. Given a depth image  $I$ , features are computed as

$$F_{\mathbf{u},\mathbf{v}}(I, \mathbf{x}) = I\left(\mathbf{x} + \frac{\mathbf{u}}{I(\mathbf{x})}\right) - I\left(\mathbf{x} + \frac{\mathbf{v}}{I(\mathbf{x})}\right) \quad (2)$$

where  $\mathbf{u}$  and  $\mathbf{v}$  are offsets relative to the pixel position  $\mathbf{x}$ , and they are normalized by the depth at  $\mathbf{x}$ ,  $I(\mathbf{x})$ .

The node data consisting of  $(I, \mathbf{x})$  pairs are split into two sets for each child as

$$C_L(\mathbf{u}, \mathbf{v}, \tau) = \{(I, \mathbf{x}) | F_{\mathbf{u},\mathbf{v}}(I, \mathbf{x}) < \tau\} \quad (3)$$

$$C_R(\mathbf{u}, \mathbf{v}, \tau) = \{(I, \mathbf{x}) | F_{\mathbf{u},\mathbf{v}}(I, \mathbf{x}) \geq \tau\}. \quad (4)$$

Since it is desired to split the data into purer children nodes, the tuple  $((\mathbf{u}, \mathbf{v}, \tau))$  that gives the maximum information gain is chosen among randomly created tuples. Maximum information gain is found using entropy. First, a candidate split is found and the total decrease in entropy that results from this split is calculated. The split score is

$$S(\mathbf{u}, \mathbf{v}, \tau) = H(C) - \sum_{s \in \{L, R\}} \frac{|C_s(\mathbf{u}, \mathbf{v}, \tau)|}{|C|} H(C_s(\mathbf{u}, \mathbf{v}, \tau)) \quad (5)$$

159 where  $H(K)$  is the Shannon entropy estimated using the normalized his-  
 160 togram of the labels in the sample set  $K$ . Then, the candidate tuple that  
 161 yields the maximum score is chosen for the particular node.

162 In classification, a pixel  $(I, \mathbf{x})$  is pushed down the tree until a leaf node  
 163 is reached. At each leaf node, a histogram represents the posterior probab-  
 164 ities  $P(c_i|I, \mathbf{x})$  for each class  $c_i$  learned during the training phase. The final  
 165 decision is made by averaging the posterior probabilities estimated by the  
 166 trees of the forest:

$$P(c_i|I, \mathbf{x}) = \frac{1}{N} \sum_{n=1}^N P_n(c_i|I, \mathbf{x}) \quad (6)$$

167 where  $N$  is the number of trees in the forest.

### 168 2.2.2. Joint Position Estimation using the Classified Pixels

169 For a given depth image, the RDF-C algorithm yields posterior probab-  
 170 ities of each pixel for each class after classification. The resulting probability  
 171 surfaces are generally multi-modal. Thus, simple averaging is not a suitable  
 172 operation. For overcoming the high impact of misclassified pixels on the cen-  
 173 troid of the pixel locations of the same class, a method that is more robust to  
 174 false positives than averaging, must be used. In this situation, mode finding  
 175 is preferred instead of averaging. A local mode finding approach, such as  
 176 mean shift, can be used.

177 First, a Gaussian kernel centered on a random point on the probability  
 178 image is placed. Then, a weighted mean of the probability image under this

179 Gaussian kernel is calculated. Weight indicates the importance of the pixel  
 180 and is an estimate of the area the pixel covers. Weights are calculated as

$$w_{I,\mathbf{x},c_i} = P(c_i|I, \mathbf{x})I(\mathbf{x})^2. \quad (7)$$

181 The newly calculated mean point is used as the starting point of the next  
 182 iteration. This is repeated until converging to a local maximum. For finding  
 183 the global maximum, the algorithm runs several times, each time starting at  
 184 a different initial point and the highest peak is selected.

### 185 *2.3. Hand Pose Estimation using Randomized Decision Forest for Regression* 186 *(RDF-R)*

187 Girshick et al. (2011) proposed a new method for the human body pose  
 188 estimation problem. This technique directly infers the joint coordinates using  
 189 random decision trees without an intermediate pixel classification represen-  
 190 tation, hence making it more robust against occlusion. This algorithm is  
 191 suitable for application to human hand pose estimation. We adapt Random-  
 192 ized Decision Forest for Regression (RDF-R) for directly inferring hand joint  
 193 positions from the depth image without the intermediary per pixel classifi-  
 194 cation phase. RDF-R can learn and estimate the joint positions even under  
 195 self-occlusions. Unlike RDF-Cs, RDF-Rs depend on the mean shift algorithm  
 196 in the training phase, as well.

#### 197 *2.3.1. Training of the Joint Positions*

198 The structure of the trees, namely the features selected in the tree nodes,  
 199 is the same as described in Section 2.2.1. Therefore, the structure is learned

200 in such a way that leaves are favored to store pixels belonging to the same  
 201 part of the hand shape.

202 RDF-Rs do not store hand part histograms at each leaf node  $l$  but a  
 203 distribution over the relative 3D offsets, called a vote, to each hand joint  $j$ .  
 204 These votes are the positions of joints relative to the pixel in question. Each  
 205 training pixel  $q$  is propagated through the tree branches until it reaches a  
 206 leaf node  $l$ . The pixel then casts a relative vote for each distinct joint  $j$ . The  
 207 relative vote can be evaluated using as:

$$\Delta_{iq \rightarrow j} = \mathbf{z}_{ij} - \mathbf{x}_{iq}, \quad (8)$$

208 where  $\mathbf{z}_{ij}$  is the ground truth joint position, and  $\mathbf{x}_{iq}$  is the 3D pixel position  
 209 for a pixel  $q$  belonging to image  $i$ . For a leaf node  $l$ , a relative vote to joint  
 210  $j$  evaluated for pixel  $q$  belonging to image  $i$  is then stored in set  $R_{lj}$ . An  
 211 example multi-modal vote distribution is illustrated in Figure 2.

