

# A Machine Vision based Gestural Interface for People with Upper Extremity Physical Impairments

Hairong Jiang, Bradley S. Duerstock, Juan P. Wachs

**Abstract**— A machine vision based gestural interface was developed to provide individuals with upper extremity physical impairments an alternative way to perform laboratory tasks that require “physical” manipulation of components. A color and depth based 3D particle filter framework was constructed with unique descriptive features for face and hands representation. This framework was integrated into an interaction model utilizing spatial and motion information to deal efficiently with occlusions and its negative effects. More specifically, the suggested method proposed solves the “false merging” and “false labeling” problems characteristic in tracking through occlusion. The same feature encoding technique was subsequently used to detect, track and recognize users’ hands. Experimental results demonstrated that the proposed approach was superior to other state of art tracking algorithms when interaction was present (97.52% accuracy). For gesture encoding, dynamic motion models were created employing the dynamic time warping (DTW) method. The gestures were classified using a Conditional Density Propagation (CONDENSATION)-based Trajectory Recognition (CTR) method. The hand trajectories were classified into different classes (commands) with a recognition accuracy of 95.9%. In addition, the new approach was validated with the “one shot learning” paradigm with comparable results to those reported in 2012. In a validation experiment, the gestures were used to control a mobile service robot and a robotic arm in a laboratory chemistry experiment. Effective control policies were selected to achieve optimal performance for the presented gestural control system through comparison of task completion time between different control modes.

**Index Terms**—Gesture recognition, particle filter, dynamic time warping, CONDENSATION, one shot learning.

## I. INTRODUCTION

Assistive technologies is about finding new ways to engage cutting-edge technologies in support of individuals with physical and/or cognitive impairments. The development of technologies relying on high usability principles, exploited new communicational channels such as eye blinking, voice,

hand gestures, sip and puff, and electromyogram (EMG) as effective control modalities [1]. These channels have led to ingenious interfaces in support of the disabled [2], such as intelligent wheelchairs systems, home medical alert systems, and assistive robotic control, to mention a few [3, 4, 5]. These interfaces offer additional degrees of mobility and control which were not possible previous to these developments, leading to a higher life quality and sense of independence.

Among all these interaction channels, hand gestures is a valuable alternative since it does not require to have the user tethered through cables or sensors, and it only requires learning a few customized gestures for a given task. In particular, upper extremity gesture control can serve as an important human computer interaction (HCI) modality for individuals with quadriplegia who lack hand fine motor skills. For instance, upper limb gesture control requires less targeting accuracy than joysticks, the mouse, and other continuous input devices. Likewise, the option to employ either continuous or discrete input control modes reduces the effort required for individuals with quadriplegics to perform navigational operations [6]. Unlike voice control, gesture control is effective in noisy environments [7]. In addition, for most of the cases, individuals with quadriplegic can only use gross motor function instead of fine motor function to perform certain tasks [8]. Apart from other common modalities, such as keyboard and joystick that require fine motor control to hit a key or move and twist a handle, upper extremity gesture control only requires gross motor function for targeting and navigational tasks [9]. Lastly, hand gesture based HCI is unencumbered because it does not require the user to directly contact or wear sensors as sip-and-puff and EMG based systems [10, 11]. While not every individual with upper extremity mobility impairments can use hand gesture control reliably, for those who are able to move their arms to some degree, gesture-based HCI can be seen as a promising alternative or complement to an existing control modality.

In our previous work [12], a prototype of gesture recognition based interface was developed for people with upper extremity mobility impairments. In the current manuscript, the tracking algorithm was greatly improved and compared with five state-of-art algorithms to demonstrate a better tracking performance. Further, more experimental results were provided with subjects with upper extremity mobility impairments and one shot learning was employed to allow instant customization of the gestural system. Face and hand tracking under frequently self-occlusion was modeled as

---

Manuscript received August 20, 2012. This work is partially funded by the National Institutes of Health through the NIH Director's Pathfinder Award to Promote Diversity in the Scientific Workforce, grant number DP4-GM096842-01.

Hairong Jiang is with School of Industrial Engineering, Purdue University, West Lafayette, IN 47907, USA (e-mail: jiang115@purdue.edu).

Bradley S. Duerstock is with School of Industrial Engineering and Weldon School of Biomedical Engineering, Purdue University, West Lafayette, IN 47907, USA (e-mail: bsd@purdue.edu).

Juan P. Wachs is with School of Industrial Engineering, Purdue University, West Lafayette, IN 47907, USA (e-mail: jpwachs@purdue.edu).

a multi-object tracking (MOT) problem. This problem is challenging since hands are non-rigid objects and their form varies among individuals, while performing a certain candidate gesture. Additionally, since the appearance of the left and right hand are similar for the same individual, trackers can focus on one hand or exchange positions when the hands are too close to each other. In this paper, an integrated approach was proposed to tackle the challenging problem of tracking under self-occlusion.

#### A. Related Work

Often, hand gesture recognition involves segmentation of the hands, tracking them through occlusion, and the classification of hand's dynamic trajectories and static pose. For vision-based real-time gesture based interfaces for assistive technologies, robustness is a critical requirement [13] for its adoption. For hand segmentation, a commonly used method is to back-project the pre-built skin color histogram model into new video frames. These methods are likely to fail in true world conditions, where illumination is uncontrolled and the background is cluttered. Adding depth information can relax at some extent the problem, by utilizing stereo vision [14] or other depth commodity sensors, such as Kinect™ [15] or Leap Motion® [16].

Face and hands tracking is a special case of MOT problem. If gestures in the lexicon only carry trajectory information, (the hand shape does not convey extra information), classical tracking approaches can be adopted. For example, CAMSHIFT [17] and CONDENSATION [18] have been shown to successfully track gestures; however they are susceptible to lose the tracked objects when occluded by new objects, or when the scene illumination changes. Another widely used technique for object tracking is particle filters [19]. Perez et al. [20] integrated color-based appearance models to a particle filter framework to enhance tracking under complex background and occlusion, and then applied the particle filter framework to multiple objects tracking. Okuma et al. [21] further extended particle filters by incorporating a boosting detector and enabling automatic initialization of potential multiple targets. One problem of these techniques is that the interaction between the tracked objects (and occlusion) was not considered part of the main framework. When the objects interact one with the other, occlusions occur frequently. Local motion information was incorporated into a color-based particle filter framework by Kristan et al. [22] to solve the self-occlusion problem through object tracking. Qu et al. [23] combined a joint state space representation with color-based particle filter and performed joint data association in a multi-object tracking scenario. All the discussed algorithms so far, attempted to solve the MOT problem; however they presented limited performance when tracking multiple non-rigid similar objects.

With the advent of Kinect™ and other 3D sensors, hand or body tracking techniques in real-time were exploited. Eichner et al. [24] presented a technique to estimate the body layout of humans by using still images. Their approach is capable of estimating upper body pose in highly uncontrolled

environment. Further, Yang et al. [25] described a method to estimate human pose from static images using body part models. By using the depth information, Shotton et al. [26] proposed a method to predict 3D positions of body joints from a single depth image. They solved the pose estimation problem through a simple per-pixel classification problem. A similar method is also used by OpenNI for human body skeleton tracking. One problem of these skeleton-based tracking methods is that they work well when users are standing with their extremities extended, but suffer sudden performance degradation for seated users with contracted limbs, as often occur with quadriplegic individuals.

Only color and depth information captured from Kinect™ were adopted for hand tracking by Oikonomidis et al. [27]. They presented a method to track the full articulation of two hands that interact with each other in an uncontrolled manner. This method is effective for static gesture recognition; however, the computation cost is excessive which affects its real-time extension for gesture tracking.

One of the most widely used techniques for gesture recognition is Hidden Markov Models (HMM) [28, 29, 30]. Common problems with HMM approach consist of finding the optimal parameters set (e.g. initial probabilities) and trajectory spotting for gesture temporal segmentation. Black and Jepson [31] proposed a CONDENSATION-based trajectory gesture recognition algorithm that can obtain less sensitive parameters set and achieve robust tracking, yet gesture temporal segmentation was not fully addressed. Alon, et al. [32] applied the dynamic time warping (DTW) approach to gesture recognition and look at sub-gestures composition to solve the temporal segmentation problem (also known as “spotting”). Interaction between hands was not specifically tackled.

