# Survey on 3D Hand Gesture Recognition

Hong Cheng, *Senior Member, IEEE*, Lu Yang, *Member, IEEE*, and Zicheng Liu, *Fellow, IEEE*

*Abstract*—Three-dimensional hand gesture recognition has attracted increasing research interests in computer vision, pattern recognition, and human–computer interaction. The emerging depth sensors greatly inspired various hand gesture recognition approaches and applications, which were severely limited in the 2D domain with conventional cameras. This paper presents a survey of some recent works on hand gesture recognition using 3D depth sensors. We first review the commercial depth sensors and public data sets that are widely used in this field. Then, we review the state-of-the-art research for 3D hand gesture recognition in four aspects: 1) 3D hand modeling; 2) static hand gesture recognition; 3) hand trajectory gesture recognition; and 4) continuous hand gesture recognition. While the emphasis is on 3D hand gesture recognition approaches, the related applications and typical systems are also briefly summarized for practitioners.

*Index Terms*—Begin–end gesture detection, depth sensor, dynamic time warping, hand gesture recognition, skeleton detection and tracking.

## I. INTRODUCTION

**H**AND gestures are elementary movements of a person's hands and are the atomic communication components representing the thoughts of a person [1]. Evolutionary anthropologists tell us that hand gestures have been used since the beginning of human history and are much older than speech [2]. Moreover, hand gestures are a natural, ubiquitous, and meaningful part of spoken language, and researchers have claimed that gesture and sound form a tightly integrated system during human cognition [3]. Inspired by human interaction, mainly by vision and sound, the use of hand gestures is one of most powerful and efficient methods of human–computer interaction (HCI) [4].

There are three basic types of sensors that are capable of sensing hand gestures: 1) mount-based sensors; 2) multitouch screen sensors; and 3) vision-based sensors. In the first case, accelerometers or gyros are used to capture the movement of hands and fingers [5]. Multitouch screen sensors [6], [7] are suitable for mobile devices, but limit the distance between users and computers. On the other hand, there are some advantages of vision-based sensors [8]–[10] since they can be

less cumbersome and uncomfortable than the mounted sensors due to no physical contact with users. Vision-based sensors also provide much larger working distance than multitouch screen sensors. However, the computational complexity is quite high for conventional vision-based hand detection and tracking [11]. To handle this challenge, colored markers or data gloves have been employed to simplify the vision tasks [12]. Although these wearable land markers circumvent the skin segmentation [13], they place additional burden on users and could feel unnatural enough to perform hand gestures, which is the fatal weakness for HCI applications [14].

During the last five years, the progress of commercial 3D depth sensing technologies [15], [16] has greatly promoted the research of hand gesture recognition [17], [18]. An efficient hand segmentation step is the primary step in hand gesture recognition approaches. The 3D depth information can be used to extract hand silhouettes for robust hand gesture recognition in a comfortable and efficient way by simply thresholding a depth map to isolate the hands. The threshold can be approximated according to the depth of the face [17]. A survey of depth-based hand segmentation techniques can be found in [19].

In general, the depth-based hand gesture recognition approaches fall into three categories: 1) static hand gesture recognition [20]; 2) hand trajectory gesture recognition [21], [22]; and 3) continuous hand gesture recognition [23], [24]. All three kinds of hand gestures can leverage the 3D hand modeling [25], [26] for fine-grained gesture (e.g., finger movement) recognition. Static gestures can represent digits, while trajectory gestures describe strokes and letters. The continuous hand gesture recognition can determine when a gesture starts and when it ends from hand motion trajectories. Hence, we can realize the online understanding of hand gestures.

Considered as a kind of human motion [27] in our daily life, hand gestures have been investigated by human motion analysis [28]–[30], which pay more attention to full-body human poses and activities [31]. Recently, there have been several surveys that emphasize general hand gesture recognition [8], [10]. Different from the above literature that gives broad and general gesture recognition reports, we mainly focus on 3D hand gesture recognition and thoroughly summarize the very recent progress in the static, trajectory, and continuous viewpoints, respectively. This survey can cover most emerging 3D hand gesture recognition approaches and systems, and reflect the research trend toward the practical HCI applications from assistive robots to wearable devices with the first person vision (FPV) technology [32], [33].

The remainder of this paper is organized as follows. Section II gives a glance at popular depth sensors that are widely accepted as basic sensing equipments in
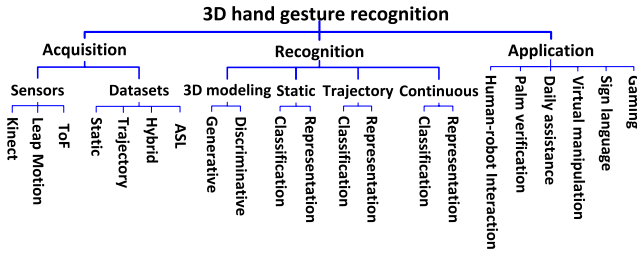
Fig. 1.　Overview of 3D hand gesture recognition techniques.

the community. Section III summarizes various 3D hand gesture data sets as benchmarks for experimental evaluations. The articulated 3D hand modeling by degrees of freedom (DoF) is reviewed in Section IV. The three core topics static hand gesture recognition, hand trajectory gesture recognition, and continuous hand gesture recognition are presented in Sections V–VII, respectively. The 3D hand gesture-based practical applications and systems are introduced in Section VIII. Finally, Section IX concludes this paper. The overview of the surveyed 3D hand gesture recognition techniques is shown in Fig. 1 with representative literatures in Table I. With the bottom-up taxonomy, all the related techniques are classified into three categories: 1) acquisition; 2) recognition; and 3) application. The acquisition works involve various 3D sensors and data sets, which are the foundations of the following gesture recognition. The key components of recognition techniques are feature representation and classification. Finally, the recognized hand gestures can be widely used in different applications. In the following sections, we will review each category in detail.

## II. Depth Sensors

Various 3D depth sensing technologies have been well reviewed by [15], [16], and [29], which mainly introduce the working mechanisms of different depth sensors.[1] In this paper, we focus on three popular depth sensors for hand gesture recognition. In the following sections, we will briefly review Kinect, leap motion, and time-of-flight (ToF) sensors (see Table II for brief comparisons). Although the conventional stereo camera such as Point Grey's Bumblebee can also sense the hand depth, its higher price restricts its applicability to the low-cost HCI application.

### A. Kinect

Since the release of Kinect 1.0 (for Xbox 360) in 2010, it has been the most popular low-cost depth sensor in the computer vision community. The associated OpenNI and Software Development Kit libraries enable the skeletal tracking of human body joints, which provide convenient information for gesture recognition [19]. Depending on whether the Kinect human skeleton [37] is used or not, the Kinect-based hand gesture recognition can be divided into two cases: 1) skeleton-based recognition and 2) depth-based recognition.

In the first case, the skeletal body joints, especially the hand palm joints, are utilized for fast hand detection and

[1]http://www.microsoft.com and https://www.leapmotion.com/

### TABLE I
### Taxonomy of Surveyed Techniques With Corresponding Publications

| Publication | Sensor | Description |
|---|---|---|
| **Acquisition** | | |
| Zhang (2012) [15] | Kinect | Introduction of Kinect |
| Nair et al. (2013) [16] | ToF | Survey of ToF sensors |
| Potter et al. (2013) [34] | Leap Motion | Leap Motion's use |
| Ren et al. (2011) [20] | Kinect | Static hand posture |
| Guyon et al. (2012) [35] | Kinect | Chalearn gesture challenge |
| Wang et al. (2012) [36] | Kinect | 3D ASL dataset |
| **Recognition** | | |
| Shotton et al. (2011) [37] | Kinect | Skeletal body parts |
| Keskin et al. (2013) [38] | Kinect | 3D hand parts |
| de La Gorce et al. (2011) [39] | 2D camera | 3D hand modeling from 2D |
| Oikonomidis et al. (2010) [26] | Multi-camera | 3D hand pose with 26 DoF |
| Oikonomidis et al. (2011) [40] | Kinect | Efficient 3D hand tracking |
| Hamer et al. (2009) [41] | Range camera | 3D object manipulating |
| Oikonomidis et al. (2011) [42] | Multi-camera | 3D hand with an object |
| Oikonomidis et al. (2012) [43] | Kinect | 3D interacting hands |
| Ballan et al. (2012) [44] | Multi-camera | Interacting with an object |
| Van den Bergh et al. (2011) [45] | Kinect | Pointing gestures |
| Cheng et al. (2013) [4] | Kinect | Image-to-class DTW |
| Xu et al. (2012) [46] | Kinect | Robot navigation |
| Zhang and Tian (2013) [47] | Kinect | 3D daily hand gestures |
| Reyes et al. (2011) [48] | Kinect | Begin-end gesture |
| Cheng et al. (2014) [49] | Kinect | Continuous hand gestures |
| **Application** | | |
| Kurakin et al. (2012) [18] | Kinect | Sign language recognition |
| Chai et al. (2009) [50] | Range camera | Virtual manipulation |
| Son and Sowmya (2013) [51] | Kinect | Daily assistance |
| McKeague et al. (2013) [52] | Kinect | Human-robot interaction |

### TABLE II
### Consumer Depth Sensors

| Sensor | Resolution | Range | Accuracy | Description |
|---|---|---|---|---|
| Kinect 1.0 | $320 \times 240$ | 0.8-4.0m | 4mm | 20 body joints |
| Leap Motion | $640 \times 240$ | 25-600mm | 0.01mm | 27 hand joints |
| Kinect 2.0 (ToF) | $512 \times 484$ | 0.8-4.5m | 1mm | 25 body joints |

tracking [53]. Due to the advantage of fast and low-cost hand extraction, hand trajectory gesture recognition has been applied to various HCI applications such as TV control [54]. In addition to skeleton-based hand trajectory gesture, the hybrid gestures [55], [56] containing both trajectory and static gestures are also recognized leveraging the Kinect skeleton system. With the power of Kinect skeleton tracking libraries, 3D hand motion-related sign language vocabularies have been successfully recognized and translated for English [57] and Chinese languages [58].

In the second case, the Kinect sensor is considered as the pure depth sensing device and all the sequential procedures are based on the Red-Green-Blue-Depth (RGB-D) data without activating the inherent skeleton system [59]. Consequently, the hand gesture recognition can be applied not only to online applications but also to the offline recorded RGB-D videos. The independence of online skeleton library is crucial for developing and evaluating hand gesture recognition approaches on public RGB-D data sets. In contrast to the skeleton provided by the official libraries that identify only the palm center, finer hand modeling such as finger joints can be accurately achieved using the RGB-D frames [40], [60], [61]. We will further discuss 3D hand modeling in Section IV. Among the works in this case, the most straightforward use of RGB-D data is the depth-based hand [45], [62], [63]/finger [64], [65] detection and tracking [66], [67]. After hand detection and tracking, either static hand gesture recognition [68]–[70]

TABLE III

SUMMARY OF 3D HAND GESTURE DATA SETS

| Dataset | Description | Availability |
|---|---|---|
| **Static Gesture Datasets** | | |
| 10-Gesture (2011) [20] | Digit gestures with wrist belt | http://eeeweba.ntu.edu.sg/computervision/people/home/renzhou/HandGesture.htm |
| ASL Finger Spelling (2011) [88] | Alphabet gestures | http://lifeprint.com/ |
| UESTC-ASL (2013) [4] | Digit gestures | http://www.uestcrobot.net/?q=download |
| **Trajectory Gesture Datasets** | | |
| ChaLearn (2012) [35] | Upper body gestures | https://www.kaggle.com/c/GestureChallenge |
| MSRC-12 (2012) [89] | Two-hand gestures | http://research.microsoft.com/en-us/um/cambridge/projects/msrc12/ |
| MSR Daily Activity (2012) [90] | Daily gestures | https://research.microsoft.com/en-us/um/people/zliu/actionrecorsrc/default.htm |
| UESTC-DHG (2013) [49] | Stroke gestures | http://www.uestcrobot.net/?q=download |
| LIRIS (2014) [91] | Annotated gestures | http://liris.cnrs.fr/voir/activities-dataset/ |
| **Hybrid Gesture Datasets** | | |
| 3DIG | Iconic gestures | http://projects.ict.usc.edu/3dig/ |
| SKIG (2009) [92] | Static and trajectory hand postures | http://lshao.staff.shef.ac.uk/data/SheffieldKinectGesture.htm |
| **ASL Datasets** | | |
| MSR Gesture3D (2012) [18, 36] | 12 ASL gestures | http://research.microsoft.com/en-us/um/people/zliu/actionrecorsrc/default.htm |
| 3D Body Part Detection Video [93] | Both one-handed and two-handed signs | Not available |

or hybrid hand gesture recognition [71], [72] can be applied.