212 We prefer to use large training sets since we want to infer joint positions of  
 213 different hand pose configurations. All the information represented by a vote  
 214 distribution of a leaf cannot be stored in memory. A consensus of relative  
 215 votes has to be reached per tree leaf for information compression. Unfortu-  
 216 nately the vote distributions are not unimodal. Representing the votes by  
 217 fitting a Gaussian is therefore not suitable. The different clusters of votes  
 218 have to be distinguished as a preliminary phase. Mean shift algorithm is a  
 219 proper candidate for the task. After selecting a suitable kernel, the mean  
 220 shift algorithm finds the number of different clusters and their means. The  
 221 percentage of relative votes belonging to a cluster is the weight of that clus-  
 222 ter. Unfortunately the training phase requires handling of great numbers of

223 votes per tree leaf. In order to learn the relative votes in a reasonable time,  
 224 the vote distribution  $R_{lj}$  is sub-sampled using reservoir sampling of Vitter  
 225 (1985). Reservoir sampling is a single-pass  $O(N)$  algorithm that facilitates  
 226 speeding up the long training phase. Sub-sampling, once a reasonable sample  
 227 size is chosen, does not affect the modes of the vote distributions, thus pro-  
 228 viding a considerable performance increase during the training phase without  
 229 compromising quality.

230 Consequently, we initialize a set  $R_{lj} = \emptyset$  for all leaf nodes  $l$  and joints  
 231  $j$ . A depth image pixel  $q$  is propagated to its respective tree leaf and casts  
 232 a vote which is stored in  $R_{lj}$ . All the reservoir sampled pixels of a leaf  $l$   
 233 cumulatively represent vote distributions for all different joints  $j$ . For all leaf  
 234 nodes  $l$  and joints  $j$  we cluster the reservoir sampled distributions  $R_{lj}$  using  
 235 mean shift and take top  $K$  weighted modes as  $V_{lj}$ .

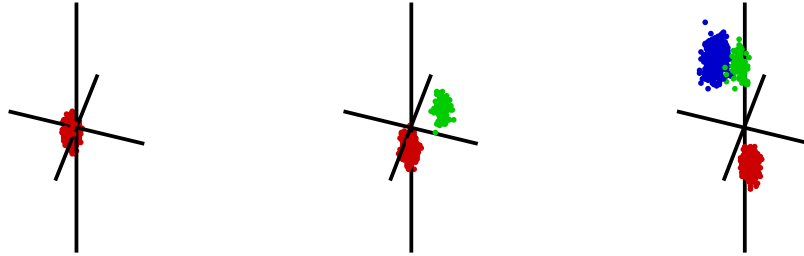


Figure 2: Sample multi-modal vote distributions for three different joints. Vote distributions may have multiple modes. Different colors indicate pixel clusters that are assigned to the same mode found by mean shift.

### 236 2.3.2. Direct Joint Position Estimation using RDF-R

237 We start by initializing  $Z_j = \emptyset$  for all joints  $j$ . All the pixels in a test  
 238 depth image are propagated to the tree leaves by starting at the root node

239 and assigning the pixel either to the left or to the right child recursively until  
 240 a leaf node  $l$  is reached. 3D pixel position of the depth image pixel is recalled  
 241 from the depth image using  $\mathbf{x}_q = (x_q, y_q, z_q)^\top$ . Each test-time depth pixel  
 242 casts its per joint vote as represented by the stored weighted relative vote  
 243 set  $V_{lj}$ . Absolute vote coordinate is evaluated using  $\mathbf{z} = \Delta_{ljk} + \mathbf{x}_q$ . The vote  
 244 is not cast if  $\|\Delta_{ljk}\|_2 \geq \lambda_j$ , where  $\lambda_j$  is a distance threshold learned for each  
 245 different hand joint. To aggregate the absolute votes  $Z_j$ , for each joint  $j$  we  
 246 define a continuous distribution over world space using a Gaussian Parzen  
 247 density estimator such as

$$P_j(z') = \sum_{(z,w) \in Z_j} w \cdot \exp \left( - \left\| \frac{z' - z}{b_j} \right\|_2^2 \right), \quad (9)$$

248 where  $b_j$  is a learned per-joint bandwidth. Running mean shift using Equa-  
 249 tion 9 produces the weighted modes as final hypotheses.

#### 250 2.4. Hierarchical Mode Selection using Geometry Constraints (RDF-R+)

251 RDF-R method outputs posterior distributions of possible joint locations.  
 252 Even though using the global modes of the multi-modal distributions seems  
 253 to be the most straightforward approach, the correct position of the joint  
 254 often corresponds to a local mode. However, considering local modes for the  
 255 joint positions results in multiple skeletal configuration candidates, instead of  
 256 a single one given by the global modes. For selecting a suitable configuration,  
 257 we introduce a hierarchical mode selection method and define a constraint  
 258 function based on our prior knowledge about the hierarchical structure of the  
 259 3D model. We penalize each candidate skeletal configuration according to  
 260 our constraints and select the one with the smallest penalty to be the hand

261 pose.

262 The hierarchy of a skeleton can be defined as

$$H = \{(c, p) : p \text{ is the parent of } c\} \quad (10)$$

263 where  $c$  and  $p$  are joints.