Recently, a new type of challenge was attracted the attention of the gesture recognition community – the “One Shot Learning” Challenge [33]. The one shot learning [34] consists of learning a gesture category by only observing one instance of that gesture, similar to how humans learn. In this context, Wu et al. [35] adopted the extended-motion-histogram image for motion feature representation and applied it to segment and classify hand gestures. Yang et al. [36] proposed discovering high level sub-actions by clustering optical flow in four dimensions (RGB-D). In our work, one shot learning provides an interesting test-bed to demonstrate the robustness of our approach, compared with the state of the art [37].

In the current paper, we also extended our method to robotic control. One advantage of using hand gestures to control robots is that it provides a natural way for navigational tasks by sending navigational information (e.g. left, right, forward and backward commands [38]).

#### B. Outline of Our Approach

In this paper, an interaction model was incorporated to the color histogram based particle filter framework to track hands through interaction and occlusion. A procedure was proposed to create dynamic motion models by DTW method and classify input gesture trajectories using the CONDENSATION

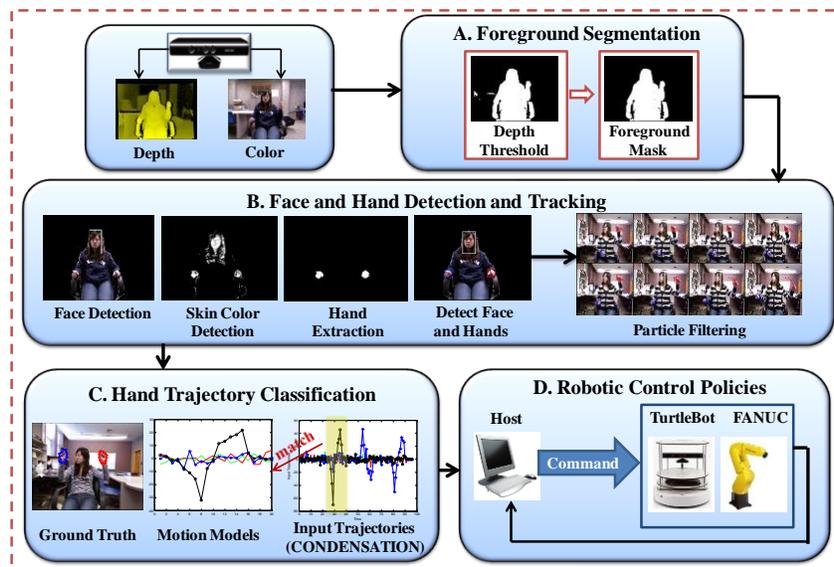


Fig. 1. System Overview.

algorithm. The system was integrated in a simplistic yet robust fashion by combining CONDENSATION algorithm with an interaction model-based particle filter, which makes it suitable for human robot interaction in assistive technologies.

The contribution of this paper is three-fold: (1) prove the effectiveness of hand gestures as an alternative modality for individuals with mobility impairments by both subjective explanation and quantified results; (2) solve the frequently hand gesture interaction and occlusion problem through integration of color and 3D spatial information as an interaction model; (3) new gestures can be created and learned through the one shot learning paradigm, leading to an almost effortless training process (a necessary attribute for subjects with severe spinal cord injuries).

The paper is organized as follows: In Section II, the architecture is presented for the gesture recognition system. In Section III, the approach suggested to track and recognize dynamic hand gestures is discussed in details. In Section IV, comparative tests and results are presented, Section V discusses and concludes the paper, and Section VI presents future work.

## II. SYSTEM ARCHITECTURE

The architecture of the proposed system is illustrated in Fig. 1. Eight gestures were selected to constitute the gesture lexicon which in turn was used to control the robots. The machine vision based gestural system included four parts: foreground segmentation, face and hand detection and tracking, hand trajectory classification, and robotic control policies. Those parts were described in the following sections.

### A. Foreground Segmentation

In foreground segmentation section, the background was ruled out from the captured frames and the whole human body was kept as the foreground.

### B. Face and Hand Detection and Tracking

Face and hand detection was to initialize the position of the

face and hands for the tracking phase. After initialization, both face and hands were tracked through video sequences by particle filter method.

### C. Hand Trajectory Classification

Hand tracking results were segmented as trajectories, compared with motion models, and decoded as commands for robotic control.

### D. Robotic Control Policies

The commands decoded by gesture recognition results were sent to control the mobile robot and the robotic arm.

## III. GESTURE RECOGNITION

### A. Foreground Segmentation

Initially, the user's body was treated as a foreground object in order to detect the user's movements. Two steps were used to segment the foreground (refer to algorithm 1 in TABLE I). In the first step, the sensed image assessed by a Kinect<sup>TM</sup> [15] sensor was thresholded using depth information. The depth value of each pixel was defined as  $D(i, j)$  with  $i$  and  $j$  indicating the horizontal and vertical coordinates of the pixel in each frame of the video sequence. An example of a depth image is shown by Fig. 2(a), where the distance between objects and the depth sensor was mapped to intensity levels. The nearer the object was to the sensor, the larger the intensity was. Two absolute depth thresholds (a low threshold  $T_{DL}$  and a high threshold  $T_{DH}$ ) were custom set by the user according to their relative distance to the depth sensor.  $T_{DL}$  was set to no less than a constant which was the minimum distance that can be registered by the depth sensor (due to its physical limitations).  $T_{DH}$  was set to be the maximum distance that can be reached by the user while seated in a wheelchair<sup>1</sup>. In this paper,  $T_{DL}$  and  $T_{DH}$  were set to be 0.4m and 2.0m to achieve an

<sup>1</sup> These values are selected since they resulted in the best performance; other thresholds can be used and the impact on the overall performance is likely to be negligible.

optimal performance for segmentation. A mask image (Fig. 2(b)) was generated by keeping the pixels with a depth value between the two thresholds while discarding the others. In the second step, the region (blob) with the largest area (denoted as  $T_{SH}$ ) was extracted from the mask image. All the remaining blobs with an area smaller than  $T_{SH}$  were discarded (Fig. 2(c)). If the extracted region contained an object that were not part of the user's body, it would be discarded in a later stage since tracking was achieved based on both color and spatial information.

TABLE I  
FOREGROUND SEGMENTATION ALGORITHM

Algorithm 1: Foreground Segmentation	
<b>Input:</b>	Low depth threshold $T_{DL}$ ; High depth threshold $T_{DH}$ ; pixel value of depth Image $D(i, j)$ ;
<b>Output:</b>	pixel value of mask image $D_1(i, j)$ ; pixel value of foreground mask image $D_2(i, j)$ .
$D_1(i, j) = \begin{cases} 1: & T_{DL} \leq D(i, j) \leq T_{DH} \\ 0: & \text{otherwise} \end{cases}$	
$T_{SH} = \max(\text{Area}(B_i)) \quad // B_i \text{ is the } i\text{th blob in the mask image } D_1$	
$D_2(i, j) = \begin{cases} 1: & D_1(i, j) \in B_i \ \& \ \text{Area}(B_i) == T_{SH} \\ 0: & \text{otherwise} \end{cases}$	

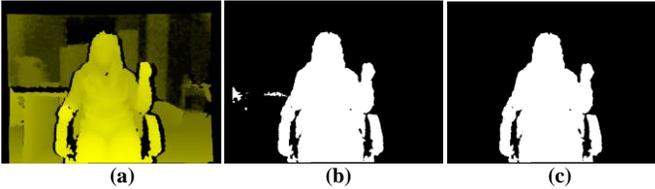


Fig. 2. Foreground Segmentation. (a) Depth image; (b) Depth threshold mask; (c) Foreground segmentation mask.

### B. Face and Hand Detection

In this section, the centroids of the face and hand regions were extracted to initialize the tracking stage. Two 3D histograms - a skin and a non-skin color histogram were created using Compaq database [39] and HSV color space to achieve higher robustness for skin color detection (referred to [12] for a detailed description). The mask image obtained from histogram back-projection is shown as in Fig. 3(a). To obtain the hand regions without the face, a face detector [40] was adopted (Fig. 3(c)) to remove the region from the target image. Two largest blobs in the target image were then selected as hand regions (Fig. 3(b)). The centroids of the hands were obtained by computing the first moment of the two blobs. This hand detection procedure was only used to provide automatic initialization to the particle filter tracking procedure. Afterwards the hands positions were continuously tracked by the particle filter.

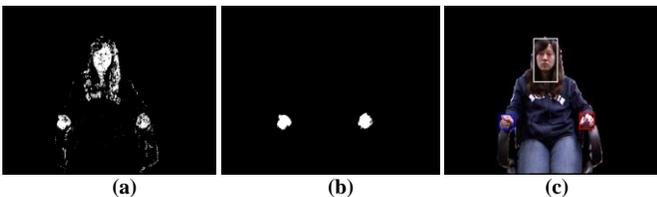


Fig. 3. Face and hand detection. (a) Skin color detection; (b) Hand extraction; (c) Face and hand localization.

