

REPORT DOCUMENTATION PAGE

AFRL-SR-BL-TR-98-

0441

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188)

1. AGENCY USE ONLY (Leave Blank)		2. REPORT DATE June 1991 2-95	3. REPORT TYPE Final
4. TITLE AND SUBTITLE Visual Recognition of American Sign Language Using Hidden Markov Models			5. FUNDING NUMBERS
6. AUTHORS Thad Eugene Starner			
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Massachusetts Institute of Technology			8. PERFORMING ORGANIZATION REPORT NUMBER
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) AFOSR/NI 110 Duncan Avenue, Room B-115 Bolling Air Force Base, DC 20332-8080			10. SPONSORING/MONITORING AGENCY REPORT NUMBER
11. SUPPLEMENTARY NOTES			
12a. DISTRIBUTION AVAILABILITY STATEMENT Approved for Public Release			12b. DISTRIBUTION CODE
13. ABSTRACT (Maximum 200 words) See attached.			
14. SUBJECT TERMS			15. NUMBER OF PAGES
			16. PRICE CODE
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT UL

DISTRIBUTION STATEMENT A
Approved for public release;
Distribution Unlimited

19980518 031

DTIC QUALITY INSPECTED 5

Visual Recognition of American Sign Language Using Hidden Markov Models

by
Thad Eugene Starner

S.B., Computer Science
S.B., Brain and Cognitive Science
Massachusetts Institute of Technology, Cambridge MA
June 1991

Submitted to the Program in Media Arts and Sciences,
School of Architecture and Planning,
in partial fulfillment of the requirements for the degree of
MASTER OF SCIENCE IN MEDIA ARTS AND SCIENCES
at the
Massachusetts Institute of Technology
February 1995

© Massachusetts Institute of Technology, 1995
All Rights Reserved

Signature of Author _____
Program in Media Arts and Sciences
20 January 1995

Certified by _____
Alex Pentland
Head, Perceptual Computing Section
Program in Media Arts and Sciences
Thesis Supervisor

Accepted by _____
Stephen A. Benton
Chairperson
Departmental Committee on Graduate Students
Program in Media Arts and Sciences

Visual Recognition of American Sign Language Using Hidden Markov Models

by
Thad Eugene Starner

Submitted to the Program in Media Arts and Sciences,
School of Architecture and Planning
on January 20, 1995
in partial fulfillment of the requirements for the degree of
Master of Science in Media Arts and Sciences

Abstract

Using hidden Markov models (HMM's), an unobtrusive single view camera system is developed that can recognize hand gestures, namely, a subset of American Sign Language (ASL). Previous systems have concentrated on finger spelling or isolated word recognition, often using tethered electronic gloves for input. We achieve high recognition rates for full sentence ASL using only visual cues.

A forty word lexicon consisting of personal pronouns, verbs, nouns, and adjectives is used to create 494 randomly constructed five word sentences that are signed by the subject to the computer. The data is separated into a 395 sentence training set and an independent 99 sentence test set. While signing, the 2D position, orientation, and eccentricity of bounding ellipses of the hands are tracked in real time with the assistance of solidly colored gloves. Simultaneous recognition and segmentation of the resultant stream of feature vectors occurs five times faster than real time on an HP 735. With a strong grammar, the system achieves an accuracy of 97%; with no grammar, an accuracy of 91% is reached (95% correct).

Thesis Supervisor: Alex Pentland
Title: Head, Perceptual Computing Section, MIT Media Lab

This work was supported in part by the USAF Laboratory Graduate Fellowship program.

Visual Recognition of American Sign Language Using Hidden Markov Models

by
Thad Eugene Starner

The following people served as readers for this thesis:

Reader: _____
Alex Pentland
Head, Perceptual Computing Section
Program in Media Arts and Sciences

Reader: _____
Nathaniel I. Durlach
Senior Research Scientist
Electrical Engineering and Computer Science

Reader: _____
Pattie Maes
Assistant Professor
Program in Media Arts and Sciences

I'd like to thank Alex Pentland, my advisor, for insisting over my objections that this project was possible within the bounds of a Master's thesis. His experienced insight into what is possible and worthwhile has guided me throughout my MIT career. I would also like to thank my readers, Nathaniel Durlach and Pattie Maes for coping with the changes over the past year.

Many thanks to the USAF Lab Graduate Fellowship program for funding my graduate studies.

Thanks also to John Makhoul, Rich Schwartz, Long Nguyen, Bruce Musicus, Paul Placeway, and numerous others in the Bolt, Barenek, and Newman Speech Group for giving me a good understanding of and intuition for modern HMM techniques. Along the same lines, thanks to Roz Picard and Alex Pentland for their lucid explanations of pattern recognition techniques which got me started in this field.

Thanks to Judy Bornstein for sharing her experience with ASL and for proofing this document.

Thanks to Mike P. Johnson, Ken Russell, and the IVE gang at the Media Lab for their contributions of code snippets, which made my work much easier. Thanks also to the Vision and Modeling Group, which has harbored me for many years as I've been "playing with the toys." It is my pleasure to work with such talented and exciting people.

A big thank-you to my parents, who have supported me without fail since I first discovered that I wanted to do research at MIT.

Last but not least, thanks to Tavenner Hall for keeping me sane through the last death throes of this thesis.

Contents

1	Introduction	6
1.1	Applications	7
1.2	Outline	7
2	Problem Description	9
2.1	Analyzing Human Body Motion	9
2.2	American Sign Language	10
2.3	Goals	11
3	Background	17
3.1	Hand Recovery	17
3.2	Machine Sign Language Recognition	18
3.3	Previous Use of Hidden Markov Models in Gesture Recognition	19
3.4	Use of HMM's for Recognizing Sign Language	20
4	Tracking and Modeling Gestures Using Hidden Markov Models	23
4.1	Hidden Markov Modeling	23
4.2	Feature Extraction Given Binarized Images of the Hands	32
4.3	Recovering the Hands from Video	33
4.4	Selecting an HMM Topology	35
4.5	Training an HMM network	37
5	Experimentation	40
6	Analysis and Discussion	43
7	Summary and Future Work	46

Chapter 1

Introduction

To date, computers have had very limited means of communicating with humans. Most common methods involve using a tethered device (keyboard, mouse, light pen, 3D tracker, etc.) that limit the user's freedom of motion. Furthermore, the expressiveness of these interactions has been very poor. With the advent of more powerful computers equipped with video, vision based interfaces are becoming more feasible. Enabling the computer to "see" its user allows for richer and more varied paradigms of man-machine interaction.

