Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

Hand Gesture Interaction with Human-Computer Dr. Dejan Chandra Gope¹ ¹ Dhaka University of Engineering and Technology *Received: 3 November 2011 Accepted: 3 December 2011 Published: 18 December 2011*

6 Abstract

Hand gestures are an important modality for human computer interaction. Compared to 7 many existing interfaces, hand gestures have the advantages of being easy to use, natural, and 8 intuitive. Successful applications of hand gesture recognition include computer games control, 9 human-robot interaction, and sign language recognition, to name a few. Vision-based 10 recognition systems can give computers the capability of understanding and responding to 11 hand gestures. The paper gives an overview of the field of hand gesture interaction with 12 Human- Computer, and describes the early stages of a project about gestural command sets, 13 an issue that has often been neglected. Currently we have built a first prototype for exploring 14 the use of pieand marking menus in gesture-based interaction. The purpose is to study if such 15 menus, with practice, could support the development of autonomous gestural command sets. 16 The scenario is remote control of home appliances, such as TV sets and DVD players, which in 17 the future could be extended to the more general scenario of ubiquitous computing in 18 everyday situations. Some early observations are reported, mainly concerning problems with 19 user fatigue and precision of gestures. Future work is discussed, such as introducing flow 20 menus for reducing fatigue, and control menus for continuous control functions. The computer 21 vision algorithms will also have to be developed further. 22

23

Index terms — Human Computer Interaction, Hand Tracking, Hand gesture, Computer Vision Based Gesture
 Recognition, HCI, Gesture Command, Marking Menu.

²⁶ 1 INTRODUCTION

27 ision-based hand gesture recognition is an active area of research in human-computer interaction (HCI), as direct use of hands is a natural means for humans to communicate with each other and more recently, with devices in 28 intelligent environments. The trend in HCI is moving towards real-time hand gesture recognition and tracking 29 for use in interacting with video games [1], remote-less control of television sets, and interacting with other 30 similar environments. Given the ubiquity of mobile devices such as smartphones and notebooks with embedded 31 cameras, a hand gesture recognition system can serve as an important way of using these camera-enabled devices 32 to interact more intuitively than traditional interfaces. The trend towards embedded, ubiquitous computing in 33 domestic environments creates a need for human-computer interaction forms that are experienced as natural, 34 35 convenient, and efficient. The traditional desktop paradigm, building on a structured office work situation, 36 and the use of keyboard, mouse and display, is no longer appropriate. Instead, natural actions in human-37 tohuman communication, such as speak and gesture, seem more appropriate for what Abowd and Mynatt [1] 38 have named everyday computing, and which should support the informal and unstructured activities of everyday life. Interaction in these situations implies that it should not be necessary to carry any equipment or to be in a 39 specific location, e.g., at a desk in front of a screen. Interfaces based on computational perception and computer 40 vision should be appropriate for accomplishing the goals of ubiquitous, everyday computing. This paper presents 41 an overview of the field of gesture-based interfaces in human-computer interaction as a background, and the first 42 stages of a project concerning the development of such interfaces. Specifically, in the project we intend to study 43

the use of hand gestures for interaction, in an approach based on computer vision. As a starting point, remote 44 control of electronic appliances in a home environment, such as TV sets and DVD players, was chosen. This is 45 an existing, common interaction situation, familiar to most. Normally it requires the use of a number of devices, 46 47 which can be a nuisance, and there are clear benefits to an appliance-free approach. In the future the application could easily be extended to a more general scenario of ubiquitous computing in everyday situations. Currently 48 we have implemented a first prototype for exploring the use of pie-and marking menus [9], [20] for gesturebased 49 interaction. Our main purpose is not menu-based interaction, but to study if such menus, with practice, could 50 support the development of an autonomous gestural command sets. The application will be described in more 51 detail later in this paper. 52

⁵³ 2 II.

54 3 RELATED WORK

Hand gesture recognition and tracking has been an important and active area of research in the field of HCI, and sign language recognition. The use of glove-based devices to measure hand location and shape, especially for virtual reality, has been actively studied. In spite of achieving high accuracy and speed in measuring hand postures, this approach is not suitable for certain applications due to the restricted hand motion caused by the attached cables.

Computer vision techniques measure hand postures and locations from a distance, providing for unrestricted 60 movement. Numerous approaches have been explored by the vision community to extract human skin regions 61 either by background subtraction or skin-complex backgrounds or real-world scenarios where the user wants 62 to use the application on-the-go. Once the image regions are identified by the system, the image regions can 63 be analyzed to estimate the hand posture. Specifically, for finger gesture recognition and tracking, a common 64 approach is to extract hand regions and then locate the fingertip to determine the pose orientation. In a 3D 65 pointing interface using image processing is presented to estimate the pose of a pointing finger gesture. This 66 67 system however, suffers from various drawbacks in real-world scenarios due to the use of a fixed threshold for 68 image binarization and the use of predetermined finger length and thickness values. Also, low-cost web cameras and infrared cameras have been used for finger detection and tracking. In finger detection is performed by fitting 69 a cone to rounded features, and in a template matching approach is used to recognize a small set of gestures. 70