### C. Face and Hand Tracking

A 3D particle filter framework based on color, depth and

spatial information was used to track the face and hands through video sequences. A detailed description of the particle filter algorithm was illustrated in [20, 41, 42]. The equation of particle filtering is given as in (1):

$$p(X_t|Z_{1:t}) = k \cdot p(Z_t|X_t) \int p(X_t|X_{t-1}) p(X_{t-1}|Z_{1:t-1}) dx_{t-1} \quad (1)$$

where  $X_t$  is the process state at time  $t$ ,  $Z_{1:t} = \{Z_1, \dots, Z_t\}$  denotes the set of observations from time 1 to  $t$ ,  $p(X_t|Z_{1:t})$  and  $p(Z_t|X_t)$  expresses the posterior and prior distribution at time  $t$ ,  $p(X_t|X_{t-1})$  is the transition probability of the system at state  $X_t$  given that the previous state was  $X_{t-1}$ , and  $k$  is a normalization factor to normalize the sum of all posterior probability to 1. In the particle filter algorithm,  $N$  weighted particles can be used to approximate the posterior as:  $p(X_{t-1}|Z_{1:t-1}) \approx \{X_{t-1}^r, \omega_{t-1}^r\}_{r=1}^N$ , where  $\omega_{t-1}^r$  denotes the weight of the particle  $r$  at time  $t-1$ . After propagation, the tracker output at time  $t$  can be approximated by the expectation of the process state:  $\hat{X}_t \approx E[X_t|Z_{1:t}] = \sum_{r=1}^N \omega_t^r X_t^r$ . Thus, (1) is converted to (2):

$$p(X_t|Z_{1:t}) \approx k \cdot p(Z_t|X_t) \sum_{r=1}^N \omega_{t-1}^r p(X_t^r|X_{t-1}^r) \quad (2)$$

The particles were initialized by using the centroids of face and hands calculated in section III. B.

The particle filter tracking process consists of three main phases: predicting, measuring and re-sampling. In the proposed system, for the predicting phase, a second order auto-regressive (AR) model (as in (3)) [20, 41] was selected to model the dynamic motion of each particle:

$$X_t^r = A_1(X_{t-1}^r - X_0^r) + A_2(X_{t-2}^r - X_0^r) + X_0^r + Bv_t \quad (3)$$

Where  $v_t \sim \mathcal{N}(0, \Sigma)$  is a Gaussian distribution with zero mean and variance matrix  $\Sigma$ ,  $X_0$  is the original particle coordinate,  $A_1$ ,  $A_2$ , and  $B$  are the optimal parameter matrices that can best match the real motion of the tracked object,  $X_t^r$  is the state of the particle  $r$  at time  $t$ . In this paper, a 3D particle filter tracking was adopted. The state of particle  $r$  at time  $t$  is written as:  $X_t^r = [x_t^r, y_t^r, z_t^r, s_t^r, x_{t-1}^r, y_{t-1}^r, z_{t-1}^r]$ , where  $s_t$  is the scale of object at time  $t$ ,  $x_t^r, y_t^r, z_t^r$  are the 3D coordinates of particle  $r$  at time  $t$ , and  $x_{t-1}^r, y_{t-1}^r, z_{t-1}^r$  are the 3D coordinates of particle  $r$  at time  $t-1$ .

For the measuring phase, the selection of the observation model determines the weight of the particles. Many appearance-based models, such as contour, edge, piece-wise, etc, were used in object tracking. Color-based pre-processing using HSV space can facilitate the extraction of the aforementioned features for face and hands tracking. As explained earlier, the initial phase of the face and hands were determined by the combination of depth-based thresholding and image processing techniques. The extracted face and hands regions were used to compute the reference HSV histogram models ( $H_f^*$ ,  $H_{hl}^*$ , and  $H_{h2}^*$ ) for tracking initialization. During the re-sampling phase, each particle, assigned in the predicting phase, was reweighted by the observation likelihood function. For every hypothesized face

or hand location of a particle  $r$ , the candidate histograms were computed as  $H_f^r, H_{h1}^r$ , and  $H_{h2}^r$ . The Bhattacharyya distance [20]  $D$  was used to measure similarity between reference and candidate histograms as (4):

$$D_i(H^*, H^r) = \left[1 - \sum \sqrt{H_i^* H_i^r}\right]^2 \quad (4)$$

where  $H_i^* = H_f^*, H_{h1}^*$ , or  $H_{h2}^*$  and  $H_i^r = H_f^r, H_{h1}^r$ , or  $H_{h2}^r$ . The observation likelihood function can be written as (5):

$$p(Z_t|X_t) \propto \exp(-\lambda_1(D_i^r)^2) \quad (5)$$

where  $\lambda_1$  measures the variance of the HSV histogram. (5) can be re-written by adding a normalization factor  $k$  to normalize the sum of all particles' weight to 1, obtaining (6):

$$p(Z_t|X_t) = k \cdot \exp(-\lambda_1(D_i^r)^2) \quad (6)$$

#### D. Hand Tracking Through Interaction and Occlusion

Color-based particle filter tracking was effective for multiple independent objects tracking when the objects did not interact or occlude each other. However, if interaction or occlusion occurs, multiple independent particle filters can be used. Standard multi-object tracking (MOT) with interaction and occlusion suffers from the “false merging” and “false labeling” problems [23]. The “false merging” problem denotes the situation that the tracker shift from the object being tracked to a different object that has higher observation likelihood. Conversely, the “false labeling” problem denotes the situation that the objects being tracked exchange their labels after interaction or occlusion occurred. In the proposed system, the face and both hands were tracked.

In this paper, two models were constructed to solve “false merging” and “false labeling” problems separately. The first model was called the “Competition Potential” (CP) model. The idea of this model comes from the Joint Markov random fields (MRF) theory [42]. The likelihood function for CP model is defined as  $\psi_1(X_{i,t}, X_{j,t})$ , which represented the pairwise interaction potential of the MRF [43]. The second model is called “Motion Consistency” (MC) model. The likelihood function for MC model is defined as  $\psi_2(X_{i,t}, X_{j,t})$ , which is based on the assumption that a particle region that has similar motion information to the previous state of that particle will have higher probability than a particle region that has distinct motion information.

For CP model, as in [43], we have  $p(X_t|X_{t-1}) \propto \prod_{i,j \in E} \psi_1(X_{i,t}, X_{j,t}) \prod_{i,j \in E} \psi_1(X_{i,t}, X_{j,t})$ . The particle filter function (2) can be rewritten as (7):

$$p(X_t|Z_{1:t}) = k \cdot p(Z_t|X_t) \prod_{i,j \in E} \psi_1(X_{i,t}, X_{j,t}) \sum_r \omega_{t-1}^r \prod_i p(X_{i,t}|X_{i,t-1}^r) \quad (7)$$

The likelihood function for CP model is then defined as:

$$\psi_{1i,t}^r(X_{i,t}, X_{j,t}) = \beta_1 \cdot \exp\left(-\frac{\lambda_2}{d(X_{i,t}^r, X_{j,t})^2}\right) \cdot \exp\left(-\lambda_3 d(X_{i,t}^r, X_{j,t-1})^2\right) \cdot \exp\left(-\frac{\lambda_4}{d_z(X_{i,t}^r - X_{j,t})^2}\right) \quad (8)$$

where  $d(X_{i,t}^r, X_{j,t})$  denotes the 2D Euclidean distance metric between two objects,  $d(X_{i,t}^r, X_{i,t-1})$  represents a distance metric between the previous and current centroid of object  $i$ ,  $d_z(X_{i,t}^r - X_{j,t})$  represents the difference of depth value between two objects, and  $\beta_1$  is a normalization factor so the sum of all particles' weight is 1.