Recently, there has been a surge in interest in recognizing human hand gestures. While there are many interesting domains, one of the most structured sets of gestures are those belonging to sign language. In sign language, each gesture already has an assigned meaning (or meanings) and strong rules of context and grammar may be applied to make recognition tractable.

Most work on sign language recognition employs expensive wired "datagloves" that the user must wear [39]. In addition, these systems mostly concentrate on finger signing, where the user spells each word with hand signs corresponding to the letters in the alphabet [10]. However, most sign does not involve finger spelling but signs that represent whole words. This allows signed conversations to proceed at about the pace of spoken conversation.

In this paper an extensible system is described that uses a single color camera to track hands in real time and recognize sentences of American Sign Language (ASL) using hidden Markov models (HMM's). The hand tracking stage of the system does not attempt to produce a fine-grain description of hand shape; studies have shown that such detailed information may be unnecessary for humans to interpret sign language [28]. Instead, the

tracking process produces only a coarse description of hand shape, orientation, and trajectory. The user is required to wear inexpensive colored gloves to facilitate the hand tracking frame rate and stability. This shape, orientation, and trajectory information is then input to an HMM for recognition of the signed words.

1.1 Applications

For many years the problem of continuous speech recognition has been a focus of research, with the goal of using speech as an interface. Similarly, if a full vocabulary recognition system for American Sign Language can be created, then ASL can be used in applications such as word processing, operating system control, etc. Perhaps the most promising application of an ASL recognition system is that of translation of ASL into written or spoken English. The translation problem involves more than recognizing signs, however; some level of grammar structure and meaning will have to be understood by the system to allow adequate translation.

Gestures are often made at points of stress in a conversation, when illustrating a motion, or when describing an object. In fact, research has been done on the language of these gestures [6]. By recognizing gestures made in conjunction with spoken language, a computer may be able to better understand the wishes of the user. If an ASL recognizer can be created, then similar technology may be applied to these conversational gestures as well.

Recently, the field of video annotation has gained popularity. Given a huge database of video footage, how are particular shots located? With a hand annotated database, a user can search the text hoping that the annotator attended what is desired. However, hand annotation is a time-consuming process. Instead, computer systems may be employed to annotate certain features of sequences. A human gesture recognition system adds another dimension to the types of features computers can automatically annotate.

1.2 Outline

Chapter 2 describes some attributes of American Sign Language and the scope of this thesis. Chapter 3 discusses previous work in related areas and develops the reasoning for choosing HMM's over other techniques. Details on the machine vision algorithms and the HMM

training and recognition methods used are provided in Chapter 4. Chapter 5 describes the experiments performed and lists the results. An analysis of the results is provided in Chapter 6, and a summary and discussion of future work is included in Chapter 7.

Chapter 2

Problem Description

When focusing the techniques of machine vision on the human body, many diverse fields are addressed. Techniques from cognition, psychophysics, dynamics, photography, and athletics can all be applied to help constrain problem domains.

2.1 Analyzing Human Body Motion

Photography has been used to help understand human body motion for over a century [25]. More recently, computers have been added to perform more complex analysis. Athletic programs may use computer tracking systems and dynamics to help maximize the amount of effort their athletes can produce. Many of these systems use hand labeled data or wired sensor systems to produce the data. Fully reconstructing the motion of the human body requires a tremendous amount of data. Therefore, both natural constraints of the human body and simplifying assumptions are often used to curtail the amount of data needed for analysis.

Several systems that address whole body systems have been developed in the past. These include gait recognition and analysis systems [27, 14, 32, 42], ballet step recognition [5], body capture [1], real-time interface systems [38, 23], and numerous others. Greater accuracy and detail can be gained by focusing attention on the body part of interest. Recent experimentation with active "focus of attention" systems is attracting interest to this topic [9]. In the case of ASL, the hands and head are of the most interest.

2.2 American Sign Language

American Sign Language is the language of choice for most deaf people in the United States. It is part of the “deaf culture” and includes its own system of puns, inside jokes, etc. However, ASL is but one of the many sign languages of the world. A speaker of ASL would have trouble understanding the Sign Language of China much as an English speaker would Chinese. ASL also uses its own grammar instead of borrowing it from English. This grammar allows more flexibility of word placement and sometimes uses redundancy for emphasis. Another variant, English Sign Language has more in common with spoken English but is not as widespread in America. ASL consists of approximately 6000 gestures of common words with finger spelling used to communicate obscure words or proper nouns. Finger spelling uses one hand and 26 gestures to communicate the 26 letters of the alphabet; however, users prefer full word signs whenever possible since this allows sign language to approach or surpass the speed of conversational English.

Conversants in ASL will often describe a person, place, or thing and then point to a place in space to temporarily store that object for later reference [35]. For example, “the man with the green sweater,” “the old grocery store,” and “the red garage” might be signed and put into various positions in space. To then say “the man with the green sweater went to the old grocery store and the red garage” would involve pointing to the location of the man and then making the sign for “walk” to the position of the store and garage in turn. For the purposes of this thesis, this particular spatial aspect of ASL will be ignored.

ASL uses facial expressions to distinguish between statements, questions, and directives. The eyebrows are raised for a question, held normal for a statement, and furrowed for a directive. While there has been work in recognizing facial gestures [11], facial features will not be used to aid recognition in the task addressed.

Traditionally, there are three components of an ASL sign:

tabular (tab): The position of the hand at the beginning of the sign.

designator (dez): The hand shape of the hand at the beginning of the sign.

signation (sig): The action of the hand(s) in the dynamic phase of the sign.

In the Stokoe ASL dictionary [37], 12 tabulars, 19 designators, and 24 signations are

distinguished. Even though both hands are used in ASL, only seven of the signations use two hands. Some signs depend on finger placement and movement to remove ambiguity, but many signs are distinct even when the fingers are ignored.

2.3 Goals

While the scope of this thesis is not to create a person independent, full lexicon system for recognizing ASL, a desired attribute of the system is extensibility towards this goal. Another goal is to allow the creation of a real-time system by guaranteeing each separate component (tracking, analysis, and recognition) runs in real-time. This demonstrates the possibility of a commercial product in the future, allows easier experimentation, and simplifies archiving of test data. "Continuous" sign language recognition of full sentences is desired to demonstrate the feasibility of recognizing complicated series of gestures. Of course, a low error rate is also a high priority.