264 Let  $j$  be a joint and  $P(\cdot|j)$  be the posterior distribution of  $j$ 's position.

265 Assuming that the distribution has  $N_j$  modes, we define a mode of this

266 distribution as  $\mathbf{x}_j^{k_j}$  where  $k_j \in \{1, \dots, N_j\}$  and  $P(\mathbf{x}_j^{k_j}|j) \geq P(\mathbf{x}_j^{k_j+1}|j)$ . With

267 this condition modes are ordered decreasingly according to their probabilities.

268 For finding the skeletal configuration, we want to select a mode for each

269  $P(\cdot|j)$ . This selection is performed by dynamic programming so that the total

270 penalty of the hierarchical model is minimized for a given penalty function

271  $f$ ,

$$(k'_1, \dots, k'_J) = \operatorname{argmin}_{k_1, \dots, k_J} \sum_{(c,p) \in H} f(\mathbf{x}_c^{k_c}, \mathbf{x}_p^{k_p}) \quad (11)$$

272 where  $J$  is the total number of joints and  $k_1, \dots, k_J$  are the mode indices

273 of the related joints. Then, the most suitable skeletal configuration can be

274 represented as

$$S' = (\mathbf{x}_1^{k'_1}, \dots, \mathbf{x}_J^{k'_J}). \quad (12)$$

275 By defining a penalty function, different hierarchical constraints can be

276 imposed. Let us define two different penalty functions

$$f_1(\mathbf{x}_c^{i_c}, \mathbf{x}_p^{i_p}) = \begin{cases} 0, & i_c = 1 \text{ and } i_p = 1 \\ 1, & \text{otherwise} \end{cases} \quad (13)$$

$$f_2(\mathbf{x}_c^{i_c}, \mathbf{x}_p^{i_p}) = (\|\mathbf{x}_c^{i_c} - \mathbf{x}_p^{i_p}\| - b_{cp})^2 \quad (14)$$

where  $b_{cp}$  is the expected length of the bone between the joints  $c$  and  $p$ .

Since  $\mathbf{x}_j^k$ s are ordered according to their probability, the global mode of the posterior distribution is  $\mathbf{x}_j^1$ . Thus,  $f_1$  behaves as the global mode finding approach used by Girshick et al. (2011).  $f_2$  tries to select the modes such that the distance between them is as close as possible to the expected bone length.

An example of the improvement made by RDF-R+ method is shown in Figure 3 which shows the vote distribution for tip of the index finger. As clearly seen, there are two different local modes in the distribution. The global mode is not the correct one to be selected. RDF-R wrongly selects the global mode whereas RDF-R+ considers the hierarchical dependencies of the joints and finds the appropriate local mode successfully.

### 3. Experiments

In Section 3.1, we give the details of the hand data generation step and introduce the datasets that we use. Section 3.2 discusses the fine-tuning of the essential training parameters, namely the forest size, the tree depth, the probe distance, the depth threshold, and the mean shift bandwidth. In Section 3.3, the fine-tuned RDFs are tested on four different datasets. We first test on the training dataset to gather information about the characteristics



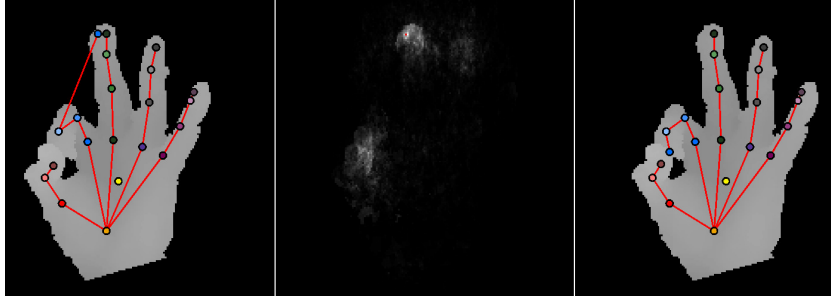


Figure 3: Improvement of RDF-R+ method over RDF-R is shown for tip of the index finger: left) incorrect fingertip estimation of RDF-R, middle) accumulated votes for index fingertip, right) correct fingertip estimation of RDF-R+.

of the methods, and then extend our tests to two different synthetic and one real test datasets.

### 3.1. Data Generation and Datasets

Training the decision trees for classification and regression requires a great amount of training data. Capturing such a big dataset and labeling different parts of the depth images is a problem on its own. In order to cope with this problem, a synthetic 3D hand mesh is modeled and a realistic skeleton is rigged. The produced skeletal object is animated. In this approach, the difficult pixel labeling problem reduces to creating a label texture that is mapped to the hand mesh. We generate the label texture such that each skeleton joint is at the center of one of the labeled parts. The variation of the human hand across different individuals is significantly less than the body. An average sized hand is selected whose length from the bottom of the palm to the tip of the middle finger is 20 cm. Each different joint of the hand is restricted in a manner to mimic the constraints of the human hand. The synthesizer program designed for this study allows for different poses to

312 be stored in keyframes along a timeline. Once an animation is designed, we  
313 animate it using linear interpolation between different keyframes.

314 In this study, we use 40 different hand poses. Poses are mainly selected  
315 from the American Sign Language alphabet. 26 poses represent the letters  
316 from A to Z, and 10 poses are for the numbers between 1 and 10. We also  
317 include four widely used hand poses, namely the closed-hand, open-hand,  
318 approval gesture, and all finger tips touching pose. For each pose, the model  
319 is rotated up to 32, 64, and 64 degrees along the x, y, and z axes, respectively.  
320 Various samples are collected for different angles with steps of six degrees.  
321 Center of the palm is always aligned with the center of the images created.  
322 We also add Gaussian depth noise to the depth image pixels with a mean  
323 of 10 mm, and a standard deviation of 5 for improving the generalization  
324 of the trained classifiers. During hand pose configuration setup, Gaussian  
325 noise is introduced to the angles of the skeleton for the unconstrained degree  
326 of freedoms with a mean of 2 degrees and a standard deviation of 1 degree.  
327 The resulting training (TRAIN) dataset consists of 29766 image samples.  
328 Rendered depth and label images are 160x160 pixels in resolution. In Figure  
329 4, we show sample frames from the TRAIN dataset.