71 4 III. HAND GESTURES FOR COMPUTER VISION

72 Gestures are expressive, meaningful body motions with the intent to convey information or interact with the 73 environment [36]. According to Cadoz [8] hand gestures serve three functional roles, semiotic, ergotic, and epistemic. The semiotic function is to communicate information, the ergotic function corresponds to the capacity 74 75 to manipulate objects in the real world, and the epistemic function allows us to learn from the environment through tactile experience. Based on this classification Quek [30] distinguishes communicative gestures, which 76 are meant for visual interpretation and where no hidden part carries information critical to understanding, from 77 manipulative gestures, which show no such constraints. Thus, it may be more appropriate to use special tools 78 for interaction, like data gloves, rather than computer vision if the intent is realistic manipulation of objects in, 79 e.g., a virtual environment. Pavlovic et al. [28] makes a similar classification, but also point out the distinction 80 81 between unintentional movements and gestures. 82 For communicative, semiotic gestures, Kendon [14] distinguishes gesticulation, gestures that accompany speech, from autonomous gestures. These can be of four different kinds: language-like gestures, pantomimes, emblems, 83

and sign languages. When moving forward in this list the association with speech diminishes, language properties
increase, spontaneity decreases and social regulation increases. Detailed descriptions and taxonomies concerning
hand gestures from the point of view of computer vision can be found in Quek [30], Pavlovic & Sharma [28] and
Turk [36].

88 Here only a brief overview will be presented.

Most work in computer vision and HCI has focused on emblems and signs because they carry more clear 89 semantic meaning, and may be more appropriate for command and control interaction [37]. It is important to 90 note, however, that they are largely symbolic, arbitrary in nature, and that universally understandable gestures 91 92 of this kind hardly exist. There is also one important exception worth mentioning. In the gesticulation category, 93 McNeill [24] defines deictic gestures as pointing gestures that refer to people, objects, or events in space and time. 94 Deictic gestures are potentially useful for all kinds of selections in humancomputer interaction, as illustrated, 95 e.g., by the early work of Bolt [4]. The deictic category itself can be further subdivided, but from a computer vision point of view all deictic gestures are performed as pointing, and the difference lies in the higher level of 96 interpretation [30]. 97

⁹⁸ In the following we limit ourselves to intentional, semiotic, hand gestures. From a computer vision point of ⁹⁹ view, we focus on the recognition of static postures and gestures involving movements of fingers, hands and arm ¹⁰⁰ with the intent to convey information to the environment.

101 5 IV. PERCEPTIVE AND MULTIMODAL USER INTER-102 FACES

The aim is to develop conversational interfaces, based on what is considered to be natural human-tohuman 103 dialog. For example, Bolt [4] suggested that in order to realize conversational computer interfaces, gesture 104 recognition will have to pick up on unintended gestures, and interpret fidgeting and other body language signs, 105 and Wexelblatt [41] argued that only the use of natural hand gestures is motivated, and that there might even be 106 added cognitive load on the user by using gestures in any other way. Two main scenarios for gestural interfaces 107 can be distinguished. One aims at developing Perceptive User Interfaces (PUI), as described by Turk [36], or 108 Perceptive Spaces, e.g., Wren [42], striving for automatic recognition of natural, human gestures integrated with 109 other human expressions, such as body movements, gaze, facial expression, and speech. 110

However, in this paper the focus is on using hand gestures given purposefully as instructions, and we restrict 111 our work to deliberate, expressive movements. This falls within the second approach to gestural interfaces, 112 113 Multimodal User Interfaces, where hand poses and specific gestures are used as commands in a command language. 114 The gestures need not be natural gestures but could be developed for the situation, or based on a standard sign 115 language. In this approach, gestures are either a replacement for other interaction tools, such as remote controls and mice, or a complement, e.g., gestures used with speech and gaze input in a multimodal interface. Oviatt et 116 al. [27] noted that there is a growing interest in designing multimodal interfaces that incorporate vision-based 117 technologies. They also contrast the passive mode of PUI with the active input mode, addressed here, and claim 118 that although passive modes may be less obtrusive, active modes generally are more reliable indicators of user 119 intent, and not as prone to error. 120

121 6 V. GESTURE-BASED APPLICATIONS IN HCI

In traditional HCI, most attempts have used some device, such as an instrumented glove, for incorporating gestures into the interface. If the goal is natural interaction in everyday situations this might not be acceptable. However, a number of applications of hand gesture recognition for HCI exist, using the untethered, unencumbered approach of computer vision. Mostly they require restricted backgrounds and camera positions, and a small set of gestures, performed with one hand. They can be classified as applications for pointing, presenting, digital desktops, and virtual workbenches and VR.

Pavlovic [28] noted that, ideally, naturalness of the interface requires that any and every gesture performed by the user should be interpretable, but that the state of the art in vision-based gesture recognition is far from providing a satisfactory solution to this problem. A major reason obviously is the complexity associated with the analysis and recognition of gestures. A number of pragmatic solutions to gesture input in HCI exist, however, such as:

133 ? use props or input devices (e.g., pen, or data glove) ? restrict the object information (e.g., silhouette of 134 the hand)? restrict the recognition situation (uniform background, restricted area)? restrict the set of gestures Pointing: A number of applications that use computer vision for pointing (deictic) gestures have been developed, 135 either in a scenario for some special kind of interaction situation, such as Put-That-There [4], or, as a replacement 136 for some input device in general, mostly the mouse. An example is Finger Mouse [31], where a down-looking 137 camera was used to create a virtual 2D mousepad above the keyboard, allowing users to perform pointing gestures 138 to control the cursor. Mouse clicks were implemented by pressing the shift key. Kjeldsen and Kender [16] used 139 a camera position below the screen, facing the user, to compute the x,y coordinates that control the cursor. 140 For window control they used a neural network to classify hand poses (point, grasp, move, menu) with a simple 141 grammar, based on pausing and retraction. They note that users had difficulties to remember the sequence 142 143 of motions and poses and that there were unexpected interface actions, because gestures were dependent on timing. O'Hagan [25] used a commercial system with a single video camera for Finger Track, which performed 144 visionbased finger tracking on top of the workspace. A pointing gesture (one finger) and a click gesture (twofingers 145 extended) could be used. A similar application, FingerMouse for controlling the mouse pointer was presented by 146 von Hardenberg and Berard [39]. The finger, moving over a virtual touchscreen, is used as mouse and selection 147 is indicated by a one sec delay in the gesture. 148