Table 2.1: ASL Vocabulary Used

<i>part of speech</i>	<i>vocabulary</i>
pronoun	I you he we you(pl) they
verb	want like lose dontwant dontlike love pack hit loan
noun	box car book table paper pants bicycle bottle can wristwatch umbrella coat pencil shoes food magazine fish mouse pill bowl
adjective	red brown black gray yellow

In this recognition system, sentences of the form "personal pronoun, verb, noun, adjective, (the same) personal pronoun" are to be recognized. This sentence structure emphasizes the need for a distinct grammar for ASL recognition and allows a large variety of meaningful sentences to be randomly generated using words from each class. Table 2.3 shows the words chosen for each class. Six personal pronouns, nine verbs, twenty nouns, and five adjectives are included making the total lexicon number forty words. The words were chosen by paging through "A Basic Course in American Sign Language" by Humphries, Padden, and O'Rourke and selecting those which would provide coherent sentences when used to generate random sentences. At first a naive eye was used to avoid ambiguities in

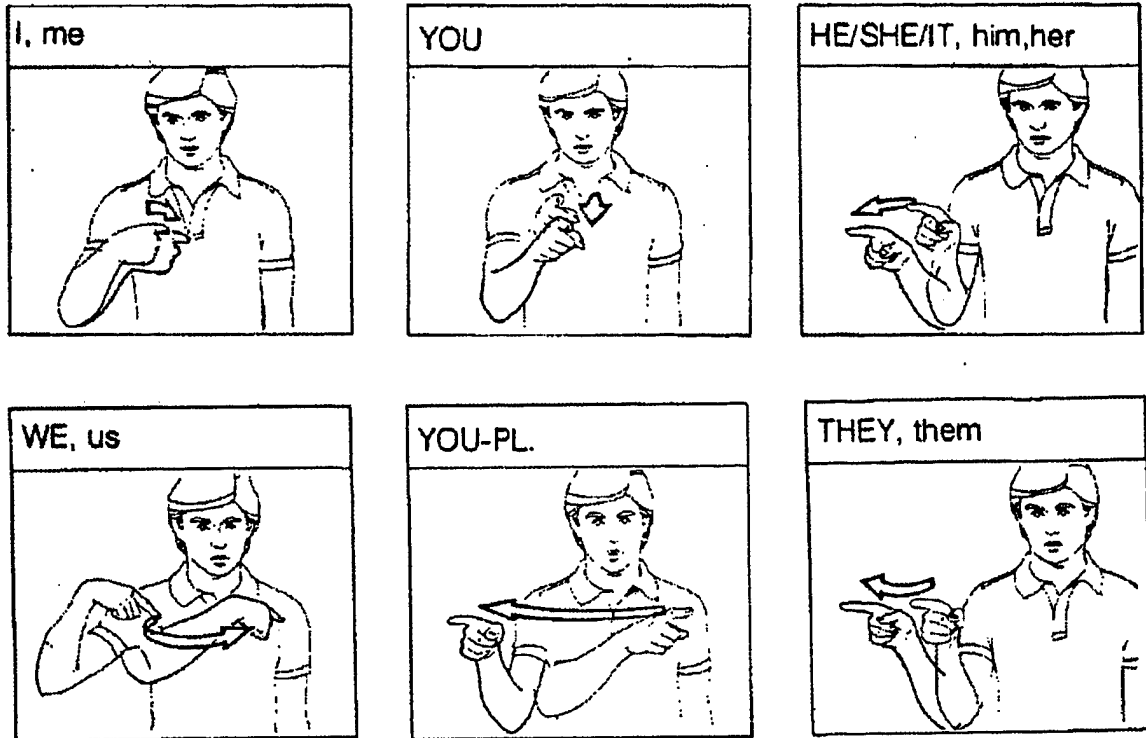


Figure 2-1: Pronouns.

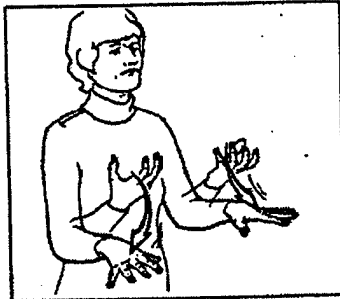
the selected signs, but this was shortly subsumed by the coherency constraint. Figures 2-1 to 2-5 illustrate the signs selected (from [17]).

The process of creating a new recognition system often limits the amount of initial training that can be collected. False starts and complications sometimes require discarding of data, making the initial process extremely frustrating for the subject. Therefore, the author learned the necessary signs to provide the database.

Recognition will occur on sentences unseen by the training process. The task is to correctly recognize the words in the given sentence in order and without inserting any additional words. Error and accuracy will be measured as in the continuous speech recognition literature, incorporating substitution, insertion, and deletion errors.



DON'T-LIKE



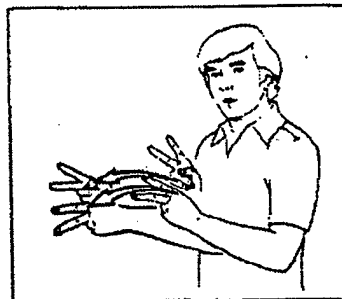
DON'T-WANT



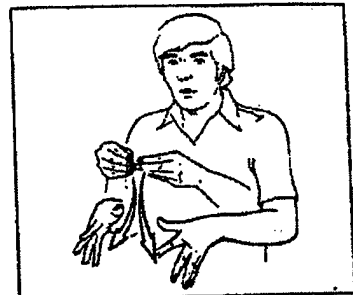
HIT



LIKE



LOAN, lend



LOSE



LOVE



PACK

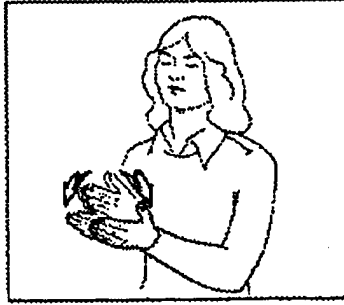


WANT, desire

Figure 2-2: Verbs.



BICYCLE



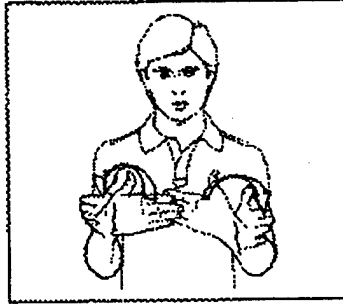
BOOK



BOTTLE



BOWL



BOX, package, room



CAR



COAT, jacket



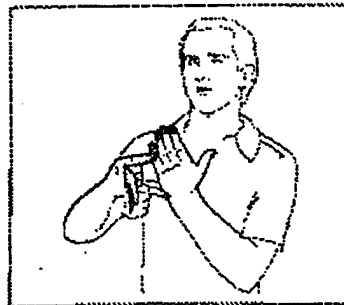
FISH



EAT, FOOD



GLASS, CAN, cup

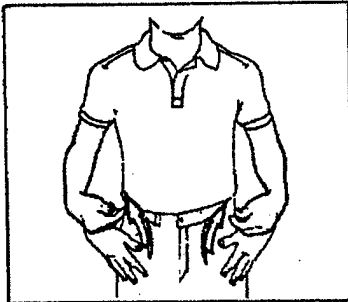


MAGAZINE, brochure,
journal, booklet



MOUSE

Figure 2-3: Nouns.