330 Training and validation are done on the TRAIN dataset using 10-fold  
331 cross validation. For testing the performance of different methods, we created  
332 two different datasets. The first one is the cropped (CROP) dataset. It is  
333 generated by retaining the center 80x80 pixels and erasing the outer pixels of  
334 each image from the TRAIN dataset. On the average, 81.09% of the pixels  
335 from each image are kept with a standard deviation of 8.91. The CROP  
336 dataset is used for testing the performance of different methods under missing



Figure 4: Sample frames from the TRAIN dataset

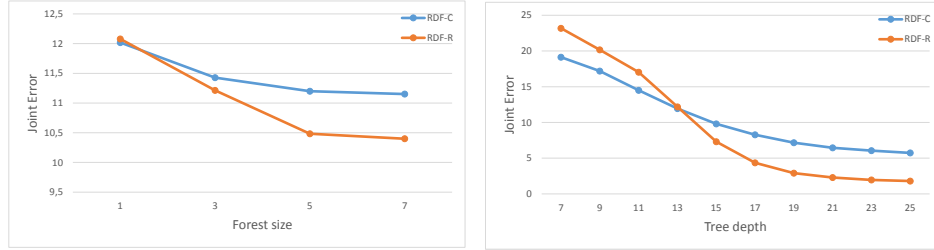
337 data conditions.

338 The second test dataset is the Rock–Paper–Scissors–Lizard–Spock (RP-  
339 SLS) dataset. The Rock–Paper–Scissors–Lizard–Spock game is invented by  
340 Kass (1995). The RPSLS dataset is a completely new dataset synthesized  
341 from 5 different well-known poses not contained in the training set. All  
342 possible transitions between pose pairs are considered and animated using  
343 3 frames per transition. The same rotation conditions are also applied to  
344 the pose animations, resulting in a dataset of 30492 images with 160x160  
345 resolution. The poses of the RPSLS dataset are not used during the training  
346 of the decision trees. This dataset is used to benchmark the generalization  
347 performance of different methods.

348 Additionally, we tested the methods on a subset of ASL Finger Spelling  
349 Dataset of Pugeault and Bowden (2011) for reporting the performance on  
350 real data. The dataset contains sample frames of 5 annotators for 24 finger  
351 spelling signs. Since 3D annotation on depth data is a tedious task, we only  
352 annotated a 55 frame subset of the dataset. The annotated subset consists  
353 of one sample for a, d, e, f, i, l, s, u, v, w, and y signs of each annotator.

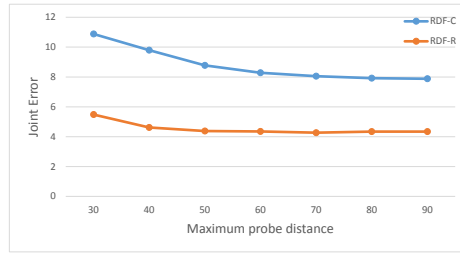
### 354 3.2. *Parameter Selection*

355 We do 10-fold cross validation for fine-tuning the essential training pa-  
356 rameters. We investigate the effects of the forest size, the tree depth, the  
357 probe distance, the depth threshold, and the mean shift bandwidth parame-  
358 ters.



(a) The effect of forest size

(b) The effect of tree depth



(c) The effect of probe distance

Figure 5: Effects of training parameters

### 3.2.1. The Effect of the Forest Size

An advantage of random forests is that the inference performance can be enhanced by combining multiple random trees. Both the generalization capability and the accuracy of the inference improves as the number of trees used is increased. There is a trade-off between inference time and the inference accuracy. Figure 5a shows the effect of number of trees on the accuracy of the system. Accuracy, recognition time, and memory usage are jointly optimized when the number of trees is selected to be 5.

### 3.2.2. The Effect of the Tree Depth

The depth of the trees is also an essential parameter. The representation capability increases as the depth of the trees increase. Unfortunately,

370 increasing the depth of the trees also increases the inference time. Selecting  
 371 an optimal depth is important for balancing the inference accuracy and the  
 372 recognition time. In addition to that, the memory requirement also increases  
 373 exponentially. If the training data complexity and size are not adequate for  
 374 the utilization of the desired tree depth, numerous empty sub-branches in  
 375 the forest structure appear. This under-utilization causes ineffective use of  
 376 the allocated memory. For instance, selecting a tree depth of 20 consumes  
 377 approximately 256 MB of memory per tree.

378 The validation performance of both RDF-C and RDF-R stabilizes after  
 379 a tree depth of 20 levels. They both slowly converge to their maximum  
 380 accuracy as we increase the tree depth. However, it is not feasible to further  
 381 increase the tree depth as the memory required increases exponentially. We  
 382 selected a tree depth of 21 for our tests as it is a good trade-off between the  
 383 accuracy, recognition time, and memory requirement.

384 Another interesting behavior is that RDF-Cs perform better than RDF-  
 385 Rs when tree depth is less than 13 as shown in Figure 5b. This behavior is  
 386 due to the fact that RDF-Cs store class label histograms whereas RDF-Rs  
 387 store 3D relative votes. For a high quality histogram, leaves should have  
 388 numerous pixels, which is the case in a shallow tree. On the other hand,  
 389 RDF-Rs work better with just enough data.