Presenting: Baudel et al. [2] used a glovebased system for controlling Microsoft PowerPointpresentations. 149 Even if the focus in this paper is on computer vision, their work should be mentioned, because it addresses the 150 question of developing gestural command sets. They suggest that command gestures should be defined according 151 to an articulatory scheme with a tense start position (e.g. all fingers outstretched), a relaxed dynamic phase (e.g. 152 a hand movement to the right) and a tense end position (e.g. all fingers bent). In a similar application, based on 153 computer vision, Lee & Kim [21] use hand movements for controlling presentations. The detection of the hand 154 155 is entirely based on skin color, which requires a controlled background. The gesture-based virtual touchscreen 156 of von Hardenberg et al. [39] included command gestures for slide changes and menu selection, in addition to general pointing gestures (see above). Hand detection relies on a time filtered background subtraction, i.e., it 157 requires a reference image. In a more advanced multimodal scenario, Kettebekov and Sharma [15] performed an 158 observational study to develop a gesture grammar for deictic gestures when presenting a weather map. Digital 159 Desks: A third kind of application aims at developing mixed reality desktops, using free hand pointing and 160 manipulation of digital objects. Kruegers VideoDesk [19] was an early desk-based system in which an overhead 161

camera and a horizontal light was used to provide hand gesture input for interactions, which were then displayed 162 on a monitor at the far end of the desk. The work was built on the early research of the VideoPlace system 163 [18]. Wellner [40] developed DigitalDesk, a more advanced digital desk system, mixing projected and electronic 164 documents on a real desktop, and using an image processing system to determine the position of the users' hands, 165 and to gather information from documents placed on the desk. Similarly, Maggioni and Kämmerer [23] explored 166 pointing gestures in vision-based virtual touchscreens for office applications, public information terminals and 167 medical applications. The detection is based on a skin segmentation step, and the approach requires controlled 168 backgrounds. More recently, Koike et al. [17] developed an augmented desk interface, EnhancedDesk, with 169 computer vision as a key technology. EnhancedDesk uses a projector for presenting information onto a physical 170 desktop, an infrared camera for detecting users arms, hands, and hand poses, and a pan-tilt camera for giving 171 detail. Users can manipulate digital information directly by using their hands and fingers. The system is reported 172 to be able to track fingertip movements in real time under any lighting condition. 173

Virtual workbenches and VR: The distinction between virtual workbenches and digital desktops is not sharp. 174 Here, a workbench is described as primarily intended for navigation and object manipulation in 3D environments. 175 As mentioned earlier, computer vision might not be suitable for these tasks. Glove-based input might be better 176 suited for intricate 3D manipulation tasks, due to the problem of occluded fingers. Recently, however, Utsumi 177 178 and Ohya [38] proposed a multipleviewpoint system for three-dimensional tracking of position, pose and shapes 179 of human hands, as a step towards replacing glove-based input. Also, many gestures for navigation and object 180 manipulation in virtual environments have a deictic component, i.e., are pointing gestures, which simplifies the problem from a computer vision point of view. Segen and Kumar [33] investigated a vision-based system for 3D 181 navigation, object manipulation and visualization. The system used stereo cameras against a plain background 182 and with stable illumination, and has been used for movement control in a 3D virtual environment, for building 183 3D scenes, and for a 2D game. Fatigue is reported as an issue, especially when the system is used for object 184 manipulation. Leibe et al. [22] experimented with 3D terrain navigation, games, and CSCW, using a FakeSpace 185 immersive workbench with infrared illuminators placed next to the camera. IR light is reflected back to the 186 camera by objects placed on the desk. A second IR camera provides a side view of users arms for recovering 187 3D pointing gestures. O'Hagan et al. [26] implemented a virtual, 3D workbench where two cameras were used 188 to provide stereo images of the users' hand. As with Segen [33], the system could be used for object and scene 189 translations, rotations, object resizing, and zoom. By combining feature-based tracking with a model-based 190 system, tracking with cluttered backgrounds and changing illumination is claimed to be possible. O'Hagan et al. 191 also point out user fatigue as a problem in this kind of application. Other examples of 3D object manipulation 192 and navigation can be found in Sato et al. [32] and Bretzner and Lindeberg [6]. 193

Finally, the work of Wren et al. regarding perceptive rooms and spaces [42] should be mentioned in this context, 194 even if it might rather be characterized as an attempt at mixed reality, multimodality and ubiquitous computing 195 in a PUI scenario. An interactive space is created in a room with constant lighting, controlled background, and a 196 large projection screen. Stereo computer vision is used to track key features of body, hand and head motion. The 197 authors point out that the possibility for users to enter the virtual environment just by stepping into the sensing 198 area is very important, not having to spend time donning equipment. Also, the importance of social context is 199 noted. Not only can the user see and hear a bystander, the bystander can easily take the users place for a few 200 seconds, without any need to "suit up", as is the case with most scenarios requiring equipment. 201

202 **7** VI.

203 8 CURRENT WORK

With the exception of Baudel et al. [2], very little attention has been paid to the selection of gestures in gesturebased interaction, and to the development on gestural command sets. Often the reason is that the gestures are deictic. However, even under circumstances when they are not, there has not been much discussion about what gestures or hand poses should be used.