### B. Leap Motion

Unlike the Kinect sensor that captures the full-body depth, the newly released leap motion focuses on the accurate 3D hand positioning. The leap motion sensor can detect hands and fingers with an accuracy of around 0.01 mm. After its release in the February of 2013, many researchers have considered it as a promising 3D sensor, particularly suitable for 3D hand gesture recognition. There have been several pioneer works using leap motion for the HCI applications. Leap motion sensor was first used to track the users' fingers to enable a hand-controlled interface in a virtual environment [73]. Meanwhile, the leap motion was also employed in the 3D molecular graphic systems [74]. Although leap motion has the potential to recognize more complex 3D hand gestures such as the Brazilian sign language [75], there still is little work in this direction. Since the accuracy [76] and the suitability of the leap motion controller for sign language recognition [34] have been well reviewed, we expect that these evaluations for leap motion can help researchers to further develop its toolkits for the 3D hand gesture recognition application.

### C. Time-of-Flight Sensors

Before the release of Kinect and leap motion, the ToF camera (e.g., Zcam [77]) had been considered as a low-cost depth measurement device [78] and was widely used in the fields of computer vision and HCI. In particular, the newly released Kinect 2.0 is also a ToF camera, which is different from Kinect 1.0 that uses a light coding technique. ToF-based 3D hand gesture recognition system can be used to control a robot by point gestures [79]. Regarding the static hand posture recognition, ToF sensors have been reported to be quite efficient and robust to the hand's orientation, size, and cluttered backgrounds [80]. Due to its accuracy and robustness for the distance measurement task, the ToF sensor has been used to generate the 3D point cloud and helped in recognizing complex hand trajectory gestures such as Polish Sign Language [81].

## III. DATA SETS

Among numerous public data sets[2,3,4] in vision field, both 2D [82], [83] and 3D [84] human action data sets have been provided for body gesture research. While those data sets collect very limited hand gestures such as hand waving, 3D hand gesture data sets have attracted more research interests for solving real-world problems [85]. A comprehensive survey for human gesture data sets can be found in [86]. In the following sections, we will focus on the emerging 3D hand gesture data sets and will categorize them from static, trajectory, and hybrid viewpoints, respectively (Table III). Furthermore, the well-known American Sign Language (ASL) data set [87] and its derived versions will also be reviewed.

### A. Static Gesture Data Sets

The static 3D gesture data sets usually capture the palm and finger postures in the RGB-D domain, which can represent basic symbols such as Arabic numerals. The 10-Gesture data set [20] was the first that collected 3D static hand gestures with a Kinect sensor. The data set was collected from ten subjects who performed ten different hand poses using the black belt on the gesturing hand's wrist. Hence, the 10-Gesture data set has 100 gesture samples in total. Each of the samples consists of a color image and a depth image. The data set was quite challenging since it was collected in uncontrolled environments with cluttered backgrounds. Besides, the subject's gesture varies in hand orientation, scale, and articulation.

Other public static hand gesture data sets include the ASL Finger Spelling data set [88] and the University of Electronic Science and Technology of China (UESTC)-ASL data set [4]. The former contains alphabet signs recorded from 4 different persons, amounting to a total of 48 000 samples. The latter has a digit sign data set with 100 samples, which is recorded from 10 subjects in different orientations, depths, and scales. In particular, the UESTC-ASL data set is more challenging than the 10-Gesture data set since no wrist belt is used to help the hand segmentation, leading to more practical but difficult gesture recognition.

[2]http://homepages.inf.ed.ac.uk/rbf/CVonline/Imagedbase.htm
[3]http://www.cvpapers.com/datasets.html
[4]http://datasets.visionbib.com/

## B. Trajectory Gesture Data Sets

The most popular and well-known trajectory gesture data set is the ChaLearn Gesture data set [35], which has been used in the series of ChaLearn Gesture Challenge. It mainly focuses on hand and arm gestures for human–machine interaction. Some gestures have full-body presentations, but most of them only have upper body. Gestures are separated by returning to a resting position. The ChaLearn Gesture data set has become the most important benchmark to evaluate 3D hand trajectory gesture recognition systems. Another 3D trajectory gesture data set is the Microsoft Research Cambridge-12 Kinect gesture data set [89], which consists of sequences of human movements, represented as body part locations, and the associated 12 gestures. It contains trajectories of 20 joints esti-mated using the Kinect pose estimation pipeline. The gestures are mainly represented by the two-hand movements. There are some other action/activity data sets such as MSR Action 3D, MSR Daily Activity data set [90], UESTC-Dynamic Hand Gesture [49], and Computer Science Laboratory for Image Processing and Information Systems human activity data set. More or less, they all provide several easy hand trajectory gestures in the RGB-D domain.

## C. Hybrid Gesture Data Sets

Sometimes, static hand pose and dynamic hand moving will together represent the gesture. The 3D Iconic Gesture data set contains many pictorial gestures, which depict entities, objects, or actions. The subjects perform iconic gestures to refer to entities through embodying their shapes. People can gesture the outline of an object either with static hand palm or dynamic trajectories using their imagination. The Sheffield Kinect Gesture (SKIG) data set [92] collects ten common categories of 3D hand gestures, which are performed with three hand postures: 1) fist; 2) index; and 3) flat. Consequently, both hand tracking and shape recognition should be processed for final gesture recognition. Although the IDIAP Hand pose/gesture data sets[5] include both static hand posture data sets [94] and hand trajectory posture data sets [95], only 2D gestures are collected without the RGB-D information.

## D. American Sign Language

ASL has been considered as the fourth most used language in the United States since the early 1970s [96]. Although the 2D ASL Lexicon Video data set has collected about 1500 lexical signs [97] in the National Center for Sign Language and Gesture Resources,[6] there still is very limited 3D ASL data sets publicly available for this research purpose. To the best of our knowledge, there are only two public 3D ASL data sets. The first one is the MSR Gesture 3D data set [18], [36], which includes 12 ASL trajectory gestures performed by 10 subjects. In this data set, the hand part has been segmented as the hand detection result. The other data set is the 3D Body Part Detection Video data set [93], which contains 1113 signs both one handed and two handed with the skeleton information generated by the Kinect sensor. The data set is expected to include most of the 3000 signs found in the Gallaudet Dictionary of ASL and serve as a benchmark for assessing future 3D hand gesture recognition algorithms.

## IV. 3D HAND MODELING

3D hand modeling and tracking estimate the articulated hand poses and motions [98]. These are key technologies in many HCI applications such as robot surgery, virtual keyboard, and sign language recognition [99]. According to the taxonomy proposed in [32], hand pose estimation approaches can be classified into two categories: 1) discriminative approaches and 2) generative approaches. Moreover, the hybrid 3D modeling that combines discriminative modeling and generative modeling has also been proposed. Next, we will review various 3D modeling approaches in the context of 3D hand pose estimation and tracking with respect to these three categories.

## A. Discriminative 3D Hand Modeling

Discriminative approaches do not constrain the 3D hand model with explicit DoF. Instead, they train the classifiers that inversely map appearance-specific hand pixel features to unknown hand parameters (part label, pose parameter, etc.). The classifiers are usually learned offline from a large set of training samples. Besides, most of those approaches employ the decision tree-based classifiers to accelerate the estimation in each frame independently.

*1) Part-Based Modeling:* Shotton *et al.* [37] proposed the human pose recognition technique, which finally extracts skeletal body parts with a Kinect sensor. They constructed the randomized decision forests with 300k synthetic training images per tree to determine the body part label for each pixel in the depth domain. Inspired by this seminal work, Keskin *et al.* [38], [61] generalized the idea of [37] to the 3D hand pose estimation using depth data. They used a 3D skinned mesh model as the hand pose generator to synthesize up to 200k training images to learn the randomized decision trees (RDTs) for hand parts. Following [37], the simple depth comparison feature for each pixel $x$ in the depth image $I$ can be computed as

$$f_{u,v}(I, x) = I\left(x + \frac{u}{I(x)}\right) - I\left(x + \frac{v}{I(x)}\right) \qquad (1)$$

where the offsets $u$ and $v$ are normalized by the pixel depth to ensure that the feature is depth invariant.

During the training step, the features of training samples are input to the RDTs to learn the optimum offsets at each leaf node. After the training of all the trees, the posterior probability $P(c|I, x)$ of part label $c$ with the depth pixel $x$ can be inferred by tree nodes. The final classification score is given by averaging all the distributions together in each part's forest. The pixel will be labeled as the specific hand part that provides the highest score. To further extract the hand skeleton, mean-shift algorithm can be used to find the centroid of each part [37], [61]. However, the part forests are trained in a large synthetic data set with considerable memory requirements.

---

[5]http://www.idiap.ch/resource/gestures/
[6]http://www.bu.edu/asllrp/ncslgr.html

*2) Efficient Part Labeling:* Based on [61], Keskin *et al.* [100] further proposed the shape classification forest to classify hand shapes rather than parts. The pixel feature computation here is exactly the same as (1). Every pixel of the input depth image will vote for a single shape label for the hand image

$$c^* = \arg\max_c \frac{1}{N} \sum_{p=1}^{N} P(c|I, x_p) \tag{2}$$

where $N$ is the number of foreground hand pixels and $c$ is the hand shape label.

With the power of shape recognition, the RDTs can be trained with a limited amount of variation, leading to the more efficient and accurate hand part estimation. To benefit from both synthetic and realistic data, Tang *et al.* [101] proposed the semisupervised transductive regression (STR) forest, which associated the sparsely labeled target domain (realistic depth images) with a fully labeled source domain (synthetic depth image). By considering the realistic synthetic discrepancies, the STR forest can improve the accuracy of hand pose estimation with a low labeling cost. The depth feature defined in (1) is simple but sensitive to background changes. In order to generate background-independent pixels features, Yao and Fu [102] fused three features: 1) shape feature; 2) surface feature; and 3) position feature from both RGB and depth cameras, which will be fed into the random forest (RF) classifier for part labeling. In contrast to [37] and [61], they employed the color glove to collect real labeled training data rather than synthetic data to circumvent the nontrivial sample rendering.

Although the discriminative approaches are based on a single frame such that there exists no track drifting issue and can be performed in real time, they require a large set of high-quality training data to train the hand part recognizers. The lack of kinematic constraint makes them less robust to articulate motions and self-occlusions.

### B. Generative 3D Hand Modeling

The generative approaches are popular among recent state-of-the-art 3D hand modeling and tracking works. They are also called model-based hand tracking approaches, which try to fit the explicit 3D DoF model to the observed hand data [103]. Note that either 2D or 3D hand images can be used to estimate the 3D hand model.