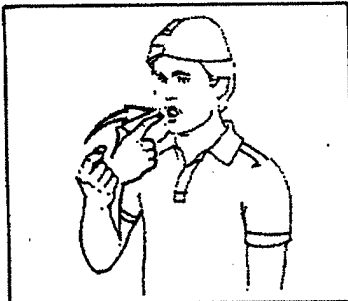
PANTS



PAPER, page



PENCIL



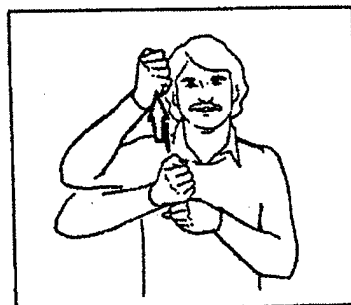
TAKE-PILL
Noun: pill



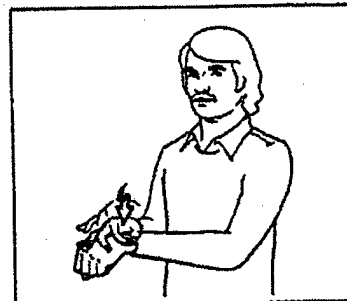
SHOES



TABLE, desk

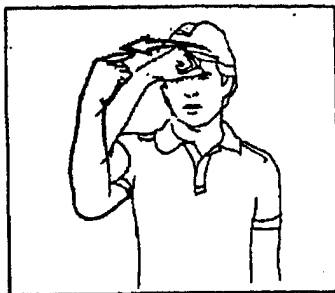


UMBRELLA

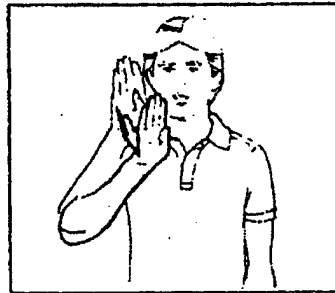


WRIST-WATCH

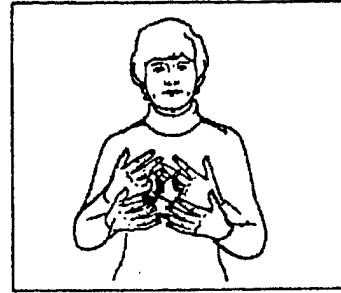
Figure 2-4: Nouns continued.



BLACK, Black-person



BROWN



GRAY



RED



YELLOW

Figure 2-5: Adjectives.

Chapter 3

Background

Visual recognition of sign language requires two main components: hand tracking and pattern recognition. Machine vision and virtual environment research have provided several tools in addressing the former, and continuous speech recognition has provided an excellent development platform for the latter. Here, recent tracking systems are surveyed, previous sign language work is reviewed, and the benefits of applying HMM technology to machine recognition of ASL are discussed.

3.1 Hand Recovery

With “multimedia” computers being packaged with video cameras, interest in human gesture recognition has grown. A large variety of interfaces have been proposed, using video driven gestures for mouse control [12], full body interactions [19, 23, 5], expression tracking [11], conducting music [24], and electronic presentation [38].

Due to their expressiveness, the hands have been a point of focus for many gesture recognition systems. Tracking the natural hand in real time using camera imagery is difficult, but successful systems have been demonstrated in controlled settings. Freeman [12] has shown a hand tracker that can be used for navigating 3D worlds. A greyscale camera tracks the hand in a small area on a table and uses hand and finger position to control the direction of a virtual graphics camera.

Rehg and Kanade have shown a two camera system that can recover the full 27 degrees of freedom in a hand [31]. While successfully demonstrating tracking with limited motion,

occlusion was not addressed. Furthermore, a simple background was required, which is often impossible when observing natural hand gestures.

Simpler, less constrained hand tracking systems have been created in a variety of environments. Krueger [19] has shown light table systems that tracked hands using single and multiple cameras aimed at the entire body. Maes *et al* [23] showed a similar system with one camera which allowed more arbitrary backgrounds to be used. Unfortunately, camera resolution often limits whole body systems to recovering just the position of the hands. When cameras are dedicated to the hands, more detail can be obtained. The "Hand Reader" by Suenaga *et al* [38] recovers 3D pointing information by dedicating two cameras, at close range, to the task. A limitation of such systems is that the working volume tends to be small. Also, the cameras for these systems tend to be obtrusive in that they are placed near the user. Longer focal length lenses may be used so that the cameras may be moved farther from the user, but then the space needed by such systems becomes prohibitive. Another solution is to monitor an entire room with a single fixed camera and use narrow field of view cameras mounted on servo platforms to direct attention to specific areas of interest [9]. This technique allows high detail and a wide range of motion, but suffers from the coupled motion problems of an active camera.

Tracking can often be simplified through using calibrated gloves or wired sensors. Dorner [10] uses a specially calibrated glove with different colors for each finger and the wrist and markers at each finger joint and tip. This, combined with Kalman filtering, simplifies occlusion problems and allows recovery of a detailed hand model through a wide range of motion. Datagloves by VPL are often used for sensing as well [24, 39, 4]. These systems offer precision, a relatively large range of motion (the sphere defined by the length of the tether), and very fast update rates at the expense of being wired to a sensing system.

3.2 Machine Sign Language Recognition

Excellent work has been done in support of machine sign language recognition. Sperling and Parish [35, 28] have done careful studies on the bandwidth necessary for a sign conversation using spatially and temporally subsampled images. Point light experiments (where "lights" are attached to significant locations on the body and just these points are used for

recognition), have been carried out by Poizner *et al* [29] and Tartter and Knowlton [41]. Tartter's 27 light experiment (13 lights per hand plus one on the nose) showed that a brief conversation in ASL was possible using only these stimuli. Poizner's experiments used only 9 lights (on the head and each shoulder, elbow, wrist, and index finger) to convey limited information about ASL signs. These experiments suggest that ASL might be recognizable even in an impoverished environment.

Most of the above experiments studied ASL in context. However, most machine recognition systems to date have studied isolated and/or static gestures. In many cases these gestures are finger spelling signs, whereas everyday ASL uses word signs for speed.