208 **9** VII.

209 10 GESTURAL COMMAND SETS

The design space for gestural commands can be characterized along three dimensions: Cognitive aspects, 210 Articulatory aspects and Technological aspects. Cognitive aspects refer to how easy commands are to learn 211 212 and to remember. It is often claimed that gestural command sets should be natural and intuitive, e.g. [4] [41], 213 mostly meaning that they should inherently make sense to the user. This might be possible for manipulative 214 gestures, but, as noted above, for communicative gestures there might not exist any shared stereotypes to build 215 on, except in very specific situations. If the aim is gestural control of devices, there is no cultural or other context for most functions. Baudel et al. [2] recommend that ease of learning should be favored and that a compromise 216 must be made between natural gestures that are immediately assimilated by the user and complex gestures that 217 give more control. They define "natural gestures" as those that involve the least effort and differ the least from 218 a rest position, i.e., that "naturalness" in part should be based on an articulatory component, according to the 219 classification used here. Articulatory aspects refer to how easy gestures are to perform, and how tiring they are 220

for the user. Gestures involving complicated hand or finger poses should be avoided, because they are difficult to articulate and might even be impossible to perform for a substantial part of the population. They are common in current computer based approaches, because they are easy to recognize by computer vision. Repetitive gestures that require the arm to be held up and moved without support are also unsuitable from an articulatory point of view because of fatigue.

Technological aspects refer to the fact that in order to be appropriate for practical use, and not only in visionary 226 scenarios and controlled laboratory situations, a command set for gestural interaction based on computer vision 227 must take into account the state-of-the art of technology, now and in the near future. For example, Sign Language 228 recognition might be desirable for a number of reasons, not least for people who need to use Sign Language for 229 communication. Although difficult to learn, once learned a Sign Language is easy to remember because of 230 its language properties, and might provide a good candidate framework for developing gestural languages for 231 interaction. Some attempts to Sign Language recognition also exist. For example, recently Starner et al. [34] 232 developed a recognition system for a subset of American Sign Language. However, Braffort [5] points out that if 233 the real aim is to deal with Sign Language, then all the different varied and complex elements of language must 234 be taken into account. This is currently far from feasible. Still, much work can be done with reduced sets of Sign 235 Language, limited to standard signs, as a first step towards a long-term objective. 236

237 Menu-based Systems for Gesture-Based Interaction: Our current work represents the first stages in a research 238 effort about computer vision based gesture interaction, primarily aimed at questions concerning gesture command 239 sets. The point of departure is cognitive, leaving articulatory aspects aside for the moment, mainly for reasons of technical feasibility. We focus on the fact that the learning curve for a gestural interface of any complexity 240 will be steeper than for a menu-based interface, because commands need to be recalled, rather than recognized. 241 As noted earlier, there are very few natural, generally understandable signs and gestures that could be used. 242 And, however desirable it might be to use some standard Sign Language it is not technically feasible, except at 243 the level of isolated signs. Using signs from Sign Language, if not the language itself, will be addressed in this 244 project in the future. Currently gestures and hand poses are kept simple, for technical reasons and for reasons 245 of articulatory simplicity. 246

As was mentioned above, menu-based systems have the cognitive advantage that commands can be recognized rather than recalled. Traditional menu-based interaction, however, is not attractive in a gesture-based scenario for everyday situations. Menu navigation would be far from the directness that gestural interaction could provide. However, by using pie-and marking menus, it might be possible to support directness, and to provide a solution for developing gestural command sets.

Pie-and Marking Menus: Pie menus were first described by Callahan et al. [9]. They are pop-up menus 252 with the alternatives arranged radially. Because the gesture to select an item is directional, users can learn to 253 make selections without looking at the menu. In principle this could be learned also with linear menus, but it 254 is much easier to move the hand without feedback in a given direction, as with a pie menu, than to a menu 255 item at a given distance, as in a linear menu. This fact can support a smooth transition between novice and 256 expert use. For an expert user, working at high speed, menus need not even be popped up. The direction of 257 the gesture is sufficient to recognize the selection. If the user hesitates at some point in the interaction, the 258 underlying menus could be popped up, always giving the opportunity to get feedback about the current selection. 259 Hierarchic marking menus [20] is a development of pie menus that allow more complex choices by the use of 260 submenus. The same principles apply: expert users could work by gesture alone, without feedback. The shape 261 of the gesture with its movements and turns can be recognized as a selection, instead of the sequence of distinct 262 choices between alternatives. A recent example can be found in Beaudouin-Lafon et al. [3]. Hierarchic Marking 263 Menus for Gesture-Based Interaction: Here the assumption is that command sets for computer vision based 264 gesture interfaces can be created from hierarchical marking menus. As to articulatory characteristics, a certain 265 hand pose, e.g., holding the hand up with all fingers outstretched, could be used for initiating a gesture and 266 activating the menu system. This would correspond to the pen-down event in a pen-based system. The gesture 267 could then be tracked by the computer vision algorithms, as the hand traverses the menu hierarchy. Finally, a 268 certain hand pose could be used to actually make the selection, e.g., the index finger and thumb outstretched, 269 corresponding to a pen-up event in pen-based interface. Put differently, the gestures in the command set would 270 consist of a start pose, a trajectory, defined by menu organization, for each possible selection, and, lastly, a 271 selection pose. Gestures ending in any other way than with the selection pose would be discarded, because either 272 they could mean that the user abandoned the gesture, or simply that tracking of the hand was lost. For a novice 273 user, this would amount to a traditional menu-selection task, where selections are made by navigating through 274 an hierarchical menu structure. This, as such, could provide for unencumbered interaction in remote control 275 situations but, as noted above, the directness of a gesture interface would be lost. The assumption here, however, 276 is that over time users will learn the gesture corresponding to each selection and no longer need visual feedback. 277 278 The interaction would develop into direct communication, using a gestural language. In addition to providing for a natural transition from novice to expert, such a gestural language makes no assumptions about naturalness or 279 semantics of gestures, because it is defined by the menu structure. In principle, if not in practice, the command 280 set is unlimited. A further advantage is that the demands put on the computer vision algorithms are reasonable. 281

Fast and stable tracking of the hand will be required, however.