### 390 *3.2.3. The Effect of the Probe Distance*

391 Probe distance is an important parameter which defines the learning  
 392 amount from spatial relations. As we increase the probe distance, even more  
 393 distant depth pixel couples are utilized for inference. With a small probe  
 394 distance value, a more localized recognizer is trained which cannot infer suc-

cessfully using correlation of distant parts or joints of the model. On the other hand, selecting a bigger than needed maximum probe distance value increases the training time. Both the RDF-C and RDF-R methods converge to their optimal performances when probe distance is selected to be 60 for our dataset as seen in Figure 5c.

#### 3.2.4. *The Effect of the Depth Threshold*

Depth threshold is similar to the probe distance. It controls the amount of learning based on depth variations. A very small value forces the learning not to depend on depth differences which produces a silhouette learner. A big value may learn the noise along the depth axis. We selected 30 as the depth threshold.

#### 3.2.5. *The Effect of the Mean Shift Bandwidth*

We used a shared bandwidth parameter for all joints. RDF-C algorithm produces less multi-modal and smoother distributions when compared to RDF-R. RDF-Rs form multi-modal distributions where the distribution peaks are clearly distinguished. Selecting an appropriate bandwidth has a greater influence on the performance of RDF-Rs. Performance of RDF-Cs do not change for reasonable values of bandwidth values whereas the performance of RDF-Rs are clearly dependent on the bandwidth parameter. We selected the bandwidth parameter to be 8.

### 3.3. *Hand Pose Estimation Test Results*

We start by testing all methods with the training dataset for demonstrating the amount of learning each method can achieve. We use CMC curves that report successful joint localization versus an acceptance threshold. For

instance, the acceptance rate at 10 mm shows the percentage of joints that are closer than 10 mm to ground truth locations. RDF-C achieves a performance rate of 76.6% at 10 mm acceptable distance threshold. For RDF-R and RDF-R+, this performance increases to 98.0% and 98.1% respectively. These results clearly show that RDF-C under-performs due to self-occlusions of the hand poses. Both regression forest based methods are robust to self-occlusion, hence, can learn all of the joints with a high accuracy. Figure 6a shows the CMC curve for the TRAIN dataset.

### 3.3.1. CROP Dataset

The hand is a limb that inherently produces self-occlusions. Moreover, depth data may be partly missing due to various other reasons. Parts of the hand may be out of sight of the camera. Similarly hand may be very close to the camera. Depth cameras such as Kinect have zero planes. The pixels closer than the zero plane are clipped. Another common cause of missing data is occlusion imposed by other objects in the environment. Depth cameras provide qualitative information where significant depth differences among neighboring regions occur. Given a depth image, segmentation algorithms are able to mark those regions occluded by other objects with a high accuracy. When those occluded regions are removed, the resulting depth image is an image where some of the valuable data is missing. The CROP dataset is specially designed for testing the inference performance of methods in the case of large amount of occlusions. The RDF-C method is by design not robust against missing data. It produces a transient state of pixel classifications where the valuable information about occluded parts are lost. RDF-R is implemented to cope with this problem. It is robust to occlusion



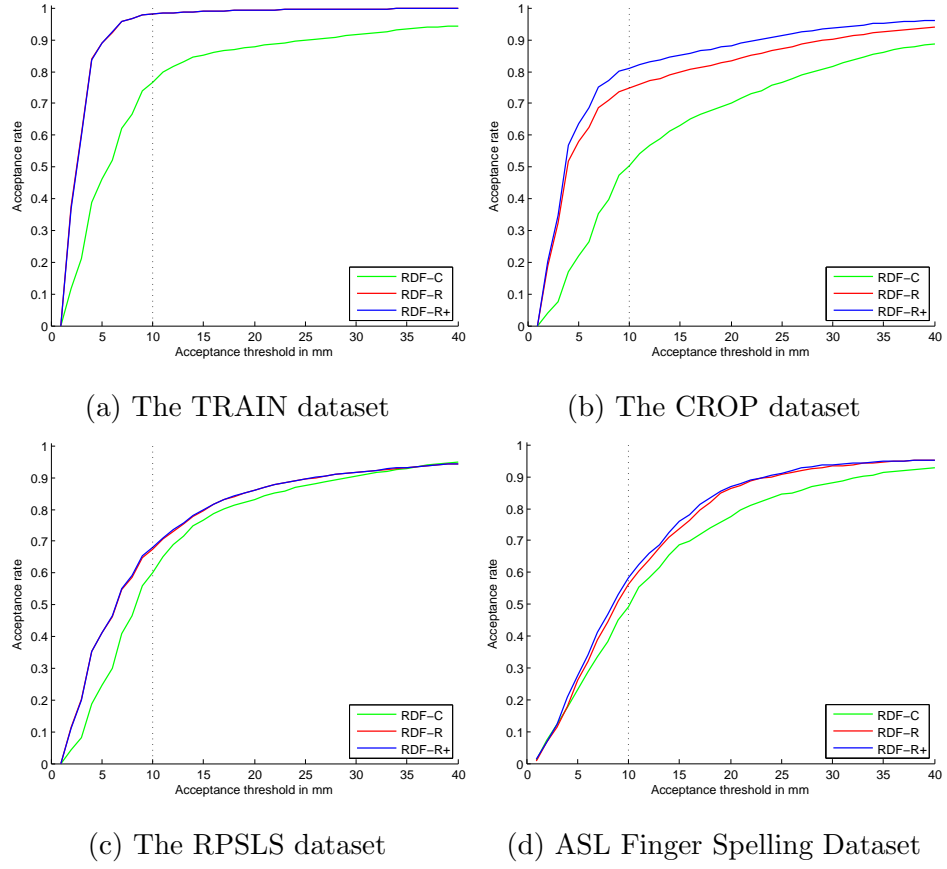


Figure 6: CMC curves for all datasets

444 by design. In addition, RDF-Rs provide multi-modal posterior distributions  
 445 that are suitable for imposing structural constraints. The prior information  
 446 provided by structural constraints of the hand is applied to create a new  
 447 algorithm, namely RDF-R+.