*1) Hand Modeling From 2D Images:* De La Gorce *et al.* [39] and de La Gorce and Paragios [104] proposed generative approaches to estimate 3D hand poses from monocular 2D images. De La Gorce and Paragios [104] modeled the hand by an articulated kinematic tree with 28 DoF. The original model was rendered by ellipsoids and polyhedra to form the hand surface model. Then, the hand silhouette could be extracted by differentiating the parameterized hand surface model. Finally, the optimal hand pose was simultaneously estimated with the foreground–background segmentation task in a maximum-likelihood framework by matching the synthesized silhouette with the observed frame. During the model fitting process, the kinematic constraints were also incorporated to prune the unrealistic parameter configurations and accelerate the optimization. De La Gorce *et al.* [39] further incorporated texture and shading information in the objective function to handle the ambiguities that cannot be solely distinguished by the hand silhouette. Compared with [104], the hand surface in [39] was more carefully synthesized by the triangulated mesh with computer graphics techniques such as shading model and texture mapping. Hence, the self-occlusion and time-varying illumination can be properly handled by the use of shading and texture. However, the 3D hand modeling from monocular 2D image requires rather high computational cost and limits its applications.

*2) Single 3D Hand Modeling:* Oikonomidis *et al.* [26] proposed to recover the 3D hand pose by matching the 26-DoF hand model to multiple camera views. The rendering of the 3D hand surface employed a sphere and a truncated cylinder as 3D basic shapes. Both hand silhouette and edge maps were computed as the observation cues to evaluate the consistency between the hypothesized pose and the observed multiple views by minimizing the objective function

$$E(h, M) = \sum_{I \in M} D(I, h, C(I)) + \lambda \cdot kc(h) \tag{3}$$

where $h$ is the hand pose hypothesis and $I$ is one of the images from the multiview set $M$. $D$ is the likelihood function that measures the dissimilarity between the projected hand surface [by camera matrix $C(I)$] and the observed image in the feature space. $kc(h)$ represents the kinematic constraint for the hand configuration.

In [26], the final optimization of (3) was carried out with the particle swarm optimization (PSO) technique and Graphics Processing Unit (GPU) implementation in real time. Their multicamera system consists of calibrated cameras with numbers ranging from two to eight. Another multicamera-based 3D hand modeling system was proposed in [105], in which different features were integrated to evaluate the observation likelihood function with particle filtering. Oikonomidis *et al.* [40] proposed an efficient generative 3D hand tracking approach using a Kinect sensor. The rendered 3D hand surface model is similar to the model of [26]. Moreover, the observation cues included both hand silhouette and depth information, leading to the more efficient objective function and simpler capturing system. With the power of Kinect sensor, the new objective function was based on the 3D structure and was robust to 2D illumination changes. Qian *et al.* [106] modeled a hand simply using spheres. A fast cost function was proposed to measure the distance between the rendered model and the observed 3D point cloud. They further combined gradient-based and stochastic optimization techniques to achieve the fast and robust hand tracking without using GPU.

*3) 3D Hand Modeling With a Manipulated Object:* All the above 3D hand tracking approaches can provide promising hand pose estimation results when the hand is solely observed in isolation. However, they may not work well when the hand is manipulating an object. To handle the strong occlusion by the manipulated object, Hamer *et al.* [41] proposed a system to recover the articulated 3D structure of the hand during object manipulation. They employed the human hand model

consisting of 27 bones and rendered each hand segment by a cylinder with a mesh approximating the skin. Both hand segmentation and depth information were used to evaluate the likelihood probability between the hand configuration and the observed RGB-D data. Moreover, the softened anatomical constraints and occlusion model were incorporated in a Markov random field, which was optimized by the belief propagation. Hamer *et al.* [107] further modeled the object-specific prior from sparse training data to improve the robustness of their 3D hand tracker. Oikonomidis *et al.* [42] extended [26] by jointly estimating the 26-DoF hand pose and the model parameters of the manipulated object. The explicit parametric 3D modeling of the interacted object effectively represents the context of the hand and improves the pose estimation both for the hand and the object.

*4) 3D Interactive Hand Modeling:* Maybe the most challenging scenario for 3D hand modeling task is the situation with interacting hands. Based on [40] and [42], Oikonomidis *et al.* [43] proposed a parametric model of the joint kinematics of two hands to track the full articulation of the strongly interacting hands. The joint model has 54 DoFs (27 for each hand) with collision constraints, and the RGB-D data are utilized with a Kinect sensor. They showed promising results of the optimized joint hand configuration. It was also demonstrated that the straightforward pose estimation for single hand tracking problem would result in a much lower accuracy in the interacting case. Ballan *et al.* [44] addressed the challenging problem of capturing the articulated motion of two hands that interact with each other and with an additional object. Regarding the observation cues, they employed multiple visual features such as edges, optical flow, salient points, and collisions to estimate the articulated pose within a single differentiable function. Compared with [43], which utilized Kinect sensor and PSO optimization, Ballan *et al.* [44] adopted a multicamera setup and the simple local optimization technique, which achieved lower pose estimation errors.

### C. Hybrid 3D Hand Modeling

Recently, the hybrid 3D hand modeling frameworks that combine discriminative and generative approaches have also been proposed. Xu and Cheng [108] presented a three-step pipeline in which the first two steps were based on the discriminative depth features of [37] with the Hough forest regression model [109]. The initial hand pose estimation was followed by the final verification step that optimized a 27-DoF hand model in the generative way. At the same time, Sridhar *et al.* [110] proposed to combine a discriminative part-based pose retrieval approach with a generative pose estimation approach based on local optimization. Both works showed that the combination of discriminative and generative ideas can achieve state-of-the-art hand modeling accuracy as well as high efficiency.

## V. STATIC HAND GESTURE RECOGNITION

After the 3D hand modeling, one may say that it is straightforward to recognize the hand shape or posture due to the rich information in the estimated hand model. However, sometimes we are mainly interested in the meaning of the global hand appearance and the detailed hand skeletal model is not necessary. In this section, we will focus on 3D approaches that were developed for recognizing the static hand shapes.

### A. Feature Representation

The hand features are usually extracted from the bounding box of the segmented hand. Both the RGB and depth information can help in detecting the hand area. Here, we assume that the region of interest (ROI) has been determined, and thus focused on the discriminative features for the ROI. Basically, the hand region's features could be roughly divided into three categories: 1) low-level features; 2) middle-level features; and 3) high-level features.

The low-level hand features can be generated either from the original 3D spatial domain or from a transformed domain. Keskin *et al.* [111] formulated the hand feature by the set of pixel-wise depth features as defined in (1) [37]. This depth feature accounts for different hand positions in a large feature space leading to the time-consuming matching. In order to significantly reduce the training time and the memory consumption, Kuznetsova *et al.* [112] proposed to use ensemble of shape function (ESF) descriptor as the feature of the segmented hand region. The ESF descriptor consists of concatenated histograms that are generated with the random points in the point cloud. In addition to the features in the original depth domain, Haarlet coefficients [45], Gabor coefficients [88], and Flusser moment invariants [113], [114] have been extracted from transformed hand images to efficiently form rotation- and intensity-invariant hand features.

In contrast to low-level hand features that are based on the global hand information, the midlevel hand features are usually based on the local patch-level descriptors [115]. The hand region can be divided into rectangle cells [116] or cylinder sectors [117]. The distributions or histograms of points in all the cells are concatenated to formulate the local descriptors. To leverage the strong shape information in the depth map, Zhang *et al.* [118] defined a 3D facet as a 3D local support surface associated with each 3D cloud point and proposed the histogram of 3D facets to represent the 3D hand shape. Not only local descriptors in regular grids, the local descriptor locates at the interest point has also been used to formulate the hand feature. Bagdanov *et al.* [119] represented the hand by the concatenation of the five Speeded Up Robust Features (SURF) descriptors for a total of 640 dimensions.

Rather than the rough hand bounding box for the above low-level and middle-level features, the detailed hand contour or hand part is necessary for extracting the high-level hand features [120]–[122]. Ren *et al.* [20] represented the hand shape by the time-series curve where each finger corresponded to a segment of the curve. The time-series curve records the relative distance between each contour vertex to a center point and reserves the topological information as well. A similar time-series curve-based hand feature was proposed in [4], in which fingerlets from those time-series curves using different finger combinations were generated. Each hand gesture was

represented by the fingerlet ensemble. The high-level hand features can also be generated in a transformed domain. Liu *et al.* [123] explored the invariance from the hand contour and utilized Fourier descriptor, edge histogram, and boundary moment invariants for the feature extraction. Besides, the number and the direction of the extracted finger tips have also been considered as features for the static hand gestures [124].

Note that the hybrid hand features can be generated by combining features in different levels [125]. The hybrid hand features can embed both local and global information of the hand region [126]. Liu *et al.* [127] adopted three types of geometric features, which are translation, rotation, and scale invariant among fingers, palm, and forearm. To evaluate different 3D hand features, Sorce *et al.* [128] compared hand mask with edges and concatenated SURF descriptors under different illuminations. Their experiments showed that edge is the worst unless some light-independent algorithms are used for the edge extraction. On the other hand, the SURF-based feature is the best one in optimal lighting conditions. Those experimental results can help researchers to design more effective hand features by the fusion of different hand descriptors.

### B. Classifiers

After the feature representation of the 3D hand, the proper classification is required to recognize the 3D hand posture. Among the static 3D hand gesture recognition works, there are four widely employed classification approaches: 1) support vector machines (SVMs); 2) neural network; 3) RFs; and 4) nearest neighbor (NN) search (template matching).

SVM is the most popular classifier for the 3D static hand gesture recognition. It has been successfully applied to single-hand [117] and double-hand poses [127]. The LIBSVM module [129] can be easily used in the implementations. The kernel of an SVM can be either linear [118] or nonlinear [119]. The linear SVM [115] has been adopted to predict the class label with the ASL data sets [88]. In more complex scenarios, the multiclass SVM [116] has successfully recognized 15 hand postures in the cluttered background with an average recognition rate of 95%. The SVM with linear kernel was also utilized with Nanyang Technological University Hand Digits data set [20] and ASL Finger Spelling data set. For nonlinear kernels, the SVM-Radial Basis Function classifiers [114] were trained to recognize four different hand postures in a 3D medical touchless interface. We find that the SVM classifier is usually trained with the middle-level features.

Neural network is still used for simple 3D static hand gesture recognition tasks. In [128], a neural network with backpropagation detected the hand pose to recognize whether it is closed or not. An exponential weighted moving average noise reduction mechanism was used to suppress the noise effects of the neural network. Meanwhile, a neural network model was constructed to recognize four gesture stages [130]. In both works, the neural network handled only fundamental and limited hand gestures.

RFs [131] have also attracted attentions due to its fast training and competitive performance for handling large data sets and feature space. Furthermore, RFs intrinsically support multiclass classification, and thus can be easily parallelized. Hence, RF-based recognition can properly handle high-dimensional low-level and middle-level hand features. The multiclass RFs have been widely applied in practical systems for recognizing ASL Fingerspelling letters [88], [111], [112]. The random decision forest has been experimentally compared with the SVM classification with various 3D hand features [126]. In general, the performance of RFs is dependent on the depth of the tree. The tradeoff of accuracy and speed always needs to be properly considered in its implementation.

NN search or template matching can be used when the high-level discriminative 3D hand features are available. The key is the definition of the distance between the features. The Euclidean distance can be directly employed for the template matching of different hand gestures [123]. However, the generated hand shape features usually contain noises caused by local distortions, pose variations, and inaccuracy depth maps. In order to robustly match the noisy hand features, finger-earth mover's distance was proposed [20] to measure the dissimilarity between time-curved hand features. To better handle the misalignment of hand features, Cheng *et al.* [4] proposed image-to-class dynamic time warping (DTW) distance to distinguish the fingerlet features of 3D hand contours. We find that NN search works quite well with the contour- or boundary-based hand features.

## VI. Hand Trajectory Gesture Recognition

In contrast to the static hand gesture recognition that works on hand shapes, the hand trajectory gesture recognition considers the sequential data of hand trajectory and explores the temporal character of hand motion. In this section, we will survey the 3D hand trajectory gesture recognition approaches.

### A. Feature Representation

The features used in hand trajectory gesture recognition can be divided into appearance-based category and tracking-based category.