²⁸³ 11 VIII. A PROTOTYPE FOR HAND GESTURE INTERAC ²⁸⁴ TION

The prototyping and experimental work is still in an early stage and only a brief overview and some early 285 impressions can be given here. Inspired by Freeman et al. [11], [12], we chose remote control of appliances in a 286 domestic environment as our first application. Freeman et al. used only one gesture to control a TV set: an open 287 hand facing the camera. An icon on a computer display followed the users hand, and by moving the icon (hand) 288 along one of two sliders, a user could control the volume or select channels. Our prototype is more intricate 289 and intended to test the hypothesis, discussed above, that hierarchical marking menus can be used to develop 290 gestural command sets. However, so far, we have only designed a first example of a hierarchic menu system for 291 292 controlling some functions of a TV, a CD player, and a lamp. The prototype has been set up in a generally accessible, open lab/demo space at CID (fig. ??). Fig. ?? : The demo space at CID. 293

²⁹⁴ 12 IX. TECHNICAL ASPECTS

The Computer Vision System: We have chosen a view-based representation of the hand, including both color 295 and shape cues. The system tracks and recognizes the hand poses based on a combination of multi-scale color 296 feature detection, view-based hierarchical hand models and particle filtering. The hand poses, or hand states, are 297 298 represented in terms of hierarchies of color image features at different scales, with qualitative inter-relations in terms of scale, position and orientation. These hierarchical models capture the coarse shape of the hand poses. 299 In each image, detection of multi-scale color features is performed. The hand states are then simultaneously 300 detected and tracked using particle filtering, with an extension of layered sampling referred to as hierarchical 301 layered sampling. The particle filtering allows for the evaluation of multiple hypotheses about the hand position, 302 state, orientation and scale, and a likelihood measure determines what hypothesis to chose. To improve the 303 performance of the system, a prior on skin color is included in the particle filtering step. In fig. 3, yellow (white) 304 ellipses show detected multi-scale features in a complex scene and the correctly detected and recognized hand 305 pose is superimposed in red (gray). A detailed description of the algorithms is given in [7]. As the coarse shape 306 of the hand is represented in the feature hierarchy, the system is able to reject other skin colored objects that 307 can be expected in the image (the face, arm, etc). The hierarchical representation can easily be further extended 308 to achieve higher discrimination to complex backgrounds, at the cost of a higher computational complexity. An 309 advantage of the approach is that it is to a large extent user and scale (distance) invariant. To some extent, 310 311 the chosen qualitative feature hierarchy also shows view invariance for rotations out of the image plane (up to 312 approx. 20-30 degrees for the chosen gestures).

313 There is a large number of works on real-time hand pose recognition in the computer vision literature. Some of the most related to our approach are, e.g., Freeman and Weissman [11] (see above) who used normalized 314 correlation of template images of hands for hand pose recognition. Though efficient, this technique can be 315 expected to be more sensitive to different users, deformations of the pose and changes in view, scale, and 316 background. Cui and Weng [10] showed promising results for hand pose recognition using an appearance based 317 method. However, the performance was far from real-time. The approach closest to ours was presented by Triesch 318 and von der Malsburg [35] representing the poses as elastic graphs with local jets of Gabor filters computed at 319 each vertex. 320

Equipment: A Dell Workstation 530 with dual 1,7 GHz Intel Xeon P4 processors running Red Hat Linux was used. The menus were shown on a 19" Trinitron monitor, placed next to the TV screen. The menu system was developed in Smalltalk. An Mvdelta 2 framegrabber, IRdeo remote IR control, and a DI-01 Data interface (X10) was used for image acquisition and to control a table lamp, a Samsung 29" TV, and a Hitachi CD player. In order to maximize speed and accuracy, gesture recognition is currently tuned to work against a uniform background within a limited area, approximately 0,5 by 0,65 m in size, at a distance of approximately 3 m from the camera, and under relatively fixed lighting conditions.

328 : An overview of the functional components and the information flow in the prototype.

329 **13** X.

330 14 MENU SYSTEM

An incomplete version with three hierarchical levels and four choices in each menu currently exists. Only a few of choices are active, however: TV on/off, Previous/Next channel, CD Play/Stop/Back/Forward, Lamp on/off. An example of a menu is shown in fig. 1. An overview of the functional components and the information flow in the prototype is presented in fig. ?? above. We have only recently begun working on the design, the arrangement, and the organization of the menus.

A hand pose with the index finger and thumboutstretched is used as the start pose for activating the menus, corresponding to pen-down in a pen-based interface. A hand with five fingers outstretched is used as the selection pose, corresponding to pen-up. Evidently, any two hand poses could be used for these purposes. Menus are activated when the start hand pose is detected by the computer vision system in the active area. The hand is tracked as long as the start pose is held. If the hand is moved over the periphery of a sector that has a submenu, the parent menu disappears, and the submenu appears. Showing the selection hand pose in an active field, e.g., TV on, makes a selection. All other ways of ending the interaction are ignored. The menus are currently shown on a computer screen, placed by the side of the TV (fig. ??). This is inconvenient, and in the future menus will

be presented in an overlay on the TV screen.