Tamura and Kawasaki demonstrated an early image processing system which could recognize 20 Japanese signs based on matching cheremes [40]. A chereme is composed of the tab, dez, and sig as discussed earlier. Charayaphan and Marble [7] demonstrated a feature set that could distinguish between the 31 isolated ASL signs in their training set. Takahashi and Kishino in [39] discuss a Dataglove-based system that could recognize 34 of the 46 Japanese kana alphabet gestures (user dependent) using a joint angle and hand orientation coding technique. From their paper, it seems the test user made each of the 46 gestures 10 times to provide data for principle component and cluster analysis. A separate test set was created from five iterations of the alphabet by the user, with each gesture well separated in time. While these systems are technically interesting, they suffer from a lack of training and have limited expandability beyond their sample domains.

3.3 Previous Use of Hidden Markov Models in Gesture Recognition

While the continuous speech recognition community adopted HMM's many years ago, these techniques are just now entering the vision community. Most early work was limited to handwriting recognition [21, 26]. More recently, He and Kundu [13] report using continuous density HMM's to classify planar shapes. Their method segmented closed shapes and exploited characteristic relations between consecutive segments for classification. The algorithm was reported to tolerate shape contour perturbation and some occlusion.

Another early effort by Yamato *et al* uses discrete HMM's to successfully recognize

image sequences of six different tennis strokes among three subjects. This experiment is significant because it used a 25x25 pixel quantized subsampled camera image as a feature vector. Even with such low-level information, the model could learn the set of motions to perform respectable recognition rates.

Schlenzig *et al* [33] also use hidden Markov models for visual gesture recognition. The gestures are limited to "hello," "good-bye," and "rotate". The authors report "intuitively" defining the HMM associated with each gesture and imply that the normal Baum-Welch re-estimation method was not implemented. However, this study shows the continuous gesture recognition capabilities of HMM's by recognizing gesture sequences.

Several vision systems have been developed with technology closely related to HMM methodology. Darrell [8] uses the dynamic time warping method to recognize gestures ("hello" and "good-bye") through time. Siskind and Morris [34] argue that event perception requires less information and may be an easier problem than object recognition. To this end they have constructed a maximum likelihood framework using methods similar to those used in hidden Markov model training to recognize the events "pick up," "put down," "drop," and "fall" from edge-detected movies of these actions. While these projects do not use HMM's *per se*, they set the stage for the idea that visual information can be modeled and used for recognition through time, much like speech recognition. In fact, Siskind and Morris make similar comparisons as those below between continuous speech recognition and vision.

3.4 Use of HMM's for Recognizing Sign Language

While the above studies show some promise for using HMM's in vision, what evidence is there that HMM's can be eventually used to address the full ASL recognition problem? The answer to that question lies in the comparison of the ASL recognition domain to the continuous speech recognition domain where HMM's have become the technology of choice.

Sign language and continuous speech share many common characteristics. Sign language can be viewed as a signal (position, shape, orientation of the hands, etc.) over time, just like speech. Silences in both speech and ASL are relatively easy to detect (hesitating between signs will be considered equivalent to stuttering or "ums" in speech), and all of the informa-

tion needed to specify a fundamental unit in both domains is given contiguously in a finite time period. The onset and offset paths of a sign depend on the temporarily neighboring signs. Correspondingly, spoken phonemes change due to coarticulation in speech. In both domains, the basic units combine to form more complex wholes (words to phrases in sign and phonemes to words in speech). Thus, language modeling can be applied to improve recognition performance for both problems.

In spite of the above similarities, sign language recognition has some basic differences from speech recognition. Unlike speech where phonemes combine to create words, the fundamental unit in much of ASL is the word itself. Thus, there is not as much support for individual word recognition in sign as there is in speech. Also, the fundamental unit in sign can switch suddenly (for example, changing into finger spelling for proper nouns). Furthermore, the grammar in ASL is significantly different than that of English speech (ASL is strongly influenced by French). However, even given these differences, there seems a strong likelihood that HMM's should also apply to sign language recognition.

Hidden Markov models have intrinsic properties which make them very attractive for ASL recognition. All that is necessary for training, except when using an optional bootstrapping process, is a data stream and its transcription (the text matching the signs). The training process can automatically align the components of the transcription to the data. Thus, no special effort is needed to label training data. The segmentation problem, as often seen in handwriting research, can be avoided altogether [36].

Recognition is also performed on a continuous data stream. Again, no explicit segmentation is necessary. The segmentation of sentences into words occurs naturally by incorporating the use of a lexicon and a language model into the recognition process. The result is a text stream that can be compared to a reference text for error calculation. Consequently sign language recognition seems an ideal machine vision application of HMM technology. The problem domain offers the benefits of scalability, well defined meanings, a pre-determined language model, a large base of users, and immediate applications for a recognizer.

Finally, HMM methodology is related to techniques that have been used previously in vision with success. Dynamic time warping, expectation maximization, Q-learning, neural nets, and several other traditional pattern recognition techniques resemble portions of the

modeling and recognition processes. The major advantage HMM's have over these other techniques is the ability to selectively, knowledgeably, and scalably tailor the model to the task at hand.

Chapter 4

Tracking and Modeling Gestures Using Hidden Markov Models

As shown in the last section, using computer vision to observe human motion is becoming a rich and diverse field. Traditionally, HMM's have been the domain of speech recognition, where a very rich system of modeling has evolved. In the following sections the vision methods used to track the hands are described, and the basics of HMM's are applied to ASL recognition.

4.1 Hidden Markov Modeling

While a substantial body of literature exists on HMM technology [43, 16, 30, 2, 20], this section modifies a traditional discussion of the algorithms so as to provide the perspective used for recognizing sign language. A simplistic example develops the fundamental theory in training and testing of a discrete HMM which is then generalized to the continuous density case used in the experiments. For broader discussion of the topic, [16] is recommended.