MC model was used to solve the “false labeling” problem. The 3D motion information was incorporated into the original likelihood function to increase the robustness of the method. We adopted a compact expression of the likelihood function similar as in [23], which integrated the magnitude and direction information of motion as (9). Instead of using 2D, 3D motion features are used to compute the motion information of the hand movement. The likelihood function for the MC model is defined as:

$$\psi_{2i,t}^r(X_{i,t}, X_{j,t}) = \beta_2 \cdot \exp(-\lambda_5(\theta_t^r)^2) \cdot \exp(-\lambda_6(A_t^r - A_{ref,t})^2) \quad (9)$$

where  $A_t^r$  and  $A_{ref,t}$  represent the norm of 3D motion vector and reference motion vector (can be computed by the difference of the current and the previous 3D position vector) of particle  $r$  at state  $t$ , respectively.,  $\theta_t^r$  is the angle between the 3D motion vector of particle  $r$  and the reference vector and  $\beta_2$  is a normalization factor to normalize the sum of all particles' weight to 1. This likelihood function assumes that a particle region that has a similar motion to the previous state will have a larger weight than one with a different motion.

When the objects' observations do not interact with each other, the approach suggested behaves as if multiple independent trackers were applied to the objects (Fig. 4(a)). However, when the objects' observations interact (e.g. partial or complete occlusion occurs), the conventional particle filter framework is extended (Fig. 4(b), (c)). The decision of when the objects interact is made based on the interdistance between the hands. When this distance is below a certain threshold, the system switches to the interaction model (Fig. 4). To find the optimal threshold, a histogram of the number of tracking frame errors at each distance is obtained (Fig. 5) and the threshold is selected so the tracking errors are minimized when the interaction model is activated. Note, a dramatic decrease of error, at the distance around 50 pixels at which hand interaction frequently occurred. The threshold  $T$  was determined according to the distribution of errors in the histogram. The extension models were given through (8) and (9). The parameters  $\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5$ , and  $\lambda_6$  in (6), (8) and (9) were optimized by utilizing a neighborhood search method [38]. The algorithm for hand tracking during interaction and occlusion is shown by TABLE II.

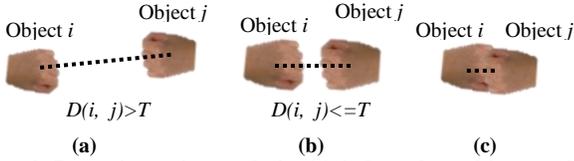


Fig. 4. Dynamic motion analysis. (a) independent object tracking; (b) interaction model added; (c) objects occlusion occurs.

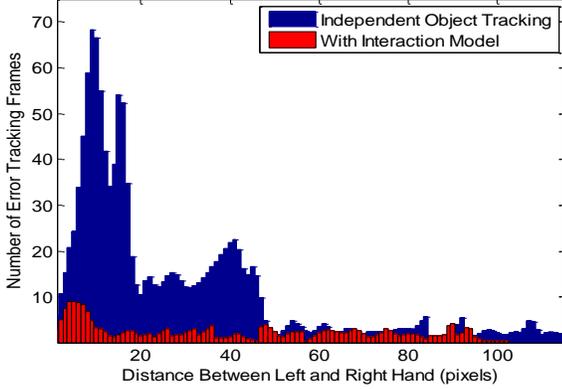


Fig. 5. Number of error tracking frames for each distance before and after the addition of interaction model.

TABLE II  
HAND TRACKING THROUGH INTERACTION AND OCCLUSION

**Algorithm 2: 3D Particle Filter tracking**

**Input:** Reference HSV histogram models  $H_f^*$ ,  $H_{h1}^*$ , and  $H_{h2}^*$ ; Optimal parameter  $\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5, \lambda_6$ .

**Output:** Centroids and the associated bounding box of the face and hands.

**1. Initialize:**

//Initialize particle states and weight for face and both hands as:

$$x_0^i = x_0^*, \quad \omega_0^i = \frac{1}{n}, \quad \text{where } i = 1, \dots, n$$

**2. Predict, Measure and Resample:**

//Select k;

**for**  $i=1,2,3$  //(1-face, 2-right hand, 3-left hand)

**for**  $r=1$  to  $N$

$$x_{i,t}^r = A_1(x_{i,t-1}^r - x_0^i) + A_2(x_{i,t-2}^r - x_{i,0}^r) + x_{i,0}^r + Bv_t$$

//Compute candidate histograms  $H^r$

$$D_i(H^*, H^r) = [1 - \sqrt{H^* H^r}]^2$$

//Calculate the weight:

$$\omega_{i,t}^r = k \exp(-\lambda D_{i,t}^2)$$

**end for**

Normalize the weights and resample the particles

$$\text{Estimate } \hat{x}_{i,t} = \sum_{r=1}^N \omega_{i,t}^r x_{i,t}^r$$

//Check interaction

**if** interactions happens for object  $i$  and  $j$

**for**  $q=1, \dots, N$

//compute interaction likelihood  $\psi_1$  and  $\psi_2$ :

Compute  $\psi_{1,i,t}^q(x_{i,t}^q, x_{j,t}^q)$  and  $\psi_{2,i,t}^q(x_{i,t}^q, x_{j,t}^q)$  using (8) and (9)

//Calculate the weight:

$$\omega_{i,t}^q = \omega_{i,t}^q \cdot \psi_{1,i,t}^q \cdot \psi_{2,i,t}^q$$

**end for**

Normalize the weights and resample the particles.

$$\text{Estimate } \hat{x}_{i,t} = \sum_{r=1}^N \omega_{i,t}^r x_{i,t}^r$$

**end if**

**end for**

**E. Gesture Lexicon**

A gesture lexicon was designed such that users with physical impairments can perform the gestures with minimal effort. These gestures were found through a series of interviews conducted with subjects with upper mobility impairments. Borg Scale [44] was used to rank the physical stress required

to perform a gesture by participants with upper mobility impairments. An eight-gesture lexicon (see Fig. 6) was then constructed by analyzing the Borg Scale results collected from the subjective rankings and selecting those corresponding to the least required effort. A detailed description of the process for gesture lexicon construction can be referred to [45].

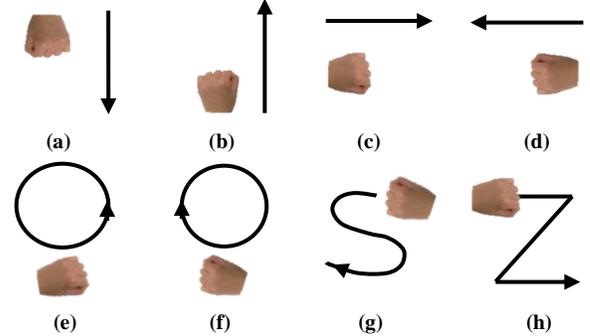


Fig. 6. Gesture lexicon. (a) upward; (b) downward; (c) rightward; (d) leftward; (e) counter-clockwise circle; (f) clockwise circle; (g) S; (h) Z.

**F. Hand Trajectory Classification**

For each frame in the video sequence, the centroids of the face and hands were obtained from the tracking stage. The motion model for each gesture trajectory was created based on the data collected from gestures performed by ten subjects. Two of the pool of ten subjects were quadriplegic due to a cervical spinal cord injury.

Even though the trajectories for each gesture performed by different subjects or the same subject in different instances may look similar, the precise duration of each sub-trajectory within the trajectory were different. To normalize the trajectories (temporal alignment), dynamic time warping (DTW) was employed [46]. The velocities' components in horizontal, vertical and depth directions of both hands were selected as the feature components for each motion model [41]. The procedure to construct the motion models is described in our previous work [12].

The CONDENSATION (Conditional Density Propagation) algorithm [31] was employed to classify hand gesture trajectories in the lexicon (as in Fig. 6). It employs a set of weighted samples instead an equation to fit the observed data. The original algorithm in [31] was extended to work for two hands. The original expression  $S_t = (\mu, \phi, \alpha, \rho)$  (the state at time  $t$ ) was extended to:

$$S_t = (\mu, \phi^i, \alpha^i, \rho^i) = (\mu, \phi^{right}, \phi^{left}, \alpha^{right}, \alpha^{left}, \rho^{right}, \rho^{left}) \quad (10)$$

where,  $\mu$  is the index of the motion models,  $\phi$  is the current phase in the model,  $\alpha$  is an amplitude scaling factor,  $\rho$  is a time dimension scaling factor, and  $i \in \{\text{right hand, left hand}\}$ .

The gestures in the lexicon (as in Fig. 6) were spotted using a rest position gesture (when the subjects put their hands on the arm rest (neural position) with no hands movement). A dynamic motion model was created for the rest position gesture. The segment between two recognized discontinuous rest position gestures is treated as a spotted gesture.