The appearance-based features rely on the local feature descriptors regardless of hand tracking or explicit motion trajectory. The features can be generated either in 2D color and depth frames or in the 3D temporal volumes. Wang *et al.* [132] concatenated body centroid, hand displacement, and relative depth level of hand to represent the hand gesture in each depth frame. Similarly, Yang and Sarkar [133] utilized groups of low-level image primitives such as region shape, proximity, or color to implicitly represent the hand as the salient object part of interest. The sequence of the frame-wise group primitives forms the hand gesture feature in temporal domain without requiring a perfect hand segmentation. The early results of the ChaLearn Gesture Challenge [35] showed that all the top ranking approaches were based on techniques making no explicit detection and tracking of humans or individual body parts. To handle the occlusion problem in depth maps, Wang *et al.* [36] proposed random occupancy pattern features that were extracted by sampling the depth space for the representation of hand action. As a special kind of human action, the hand trajectory gesture can

also be represented by concatenating Histogram of Oriented Gradients-based descriptor [134], [135]. All the above features require no explicit hand tracking but discriminative appearance feature instead. The trajectory or motion of the hand can be implicitly embedded in those appearance cues.

The tracking-based features require the tracking of the hand with hand's centroid position or the body skeleton. Hand positions, velocity, acceleration, and chain code have been widely used to formulate trajectory features [136]. The hand position can be either determined by hand segmentation [137] or the body skeleton generated by NITE middleware with the Kinect sensor [138]. In addition to the palm position, the orientation or angle of the hand centroid in a 3D hand gesture trajectory can also be used [46]. Miranda *et al.* [139] described the pose in each frame using a tailored angular representation of the skeleton joints. The key poses were identified with those descriptors. The gesture was represented as the sequence of key poses. For velocity-based trajectory feature, Ren and O'Neill [140] used the speed and direction of hand motion for the menu selection gesture. Similarly, Wu *et al.* [141] utilized the incremental changes of the 3D coordinates in a unit time as the features. Those features rely on the valid hand tracking results and could fail when serious occlusion occurs.

More complex hand features have also been formulated based on the segmented and tracked hands. Zhang and Tian [47] proposed a novel edge enhanced depth motion map together with histogram of gradient descriptor to generate the vector representation of the hand trajectory gestures from the depth video. To ensure the generation of valid hand trajectory, Wang *et al.* [142] utilized the potential active region. The hand gesture was represented by the segmented series of movements in the form of motion history images. Besides, more discriminative features were also proposed by projecting the motion trajectory to the higher dimensional feature space [143]. Beh *et al.* [144] further composed the hand motion trajectory as a unique series of straight and curved segments. They proposed an automated process of segmenting gesture trajectories based on a simple set of threshold values in the angular change measure. The strokes and segments of hand trajectories played a crucial role in those approaches.

### B. Classifiers

There are four main classifying approaches for recognizing hand trajectory gestures. Similar to the static hand gesture recognition, the hand trajectory gesture recognition can be carried out by the traditional classifiers such as SVM and NN matching. Considered as the sequential data, the hand trajectory features can also be handled by hidden Markov model (HMM) and DTW.

SVMs have been trained to classify the ASL trajectory gestures in the MSR Gesture3D data set [36], [47]. Five individual SVMs integrated by a second-stage linear SVM were utilized in the naturalistic CVRR-HANDS 3D data set [135]. For NN search, a matching score between the model sequence and the input sequence was adopted as the

measure to distinguish 39 different ASL signs and 7 hand actions in a two-view hand gesture data set [133]. The maximum correlation coefficient was also used for the gesture recognition on Chalearn data set [134]. For complex gesture features [143], the straightforward NN (1-NN) with the Euclidean distance has also shown the powerful discriminative potential.

HMM has been widely used for the analysis of sequential data. It is very suitable to classify the tracking-based gestures such as trajectory-based signed digits [138]. HMM was also applied to the robotic navigation [46] and sign language recognition [141], [144]. Compared with HMM, the main advantage of DTW is that it can automatically align the sequences that have different lengths and return the proper distance. Hence, DTW distance can be combined with $k$-NN classifiers for robust signed digit recognition [136], [137].

## VII. Continuous Hand Gesture Recognition

For general hand trajectory gesture recognition, we usually assume that the hand gesture video clip has been well segmented in the temporal domain. The focus is to recognize the single meaningful hand gesture in the video. However, the practical HCI applications require the continuous hand gesture recognition in video streams. This means that both the spatial segmentation and temporal segmentation are necessary to recognize the sequential hand gestures.

Following [23], we can roughly divide continuous hand gestures into direct approaches and indirect approaches depending on whether the explicit temporal segmentation proceeds or not. For 2D continuous hand gesture recognition, both direct [145] and indirect [23], [146], [147] approaches have been well investigated. Besides, continuous hand gesture recognition approaches that rely on the 3D accelerometer sensor [148] and touch screen [149] have also been proposed. Here, we will mainly emphasize on 3D continuous hand gesture recognition approaches with the depth data.

For direct 3D continuous hand gesture recognition, it is convenient to use low-level motion features such as velocity and trajectory to detect abrupt changes for the spotting. Elmezain *et al.* [150] proposed a system to recognize continuous digit gestures in real time using HMM. They first generated orientation dynamic features from spatiotemporal trajectories in depth domain, and then quantized them to codewords. The segmentation of the continuous gestures was based on the detection of the zero codeword, which actually detected the static velocity and the endpoint of the gesture. Kristensson *et al.* [151] presented a markerless gesture interface with a Kinect sensor. They defined the input zone for the gesture delimitation purpose. When the depth of the hand was below a set threshold, the hand was defined to be within the input zone. By this zoning technique, the beginning and the end of the gesture can be determined.

The main limitations of direct approaches are twofold.

1) They require the presegmentation of gestures, which may delay the recognition result.
2) The requirement of begin/end signals (e.g., gesture interval) makes them not flexible in the HCI applications.

Hence, the indirect approaches have received increasing interests due to its natural applicability and fast recognition for continuous hand gestures.

Among the indirect continuous hand gesture recognition approaches, DTW is widely used for finding the matched gesture segment in the temporal domain. Lichtenauer *et al.* [152] claimed that time warping and classification should be separated because of conflicting likelihood modeling demands. They proposed to use statistical DTW only for time warping while classifying the warped with different statistical classifiers. The proposed hybrid approaches provided a significant improvement over HMM on the 3D Dutch Sign Language data set. To improve the recognition accuracy, Keskin *et al.* [184] proposed a DTW-based preclustering technique for 3D digit recognition using graphical models. To leverage the discriminativity of each hand part, Arici *et al.* [153] proposed a weighted DTW approach that weighted joints by optimizing a discriminant ratio to improve the 3D hand-arm gesture recognition.

In the context of begin–end hand gesture recognition, Reyes *et al.* [48] presented a begin–end gesture recognition approach using feature weighting in the DTW framework. The feature weighting approach was proposed to improve the cost distance computation in the conventional begin–end DTW algorithm. They associated a discriminatory weight with each joint of the skeletal model depending on its participation in a particular gesture. DTW can also be applied to segment the gestures in the preprocessing step. Hernandez-Vela *et al.* [154] extended the Bag of Words [155] model to bag-of-visual-and-depth-words model for gesture recognition. The state-of-the-art RGB and depth features were fused to generate the final gesture feature, which was integrated in a continuous gesture recognition pipeline, and the DTW algorithm was used to perform the begin–end segmentation of gestures. Besides, Bhuyan *et al.* [156] proposed to use a novel set of features with conditional random fields for gesture spotting to distinguish meaningful gestures from unintentional movements. Although their original approach works for 2D continuous hand gesture recognition with skin-color-based hand detection, it could be utilized for 3D scenario with hand tracking in the RGB-D space.

One of the most challenging practical problems for the continuous gesture recognition system is the subgesture problem when some gestures are similar to parts of other longer gestures. Subgesture problem is very common in real-world gesture data sets (e.g., ASL and ChaLearn data sets). Recently, Cheng *et al.* [49] have proposed a windowed DTW (WDTW) approach for 3D continuous hand trajectory gesture recognition. In their work, a parameterized searching window was introduced in the cost matrix of traditional DTW approaches to detect the beginning and end of the specific gesture from an infinite trajectory gesture sequence. Hence, the continuous gesture recognition can be formulated into online parameter estimation of the searching window. Compared with the feature weighting DTW approach in [154], the WDTW approach can handle the subgesture problem by moving the searching window forward and judging whether this is a single gesture after some delay time.

## VIII. APPLICATIONS AND TYPICAL SYSTEMS

The 3D hand gesture recognition approaches are mainly used in six application domains: 1) sign language recognition; 2) virtual manipulation; 3) daily assistance; 4) palm verification; 5) gaming; and 6) the emerging human–robot interaction (HRI).

*Sign language recognition* with hand gestures has been investigated for the ASL as mentioned in the previous sections. For disable people who are deaf mute, a recognition system of sign language can greatly help them to keep in touch with others. Sun *et al.* [157] employed a novel two-Kinect system to collect the Japanese Sign Language (JSL) data. Two Kinect sensors were located perpendicularly to each other, and the Point Cloud Library was used to recognize the JSL gestures. Considering personal experiences, a fully automatic hand gesture recognition system can make people feel natural to use it. To approach this goal, Kurakin *et al.* [18] proposed a real-time system for hand trajectory gesture recognition. The system was designed for practical ASL applications, fully automatic, and robust to variations in speed and style as well as in hand orientations. This system is the first data-driven system that is capable of automatic hand gesture recognition. To further improve user experiences, Suau *et al.* [158] presented an online hand-based touchless interaction system named intAIRact. The intAIRact is user friendly and highly configurable with easy-to-learn hand gestures. Consequently, a user can trigger a large number of events by remembering nine hand gestures and combining them with simple translations and rotations.

*Virtual manipulation* is the popular application of 3D hand gesture recognition due to its natural user interface for HCI tasks. Fraunhofer FIT developed the 3D gesture-based interaction system,[7] which is the noncontact gesture and finger recognition system. The system detects hand and finger positions in real time and translates these into appropriate interaction commands. For picture gallery browsing, Chai *et al.* [50] leveraged accurate hand segmentation to greatly improve the interaction usability. To approach user interfaces in dark environments, Lee *et al.* [159] presented an approach for tracking hand rotation and various grasping gestures through an infrared camera. The 3D hand gesture recognition has also been combined with the 3D stereoscopic display to provide the immersive HCI experience in a virtual reality system with stereo vision [160]. In this system, a user can manipulate the virtual object in the 3D stereoscopy scene. The Kinect sensor was used to track the user's hand and render the virtual objects according to the user's viewpoint. In particular, several recent systems involving Kinect-based touchless interaction with surgical images have been developed for operating theater practices [161]. The pioneer works include medical image manipulation systems at Sunnybrook Hospital in Toronto [162] and Guy's and St. Thomas' Hospital in London [161]. In these applications, 3D hand gesture vocabularies were carefully designed to realize the touchless image manipulation in sterile environments.

*Daily assistance* using hand gestures can help older people to perform activities of daily living such as handwashing [163].

---

[7]http://www.fit.fraunhofer.de/en/fb/cscw/projects/3d-multi-touch.html

The important senior gestures such as eating and drinking can also be monitored [164]. To make the smart home that facilitates aging in place, Yanik *et al.* [165] used Kinect depth data and growing neural gas algorithm for gesture-based robot control. It was the initial effort toward the goal of assistive robot in which the response of the robot converged the user's desired response. For single-hand driving system, Son and Sowmya [51] proposed a Kinect-based system that can help those people who have difficulties in moving one of their arms, to drive and control the vehicles with only one hand. It also has the potential to be applied to any wheeled vehicles. Regarding the typical systems, CMU's assistive robot Home Exploring Robotic Butler[8] can help the user to accomplish the seamless teleoperation by predicting the user's manipulation intent [166]. Since older people also experience challenges in maintaining their home, care robots are considered as one option to support. Recently, Fischinger *et al.* [167] have released a prototype care robot Hobbit in which the 3D hand gesture interface was developed for the interaction of older people with the robot.