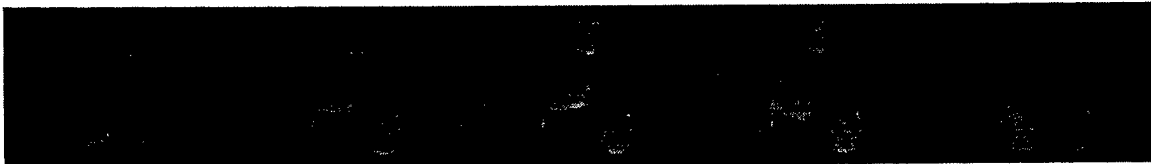
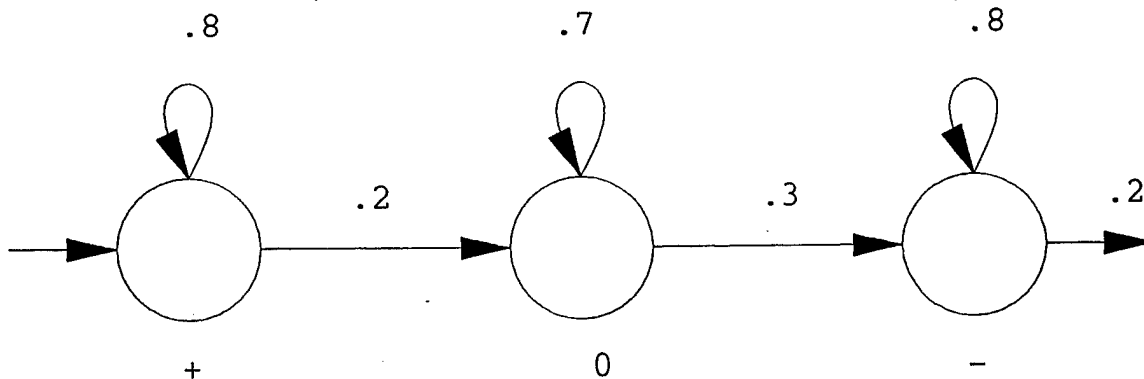
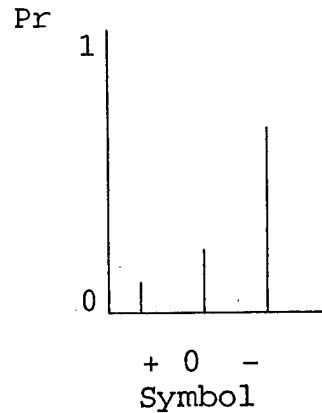
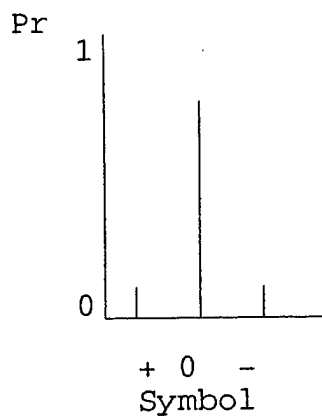
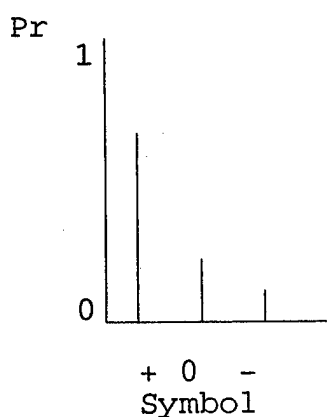
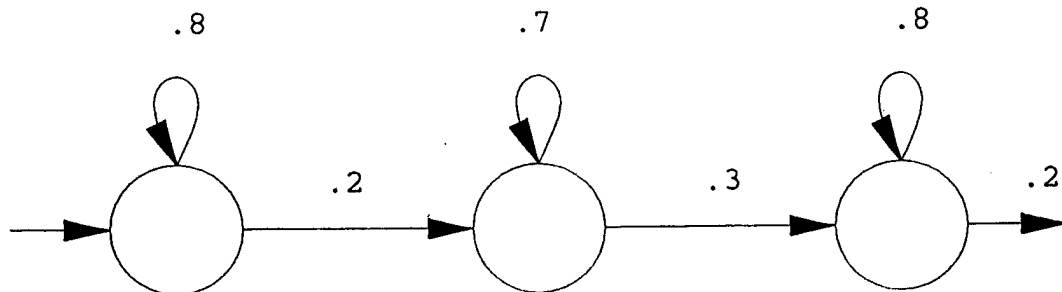


Figure 4-1: The ASL sign for "I."



Generated sequence: +++++|000|-----

Figure 4-2: A sample HMM topology with transition probabilities. Note that generated sequence can be divided into the states which produced each section.



Generated sequence: +0++++00-0---+---

Figure 4-3: The same 3 state topology with output probabilities added. Now the state sequence can no longer be precisely recovered.

A time domain process demonstrates a Markov property if the conditional probability density of the current event, given all present and past events, depends only on the j th most recent events. If the current event depends solely on the most recent past event, then the process is a first order Markov process. While the order of words in American Sign Language is not truly a first order Markov process, it is a useful assumption when considering the positions and orientations of the hands of the signer through time. Consider the hand actions for the isolated gesture "I" (Figure 4-1). These actions might be separated into the onset movement of the hand to the chest, pointing at the chest, and the offset movement of the hand back to a rest position. If the y movement of the right hand is used as a feature in recognizing these actions (say, at 5 frames/sec), the onset, pointing, and offset actions would correspond to a predominantly positive motion from the floor (+), a relatively stable period (0), and a predominantly negative motion respectively (-). Consider the onset, pointing, and offset actions to be "states" in a Markov model of the motions involved in the sign "I" (see Figure 4-2). Note that each state might transition back to itself. This corresponds to the inherent time aspect of the separate actions of the gesture. For example, the onset action may take a second, corresponding to four transitions (5 frames/sec) of the onset state back to itself. Similarly, the pointing and offset actions may take varying amounts of time to complete. This is reflected in the transition probabilities shown in Figure 4-2. In this simple example, associated with each state are the data of positive, negative, or no y movement. Thus, the model can be used to generate appropriate y movement for the gesture "I". However, in real life the motions for the gesture "I" are not deterministic. Due to possible indecision and wavering of the hand the actual y motion of the right hand might consist of varying proportions of positive, negative, and no motion during any given action. To reflect this, the model is changed to that in Figure 4-3 where the output of each state is a discrete probability distribution (output probability) of the three possible y classes. Now, if the model is used to generate appropriate y motion it is not obvious from which state each y movement is generated. This attribute is the reason for the word "hidden" in hidden Markov models; the state sequence is not obvious from the generated set of motions. This example demonstrates one other assumption for first order HMM's, that the output probability depends only on the current state and not on how or when the state was entered.

The previous paragraph described how to use HMM's to generate a sequence of symbols

that statistically resembles the data that the modeled process might produce. For the rest of this section, this idea is turned around to show how HMM's can be used to recognize a similar string of symbols (which will now be called observations). In addition, the algorithms needed to train a set of HMM's are described.

In order to proceed, a clear standard for notation is needed. Below is a list of symbols that will be used in this discussion. The meaning for some of these variables will become clearer in context, but the reader is urged to gain some familiarity with them before continuing.

T : the number of observations.

N : number of states in the HMM.

L : distinct number of possible observations (in this example three: +, 0, -).

s : a state. For convenience (and with regard to convention in the HMM literature), state i at time t will be denoted as $s_t = i$.