448 Figure 6b shows the CMC curve for the CROP dataset. The performance  
 449 plots clearly demonstrate the strength of RDF-R over RDF-C. 50.4% of all  
 450 joints recognized by the RDF-C method are in a neighborhood of 10 mm  
 451 of ground-truth coordinates. RDF-R method enhances this performance to  
 452 74.8%, which is a very significant improvement. Applying skeletal constraints  
 453 still improves the results. RDF-R+ performs significantly better than RDF-  
 454 R, increasing the performance to 81.3%.

### 455 3.3.2. *Rock-Paper-Scissors-Lizard-Spock (RPSLS) Dataset*

456 The RPSLS dataset is a difficult dataset to recognize by the trees trained  
 457 with the TRAIN dataset. It is chosen to evaluate the sensitivity of the  
 458 methods against extreme situations. The poses used during the creation  
 459 of the dataset are either completely new poses or very similar but different  
 460 poses. In either case the exact poses are not included in the training dataset.

461 Training dataset includes a thumb-up pose which means alright in English  
 462 body language. This pose is similar to the rock-pose of RPSLS dataset. The  
 463 open-hand pose of training dataset is similar to the Spock-pose of RPSLS  
 464 dataset. The open-hand pose is not regarded as similar to the paper pose  
 465 due to their different alignments around z-axis. Training set rotates the  
 466 poses around z-axis only 32 degrees which cannot cover 90 degree difference  
 467 between open-hand and paper poses. None of the poses in RPSLS dataset is  
 468 included in training. During this test, we check for the generalization limits

469 of different methods.

470 Figure 6c shows the CMC curve for the RPSLS dataset. The performance  
471 plots demonstrate the similar relative performance improvements between  
472 RDF-C, RDF-R. RDF-R+, however, behaves similar to RDF-R. 60.2% of all  
473 joints recognized by the RDF-C method are in a neighborhood of 10 mm  
474 of ground-truth coordinates. RDF-R method increases the performance to  
475 67.5%. An interesting result is that the performance of RDF-C is better on  
476 RPSLS dataset compared to the CROP dataset with their respective values of  
477 60.2% and 50.4%. This is caused due to successful recognition rate on poses  
478 that are similar to the training poses and absence of cropping. Inferring in  
479 case of missing data is where RDF-C is weak. Applying skeletal constraints  
480 this time improves the results not so significantly, as its performance is 67.7%.  
481 When we examine the multi-modal posterior distributions of different joints,  
482 we see that the joint configurations are either detected with a high confidence  
483 or not detected at all with a high variance. Constraints cannot improve  
484 the performance significantly since similar poses are already detected well  
485 enough.

### 486 3.3.3. *ASL Finger Spelling Dataset (SURREY)*

487 Figure 6d shows the CMC curve for the ASL Finger Spelling Dataset. The  
488 performance of all algorithms are lower than those for the CROP and RPSLS  
489 datasets. This performance degradation is due to several factors: The first,  
490 is the significantly different characteristics of the training and test datasets:  
491 The test dataset contains many different unseen poses, and variations. The  
492 second is the presence of 5 different subjects, with different hand geometries.  
493 To see which of these factors weighs more in performance degradation, we

494 have looked at performance on different shapes and on different subjects.  
 495 Instead of giving CMC curves, we provide performance figures at the 10 mm  
 496 acceptance threshold in Table 1 and Table 2, for different shapes and different  
 497 subjects, respectively.

Table 1: Percentage of joints that are closer than 10 mm to the ground truth for each sign on ASL Finger Spelling Dataset

ASL Letter	RDF-C	RDF-R	RDF-R+
a	46.3	44.2	<b>48.4</b>
d	49.5	<b>61.1</b>	<b>61.1</b>
e	41.1	<b>53.7</b>	52.6
f	48.4	47.4	<b>53.7</b>
i	41.1	52.6	<b>58.9</b>
l	54.7	65.3	<b>67.4</b>
s	35.8	<b>47.4</b>	45.3
u	55.8	<b>67.4</b>	63.2
v	58.9	58.9	<b>62.1</b>
w	51.6	64.2	<b>65.3</b>
y	58.9	55.8	<b>62.1</b>

498 As observed from Table 1 and Table 2, the performance varies among  
 499 different hand shapes and subjects. For example, the ASL letters a and s  
 500 perform the worst: Upon inspection, it is seen that these ASL letters have  
 501 been performed differently than they are rendered in the training database.  
 502 If one excludes these shapes from the test set, performance increases by  
 503 2.5%. Different subjects, on the other hand, affect performance; somewhat

Table 2: Percentage of joints that are closer than 10 mm to the ground truth for each subject on ASL Finger Spelling Dataset

Subject	RDF-C	RDF-R	RDF-R+
Subject 1	52.2	56.9	<b>61.7</b>
Subject 2	<b>50.2</b>	49.3	47.8
Subject 3	46.4	51.2	<b>56.5</b>
Subject 4	45.9	<b>57.9</b>	56.9
Subject 5	51.7	65.6	<b>67.9</b>

less. The hand size of the subject is an important factor. For instance, RDF-C performs better for subject 2. When subject 2’s live samples are examined, it is seen that she has a considerably bigger palm and shorter fingers than the synthetic model used in training. In some rare cases RDF-R+ algorithm decreases the performance of RDF-R. It is due to imposing bone length constraints which are not very compatible with the test data that is estimated. It is apparent that the system would further benefit from more rigorous training, with poses closer to those in the test set, with different hand shapes, and with different hand sizes.