*Palm verification* is a key biometric technology for many security applications. The shape of the hand can be easily captured in a relatively user-friendly manner using 2D/3D cameras thus acceptable by the public. Amayeh *et al.* [168] proposed a component-based approach to hand-based verification and identification. Their approach utilized a 2D camera plus a lighting table to decomposite the hand silhouette into different regions corresponding to the back of the palm and the fingers. To identify the user's hand shape in free pose, Kanhangad *et al.* [169] proposed a contactless and pose invariant biometric identification system that utilized a 3D digitizer to simultaneously acquire intensity and range images of the user's hand. Their approaches determined the orientation of the hand in 3D space and normalized the pose of the acquired RGB-D images to register the ROIs for the matching-based identification.

*Gaming* with interactive hand gestures has been significantly promoted by the next-generation game consoles. Sony's Play Station and Nintendo's Wii are equipped with handy controllers (PS Move and Wiimote, respectively) for the player's hand tracking. Roccetti *et al.* [170] addressed the problem of differentiating the design of a gesture-based interface for a console from the problem of designing it for a public space setting. The most well-known somatic games involve Microsoft's Xbox with Kinect. There have been quite a few body-sensing games released in the consumer market.

*HRI* is the most important function for the emerging social robot that is able to interact with people using the natural gestures [171]. Hand gesture-based interface offers a way to enable human to interact with robots more easily and efficiently [172]. Yin and Xie [173] implemented a posture recognition system on a real humanoid service robot HARO-1, and the experimental results demonstrated the effectiveness and robustness of the system. Among the set of gestures intuitively performed by humans, pointing gestures are especially important for interaction with robots.

Nickel and Stiefelhagen [174] presented an approach for recognizing pointing gestures in the context of HRI. The stereo camera was employed for the hand detection. The system aims at run-on gesture recognition in real time and allows for the robot ego-motion. To handle the challenge of hand tracking in crowded and dynamic environments for HRI applications, McKeague *et al.* [52] proposed a sensor fusion-based hand tracking algorithm for crowded environments. It significantly improved the accuracy of existing hand detectors, based on the RGB-D information. To support robust and efficient HRI for socially assistive robots, Michel *et al.* [175] proposed a 3D vision-based gesture recognition approach, which considered gestural vocabulary in the context of human–robot dialog.

In particular, to build an open and effective system for empowering natural modalities of HCI, Pedersoli *et al.* [176] developed the first open source package for Kinect, which targeted both static and trajectory-based hand gestures in a unified framework. The FPV technology [33] also has the potential to promote the emerging pervasive wearable devices (e.g., Google Glasses, Vuzix SmartGlass, Lenovo New Glass, and Microsoft HoloLens) by hand gestures. Some pioneer works [177]–[179] in this direction have investigated hand detection and segmentation tasks with the ego-centric vision. Li and Kitani [179], [180] addressed the task of pixel-level hand detection with challenges of illumination changes and camera motion. The sequential classifier [181] and the super-pixel classification [178], [182] techniques were also employed to explore the temporal and spatial coherence of hand gestures in first-person views. The real-time wearable hand gesture recognition [32] can be used for the HRI event detection (e.g., object manipulation) and the intention understanding. More hand gesture systems and applications have been well summarized in [183].

## IX. CONCLUSION

In this paper, we have given a comprehensive survey of the emerging progress on 3D hand gesture recognition. We discussed a variety of 3D hand gesture recognition aspects along with their applications. The survey reviewed important progress made on 3D depth sensors, data sets, 3D hand modeling, static hand gesture recognition, hand trajectory gesture recognition, continuous hand gesture recognition, related applications, and typical systems. One of the major challenges is the online recognition of 3D hand gestures. The absence of explicit begin/end hints in practical scenarios will degrade the performance of traditional static and trajectory approaches. Hence, the continuous hand gesture recognition will attract more attention due to its applicability. Other challenges include different meanings of similar hand gestures, which need to be distinguished by fine-grained gesture recognition. In the future, we expect that more tiny finger gestures will be well recognized, leveraging more accurate depth sensors. Finally, we envision the booming of intelligent products using 3D hand gesture recognition for HRI purposes. Most of the existing 3D hand gesture recognition works employ depth sensors with a fixed position. However, users may move freely and disappear during the interaction with robots. We believe that there will be much research room for interactive hand gesture recognition for HRI.

[8]http://www.cmu.edu/herb-robot/

## REFERENCES

[1] J. K. Aggarwal and M. S. Ryoo, "Human activity analysis: A review," *ACM Comput. Surv.*, vol. 43, no. 3, pp. 16:1–16:43, Apr. 2011.

[2] R. Lefevre, *Rude Hand Gestures of the World: A Guide to Offending Without Words*. San Francisco, CA, USA: Chronicle Books, 2011.

[3] S. D. Kelly, S. M. Manning, and S. Rodak, "Gesture gives a hand to language and learning: Perspectives from cognitive neuroscience, developmental psychology and education," *Lang. Linguistics Compass*, vol. 2, no. 4, pp. 569–588, Jul. 2008.

[4] H. Cheng, Z. Dai, and Z. Liu, "Image-to-class dynamic time warping for 3D hand gesture recognition," in *Proc. IEEE ICME*, Jul. 2013, pp. 1–6.

[5] X. Zhang, X. Chen, Y. Li, V. Lantz, K. Wang, and J. Yang, "A framework for hand gesture recognition based on accelerometer and EMG sensors," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 41, no. 6, pp. 1064–1076, Nov. 2011.

[6] D. Kammer, J. Wojdziak, M. Keck, R. Groh, and S. Taranko, "Towards a formalization of multi-touch gestures," in *Proc. ACM Int. Conf. Interact. Tabletops Surf.*, 2010, pp. 49–58.

[7] E. Hoggan, J. Williamson, A. Oulasvirta, M. Nacenta, P. O. Kristensson, and A. Lehtiö, "Multi-touch rotation gestures: Performance and ergonomics," in *Proc. ACM SIGCHI*, 2013, pp. 3047–3050.

[8] X. Zabulis, H. Baltzakis, and A. Argyros, "Vision-based hand gesture recognition for human-computer interaction," in *The Universal Access Handbook*. Boca Raton, FL, USA: CRC Press, 2009.

[9] P. Buehler, M. Everingham, D. P. Huttenlocher, and A. Zisserman, "Upper body detection and tracking in extended signing sequences," *Int. J. Comput. Vis.*, vol. 95, no. 2, pp. 180–197, Nov. 2011.

[10] S. S. Rautaray and A. Agrawal, "Vision based hand gesture recognition for human computer interaction: A survey," *Artif. Intell. Rev.*, vol. 43, no. 1, pp. 1–54, Jan. 2012.

[11] X. Suau, J. Ruiz-Hidalgo, and J. R. Casas, "Real-time head and hand tracking based on 2.5D data," *IEEE Trans. Multimedia*, vol. 14, no. 3, pp. 575–585, Jun. 2012.

[12] R. Y. Wang and J. Popović, "Real-time hand-tracking with a color glove," in *Proc. ACM SIGGRAPH*, 2009, p. 63.

[13] S. L. Phung, A. Bouzerdoum, and D. Chai, Sr., "Skin segmentation using color pixel classification: Analysis and comparison," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 1, pp. 148–154, Jan. 2005.

[14] R. Yang, S. Sarkar, and B. Loeding, "Handling movement epenthesis and hand segmentation ambiguities in continuous sign language recognition using nested dynamic programming," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 3, pp. 462–477, Mar. 2010.

[15] Z. Zhang, "Microsoft Kinect sensor and its effect," *IEEE MultiMedia*, vol. 19, no. 2, pp. 4–10, Feb. 2012.

[16] R. Nair *et al.*, "A survey on time-of-flight stereo fusion," in *Time-of-Flight and Depth Imaging. Sensors, Algorithms, and Applications* (Lecture Notes in Computer Science), vol. 8200. New York, NY, USA: Springer-Verlag, 2013, pp. 105–127.

[17] M. Van den Bergh and L. Van Gool, "Combining RGB and ToF cameras for real-time 3D hand gesture interaction," in *Proc. IEEE Workshop Appl. Comput. Vis.*, Jan. 2011, pp. 66–72.

[18] A. Kurakin, Z. Zhang, and Z. Liu, "A real time system for dynamic hand gesture recognition with a depth sensor," in *Proc. 20th Eur. Signal Process. Conf.*, Aug. 2012, pp. 1975–1979.

[19] J. Suarez and R. R. Murphy, "Hand gesture recognition with depth images: A review," in *Proc. IEEE RO-MAN*, Sep. 2012, pp. 411–417.

[20] Z. Ren, J. Yuan, and Z. Zhang, "Robust hand gesture recognition based on finger-earth mover's distance with a commodity depth camera," in *Proc. ACM MM*, 2011, pp. 1093–1096.

[21] W.-G. Liou, C.-Y. Hsieh, and W.-Y. Lin, "Trajectory-based sign language recognition using discriminant analysis in higher-dimensional feature space," in *Proc. IEEE ICME*, Jul. 2011, pp. 1–4.

[22] C. Tran and M. M. Trivedi, "3-D posture and gesture recognition for interactivity in smart spaces," *IEEE Trans. Ind. Informat.*, vol. 8, no. 1, pp. 178–187, Feb. 2012.

[23] J. Alon, V. Athitsos, Q. Yuan, and S. Sclaroff, "A unified framework for gesture recognition and spatiotemporal gesture segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 9, pp. 1685–1699, Sep. 2009.

[24] S. Sarkar, B. Loeding, R. Yang, S. Nayak, and A. Parashar, "Segmentation-robust representations, matching, and modeling for sign language," in *Proc. IEEE CVPR Workshops*, Jun. 2011, pp. 13–19.

[25] J. Davis and M. Shah, "Toward 3-D gesture recognition," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 13, no. 3, pp. 381–393, 1999.

[26] I. Oikonomidis, N. Kyriazis, and A. A. Argyros, "Markerless and efficient 26-DOF hand pose recovery," in *Proc. ACCV*, 2010, pp. 744–757.

[27] R. Poppe, "Vision-based human motion analysis: An overview," *Comput. Vis. Image Understand.*, vol. 108, nos. 1–2, pp. 4–18, Oct./Nov. 2007.

[28] M. Ye, Q. Zhang, L. Wang, J. Zhu, R. Yang, and J. Gall, "A survey on human motion analysis from depth data," in *Time-of-Flight and Depth Imaging. Sensors, Algorithms, and Applications* (Lecture Notes in Computer Science), vol. 8200. Berlin, Germany: Springer-Verlag, 2013, pp. 149–187.

[29] L. Chen, H. Wei, and J. Ferryman, "A survey of human motion analysis using depth imagery," *Pattern Recognit. Lett.*, vol. 34, no. 15, pp. 1995–2006, Nov. 2013.

[30] J. Han, L. Shao, D. Xu, and J. Shotton, "Enhanced computer vision with Microsoft Kinect sensor: A review," *IEEE Trans. Cybern.*, vol. 43, no. 5, pp. 1318–1334, Oct. 2013.

[31] Y. Song, J. Tang, F. Liu, and S. Yan, "Body surface context: A new robust feature for action recognition from depth videos," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 6, pp. 952–964, Jun. 2014.

[32] W. W. Mayol and D. W. Murray, "Wearable hand activity recognition for event summarization," in *Proc. IEEE Int. Symp. Wearable Comput.*, Oct. 2005, pp. 122–129.

[33] A. Betancourt, P. Morerio, C. S. Regazzoni, and M. Rauterberg, "The evolution of first person vision methods: A survey," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 25, no. 5, pp. 744–760, May 2015.

[34] L. E. Potter, J. Araullo, and L. Carter, "The leap motion controller: A view on sign language," in *Proc. Austral. Comput.-Human Interact. Conf.*, 2013, pp. 175–178.