TABLE III  
HAND TRACKING PERFORMANCE

Method	False Merging (6080 Frames)	False Labeling (157 Interactions)	Tracking Accuracy (%)	Particle Number	Number of Body Parts
MCMC [41]	568	33	74.87	100	--
MI [23]	323	20	82.58	100	--
ETH (color) [24]	11	4	74.80	--	6
ETH (depth) [24]	0	3	58.25	--	6
Body Part (color) [25]	72	33	64.80	--	26
Body Part (depth) [25]	220	14	27.58	--	26
Kinect Skeleton [48]	0	3	73.87	--	16
CPMC(Proposed)	1	4	97.52	100	--

### G. Gesture Customization (One Shot Learning)

One of the objectives of our prototype follows the “came as you are” paradigm [13], where new gestures can be learned by the system automatically or by observing only one instance of it. The reason for this is to reduce the level of effort involved in the training phase of the system. In this section, we validated our approach in the context of one shot learning to assess the ability of the system to generalize learning from very few observations. A Savitzky-Golay smoothing filter [47] was added to smooth the 3D trajectories during the creation of the motion models. Two new gestures (see Fig. 7) were added to the lexicon to offer another degree of navigational control.

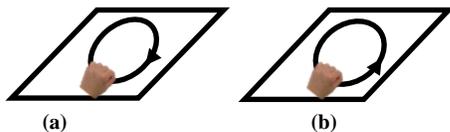


Fig. 7. Extended lexicon. (a) clockwise circle in horizontal plane; (b) counter-clockwise circle in horizontal plane.

## IV. EXPERIMENTS AND RESULTS

### A. Experiment 1: Hand Tracking Performance

A dataset of 16 videos (4 subjects x 4 activities) was used to evaluate the proposed tracking algorithm. The videos were captured with a Kinect™ camera at 30Hz using an image size of 640×480. Among the 4 subjects, three were able-bodied individuals and one was an individual with Cervical-6 level quadriplegia. The 4 activities performed by the subjects were: (a) “holding a cup”, (b) “clapping hands”, (c) “moving one hand up and down” (to occlude the other hand), (d) “rotating two hands forward and backward” (to occlude each other). The total number of frames for all the videos was 6080, while the total number of interactions between the two hands was 157. The total number of frames of each video corresponding to each of the four activities was: 930, 2230, 1320, and 1600, respectively (two sample sequences are shown in Fig. 17 and 18 in appendix). The ground truth position of the left and right hands in each video was provided by manually hand labeling.

The local likelihood  $p(z_t^i | x_t^i)$  was calculated using the 3D color histograms and two interaction models as the algorithm mentioned in TABLE II. The performance of the proposed method – competition potential and motion consistency (CPMC) – was compared to other existing methods such as: (i) Markov Chain Monte Carlo (MCMC) based particle filter

tracking [43], (ii) Magnetic-Inertial based particle filter tracking [23], (iii) ETH skeleton tracking based on color or depth frames [24], (iv) Body-parts tracking based on color or depth frames [25], and (v) Kinect™ OpenNI SDK skeleton tracking [48]. For the method (i), (ii), and the proposed method, 100 particles were used for face and each hand tracking. For methods (iii), (iv) and (v), the number of body parts (segments of a human body, i.e. hand, head, leg, and part of the arms) being tracked were: 6, 26, and 16, respectively. For the method (iii), only the upper body parts were tracked. 6 parts were used. Since the focus was on hand tracking, the results of two hands interaction for all the activities are shown as in TABLE III.

The tracking performance of these algorithms was evaluated by employing three metrics: false merging, false labeling and tracking accuracy. The “false merging” is defined as the situation where the tracker of one hand occupies 80% of the area of the other hand. The “false labeling” is defined as the situation where the trackers of both hands change positions during/after interaction or occlusion. The “tracking accuracy” is defined by (11):

$$\text{Tracking Accuracy} = \frac{\text{total number of (true positives + true negatives)}}{\text{total number of tracked frames}} \quad (11)$$

where a true positive is defined as the situation whereas a target object is present and the tracker was able to find it. True negatives are instances where the target object is not present and the tracker also agreed that the object was absent [47]. TABLE III shows that the proposed algorithm (CPMC) exhibits the best performance for the interaction and occlusion conditions among the three particle based methods. There is a marginal decrease in the algorithm speed. When there is no interaction or occlusion occurring, CPMC has the same speed as MCMC and MI approaches. Comparing to the skeleton tracking [24, 48] and body part tracking method proposed by [25], the proposed method obtained higher tracking accuracy. Since our targeting user group is individuals with upper mobility impairments, two challenges exist for hand tracking in the proposed system that could not be tackled very well by skeleton based or other articulated pose tracking method. One challenge is that the users with upper extremity mobility impairments need to sit most of the time on a wheelchair and

perform limited space hand gestures. Their hands could be very close to the body when they perform the gestures. The tracking method based purely on depth information could easily lose track when the hands are so close to the torso [24, 25]. The second challenge is that most of the wheelchairs have armrest, which can be easily confused with human arms and hands. This can explain why the skeleton based tracking method ((iii) and (v)), and the articulated pose tracking method (iv) does not work well for our dataset.

The values of the performance metrics “false merging” and “tracking accuracy” versus the distance between left and right hands are shown in Fig. 8, and Fig. 9. From Fig. 8, we can see that the proposed approach outperforms method (i), (ii), (iv), and method (iii) (with color images). Additionally, the result of the proposed approach (1 false merging frame) were very close to the results of method (iii) (with depth images) and method (v) (no false merging occurred). In Fig. 9, the total number of false positive and false negative frames versus the hands’ interdistance was presented for each method. This figure showed that the proposed approach outperformed all the other state of art algorithms, since it displayed the fewest number of false positive and false negative frames among all the algorithms for nearly all distances.

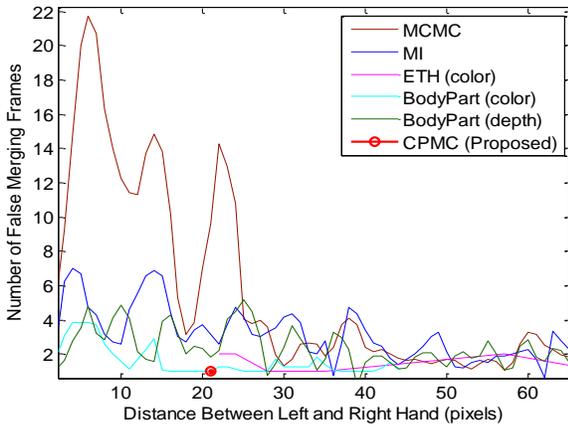


Fig. 8. Number of “false merging” occurred vs. hand distance.

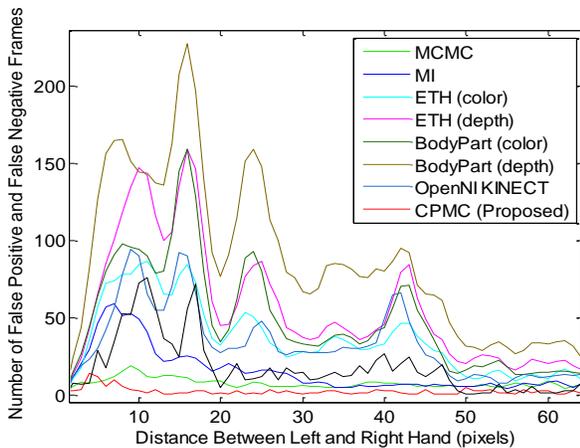


Fig. 9. Number of false positive and false negative vs. hand distance.

### B. Experiment 2: Gesture Recognition Performance

The motion models were constructed using the DTW

algorithm. The velocities (3D directions) of right and left hands were used as the main feature components. The gesture lexicon in Fig. 6 was adopted, and those gestures were used to create spatio-temporal trajectories that later were classified by the gesture-based recognition system.

The system was validated by eight able-bodied subjects and two subjects with quadriplegia due to cervical spinal cord injuries aged around 24-40. The ten subjects performed all the gestures in the lexicon each ten times (8 gestures x 10 subjects x 10 repetitions). Ten sessions were used for cross validation for each gesture (k-fold with k=10). In each session, 720 observations (8 gestures x 9 subjects x 10 repetitions) were used for training and 80 gestures (8 gestures x 1 subject x 10 repetitions) were used for testing. This cross validation resulted in an average accuracy of 95.9%. A confusion matrix was computed and shown by Fig. 10 (with a temporal window size of  $w=19$ ). Confusions occurred when the subjects performed a gesture mistakenly in a single direction or not enough motion was exhibited as expected in other directions. Other cases of misclassification occurred when two gestures shared similar sub-trajectories (i.e., counter clock and S gestures).