S : the set of states. S_I and S_F will be used to denote the set of initial and final states respectively.

O_t : an observation at time t .

O : an observation sequence O_1, O_2, \dots, O_T .

v : a particular type of observation. For example, v_1 , v_2 , and v_3 might represent +, 0, and - y motion respectively.

a : state transition probability. a_{ij} represents the transition probability from state i to state j .

A : the set of state transition probabilities.

b : state output probability. $b_j(k)$ represents the probability of generating some discrete symbol v_k in state j .

B : the set of state output probabilities.

π : initial state distribution.

λ : a convenience variable representing a particular hidden Markov model. λ consists of A , B , and π .

α : the “forward variable,” a convenience variable. $\alpha_t(i)$ is the probability of the partial observation sequence to time t and state i , which is reached at time t , given the model λ . In notation, $\alpha_t(i) = Pr(O_1, O_2, \dots, O_t, s_t = i | \lambda)$.

β : the “backward variable,” a convenience variable. Similar to the forward variable, $\beta_t(i) = Pr(O_{t+1}, O_{t+2}, \dots, O_T | s_t = i, \lambda)$, or the probability of the partial observation sequence from $t + 1$ to the final observation T , given state i at time t and the model λ .

γ : generally used for *a posteriori* probabilities. $\gamma_t(i, j)$ will be defined as the probability of a path being in state i at time t and making a transition to state j at time $t + 1$, given the observation sequence and the particular model. In other words, $\gamma_t(i, j) = Pr(s_t = i, s_{t+1} = j | \mathbf{O}, \lambda)$. $\gamma_t(i)$ will be defined as the posterior probability of being in state i at time t given the observation sequence and the model, or $\gamma_t(i) = Pr(s_t = i | \mathbf{O}, \lambda)$.

While a few other variables will be introduced in the description of the following algorithms, the above variables are typically found in descriptions of work in the field.

There are three key problems in HMM use. These are the evaluation problem, the estimation problem, and the decoding problem. The evaluation problem is that given an observation sequence and a model, what is the probability that the observed sequence was generated by the model ($Pr(\mathbf{O} | \lambda)$)? If this can be evaluated for all competing models for an observation sequence, then the model with the highest probability can be chosen for recognition.

$Pr(\mathbf{O} | \lambda)$ can be calculated several ways. The naive way is to sum the probability over all the possible state sequences in a model for the observation sequence:

$$Pr(\mathbf{O} | \lambda) = \sum_{all\ S} \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t}(O_t)$$

The initial distribution π_{s_1} is absorbed into the notation for $a_{s_0s_1}$ for simplicity in this discussion. The above equation can be better understood by ignoring the outside sum and product and setting $t = 1$. Assuming a particular state sequence through the model and

the observation sequence, the inner product is the probability of transitioning to the state at time 1 (in this case, from the initial state) times the probability of observation 1 being output from this state. By multiplying over all times 1 through T , the probability that the state sequence S and the observation sequence \mathbf{O} occur together is obtained. Summing this probability for all possible state sequences S produces $Pr(\mathbf{O}|\lambda)$. However, this method is exponential in time, so the more efficient forward-backward algorithm is used in practice.

The forward variable has already been defined above. Here its inductive calculation, called the forward algorithm, is shown (from [16]).

- $\alpha_1(i) = \pi_i b_i(O_1)$, for all states i (if $i \in S_I$, $\pi_i = \frac{1}{n_I}$; otherwise $\pi_i = 0$)
- Calculating $\alpha()$ along the time axis, for $t = 2, \dots, T$, and all states j , compute

$$\alpha_t(j) = \left[\sum_i \alpha_{t-1}(i) a_{ij} \right] b_j(O_t)$$

- Final probability is given by

$$Pr(\mathbf{O}|\lambda) = \sum_{i \in S_F} \alpha_T(i)$$

The first step initializes the forward variable with the initial probability for all states, while the second step inductively steps the forward variable through time. The final step gives the desired result $Pr(\mathbf{O}|\lambda)$, and it can be shown by constructing a lattice of states and transitions through time that the computation is only order $O(N^2T)$.

Another way of computing $Pr(\mathbf{O}|\lambda)$ is through use of the backward variable β , as already defined above, in a similar manner.

- $\beta_T(i) = \frac{1}{N_F}$, for all states $i \in S_F$; otherwise $\beta_T(i) = 0$
- Calculating $\beta()$ along the time axis, for $t = T-1, T-2, \dots, 1$ and all states j , compute:

$$\beta_t(j) = \sum_i a_{ji} b_i(O_{t+1}) \beta_{t+1}(i)$$

- Final probability is given by:

$$Pr(\mathbf{O}|\lambda) = \sum_{i \in S_T} \pi_i b_i(O_1) \beta_1(i)$$

The estimation problem concerns how to adjust λ to maximize $Pr(\mathbf{O}|\lambda)$ given an observation sequence \mathbf{O} . Given an initial model, which can have flat probabilities, the forward-backward algorithm allows us to evaluate this probability. All that remains is to find a method to improve the initial model. Unfortunately, an analytical solution is not known, but an iterative technique can be employed.

Referring back to the definitions and the gesture "I" example earlier, it can be seen that

$$\sum_{t \in O_t=+} \gamma_t(1)$$

is the expected number of times "+" should occur (given the observation and model) during a typical onset action of the gesture "I." The total expected motions (+, 0, -) given the model and observation during the onset action is

$$\sum_{t=1}^T \gamma_t(1)$$

Note that now, from the actual evidence and these two calculations, a new estimate for the respective output probability (for up in the onset gesture) can be assigned. Generalizing for a new estimate of a given output probability,

$$\bar{b}_j(k) = \frac{\sum_{t \in O_t=v_k} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$

Similarly, the evidence can be used to develop a new estimate of the probability of a state transition. Thus,

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \gamma_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

Initial state probabilities can also be re-estimated through the formula

$$\bar{\pi}_i = \gamma_1(i)$$

Thus all the components of λ , namely A , B , and π can be re-estimated. Since either the forward or backward algorithm can be used to evaluate $Pr(\mathbf{O}|\bar{\lambda})$ versus the previous estimation, the above technique can be used iteratively to converge the model to some limit. While the technique described only handles a single observation sequence, it is easy to extend to a set of observation sequences. A more formal discussion can be found in [16, 2, 43].