[35] I. Guyon, V. Athitsos, P. Jangyodsuk, B. Hamner, and H. J. Escalante, "ChaLearn gesture challenge: Design and first results," in *Proc. IEEE CVPR Workshops*, Jun. 2012, pp. 1–6.

[36] J. Wang, Z. Liu, J. Chorowski, Z. Chen, and Y. Wu, "Robust 3D action recognition with random occupancy patterns," in *Proc. ECCV*, 2012, pp. 872–885.

[37] J. Shotton *et al.*, "Real-time human pose recognition in parts from single depth images," in *Proc. IEEE CVPR*, Jun. 2011, pp. 1297–1304.

[38] C. Keskin, F. Kıraç, Y. Kara, and L. Akarun, "Real time hand pose estimation using depth sensors," in *Proc. Consum. Depth Cameras Comput. Vis.*, 2013, pp. 119–137.

[39] M. de La Gorce, D. J. Fleet, and N. Paragios, "Model-based 3D hand pose estimation from monocular video," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 9, pp. 1793–1805, Sep. 2011.

[40] I. Oikonomidis, N. Kyriazis, and A. A. Argyros, "Efficient model-based 3D tracking of hand articulations using Kinect," in *Proc. BMVC*, 2011, pp. 1–3.

[41] H. Hamer, K. Schindler, E. Koller-Meier, and L. Van Gool, "Tracking a hand manipulating an object," in *Proc. IEEE ICCV*, Sep./Oct. 2009, pp. 1475–1482.

[42] I. Oikonomidis, N. Kyriazis, and A. A. Argyros, "Full DOF tracking of a hand interacting with an object by modeling occlusions and physical constraints," in *Proc. IEEE ICCV*, Nov. 2011, pp. 2088–2095.

[43] I. Oikonomidis, N. Kyriazis, and A. A. Argyros, "Tracking the articulated motion of two strongly interacting hands," in *Proc. IEEE CVPR*, Jun. 2012, pp. 1862–1869.

[44] L. Ballan, A. Taneja, J. Gall, L. Van Gool, and M. Pollefeys, "Motion capture of hands in action using discriminative salient points," in *Proc. ECCV*, 2012, pp. 640–653.

[45] M. Van den Bergh *et al.*, "Real-time 3D hand gesture interaction with a robot for understanding directions from humans," in *Proc. IEEE RO-MAN*, Jul./Aug. 2011, pp. 357–362.

[46] D. Xu, Y.-L. Chen, C. Lin, X. Kong, and X. Wu, "Real-time dynamic gesture recognition system based on depth perception for robot navigation," in *Proc. IEEE Int. Conf. Robot. Biomimetics*, Dec. 2012, pp. 689–694.

[47] C. Zhang and Y. Tian, "Edge enhanced depth motion map for dynamic hand gesture recognition," in *Proc. IEEE CVPR Workshops*, Jun. 2013, pp. 500–505.

[48] M. Reyes, G. Dominguez, and S. Escalera, "Featureweighting in dynamic timewarping for gesture recognition in depth data," in *Proc. IEEE ICCV Workshops*, Nov. 2011, pp. 1182–1188.

[49] H. Cheng, J. Luo, and X. Chen, "A windowed dynamic time warping approach for 3D continuous hand gesture recognition," in *Proc. IEEE ICME*, Jul. 2014, pp. 1–6.

[50] X. Chai, Y. Fang, and K. Wang, "Robust hand gesture analysis and application in gallery browsing," in *Proc. IEEE ICME*, Jun./Jul. 2009, pp. 938–941.

[51] J. P. Son and A. Sowmya, "Single-handed driving system with Kinect," in *Proc. HCI*, 2013, pp. 631–639.

[52] S. McKeague, J. Liu, and G.-Z. Yang, "An asynchronous RGB-D sensor fusion framework using Monte–Carlo methods for hand tracking on a mobile robot in crowded environments," in *Proc. Int. Conf. Soc. Robot.*, 2013, pp. 491–500.

[53] K. Lai, J. Konrad, and P. Ishwar, "A gesture-driven computer interface using Kinect," in *Proc. IEEE Southwest Symp. Image Anal. Interpretation*, Apr. 2012, pp. 185–188.

[54] R.-D. Vatavu, "User-defined gestures for free-hand TV control," in *Proc. Eur. Conf. Interact. TV Video*, 2012, pp. 45–48.

[55] L. Gallo, A. P. Placitelli, and M. Ciampi, "Controller-free exploration of medical image data: Experiencing the Kinect," in *Proc. Int. Symp. Comput.-Based Med. Syst.*, Jun. 2011, pp. 1–6.

[56] M. Caputo, K. Denker, B. Dums, and G. Umlauf, "3D hand gesture recognition based on sensor fusion of commodity hardware," in *Proc. Mensch Comput.*, 2012, pp. 293–302.

[57] Z. Zafrulla, H. Brashear, T. Starner, H. Hamilton, and P. Presti, "American sign language recognition with the Kinect," in *Proc. Int. Conf. Multimodal Inter.*, 2011, pp. 279–286.

[58] X. Chai *et al.*, "Sign language recognition and translation with Kinect," in *Proc. IEEE Int. Conf. Autom. Face Gesture Recognit.*, Apr. 2013.

[59] M. Tang, "Recognizing hand gestures with Microsoft's Kinect," Dept. Elect. Eng., Stanford Univ., Stanford, CA, USA, Tech. Rep., 2011.

[60] W. Zhao, J. Chai, and Y.-Q. Xu, "Combining marker-based mocap and RGB-D camera for acquiring high-fidelity hand motion data," in *Proc. ACM SIGGRAPH Eurograph. Symp. Comput. Animation*, 2012, pp. 33–42.

[61] C. Keskin, F. Kıraç, Y. E. Kara, and L. Akarun, "Real time hand pose estimation using depth sensors," in *Proc. IEEE ICCV Workshops*, Nov. 2011, pp. 1228–1234.

[62] K. K. Biswas and S. K. Basu, "Gesture recognition using Microsoft Kinect," in *Proc. Int. Conf. Autom., Robot. Appl.*, Dec. 2011, pp. 100–103.

[63] Y. Wen, C. Hu, G. Yu, and C. Wang, "A robust method of detecting hand gestures using depth sensors," in *Proc. IEEE Int. Workshop Haptic Audio Vis. Environ. Games*, Oct. 2012, pp. 72–77.

[64] J. L. Raheja, A. Chaudhary, and K. Singal, "Tracking of fingertips and centers of palm using KINECT," in *Proc. Int. Conf. Comput. Intell., Modelling Simulation*, Sep. 2011, pp. 248–252.

[65] G. Du, P. Zhang, J. Mai, and Z. Li, "Markerless Kinect-based hand tracking for robot teleoperation," *Int. J. Adv. Robot. Syst.*, vol. 9, no. 36, pp. 1–10, 2012.

[66] V. Frati and D. Prattichizzo, "Using Kinect for hand tracking and rendering in wearable haptics," in *Proc. IEEE World Haptics Conf.*, Jun. 2011, pp. 317–321.

[67] U. Lee and J. Tanaka, "Hand controller: Image manipulation interface using fingertips and palm tracking with Kinect depth data," in *Proc. Asia Pacific Conf. Comput. Human Interact.*, 2012, pp. 705–706.

[68] Z. Ren, J. Meng, J. Yuan, and Z. Zhang, "Robust hand gesture recognition with Kinect sensor," in *Proc. ACM MM*, 2011, pp. 759–760.

[69] Y. Li, "Hand gesture recognition using Kinect," in *Proc. IEEE Int. Conf. Softw. Eng. Service Sci.*, Jun. 2012, pp. 196–199.

[70] Z. Ren, J. Yuan, J. Meng, and Z. Zhang, "Robust part-based hand gesture recognition using Kinect sensor," *IEEE Trans. Multimedia*, vol. 15, no. 5, pp. 1110–1120, Aug. 2013.

[71] M. F. Shiratuddin and K. W. Wong, "Non-contact multi-hand gestures interaction techniques for architectural design in a virtual environment," in *Proc. Int. Conf. Inf. Technol. Multimedia*, Nov. 2011, pp. 1–6.

[72] F. Pedersoli, N. Adami, S. Benini, and R. Leonardi, "XKin— Extendable hand pose and gesture recognition library for Kinect," in *Proc. ACM MM*, 2012, pp. 1465–1468.

[73] H. Regenbrecht, J. Collins, and S. Hoermann, "A leap-supported, hybrid AR interface approach," in *Proc. Austral. Comput.-Human Interact. Conf.*, 2013, pp. 281–284.

[74] K. Sabir, C. Stolte, B. Tabor, and S. I. O'Donoghue, "The molecular control toolkit: Controlling 3D molecular graphics via gesture and voice," in *Proc. IEEE Symp. Biol. Data Visualizat.*, Oct. 2013, pp. 49–56.

[75] A. J. Porfirio, K. Lais Wiggers, L. E. S. Oliveira, and D. Weingaertner, "LIBRAS sign language hand configuration recognition based on 3D meshes," in *Proc. IEEE SMC*, Oct. 2013, pp. 1588–1593.

[76] F. Weichert, D. Bachmann, B. Rudak, and D. Fisseler, "Analysis of the accuracy and robustness of the leap motion controller," *Sensors*, vol. 13, no. 5, pp. 6380–6393, 2013.

[77] G. J. Iddan and G. Yahav, "Three-dimensional imaging in the studio and elsewhere," *Proc. SPIE*, vol. 4298, no. 1, pp. 48–55, Apr. 2001.

[78] A. Kolb, E. Barth, and R. Koch, "ToF-sensors: New dimensions for realism and interactivity," in *Proc. IEEE CVPR Workshops*, Jun. 2008, pp. 1–6.

[79] D. Droeschel, J. Stückler, and S. Behnke, "Learning to interpret pointing gestures with a time-of-flight camera," in *Proc. Int. Conf. Human-Robot Interact.*, Mar. 2011, pp. 481–488.

[80] S. Oprisescu, C. Rasche, and B. Su, "Automatic static hand gesture recognition using ToF cameras," in *Proc. 20th Eur. Signal Process. Conf.*, Aug. 2012, pp. 2748–2751.

[81] T. Kapuscinski, M. Oszust, and M. Wysocki, "Recognition of signed dynamic expressions observed by ToF camera," in *Proc. Signal Process., Algorithms, Archit., Arrangements, Appl.*, Sep. 2013, pp. 291–296.

[82] O. Kliper-Gross, T. Hassner, and L. Wolf, "The action similarity labeling challenge," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 3, pp. 615–621, Mar. 2012.

[83] T. Hassner, "A critical review of action recognition benchmarks," in *Proc. IEEE CVPR Workshops*, Jun. 2013, pp. 245–250.

[84] N. Gkalelis, H. Kim, A. Hilton, N. Nikolaidis, and I. Pitas, "The i3DPost multi-view and 3D human action/interaction database," in *Proc. Conf. Vis. Media Prod.*, Nov. 2009, pp. 159–168.

[85] Y. Song, D. Demirdjian, and R. Davis, "Tracking body and hands for gesture recognition: NATOPS aircraft handling signals database," in *Proc. IEEE Int. Conf. Autom. Face Gesture Recognit. Workshops*, Mar. 2011, pp. 500–506.

[86] S. Ruffieux, D. Lalanne, E. Mugellini, and O. A. Khaled, "A survey of datasets for human gesture recognition," in *Proc. 16th Int. Conf. Human-Comput. Interact.*, 2014, pp. 337–348.

[87] A. M. Martinez, R. B. Wilbur, R. Shay, and A. C. Kak, "Purdue RVL-SLLL ASL database for automatic recognition of American sign language," in *Proc. IEEE Int. Conf. Multimodal Inter.*, 2002, pp. 167–172.

[88] N. Pugeault and R. Bowden, "Spelling it out: Real-time ASL fingerspelling recognition," in *Proc. IEEE ICCV Workshops*, Nov. 2011, pp. 1114–1119.