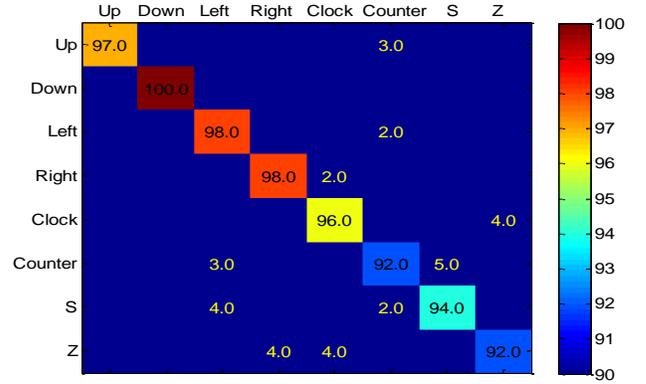


Fig. 10. Confusion matrix with window size of  $w=19$ .

The recognition performance for the CONDENSATION algorithm with our training procedures (CONDENSE) was compared to four other existing state-of-the-art recognition algorithms: (i) Basic motion [50]; (ii) Motion based PCA [51]; (iii) Dynamic time warping (DTW) [52], and (iv) Hidden Markov Model (HMM) [28]. After applying each gesture recognition method to our data set, the results shown in TABLE IV were obtained. The confusion matrices for the different methods are shown in Fig. 11, Fig. 12, Fig. 13 and Fig. 14 respectively. Method (i), (ii), and (iii) used motion information to recognize hand gestures, while (iv) and the CONDENSATION method recognized hand gestures by extracting and classifying hand trajectories. The comparison results demonstrate a high recognition accuracy for the trajectories classification based method. HMM based recognition method can get comparable results as the method used in our paper.

TABLE IV  
GESTURE RECOGNITION PERFORMANCE

Method	Basic	PCA	DTW	HMM	CONDENSE
Accuracy (%)	54.6	66.0	67.4	94.2	95.9

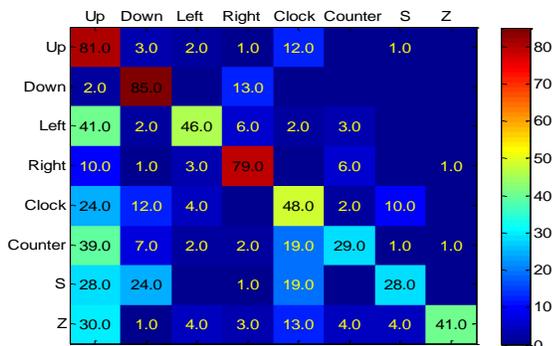


Fig. 11. Confusion matrix for Basic motion method.

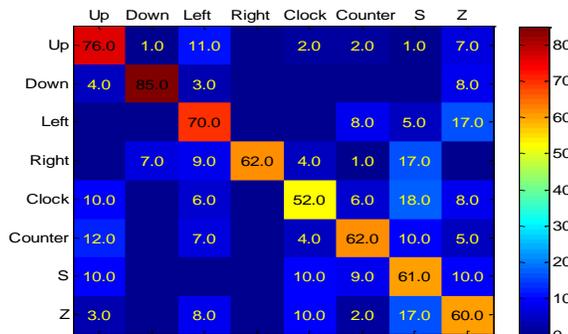


Fig. 12. Confusion matrix for PCA method (with 12 principal components).

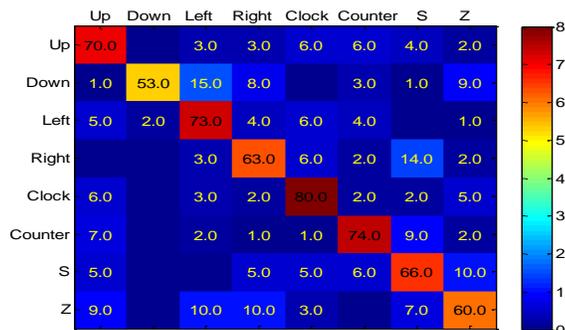


Fig. 13. Confusion matrix for DTW method.

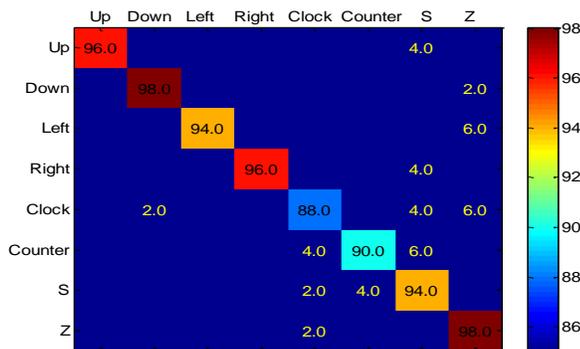


Fig. 14. Confusion matrix for HMM method.

### C. Experiment 3: One Shot Learning Performance

One instance (one repetition by a subject) for each gesture in the lexicon was used for training and remaining observations were used for testing. Ten sessions were used for cross validation for each gesture (k-fold with k=10). In each session, 10 observation (10 gestures x 1 subjects x 1 repetitions) was used for training and 740 gestures (8 gestures

x 9 subject x 10 repetitions and 2 gestures x 1 subjects x 10 repetitions) were used for testing. This cross validation resulted in an average accuracy of 82.78%. A confusion matrix was computed and is shown by Fig. 15 (with a temporal window size of w=19). The recognition accuracy found is comparable to those reported in the ChaLearn Competition [33] in 2012 (fourth place in the competition).

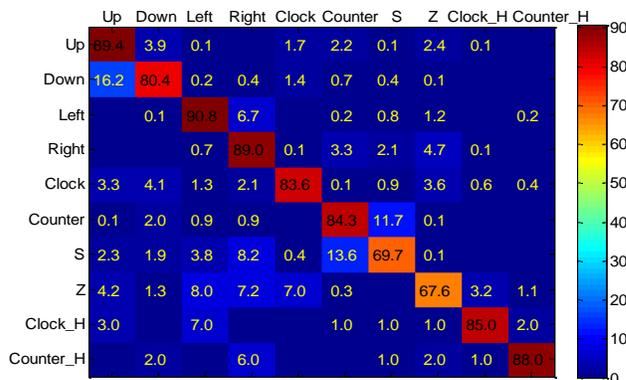


Fig. 15. Confusion matrix for one shot learning (window size w=19) .

### D. Experiment 4: Robotic Control Performance

Since the first remotely driven robotic arm developed by Goldberg et al. (1995) [53] for gardening tasks (Telegarden), there has been an extensive wave of remote labs enabling users the ability to perform lab experiments without the need to physically attend them. Some examples include a work-cell with a 6 axes robot for conducting experiments remotely [54]; a LEGO mobile robotic platform for experimenting with autonomous navigation [55]; and a remote laboratory on robotics was developed at University of Siena called TeleTab [56].

A chemistry laboratory based experiment was performed by five subjects including two individuals with quadriplegia due to a cervical spinal cord injury and three able-bodied individuals. In the laboratory case study experiment, a mobile robot was controlled by the gesture algorithm to transport a beaker to a position near a robotic arm. The robotic arm was activated by the operator to add a reagent to the beaker and then, the mobile robot was brought back to its original position. The gestures (a)-(h) (from the lexicon in Fig. 6) were used and mapped to the commands: ‘change mode’, ‘robotic arm action’, ‘go forward’, ‘go backward’, ‘turn left’, ‘turn right’, ‘stop’ and ‘enable robotic arm’. The two robots were controlled by three modes; discrete, continuous and hybrid mode (discrete plus continuous mode). In discrete mode, for each issued command, the mobile robot moved a fixed increment of distance. While in continuous mode, the mobile robot responded to a given command, until the ‘stop’ command was issued. To switch between the discrete and continuous mode one distinctive gesture (‘upward’) was used. In the experiment, the discrete, continuous and hybrid (continuous plus discrete) control modes were each tested five times by all subjects. The resulting average task completion times were 241.8, 134.7 and 169.6 seconds, for the discrete, continuous and hybrid mode, respectively (Fig. 16). From the results, the completion time of discrete mode took longer time

than the continuous. Continuous and hybrid modes require commands to be issued only when the robot needs to change directions or stop, therefore fewer operations were required for continuous and hybrid modes than for discrete mode for the task observed.