Overall live data performance of RDF-C, RDF-R, and RDF-R+ at 10 mm acceptance threshold are 49.3%, 56.2%, and 58.2% respectively. Comparing the different algorithms, we observe that RDF-R algorithms perform significantly better than RDF-C; and RDF-R+ has a 2% advantage over RDF-R.

Figure 7 illustrates sample problematic cases for different methods. The average joint estimation performance rates of methods for all four datasets

520 are shown in Table 3.

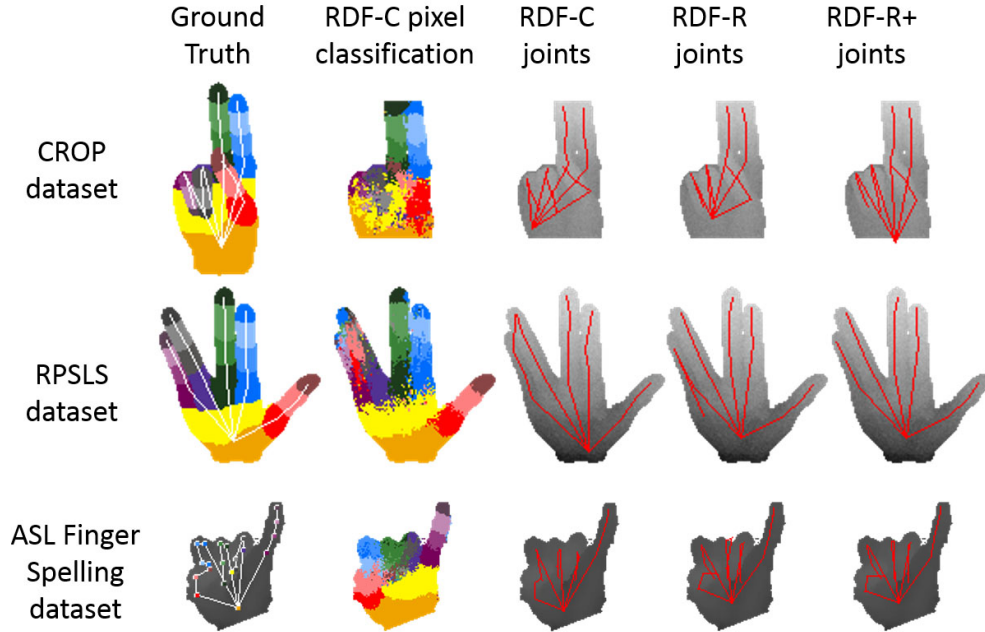


Figure 7: Joint estimation illustrations on test datasets

## 521 4. Conclusions

522 We have demonstrated an implementation of regression forests for esti-  
523 mating the articulated 3D pose of the human hand. Previous attempts at  
524 articulated hand pose estimation used RDF-Cs. We have adapted RDF-  
525 Rs to this problem and implemented an improved hierarchically constrained  
526 version for further enhancing the robustness against heavy occlusion by im-  
527 plementing an algorithm that exploits the prior knowledge about the hierar-  
528 chy of human hand. Considering the hierarchical dependencies of the joints  
529 improved the joint position accuracy significantly. Skeletal constraints are

Table 3: Acceptance rates for the threshold of 10 mm. In all datasets, RDF-R+ method outperforms the other methods.

Dataset	RDF-C	RDF-R	RDF-R+
TRAIN	76.6	98.0	<b>98.1</b>
CROP	50.4	74.8	<b>81.3</b>
RPSLS	60.2	67.5	<b>67.7</b>
SURREY	48.7	55.8	<b>57.6</b>

530 exhaustively evaluated using dynamic programming in real-time. Tests with  
 531 real data have shown us that although performance is lower with tests on  
 532 real data, the results are consistent and performance is still acceptable. In  
 533 order to improve performance, more rigorous training with i) more poses,  
 534 ii) different hand geometries and, iii) real data is left as future work. The  
 535 inference algorithms altogether run with an approximate speed of 200 FPS  
 536 on a conventional notebook computer (Core i7 Quad 2.7 Ghz). Moreover  
 537 the approach only uses a single depth image for inference. Temporal infor-  
 538 mation can still be utilized for extra performance in future studies. Being  
 539 able to detect the hand configuration without using a prior calibration step is  
 540 important for commercial applications. Although this method works with a  
 541 high accuracy, it can also be used as an initialization and/or observation step  
 542 for a temporal domain tracker. For future studies, other skeletal constraints  
 543 can be used and combined. Distances between all different joint pairs can be  
 544 learned from the dataset for applying more restrictive hand configurations.  
 545 Posterior distribution of joints can also be used as an observation step of a  
 546 particle filter that fits a skeleton with a fast local search.

## 547 **References**

- 548 Athitsos, V., Sclaroff, S., 2003. Estimating 3D hand pose from a  
549 cluttered image. 2003 IEEE Computer Society Conference on Com-  
550 puter Vision and Pattern Recognition, 2003. Proceedings. , II-432-  
551 9doi:10.1109/CVPR.2003.1211500.
- 552 Billiet, L., Oramas M., J., Hoffmann, M., Meert, W., Antanas, L., 2013.  
553 Rule-based hand posture recognition using qualitative finger configura-  
554 tions acquired with the Kinect, in: Proceedings of the 2nd International  
555 Conference on Pattern Recognition Applications and Methods, pp. 1-4.
- 556 de Campos, T., Murray, D., 2006. Regression-based Hand Pose Estima-  
557 tion from Multiple Cameras. 2006 IEEE Computer Society Conference  
558 on Computer Vision and Pattern Recognition - Volume 1 (CVPR'06) ,  
559 782-789doi:10.1109/CVPR.2006.252.
- 560 de La Gorce, M., Fleet, D.J., Paragios, N., 2011. Model-Based 3D Hand Pose  
561 Estimation from Monocular Video. IEEE transactions on pattern analysis  
562 and machine intelligence , 1-14doi:10.1109/TPAMI.2011.33.
- 563 Doliotis, P., Athitsos, V., 2012. Hand Shape and 3D Pose Estimation Using  
564 Depth Data from a Single Cluttered Frame, in: International Symposium  
565 on Visual Computing, pp. 148-158.
- 566 Erol, A., Bebis, G., Nicolescu, M., Boyle, R.D., Twombly, X., 2007. Vision-  
567 based hand pose estimation: A review. Computer Vision and Image Un-  
568 derstanding 108, 52-73. doi:10.1016/j.cviu.2006.10.012.