[89] S. Fothergill, H. M. Mentis, P. Kohli, and S. Nowozin, "Instructing people for training gestural interactive systems," in *Proc. ACM CHI*, 2012, pp. 1737–1746.

[90] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Mining actionlet ensemble for action recognition with depth cameras," in *Proc. IEEE CVPR*, Jun. 2012, pp. 1290–1297.

[91] C. Wolf *et al.*, "Evaluation of video activity localizations integrating quality and quantity measurements," *Comput. Vis. Image Understand.*, vol. 127, pp. 14–30, Oct. 2014.

[92] L. Liu and L. Shao, "Learning discriminative representations from RGB-D video data," in *Proc. Int. Joint Conf. Artif. Intell.*, 2013, pp. 1493–1500.

[93] C. Conly, P. Doliotis, P. Jangyodsuk, R. Alonzo, and V. Athitsos, "Toward a 3D body part detection video dataset and hand tracking benchmark," in *Proc. Int. Conf. Pervasive Technol. Rel. Assist. Environ.*, 2013, pp. 1–2.

[94] J. Triesch and C. von der Malsburg, "A system for person-independent hand posture recognition against complex backgrounds," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 12, pp. 1449–1453, Dec. 2001.

[95] S. Marcel, O. Bernier, J. E. Viallet, and D. Collobert, "Hand gesture recognition using input-output hidden Markov models," in *Proc. IEEE Int. Conf. Autom. Face Gesture Recognit.*, Mar. 2000, pp. 456–461.

[96] R. E. Mitchell, T. A. Young, B. Bachleda, and M. A. Karchmer, "How many people use ASL in the United States? Why estimates need updating," *Sign Lang. Stud.*, vol. 6, no. 3, pp. 306–335, 2006.

[97] A. Thangali, J. P. Nash, S. Sclaroff, and C. Neidle, "Exploiting phonological constraints for handshape inference in ASL video," in *Proc. IEEE CVPR*, Jun. 2011, pp. 521–528.

[98] A. Erol, G. Bebis, M. Nicolescu, R. D. Boyle, and X. Twombly, "Vision-based hand pose estimation: A review," *Comput. Vis. Image Understand.*, vol. 108, nos. 1–2, pp. 52–73, Oct./Nov. 2007.

[99] D. Mohr and G. Zachmann, "A survey of vision-based markerless hand tracking approaches," to be published.

[100] C. Keskin, F. Kıraç, Y. E. Kara, and L. Akarun, "Hand pose estimation and hand shape classification using multi-layered randomized decision forests," in *Proc. ECCV*, 2012, pp. 852–863.

[101] D. Tang, T.-H. Yu, and T.-K. Kim, "Real-time articulated hand pose estimation using semi-supervised transductive regression forests," in *Proc. IEEE ICCV*, Dec. 2013, pp. 3224–3231.

[102] Y. Yao and Y. Fu, "Real-time hand pose estimation from RGB-D sensor," in *Proc. IEEE ICME*, Jul. 2012, pp. 705–710.

[103] Y. Wu, J. Y. Lin, and T. S. Huang, "Capturing natural hand articulation," in *Proc. 8th IEEE ICCV*, Jul. 2001, pp. 426–432.

[104] M. de La Gorce and N. Paragios, "A variational approach to monocular hand-pose estimation," *Comput. Vis. Image Understand.*, vol. 114, no. 3, pp. 363–372, 2010.

[105] M.-F. Ho, C.-Y. Tseng, C.-C. Lien, and C.-L. Huang, "A multi-view vision-based hand motion capturing system," *Pattern Recognit.*, vol. 44, no. 2, pp. 443–453, 2011.

[106] C. Qian, X. Sun, Y. Wei, X. Tang, and J. Sun, "Realtime and robust hand tracking from depth," in *Proc. IEEE Conf. CVPR*, Jun. 2014, pp. 1106–1113.

[107] H. Hamer, J. Gall, T. Weise, and L. Van Gool, "An object-dependent hand pose prior from sparse training data," in *Proc. IEEE Conf. CVPR*, Jun. 2010, pp. 671–678.

[108] C. Xu and L. Cheng, "Efficient hand pose estimation from a single depth image," in *Proc. IEEE ICCV*, Dec. 2013, pp. 3456–3462.

[109] R. Girshick, J. Shotton, P. Kohli, A. Criminisi, and A. Fitzgibbon, "Efficient regression of general-activity human poses from depth images," in *Proc. IEEE ICCV*, Nov. 2011, pp. 415–422.

[110] S. Sridhar, A. Oulasvirta, and C. Theobalt, "Interactive markerless articulated hand motion tracking using RGB and depth data," in *Proc. IEEE ICCV*, Dec. 2013, pp. 2456–2463.

[111] C. Keskin, F. Kıraç, Y. E. Kara, and L. Akarun, "Randomized decision forests for static and dynamic hand shape classification," in *Proc. IEEE Comput. Soc. Conf. CVPR Workshops*, Jun. 2012, pp. 31–36.

[112] A. Kuznetsova, L. Leal-Taixe, and B. Rosenhahn, "Real-time sign language recognition using a consumer depth camera," in *Proc. IEEE ICCV Workshops*, Dec. 2013, pp. 83–90.

[113] L. Gallo and A. P. Placitelli, "View-independent hand posture recognition from single depth images using PCA and Flusser moments," in *Proc. 8th Int. Conf. Signal Image Technol. Internet Syst.*, Nov. 2012, pp. 898–904.

[114] L. Gallo, "Hand shape classification using depth data for unconstrained 3D interaction," *J. Ambient Intell. Smart Environ.*, vol. 6, no. 1, pp. 93–105, Jan. 2014.

[115] X. Zhu and K. K. Wong, "Single-frame hand gesture recognition using color and depth kernel descriptors," in *Proc. ICPR*, Nov. 2012, pp. 2989–2992.

[116] H. Wang, Q. Wang, and X. Chen, "Hand posture recognition from disparity cost map," in *Proc. 11th ACCV*, 2012, pp. 722–733.

[117] P. Suryanarayan, A. Subramanian, and D. Mandalapu, "Dynamic hand pose recognition using depth data," in *Proc. 20th ICPR*, Aug. 2010, pp. 3105–3108.

[118] C. Zhang, X. Yang, and Y. Tian, "Histogram of 3D facets: A characteristic descriptor for hand gesture recognition," in *Proc. 10th IEEE Int. Conf. Workshops Autom. Face Gesture Recognit.*, Apr. 2013, pp. 1–8.

[119] A. D. Bagdanov, A. Del Bimbo, L. Seidenari, and L. Usai, "Real-time hand status recognition from RGB-D imagery," in *Proc. 21st ICPR*, Nov. 2012, pp. 2456–2459.

[120] M. Flasiński and S. Myśliński, "On the use of graph parsing for recognition of isolated hand postures of Polish sign language," *Pattern Recognit.*, vol. 43, no. 6, pp. 2249–2264, 2010.

[121] S. P. Priyal and P. K. Bora, "A robust static hand gesture recognition system using geometry based normalizations and Krawtchouk moments," *Pattern Recognit.*, vol. 46, no. 8, pp. 2202–2219, 2013.

[122] Y.-T. Li and J. P. Wachs, "HEGM: A hierarchical elastic graph matching for hand gesture recognition," *Pattern Recognit.*, vol. 47, no. 1, pp. 80–88, 2014.

[123] Y. Liu *et al.*, "Image processing and recognition of multiple static hand gestures for human–computer interaction," in *Proc. 7th ICIG*, Jul. 2013, pp. 465–470.

[124] S. Qin, X. Zhu, H. Yu, S. Ge, Y. Yang, and Y. Jiang, "Real-time markerless hand gesture recognition with depth camera," in *Proc. 13th PCM*, 2012, pp. 186–197.

[125] D. Uebersax, J. Gall, M. Van den Bergh, and L. Van Gool, "Real-time sign language letter and word recognition from depth data," in *Proc. IEEE ICCV Workshops*, Nov. 2011, pp. 383–390.

[126] D. Minnen and Z. Zafrulla, "Towards robust cross-user hand tracking and shape recognition," in *Proc. IEEE ICCV Workshops*, Nov. 2011, pp. 1235–1241.

[127] L. Liu, J. Xing, H. Ai, and X. Ruan, "Hand posture recognition using finger geometric feature," in *Proc. 21st ICPR*, Nov. 2012, pp. 565–568.

[128] S. Sorce, V. Gentile, and A. Gentile, "Real-time hand pose recognition based on a neural network using Microsoft Kinect," in *Proc. 8th IEEE Int. Conf. Broadband Wireless Comput., Commun., Appl.*, Oct. 2013, pp. 344–350.

[129] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, 2011, Art. ID 27.

[130] M. Asad and C. Abhayaratne, "Kinect depth stream pre-processing for hand gesture recognition," in *Proc. 20th IEEE ICIP*, Sep. 2013, pp. 3735–3739.

[131] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.

[132] Y. Wang, T. Yu, L. Shi, and Z. Li, "Using human body gestures as inputs for gaming via depth analysis," in *Proc. IEEE ICME*, Apr./Jun. 2008, pp. 993–996.

[133] R. Yang and S. Sarkar, "Coupled grouping and matching for sign and gesture recognition," *Comput. Vis. Image Understand.*, vol. 113, no. 6, pp. 663–681, 2009.

[134] B. Liang, "Gesture recognition using depth images," in *Proc. 15th ACM ICMI*, 2013, pp. 353–356.

[135] E. Ohn-Bar and M. M. Trivedi, "The power is in your hands: 3D analysis of hand gestures in naturalistic video," in *Proc. IEEE Conf. CVPR Workshops*, Jun. 2013, pp. 912–917.

[136] M.-K. Sohn, S.-H. Lee, D.-J. Kim, B. Kim, and H. Kim, "A comparison of 3D hand gesture recognition using dynamic time warping," in *Proc. 27th Conf. Image Vis. Comput. New Zealand*, 2012, pp. 418–422.

[137] P. Doliotis, A. Stefan, C. McMurrough, D. Eckhard, and V. Athitsos, "Comparing gesture recognition accuracy using color and depth information," in *Proc. 4th Int. Conf. Pervasive Technol. Rel. Assistive Environ.*, 2011, Art. ID 20.

[138] H.-M. Zhu and C.-M. Pun, "Real-time hand gesture recognition from depth image sequences," in *Proc. 9th Int. Conf. Comput. Graph., Imag., Visualizat.*, Jul. 2012, pp. 49–52.

[139] L. Miranda, T. Vieira, D. Martinez, T. Lewiner, A. W. Vieira, and M. F. M. Campos, "Real-time gesture recognition from depth data through key poses learning and decision forests," in *Proc. 25th SIBGRAPI Conf. Graph., Patterns, Images*, Aug. 2012, pp. 268–275.

[140] G. Ren and E. O'Neill, "3D selection with freehand gesture," *Comput. Graph.*, vol. 37, no. 3, pp. 101–120, 2013.

[141] Y.-K. Wu, H.-C. Wang, L.-C. Chang, and K.-C. Li, "Using HMMs and depth information for signer-independent sign language recognition," in *Proc. 7th Int. Workshop Multi-Discipl. Trends Artif. Intell.*, 2013, pp. 79–86.

[142] H. Wang, J. Fu, Y. Lu, X. Chen, and S. Li, "Depth sensor assisted real-time gesture recognition for interactive presentation," *J. Vis. Commun. Image Represent.*, vol. 24, no. 8, pp. 1458–1468, 2013.

[143] W.-Y. Lin and C.-Y. Hsieh, "Kernel-based representation for 2D/3D motion trajectory retrieval and classification," *Pattern Recognit.*, vol. 46, no. 3, pp. 662–670, 2013.

[144] J. Beh, D. Han, and H. Ko, "Rule-based trajectory segmentation for modeling hand motion trajectory," *Pattern Recognit.*, vol. 47, no. 4, pp. 1586–1601, 2014.