345 XI.

³⁴⁶ 15 RESULTS AND DISCUSSION

Menu-based systems are more complex, and there is simply more to learn at the outset. However, learning the 347 principles for using the menus was not a main issue, and the principles are the same no matter the number of 348 choices in the menu system. There are major drawbacks with using static hand poses for direct control as in 349 the earlier prototype. First, the number of usable poses is limited. Second, many people have difficulties using 350 finger poses. Third, the association of poses to functions is arbitrary, and difficult to remember. There are 351 also culturally specific hand poses (emblems) that have to be avoided. We have not yet been able to bring the 352 technical performance (speed and accuracy) of the menu-based system to a level where true gesturebased control 353 without feedback can be accomplished. However, observations with the current system, as it is, indicate that 354 gesture-based control with simple, singlelevel pie menus is feasible, but that gestures based on hierarchical menus 355 create some problems. It is difficult for users to make the gestures for multiple-level selections sufficiently distinct, 356 based on feedback only from the proprioceptive system of the arm. Thus, computer algorithms for recognition 357 of fuzzy gestures might also be required. Another solution could be to Global Journal of Computer Science and 358 Technology Volume XI Issue XXIII Version I 9 reduce the number of choices at each level. The current setup, 359 with subjects seated facing the TV and making gestures with one arm and hand held out by the side of the body 360 without support, is not suitable from an articulatory point of view. It is inconvenient and fatigue quickly sets in. 361 This is also a consequence of the fact that gestures have to cover a relatively large area if the hierarchy is deep. 362 Also, the gesture might end up outside of the recognition area. The problem of fatigue is known from earlier 363 attempts with gesture-based interfaces and must be addressed. In the current application much could be gained 364 by providing support for the arm, by making gestures smaller, and by making © 2011 Global Journals Inc. (US) 365 Fig. ?? XII. 366

³⁶⁷ 16 FUTURE WORK

As to the computer vision algorithms there is ongoing work to increase the speed and performance of the system, 368 to acquire more position independence for recognition of gestures, to increase the tolerance for varying lighting 369 conditions, and to increase recognition performance with complex backgrounds. The main effort, however, is 370 currently aimed at the design and organization of menus. Recently we have begun development of Flow Menus, a 371 version of hierarchical marking menus in which successive levels of the hierarchy are shown in the same position 372 [13]. In our application this would greatly reduce the area which the gestures have to cover when the hierarchy 373 is deep. An additional problem we faced is that not all kinds of functions, e.g., increasing sound volume, are 374 suitable for standard pie menus. Thus, we are working on including a version of control menus [29] into the 375 hierarchy. With control menus, repeated control signals are sent as long as the hand is kept within the menu 376 item in a selection pose. 377

We have started to implement a Hidden Markov Model for gesture learning and recognition in hopes to be able to create better and more natural gestures. The gestures currently implemented all use a heuristic approach. HMMs have been used extensively for gesture recognition in pen computing [DT04] and in vision [ER98] before. Using a type of machine learning instead of heuristics for a gesture recognizer is no more difficult to have interact with our system.

We are also considering a different scenario in which a few gestures (hand poses or deictic gestures) are used for direct control of common functions, such as controlling the sound level or lighting, and menubased gestures are used for more complex selections. In this situation it seems attractive to investigate if signs from Sign Language could be used for the static hand poses and poses for menu control.

387 **17 XIII.**

388 18 CONCLUSIONS

Human-computer interaction is still in its infancy. Visual interpretation of hand gestures would allow the 389 development of potentially natural interfaces to computer controlledenvironments. In response to this potential, 390 391 thenumber of different approaches to videobased hand gesturerecognition has grown tremendously in recent 392 years. Thus there is a growing need for systematization and analysis of many aspects of gestural interaction. 393 Several simple HCI systems have been proposed that demonstrate the potential of visionbased gestural interfaces. 394 However, from a practical standpoint, the development of such systems is in its infancy. Though most current systems employ hand gestures for the manipulation of objects, the complexity of the interpretation of gestures 395 dictates the achievable solution. For example, the gestures used to convey manipulative actions today are usually 396 of the communicative type. Further, hand gestures for HCI are mostly restricted to single-handed and produced 397 only by a single user in the system. This consequently downgrades the effectiveness of the interaction. We 398 suggest several directions of research for raising these limitations toward gestural HCI. For example, integration 399

- 400 of hand gestures with speech, gaze and other naturally related modes of communication in a multimodal interface.
- 401 However, substantial research effort that connects advances in computer vision with the basic study of humancomputer interaction will be needed in the future to develop an effective andnatural hand gesture interface. ¹



Figure 1: Fig. 1:

 $\begin{array}{c}
402 \\
403 \\
2 3 4 5 6 7
\end{array}$

⁵DecemberHand Gesture Interaction with Human-Computer

 $^{^1 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 4 2011 December Hand Gesture Interaction with Human-Computer

 $^{^2 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 5 2011 December Hand Gesture Interaction with Human-Computer

 $^{^3 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 6 2011 December Hand Gesture Interaction with Human-Computer

 $^{^4 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 8 2011 December Hand Gesture Interaction with Human-Computer

 $^{^6 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 10 2011 December Hand Gesture Interaction with Human-Computer

 $^{^7 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I 12 2011 December Hand Gesture Interaction with Human-Computer



Figure 2: Fig. 3:



Figure 3:

18 CONCLUSIONS

404 .1 ACKNOWLEDGMENTS

- 405 We thank Björn Eiderbäck, CID, who performed the Smalltalk programming for the menu system. Olle Sundblad,
- 406 CID did Java programming for the Application control server.
- 407 [Beaudoin-Lafon et al.], M Beaudoin-Lafon, W Mackay, P Andersen, P Janecek, M Jensen, M Lassen, K
 408 Lund, K Mortensen, S Munck, K Ravn, A Ratzer, Christensen, Jensen.
- [Baudel et al. ()], T Baudel, M Beaudouin-Lafon, Charade. Communications of the ACM 1993. 36 (7) p. .
- 410 [Segen and Kumar (2000)], J Segen, S Kumar. No Mouse! Communications of the ACM 2000. July 2. 43 (7).
- 411 [Triesch et al. ()] 'A system for person-independent hand posture recognition against complex backgrounds'. J
- Triesch , C Von , Malsburg . *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2001. 23 p.
 12.
- 414 [Wexelblatt (1995)] 'An Approach to Natural Gesture in Virtual Environments'. A Wexelblatt . ACM ToCHI
 415 1995. September. 2 (3) p. .
- [Callahan et al. ()] 'An Empirical Comparision of Pie vs. Linear Menus'. J Callahan , D Hopkins , M Weiser , B
 Shneiderman . *Proceedings of CHI'88*, (CHI'88) 1988. p. .
- [Lee and Kim ()] 'An HMM-based threshold model approach for gesture recognition'. H.-K Lee , J H Kim . In
 IEEE Trans. on Pattern Analysis and Machine Intelligence 1999. 21 p. 10.
- 420 [Krueger ()] Artificial Reality II, M Krueger . 1991. Addison-Wesley.
- 421 [Automatic Face-and Gesture-Recognition] Automatic Face-and Gesture-Recognition, Zurich, Switzerland. p. .
- Hardenberg and Berard ()] 'Bare-Hand Human-Computer Interaction'. Von Hardenberg , C Berard , F . Proc.
 of ACM Workshop on Perceptive User Interfaces, (of ACM Workshop on Perceptive User InterfacesOrlando,
 Florida) 2001.
- 425 [Cadoz ()] C Cadoz . Les Réalites Virtuelles. Dominos, Flammarion, 1994.
- [Abowd and Mynatt ()] 'Charting Past, Present, and Future Research in Ubiquitous Computing'. G D Abowd ,
 E D Mynatt . ACM ToCHI 2000. 7 (1) p. .
- [Wren et al. (1999)] Combining Audio and video in Perceptive Spaces. 1st International Workshop on Managing
 Interactions in Smart Environments, C R Wren, S Basu, F Sparacino, A Pentland. 1999. Dec. 13-14.
- ⁴³⁰ Dublin, Ireland.
 ⁴³¹ [Freeman et al. (1998)] 'Computer Vision for Interactive Computer Graphics'. W T Freeman , D Anderson , P
- Beardsley, C Dodge, H Kage, K Kyuma, Y Miyake, M Roth, K Tanaka, C Weissman, W Yerazunis.
 IEEE Computer Graphics and Applications 1998. May-June. p. .
- (Pook et al. ()] 'Control Menus: Execution and Control in a Single Interactor'. S Pook , E Lecolinet , G Vaysseix
 , E Barillot . *Extended Abstracts of CHI2000*, 2000. p. .
- [K ()] CPN/Tools: Revisiting the Desktop Metaphor with Post-WIMP Interaction Techniques, K. 2001. p. .
 (Extended Abstracts from CHI2001)
- ⁴³⁸ [Oviatt et al. ()] 'Designing the User Interface for Multimodal Speech and Pen-Based Gesture Applications:
 ⁴³⁹ State-of-the-Art Systems and Future Research Directions'. , Oviatt , S Cohen , P Wu , L Vergo , J Duncan , L Suhm , B Bers , J Holzman , T Winograd , T Landay , J Larson , J Ferro , D . *Human-Computer*
- Interaction 2000. 15 p. .
- [European Conference on Computer Vision] European Conference on Computer Vision, (Berlin) Springer Verlag.
 1406 p. .
- 444 [Quek ()] 'Eyes in the Interface'. F Quek . International Journal of Image and Vision Computing 1995. 13 (6) p. 445 .
- [Quek et al. ()] 'FingerMouse: A Freehand Pointing Interface'. F Quek, T Mysliwiec, M Zhao. Proceedings of
 the International Workshop on © 2011 Global Journals Inc. (US), (the International Workshop on © 2011
- 448 Global Journals Inc. (US)) 1995.
- [Guimbretière and Winograd ()] 'FlowMenu: combining Command, Text and Data Entry'. F Guimbretière , T
 Winograd . Proceedings of UIST'2000, (UIST'2000) 2000. p. .
- 451 [Maggioni and Kämmerer ()] 'Gesture Computer -History, Design and Applications'. C Maggioni , B Kämmerer .
- 452 Computer Vision for Human-Computer Interaction, Pentland Cipolla (ed.) 1998. 1998. Cambridge University
 453 Press. p. 2351.
- [Turk ()] 'Gesture Recognition'. M Turk . Handbook of Virtual Environments. Design, Implementation, and
 Applications. Lawrence-Erlbaum Assoc, K Stanney (ed.) 2002.
- 456 [Global Journal of Computer Science and Technology Volume XI Issue XXIII Version I] Global Journal of
 457 Computer Science and Technology Volume XI Issue XXIII Version I, 11 p. 2011.