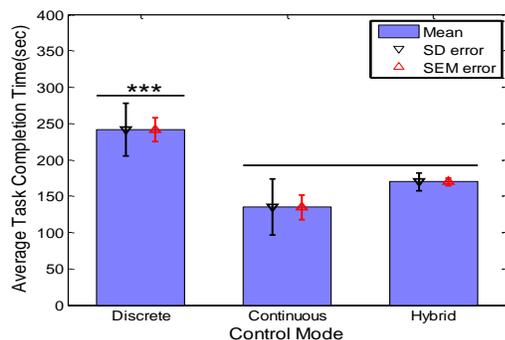


Fig. 16. Average task completion time, unpaired t-test,  $p < 0.001$ .

## V. DISCUSSION & CONCLUSION

A machine vision based gestural interface was developed for individuals with upper extremity physical impairments. Since skin and non-skin color histogram models were used to initialize the face and hands' centroid, the performance of the system may be affected when the users wear short sleeves. In addition, it is expected that users will be seated within the working distance to range specified by the Kinect sensor. An interaction model was incorporated into the color-based particle filter framework for hand tracking. When there was no interaction between the face and hands, multiple independent particle filters tracked the users' movements. When interaction was present, the multiple independent particle filter trackers were combined with an interaction model to solve "false merging" and "false labeling" problems. A comparison between our proposed approach (CPMC), and five state of the art algorithms demonstrated that our approach can achieve robust performance for hand tracking through interaction and occlusion conditions. The proposed tracking strategy can obtain significantly better performance than the other three methods for both "false merge" and "false labeling" problems in hand tracking through interaction and occlusion. Yet, improvements are still need for "false labeling" solving. A training procedure was proposed to obtain motion models for each gesture in the lexicon. The CONDENSATION algorithm with the proposed training procedure was used and compared with 4 other recognition algorithms to classify bimanual gestures. Results showed that HMM based recognition methods may deliver comparable results to our method. Thus, higher recognition could be achieved by using trajectories classification based method. The gesture recognition algorithm designed was found to reach a recognition accuracy of 95.8%. One shot learning was applied in this paper to customize gestures and reduce the number of repetitions required to train/teach the system to a minimum (one observation). The results obtained were comparable to the state of art one shot gesture recognition algorithms presented in the ChaLearn Challenge [33].

A laboratory task experiment was conducted, a typical

biomedical lab procedure with the help of two robots, which were controlled through a gestural interface. Subjects with upper extremity physical impairments can successfully use the machine vision based gestural interface to control the two robots. It was found that the proposed gestural interface was robust enough to support the completion of this task for subjects with upper extremity mobility impairments. In addition, three modes of operation were compared: discrete, continuous and hybrid. Results showed that the continuous mode required the least average task completion time, while the discrete control mode requires the most. Therefore, the authors recommend to use continuous control mode in general, and to use discrete mode only when the robot is very near to the target, for precise location and manipulation.

## VI. FUTURE WORK

Future work for this paper may include: (1) develop more effective and robust algorithms to solve "false merge" and "false labeling" problems of hand tracking through interaction and occlusion. (2) extend the laboratory task to increase the pool of participating users. Ideally, users with physical impairments can participate and provide feedback about the usability, learning and adaptability to the interface suggested.

## APPENDIX

The video sequences of two activities for hand tracking through interaction are shown as in Fig. 17, and 18.

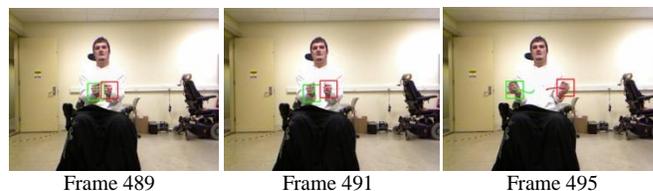


Fig. 17. Hand tracking sequence for "clapping hands" activity.



Fig. 18. Hand tracking sequence for "moving one hand up and down" activity.

## ACKNOWLEDGMENT

This work is partially funded by the National Institutes of Health through the NIH Director's Pathfinder Award to Promote Diversity in the Scientific Workforce, grant number DP4-GM096842-01.

## REFERENCES

- [1] Ahsan, M. R., EMG Signal Classification for Human Computer Interaction: A Review. *European Journal of Scientific Research*, pp. 480-501 (2009).
- [2] Jacko, J. A., Human-Computer Interaction Design and Development Approaches. In: *14th HCI International Conference*, pp. 169-180 (2011).
- [3] Moon, I. H., Lee, M., Ryu, J. C., and Mun, M., Intelligent robotic wheelchair with EMG-, gesture-, and voice-based interface. In: *Intelligent Robots and Systems*, pp. 3453-3458 (2003).
- [4] Walters, M., Marcos, S., Syrdal, D. S., and Dautenhahn, K., An Interactive Game with a Robot: People's Perceptions of Robot Faces and a Gesture-Based User Interface. In: *6th international conference on advances in computer-human interactions*, pp. 123-128 (2013).
- [5] Brdiczka, O., Langet, M., Maisonnasse, J., Crowley, J. L., Detection Human Behavior Models from Multimodal Observation in a Smart Home. *IEEE Transactions on Automation Science and Engineering*, pp. 588-597 (2009).
- [6] Cook, M.C. and Miller Polgar, J. (Eds.) (2008) *Cook & Hussey's Assistive Technologies: Principles and Practice 3<sup>rd</sup> Edition*, Missouri, Mosby Elsevier.
- [7] Murthy, G. R. S. And Jadon, R. S., A Review of Vision based Hand Gesture Recognition. *Internation Journal of Information Technology and Knowledge Management*, vol 2, No. 2, pp. 405-410 (2009).
- [8] Debuse, D., Gibb, C., and Chandler, C., Effects of hippotherapy on people with cerebral palsy from the users' perspective: A qualitative study, vol. 25, No. 3, pp. 174-192 (2009).
- [9] Sterba, J. A., Rogers, B. T., France, A. P., and Vokes, D. A., Horseback riding in children with cerebral palsy: effect on gross motor function. *Developmental Medicine & Child Neurology*, vol 44, No. 5, pp. 301-308 (2002).
- [10] Kitto, K. L., Development of a low-cost sip and puff mouse. In: *Proceedings of 16th Annual Conference of RESNA*. pp. 452-454 (1993).
- [11] Yin, Y. H., Fan, Y. J., and Xu, L. D., EMG and EPP-Integrated Human-Machine Interface Between the Paralyzed and Rehabilitation Exoskeleton. *IEEE Transactions on Information Technology in Biomedicine*, vol 16, No. 4, pp. 542-549 (2012).
- [12] Jiang, H., Wachs, J. P., Duerstock, B. S., Facilitated gesture recognition based interfaces for people with upper extremity physical impairments. In: *Proceedings in Pattern Recognition, Image Analysis, Computer Vision, and Applications. Lecture Notes in Computer Science*, pp. 228-235 (2012).
- [13] Wachs, J., Kölsch, M., Stern, H. and Edan, Y., Vision-Based Hand Gesture Applications: Challenges and Innovations. *Communication of the ACM, Cover Article*, vol 54, No. 2, pp: 60-71 (2011).
- [14] Li, Z., and Jarvis, R., A multi-modal gesture recognition system in a Human-robot interaction scenario. In: *Proceedings of the IEEE International workshop on robotic and sensors environments*, pp. 41-46 (2009).
- [15] Suma, E. A., Lange, B., Rizzo, A., Krum, DM., and Bolas, M., FFAST: The Flexible Action and Articulated Skeleton Toolkit. *IEEE Virtual Reality*, pp. 247-248 (2011).
- [16] Leap Motion: <https://www.leapmotion.com/>.
- [17] Bradski, G.R., Computer vision face tracking as a component of a perceptual user interface. In: *Workshop on applications of computer vision*, pp. 214-219 (1998).
- [18] Isard, M., and Black, A., CONDENSATION: Conditional density propagation for visual tracking. *J. International Journal of Computer Vision*, Vol. 29, pp. 5-28 (1998)
- [19] Maskell, S., Gordon, N., and Clapp, T., A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE transaction on signal processing*, pp. 174-188 (2002).
- [20] Perez, P., Hue, C., Vermaak, J., and Gangnet, M., Color-based probabilistic tracking. *LNCS, Springer, Heidelberg*, vol. 2350, pp. 661-675 (2002).
- [21] Okuma, K., Taleghani, A., Freitas, N., Little, J. J., and Lowe, D. G., A boosted particle filter: multitarget detection and tracking. *European Conference on Computer Vision* (2004).
- [22] Kristan, M., Pers, J., Kovacic, S., and Leonardis, A., A local-motion-based probabilistic model for visual tracking. *Pattern Recognition*. vol. 42, No. 9, pp. 2160-2168 (2009).
- [23] Qu, W., Schonfeld, D., and Mohamed, M., Real-time distributed multi-object tracking using multiple interactive trackers and a magnetic-inertia potential model. *IEEE transactions on Multimedia*, vol. 9, No. 3, pp. 511-519 (2007).
- [24] Eichner, M., Marin-Jimenez, M., Zisserman, A., and Ferrari, V., Articulated Human Pose Estimation and Search in (Almost) Unconstrained Still. *ETH Technical Report* (2010).
- [25] Yang, Y., and Ramanan, D., Articulated Pose Estimation with Flexible Mixture of Parts. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1385-1392 (2011).
- [26] Shotton, J., Fitzgibbon, A. W., Cook, M., Sharp, T., Finocchio, M., Moore, R., Kipman, A., and Blake, A., Real-time human pose recognition in parts from single depth images. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1297-1304 (2011).
- [27] Oikonomidis, I., Kyriazis, N., and Argyros, A. A., Tracking the articulated motion of two strongly interacting hands. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1862-1869 (2012).
- [28] Bilal, S., Akmeliawati, R., Shafie, A. A., and Salami, M. J. E.: Hidden Markov Model for human to computer interaction: a study on human hand gesture recognition. *Artificial Intelligence* (2011).
- [29] Miners, B. W., Basir, O. A., and Kamel, M. S., Understanding hand gestures using approximate graph matching. *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol. 35, no. 2, pp. 239-248 (2005).
- [30] Xu, Z., Xiang, C., Li, Y., Lantz, V., Wang, K., and Yang, J., A Framework for Hand Gesture Recognition Based on Accelerometer and EMG Sensors, *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol.41, No.6, pp.1064-1076 (2011).
- [31] Black, M. J., and Jepson, A. D., A probabilistic framework for matching temporal trajectories: CONDENSATION-based recognition of gesture and expressions. *ECCV*. 1998.
- [32] Alon, J., Athitsos, V., and Yuan, W., A Unified Framework for gesture Recognition and Spatiotemporal Gesture Segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, pp. 1685-1699 (2009).
- [33] Guyon, I., Athitsos, V., Jangyodsuk, P., Escalante, H. J., and Hamner,