[145] H. Li and M. Greenspan, "Model-based segmentation and recognition of dynamic gestures in continuous video streams," *Pattern Recognit.*, vol. 44, no. 8, pp. 1614–1628, 2011.

[146] L. Morency, A. Quattoni, and T. Darrell, "Latent-dynamic discriminative models for continuous gesture recognition," in *Proc. IEEE Conf. CVPR*, Jun. 2007, pp. 1–8.

[147] H.-I. Suk, B.-K. Sin, and S.-W. Lee, "Hand gesture recognition based on dynamic Bayesian network framework," *Pattern Recognit.*, vol. 43, no. 9, pp. 3059–3072, 2010.

[148] G. Raffa, J. Lee, L. Nachman, and J. Song, "Don't slow me down: Bringing energy efficiency to continuous gesture recognition," in *Proc. Int. Symp. Wearable Comput.*, Oct. 2010, pp. 1–8.

[149] P. O. Kristensson and L. C. Denby, "Continuous recognition and visualization of pen strokes and touch-screen gestures," in *Proc. 8th Eurograph. Symp. Sketch-Based Interf. Modelling*, 2011, pp. 95–102.
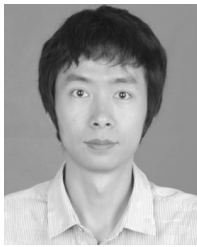
[150] M. Elmezain, A. Al-Hamadi, J. Appenrodt, and B. Michaelis, "A hidden Markov model-based continuous gesture recognition system for hand motion trajectory," in *Proc. 19th ICPR*, Dec. 2008, pp. 1–4.

[151] P. O. Kristensson, T. Nicholson, and A. Quigley, "Continuous recognition of one-handed and two-handed gestures using 3D full-body motion tracking sensors," in *Proc. ACM Int. Conf. Intell. User Interf.*, 2012, pp. 89–92.

[152] J. F. Lichtenauer, E. A. Hendriks, and M. J. T. Reinders, "Sign language recognition by combining statistical DTW and independent classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 11, pp. 2040–2046, Nov. 2008.

[153] T. Arici, S. Celebi, A. S. Aydin, and T. T. Temiz, "Robust gesture recognition using feature pre-processing and weighted dynamic time warping," *Multimedia Tools Appl.*, vol. 72, no. 3, pp. 3045–3062, 2014.

[154] A. Hernandez-Vela *et al.*, "BoVDW: Bag-of-visual-and-depth-words for gesture recognition," in *Proc. ICPR*, Nov. 2012, pp. 449–452.

[155] J. C. Niebles, H. Wang, and L. Fei-Fei, "Unsupervised learning of human action categories using spatial-temporal words," *Int. J. Comput. Vis.*, vol. 79, no. 3, pp. 299–318, 2008.

[156] M. K. Bhuyan, D. A. Kumar, K. F. MacDorman, and Y. Iwahori, "A novel set of features for continuous hand gesture recognition," *J. Multimodal User Interf.*, vol. 8, no. 4, pp. 333–343, 2014.

[157] Y. Sun, N. Kuwahara, and K. Morimoto, "Development of recognition system of Japanese sign language using 3D image sensor," in *Proc. Int. Conf. HCI*, 2013, pp. 286–290.

[158] X. Suau, M. Alcoverro, A. Lopez-Mendez, J. Ruiz-Hidalgo, and J. Casas, "INTAIRACT: Joint hand gesture and fingertip classification for touchless interaction," in *Proc. ECCV*, 2012, pp. 602–606.

[159] C.-S. Lee, S. Chun, and S. W. Park, "Tracking hand rotation and various grasping gestures from an IR camera using extended cylindrical manifold embedding," *Comput. Vis. Image Understand.*, vol. 117, no. 12, pp. 1711–1723, 2013.

[160] V. T. Hoang, A. N. Hoang, and D. Kim, "Real-time stereo rendering technique for virtual reality system based on the interactions with human view and hand gestures," in *Proc. 5th Int. Conf. Virt. Augmented Mixed Reality*, 2013, pp. 103–110.

[161] K. O'Hara *et al.*, "Touchless interaction in surgery," *Commun. ACM*, vol. 57, no. 1, pp. 70–77, 2014.

[162] M. Strickland, J. Tremaine, G. Brigley, and C. Law, "Using a depth-sensing infrared camera system to access and manipulate medical imaging from within the sterile operating field," *Can. J. Surgery*, vol. 56, no. 3, pp. E1–E6, 2013.

[163] J. Hoey, P. Poupart, A. von Bertoldi, T. Craig, C. Boutilier, and A. Mihailidis, "Automated handwashing assistance for persons with dementia using video and a partially observable Markov decision process," *Comput. Vis. Image Understand.*, vol. 114, no. 5, pp. 503–519, 2010.

[164] H. Cheng, Z. Liu, Y. Zhao, G. Ye, and X. Sun, "Real world activity summary for senior home monitoring," *Multimedia Tools Appl.*, vol. 70, no. 1, pp. 177–197, 2014.

[165] P. M. Yanik *et al.*, "Use of Kinect depth data and growing neural gas for gesture based robot control," in *Proc. 6th Int. Conf. Pervasive Comput. Technol. Healthcare*, May 2012, pp. 283–290.

[166] A. Dragan and S. Srinivasa, "Formalizing assistive teleoperation," in *Proc. Robot., Sci. Syst.*, 2012, pp. 1–8.

[167] D. Fischinger, "Hobbit, a care robot supporting independent living at home: First prototype and lessons learned," *Robot. Auto. Syst.*, vol. 75, pp. 60–78, Jan. 2016.

[168] G. Amayeh, G. Bebis, A. Erol, and M. Nicolescu, "Hand-based verification and identification using palm–finger segmentation and fusion," *Comput. Vis. Image Understand.*, vol. 113, no. 4, pp. 477–501, 2009.

[169] V. Kanhangad, A. Kumar, and D. Zhang, "Contactless and pose invariant biometric identification using hand surface," *IEEE Trans. Image Process.*, vol. 20, no. 5, pp. 1415–1424, May 2011.

[170] M. Roccetti, G. Marfia, and A. Semeraro, "Playing into the wild: A gesture-based interface for gaming in public spaces," *J. Vis. Commun. Image Represent.*, vol. 23, no. 3, pp. 426–440, 2012.

[171] S. Waldherr, R. Romero, and S. Thrun, "A gesture based interface for human–robot interaction," *Auto. Robots*, vol. 9, no. 2, pp. 151–173, 2000.

[172] X. Yin and X. Zhu, "Hand posture recognition in gesture-based human–robot interaction," in *Proc. IEEE Conf. Ind. Electron. Appl.*, May 2006, pp. 1–6.

[173] X. Yin and M. Xie, "Finger identification and hand posture recognition for human–robot interaction," *Image Vis. Comput.*, vol. 25, no. 8, pp. 1291–1300, 2007.

[174] K. Nickel and R. Stiefelhagen, "Visual recognition of pointing gestures for human–robot interaction," *Image Vis. Comput.*, vol. 25, no. 12, pp. 1875–1884, 2007.

[175] D. Michel, K. Papoutsakis, and A. A. Argyros, "Gesture recognition for the perceptual support of assistive robots," in *Proc. Int. Symp. Vis. Comput.*, 2014.

[176] F. Pedersoli, S. Benini, N. Adami, and R. Leonardi, "XKin: An open source framework for hand pose and gesture recognition using Kinect," *Vis. Comput.*, vol. 30, no. 10, pp. 1107–1122, 2014.

[177] P. Morerio, L. Marcenaro, and C. S. Regazzoni, "Hand detection in first person vision," in *Proc. 16th Int. Conf. Inf. Fusion*, Jul. 2013, pp. 1502–1507.

[178] G. Serra, M. Camurri, L. Baraldi, M. Benedetti, and R. Cucchiara, "Hand segmentation for gesture recognition in ego-vision," in *Proc. 3rd ACM IMMPD*, 2013, pp. 31–36.

[179] C. Li and K. M. Kitani, "Pixel-level hand detection in ego-centric videos," in *Proc. IEEE Conf. CVPR*, Jun. 2013, pp. 3570–3577.

[180] C. Li and K. M. Kitani, "Model recommendation with virtual probes for egocentric hand detection," in *Proc. IEEE ICCV*, Dec. 2013, pp. 2624–2631.

[181] A. Betancourt, M. M. López, C. S. Regazzoni, and M. Rauterberg, "A sequential classifier for hand detection in the framework of egocentric vision," in *Proc. IEEE Conf. CVPR Workshops*, Jun. 2014, pp. 600–605.

[182] L. Baraldi, F. Paci, G. Serra, L. Benini, and R. Cucchiara, "Gesture recognition in ego-centric videos using dense trajectories and hand segmentation," in *Proc. IEEE Conf. CVPR Workshops*, Jun. 2014, pp. 702–707.

[183] F. Coleca, T. Martinetz, and E. Barth, "Gesture interfaces with depth sensors," in *Proc. Time-Flight Depth Imag.*, 2013, pp. 207–227.

[184] C. Keskin, A. T. Cemgil, and L. Akarun, "DTW based clustering to improve hand gesture recognition," in *Proc. Human Behavior Understanding*, Nov. 2011, pp. 72–81.

**Hong Cheng** (M'06–SM'14) received the Ph.D. degree in pattern recognition and intelligent systems from Xi'an Jiaotong University, Xi'an, China, in 2003.

He was a Post-Doctoral Researcher with the School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA, from 2006 to 2009. He has been an Associate Professor with Xi'an Jiaotong University since 2005. He is currently a Full Professor with the School of Automation Engineering, University of Electronic Science and Technology of China, Chengdu, China, where he is also the Executive Director of the Center for Robotics. He has over 80 academic publications including two books. His research interests include computer vision, machine learning, and robotics.

Dr. Cheng is a Senior Member of the Association for Computing Machinery and an Associate Editor of *IEEE Computational Intelligence Magazine*. He is a Reviewer of many important journals and conferences, such as IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS, Machine Vision and Application, the Conference on Computer Vision and Pattern Recognition, the International Conference on Computer Vision, International Conference on Intelligent Transportation Systems, the Israel Vacuum Society Conference, and the Asian Conference on Computer Vision. He serves as the General Chair of Vision And Learning Seminar 2015, the Finance Chair of the IEEE International Conference on Multimedia and Expo in 2014, the Local Arrangement Chair of USA–Sino Summer School in Vision, Learning and Pattern Recognition 2012, and the Registration Chair of the IEEE International Conference on Vehicular Electronics and Safety in 2005.

**Lu Yang** (M'11) received the bachelor's degree from University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2005, and the Ph.D. degree in electrical engineering and computer science from Nagoya University, Nagoya, Japan, in 2011.

He was with the joint Doctor–Master Program, School of Automation Engineering, UESTC, from 2005 to 2008. He was sponsored by the China Scholarship Council to study at Nagoya University from 2008 to 2011. Then, he joined the School of Automation Engineering, UESTC, in 2011, as a Lecturer, where he has been an Associate Professor with the Center for Robotics since 2013. His research interests include computer vision, pattern recognition, machine learning, and 3D multimedia applications.

Dr. Yang is the Committee Member of the IEEE International Conference on Multimedia and Expo in 2014, Vision And Learning Seminar 2015, the China Summit and International Conference on Signal and Information Processing in 2015, and the Pacific-Rim Conference on Multimedia in 2015.

**Zicheng Liu** (SM'05–F'15) received the B.S. degree in mathematics from Huazhong Normal University, Wuhan, China, in 1984; the M.S. degree in operation research from Institute of Applied Mathematics, Chinese Academy of Sciences, Beijing, China, in 1989; and the Ph.D. degree in computer science from Princeton University, Princeton, NJ, USA, in 1996.

He was with Silicon Graphics, Inc., Mountain View, CA, USA. He is currently a Principle Researcher with Microsoft Research, Redmond, MA, USA. His research interests include human activity recognition, 3D face modeling and animation, and multimedia signal processing.