- 569 Gallo, L., Placitelli, A.P., 2012. View-independent Hand Posture Recognition  
570 from Single Depth Images Using PCA and Flusser Moments, in: 2012  
571 Eighth International Conference on Signal Image Technology and Internet  
572 Based Systems, Ieee. pp. 898–904. doi:10.1109/SITIS.2012.133.
- 573 Girshick, R., Shotton, J., Kohli, P., Criminisi, A., Fitzgibbon, A., 2011.  
574 Efficient Regression of General-Activity Human Poses from Depth Images,  
575 in: 2011 International Conference on Computer Vision, Ieee. pp. 415–422.  
576 doi:10.1109/ICCV.2011.6126270.
- 577 Kass, S., 1995. Rock paper scissors spock lizard. URL:  
578 <http://www.samkass.com/theories/RPSSL.html>.
- 579 Keskin, C., Kirac, F., Kara, Y., Akarun, L., 2011. Real-time hand pose  
580 estimation using depth sensors, in: Proceedings of Thirteenth IEEE Inter-  
581 national Conference on Computer Vision Workshops. ICCV 2011, IEEE  
582 Comput. Soc. pp. 1228—1234.
- 583 Lee, C., Chun, S., Park, S., 2012. Articulated Hand Configuration and  
584 Rotation Estimation using Extended Torus Manifold Embedding, in: 21st  
585 International Conference on Pattern Recognition (ICPR 2012), pp. 441–  
586 444.
- 587 Liang, H., Yuan, J., Thalmann, D., 2012. Hand Pose Estimation by Com-  
588 bining Fingertip Tracking and Articulated ICP, in: Proceedings of the  
589 11th ACM SIGGRAPH International Conference on Virtual-Reality Con-  
590 tinuum and its Applications in Industry - VRCAI '12, ACM Press, New  
591 York, New York, USA. pp. 87–90. doi:10.1145/2407516.2407543.

592 Liu, X., Fujimura, K., 2004. Hand gesture recognition using depth data. Sixth  
593 IEEE International Conference on Automatic Face and Gesture Recogni-  
594 tion, 2004. Proceedings. , 529–534doi:10.1109/AFGR.2004.1301587.

595 Malassiotis, S., Srintzis, M., 2008. Real-time hand posture recogni-  
596 tion using range data. Image and Vision Computing 26, 1027–1037.  
597 doi:10.1016/j.imavis.2007.11.007.

598 Mo, Z., Neumann, U., 2006. Real-time Hand Pose Recognition Using Low-  
599 Resolution Depth Images. 2006 IEEE Computer Society Conference on  
600 Computer Vision and Pattern Recognition - Volume 2 (CVPR’06) , 1499–  
601 1505doi:10.1109/CVPR.2006.237.

602 Oikonomidis, I., Kyriazis, N., Argyros, A., 2011. Markerless and efficient  
603 26-DOF hand pose recovery, in: Proceedings of the 10th Asian conference  
604 on Computer vision-Volume Part III, Springer. pp. 744–757.

605 Pugeault, N., Bowden, R., 2011. Spelling it out: Real-time  
606 ASL fingerspelling recognition. 2011 IEEE International Confer-  
607 ence on Computer Vision Workshops (ICCV Workshops) , 1114–  
608 1119doi:10.1109/ICCVW.2011.6130290.

609 Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore, R.,  
610 Kipman, A., Blake, A., 2011. Real-time human pose recognition in parts  
611 from single depth images, in: CVPR.

612 Shotton, J., Sharp, T., Kipman, A., Fitzgibbon, A., Finocchio, M., Blake,  
613 A., Cook, M., Moore, R., 2013. Real-time human pose recognition in parts  
614 from single depth images. Communications of the ACM 56, 116–124.

- 615 Suarez, J., Murphy, R.R., 2012. Hand Gesture Recognition with Depth  
616 Images: A Review, in: 2012 IEEE RO-MAN: The 21st IEEE International  
617 Symposium on Robot and Human Interactive Communication, Ieee. pp.  
618 411–417. doi:10.1109/ROMAN.2012.6343787.
- 619 Suryanarayan, P., Subramanian, A., Mandalapu, D., 2010. Dynamic Hand  
620 Pose Recognition Using Depth Data. 2010 20th International Conference  
621 on Pattern Recognition , 3105–3108doi:10.1109/ICPR.2010.760.
- 622 Thippur, A., Ek, C., Kjellström, H., 2013. Inferring Hand Pose: A Com-  
623 parative Study of Visual Shape Features, in: International Conference on  
624 Face and Gesture.
- 625 Vitter, J.S., 1985. Random sampling with a reservoir. ACM Transactions on  
626 Mathematical Software 11, 37–57. doi:10.1145/3147.3165.
- 627 Yao, Y., Fu, Y., 2012. Real-Time Hand Pose Estimation from RGB-D Sensor,  
628 in: 2012 IEEE International Conference on Multimedia and Expo, Ieee. pp.  
629 705–710. doi:10.1109/ICME.2012.48.