- B., Results and Analysis of the ChaLearn Gesture Challenge (2012).
- [34] Fei-Fei, L., Fergus, R., and Perona, P., One-shot learning of object categories, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, No. 4, pp. 594-611 (2006).
- [35] Wu, D., Zhu, F., and Shao, L., One shot learning gesture recognition from rgbd images. In: *CVPR2012 workshop on gesture recognition* (2012).
- [36] Yang, Y., Saleemi, I., and Shah, M., Discovering Motion Primitives for Unsupervised Grouping and One-shot Learning of Human Actions, Gestures, and Expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2012).
- [37] Guyon, I., Athitsos, V., Jangyodsuk, P., Hamner, B., and Escalante, H. J., ChaLearn gesture challenge: Design and first results, *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1-6 (2012).
- [38] Wachs, J., Stern, H., and Edan, Y., Cluster Labeling and Parameter Estimation for the Automated Setup of a Hand-Gesture Recognition System, *IEEE Transactions on Systems, Man and Cybernetics*, vol. 35, no. 6, pp. 932-944 (2005).
- [39] Jones, M. J., and Rehg, J. M., Statistical color models with application to skin detection, In: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, pp. 81-96 (1999).
- [40] Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. In: *International Conference on Computer Vision and Pattern Recognition*, pp. 511-518 (2001).
- [41] Hess, R., Fern, A.: Discriminatively Trained Particle Filters for Complex Multi-Object Tracking. In: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 240-247 (2009).
- [42] Yu, T. and Wu, Y.. Collaborative tracking of multiple targets. *Computer society conference on computer vision and pattern recognition (CVPR)*. 2004.
- [43] Khan, Z., Balch, T. and Dellaert, F., An MCMC-based particle filter for tracking multiple interacting targets. Technical report. 2003.
- [44] Borg, G., Psychophysical bases of perceived exertion. *Medicine & Science in Sports & Exercise*, vol. 14(5), 377-383 (1982).
- [45] Jiang, H., Duerstock, B. S., and Wachs, J. P. Integrated gesture recognition based interface for people with upper extremity mobility impairments, In: *Proc. of the 4th Internaiton Conference on Applied Human Factors and Ergonomics* (2012).
- [46] Aach, J., Church, G.M.: Alignment gene expression time series with time warping algorithms, *J: Bioinformatics*, vol. 17, no. 6, pp. 495-508, Oxford University Press (2001).
- [47] Savitzky, A., and Golay, M. J. E., Smoothing and diferentiation of data by simplified leart squares procedures, *Analytical Chemistry*, vol. 36, no. 8, pp. 1627-1639 (1964).
- [48] OpenNI: <http://www.openni.org/>
- [49] Pingali, G., and Segen, J., Performance evaluation of people tracking systems, In: *Proc. IEEE Workshop on Application of Computer vision*, pp. 33-38 (1996).
- [50] Guyon, I., Athitsos, V., Jangyodsuk, P., and Escalante, H. J., The ChaLearn Gesture Dataset (2011).
- [51] Escalante, H. J., and Guyon, I., Principal motion: Principal motion: Pca-based reconstruction of motion histograms. Technical report, ChaLearn Technical Memorandum (2012).
- [52] Sohn, M. K., Lee, S. H., Kim, D. J., Kim, B., and Kim, H., A comparison of 3D hand gesture recognition using dynamic time warping. In: *Proceedings of the 27th Conference on Image and Vison Computing*, pp. 418-422 (2012).
- [53] Goldberg, K. *The Robot in the Garden: Telerobotics and Telepistemology in the Age of the Internet*. Mit Press. 2000.
- [54] Safaric, M. Truntic, D. Hercog, and G. Pacnik. Con-trol and robotics remote laboratory for engineering ed-ucation. *International Journal of Online Engineering (iJOE)*, 1(1), 2005.
- [55] Carusi, F., Casini, M., Prattichizzo, D., and Vicino, A. Distance learning in robotics and automation by remote control of LEGO mobile robots. In *Proc. Int. Conf. on Robotics and Automation*, pages 1820–1825, New Orleans, USA, Aprile 2004.
- [56] Casini, M., Chinello, F., Prattichizzo, D., & Vicino, A. (2008, July). RACT: A remote lab for robotics experiments. In *Proceedings of the 17th IFAC World Congress*. Seoul (Korea).

## BIOGRAPHIES



**Hairong Jiang** received her B.S. degree in control science and engineering from Harbin Institute of Technology, in 2008. Her M.S. degree in control science and engineering from Harbin Institute of Technology Shenzhen graduate school, in 2010. She is currently working toward her Ph.D. degree at the School of Industrial Engineering, Purdue University, USA. Her primary research interests include gesture recognition and assistive technology.



**Bradley S. Duerstock** received a B.S. degree in Biomedical Engineering at the School of Interdisciplinary Engineering in 1994 and Ph.D. degree in Neurobiology at the College of Veterinary Medicine in 1999 from Purdue University, West Lafayette, IN, USA. He was a postdoctoral research associate at the Center for Paralysis Research at Purdue University. He is an associate professor of Engineering Practice in the Weldon School of Biomedical Engineering and School of Industrial Engineering at Purdue University. He is the Director of Institute for Accessible Science. His research interests focus on the restoration of functional impairment through repair of central nervous system damage or development of assistive technologies and accessible design.



**Juan P. Wachs** is an Assistant Professor in the School of Industrial Engineering at Purdue University. He is the director of the Intelligent Systems and Assistive Technologies Lab (ISAT) and he is affiliated with the Regenstrief Center for Healthcare Engineering. He completed a postdoctoral training at the Naval Postgraduate School's MOVES Institute in the area of computer vision, under a National Research Council Fellowship from the National Academics of Sciences and he was awarded the Air Force Young Investigator Award 2013. His research interests include machine and computer vision, robotics, teleoperations, human robot interaction, and assistive technologies. Juan Wachs is a member of IEEE and the Operation Research Society of Israel (ORSIS). He has published in journals including *IEEE Trans. Systems, Man, and Cybernetics*, *Journal of American Medical Informatics*, *Communications of the ACM*, and the *Journal of Robotic Surgery*. He received his M. Sc. and Ph.D. in Industrial Engineering and Management from the Ben-Gurion University of the Negev.