- 458 [Bretzner et al. (2002)] 'Hand Gesture Recognition using Multi-Scale Colour Features, Hierarchical Models and
- Particle Filtering'. L Bretzner, I Laptev, T Lindeberg. the 5th International Conference on Automatic Face
 and Gesture Recognition, (Washington, D.C) 2002. May 2002. (To appear in)
- [Cui and Weng ()] 'Hand sign recognition from intensity image sequences with complex background'. Y Cui , J
 Weng . Proc. IEEE Conference on Computer Vision and Pattern Recognition, (IEEE Conference on Computer
 Vision and Pattern Recognition) 1996. p. .
- [Koike et al. ()] 'Integrating Paper and Digital Information on EnhancedDesk: A Method for Realtime Finger
 Tracking on an Augmented Desk System'. H Koike, Y Sato, Y Kobauashi. ACM ToCHI 2001. 8 (4) p.
- [Leibe et al. (2001)] 'Integration of Wireless Gesture Tracking, Object tracking and 3D Reconstruction in the
 Perceptive Workbench'. B Leibe, D Minnen, J Weeks, T Starner. Proceedings of 2nd International Workshop
 on Computer Vision Systems, (2nd International Workshop on Computer Vision SystemsVancouver, BC,
 Canada) 2001. ICVS 2001. July 2001.
- [Utsumi and Ohya ()] 'Multiple-Hand-Gesture Tracking Using Multiple Cameras'. A Utsumi , J Ohya . Proc.
 IEEE Conference on Computer Vision and Pattern Recognition, (IEEE Conference on Computer Vision and
- Pattern Recognition) 1999. p. .
- [O'hagan and Zelinsky ()] R O'hagan , A Zelinsky . Finger Track A Robust and Real-Time Gesture Interface.
 Australian Joint Conference on AI, (Perth) 1997.
- [Bolt ()] 'Put-that-there: Voice and Gesture in the graphics interface'. R A Bolt . Computer Graphics 1980. 14
 (3) p. .
- 479 [Sato et al. ()] 'Real-time input of 3D pose and gestures of a user's hand and its applications for HCI'. Y Sato
- 480 , M Saito , H Koike . Proc. 2001 IEEE Virtual Reality Conference, (2001 IEEE Virtual Reality Conference)
 481 2001. IEEE VR2001. p. .
- 482 [Starner et al. ()] 'Realtime American Sign Language recognition using desk and wearable computer-based video'.
- 483 T Starner, J Weaver, A Pentland. *IEEE Transactions on Pattern. Analysis and Machine. Intelligence* 1998.
- ⁴⁸⁴ [Braffort ()] 'Research on Computer Science and Sign Language: Ethical Aspects'. A Braffort . Lecture Notes in
 ⁴⁸⁵ Artificial Intelligence Roy, D. & Panayi, (ed.) 2001. Springer-Verlag.
- 486 [Mcneill ()] 'So you think gestures are nonverbal'. D Mcneill . Psychological Review 1985. 92 (3) p. .
- [Freeman and Weissman ()] 'Television Control by Hand Gestures'. W T Freeman , C D Weissman . 1st Intl.
 Conf. on Automatic Face and Gesture Recognition, 1994.
- [Kendon (ed.) ()] The biological Foundation of Gestures: Motor and Semiotic Aspects, A Kendon . Nespoulous,
 J.-L., Peron, P. & Lecours, A.R. (ed.) 1986. Lawrence -Erlbaum. p. . (Current issues in the study of gesture)
- ⁴⁹¹ [Wellner (1991)] 'The DigitalDesk Calculator: Ta ctile Manipulation on a Desk Top Display'. P Wellner .
 ⁴⁹² Proceedings of UIST'91, (UIST'91) 1991. Nov. 11-13. p. .
- ⁴⁹³ [Kurtenbach and Buxton ()] 'The Limits of Expert Performance Using Hierarchic Marking Menus'. G Kurtenbach
 ⁴⁹⁴ , W Buxton . *Proceedings of CHI'94*, (CHI'94) 1994. p. .
- [Kettebekov and Sharma ()] 'Toward Natural Gesture/Speech Control of a Large Display'. S & Kettebekov , R
 Sharma . Proceedings of EHCT'01, Lecture Notes in Computer Science (EHCT'01) 2001. Springer Verlag.
- ⁴⁹⁷ [Kjeldsen and Kender ()] 'Toward the use of gesture in traditional user interfaces'. R Kjeldsen , J Kender .
 ⁴⁹⁸ Proc. of IEEE International Conference on Automatic Face and Gesture Recognition, (of IEEE International Conference on Automatic Face and Gesture Recognition) 1996. p. .
- [Bretzner and Lindeberg ()] 'Use Your Hand as a 3-D Mouse or Relative Orientation from Extended Sequences
 of Sparse Point and Line Correspondances Using the Affine Trifocal Tensor'. L Bretzner, T Lindeberg. Proc.
 5th, H Burkhardt, B Neumann (ed.) (5th) 1998.
- [Krueger et al. ()] 'Videoplace -an artificial reality'. M W Krueger , T Gionfriddo , K Hinrichsen . Proceedings
 of CHI'85, (CHI'85) 1985. p. .
- [O'hagan et al. ()] 'Visual Gesture Interfaces for Virtual Environments'. R G O'hagan , A Zelinsky , S Rougeaux
 Interacting with Computers 2002. 14 p. .
- 507 [Pavlovic et al. ()] 'Visual Interpretation of Hand Gestures for Human-Computer Interaction: A Review'. V I
- Pavlovic, R Sharma, T S Huang. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1997.
 19 p. .