While the estimation and evaluation processes described above are sufficient for the development of an HMM system, the Viterbi algorithm provides a quick means of evaluating a set of HMM's in practice as well as providing a solution for the decoding problem. In decoding, the goal is to recover the state sequence given an observation sequence. The Viterbi algorithm can be viewed as a special form of the forward-backward algorithm where only the maximum path at each time step is taken instead of all paths. This optimization reduces computational load and additionally allows the recovery of the most likely state sequence. The steps to the Viterbi are

- Initialization. For all states i , $\delta_1(i) = \pi_i b_i(O_1)$; $\psi_1(i) = 0$
- Recursion. From $t = 2$ to T and for all states j , $\delta_t(j) = \text{Max}_i[\delta_{t-1}(i) a_{ij}] b_j(O_t)$;
 $\psi_t(j) = \text{argmax}_i[\delta_{t-1}(i) a_{ij}]$
- Termination. $P = \text{Max}_{s \in S_F}[\delta_T(s)]$; $s_T = \text{argmax}_{s \in S_F}[\delta_T(s)]$
- Recovering the state sequence. From $t = T - 1$ to 1 , $s_t = \psi_{t+1}(s_{t+1})$

In many HMM system implementations, the Viterbi algorithm is used for evaluation at recognition time. Note that since Viterbi only guarantees the maximum of $Pr(\mathbf{O}, S|\lambda)$ over all S (as a result of the first order Markov assumption) instead of the *sum* over all possible state sequences, the resultant scores are only an approximation. For example, if there are two mostly disjoint state sequences through one model with medium probability and one state sequence through a second model with high probability, the Viterbi algorithm would

favor the second HMM over the first. However, [30] shows that the probabilities obtained from both methods may be typically very close.

In practice, the Viterbi algorithm may be modified with a limit on the lowest numerical value of the probability of the state sequence, which in effect causes a beam search of the space. While this modification no longer guarantees optimality, a considerable speed increase may be obtained. Furthermore, to aid in estimation, the Baum-Welch algorithm may be manipulated so that parts of the model are held constant while other parts are trained.

To date, the example using the y motion in the gesture "I" assumed quantization of the motion into three classes +,0, and -. It is easy to see how, instead of quantizing, the actual probability density for y motion might be used. However, the above algorithms must be modified to accept continuous densities. The efforts of Baum, Petrie, Liporace, and Juang [3, 2, 22, 18] showed how to generalize the Baum-Welch, Viterbi, and forward-backward algorithms to handle a variety of characteristic densities. In this context, however, the densities will be assumed to be Gaussian. Specifically,

$$b_j(O_t) = \frac{1}{\sqrt{(2\pi)^n |\sigma_j|}} e^{\frac{1}{2}(O_t - \mu_j)' \sigma_j^{-1} (O_t - \mu_j)}$$

Initial estimations of μ and σ may be gotten by dividing the evidence evenly among the states of the model and calculating the mean and variance in the normal way.

$$\mu_j = \frac{1}{T} \sum_{t=1}^T O_t$$

$$\sigma_j = \frac{1}{T} \sum_{t=1}^T (O_t - \mu_j)(O_t - \mu_j)'$$

Whereas flat densities were used for the initialization step before, here the evidence is used. Now all that is needed is a way to provide new estimates for the output probability. We wish to weight the influence of a particular observation for each state based on the likelihood of that observation occurring in that state. Adapting the solution from the discrete case yields

$$\bar{\mu}_j = \frac{\sum_{t=1}^T \gamma_t(j) O_t}{\sum_{t=1}^T \gamma_t(j)}$$

and

$$\bar{\sigma}_j = \frac{\sum_{t=1}^T \gamma_t(j) (O_t - \bar{\mu}_j)(O_t - \bar{\mu}_j)^t}{\sum_{t=1}^T \gamma_t(j)}$$

In practice, μ_j is used to calculate $\bar{\sigma}_j$ instead of the re-estimated $\bar{\mu}_j$ for convenience. While this is not strictly proper, the values are approximately equal in contiguous iterations [16] and seem not to make an empirical difference [43]. Since only one stream of data is being used and only one mixture (Gaussian density) is being assumed, the algorithms above can proceed normally incorporating these changes for the continuous density case.

4.2 Feature Extraction Given Binarized Images of the Hands

In the previous discussion, the y motion of the right hand, a scalar quantity, was used to demonstrate the mathematics behind continuous density HMM's. However, there is no factor excluding the use of vectors (in fact, the equations are written for vector form). Feature vectors will simply result in multi-dimensional Gaussian densities. Given this freedom, what features should be used to recognize sign language?

Previous experience has shown that starting simple and evolving the feature set is often best [36]. Since finger spelling is not being allowed and there are few ambiguities in the test vocabulary based on individual finger motion, a relatively coarse tracking system may be used. Based on previous work [23], it was assumed that a system could be designed to separate the hands from the rest of the scene (explained in the next section). Traditional vision algorithms could then be applied to the binarized result. Besides the position of the hands, some concept of the shape of the hand and the angle of the hand relative to horizontal seemed necessary. Thus, an eight element feature vector consisting of each hand's x and y position, angle of axis of least inertia, and eccentricity of bounding ellipse was decided upon. The eccentricity of the bounded ellipse was found by determining the ratio of the square roots of the eigenvalues that correspond to the matrix

$$\begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix}$$

where a, b, and c are defined as

$$a = \int \int_{I'} (x')^2 dx' dy'$$

$$b = \int \int_{I'} x' y' dx' dy'$$

$$c = \int \int_{I'} (y')^2 dx' dy'$$

(x' and y' are the x and y coordinates normalized to the centroid)

The axis of least inertia is then determined by the major axis of the bounding ellipse, which corresponds to the primary eigenvector of the matrix [15]. Figure 4-4 demonstrates bounding ellipses fitted to the images of the hands. Note the 180 degree ambiguity in the angle of the ellipses. To address this problem, the angles were only allowed to range from -90 to +90 degrees.

4.3 Recovering the Hands from Video

Since real-time recognition is a goal, several compromises were made. The subject wears distinctly colored gloves on each hand (a yellow glove for the right hand and an orange glove for the left) and sits in a chair before the camera (see Figure 4-5). Figure 4-6 shows the view from the camera's perspective and gives an impression of the quality of video that is used. Color NTSC composite video is captured and analyzed at a constant 5 frames per second at 320 by 243 pixel resolution on a Silicon Graphics Indigo 2 with Galileo video board. To initially find each hand, the algorithm scans the image until it finds a pixel of the appropriate color. Given this pixel as a seed, the region is grown by checking the eight nearest neighbors for the appropriate color. Each pixel checked is considered to be part of the hand. This, in effect, performs a simple morphological dilation upon the resultant image that helps to prevent edge and lighting aberrations. The centroid is calculated as

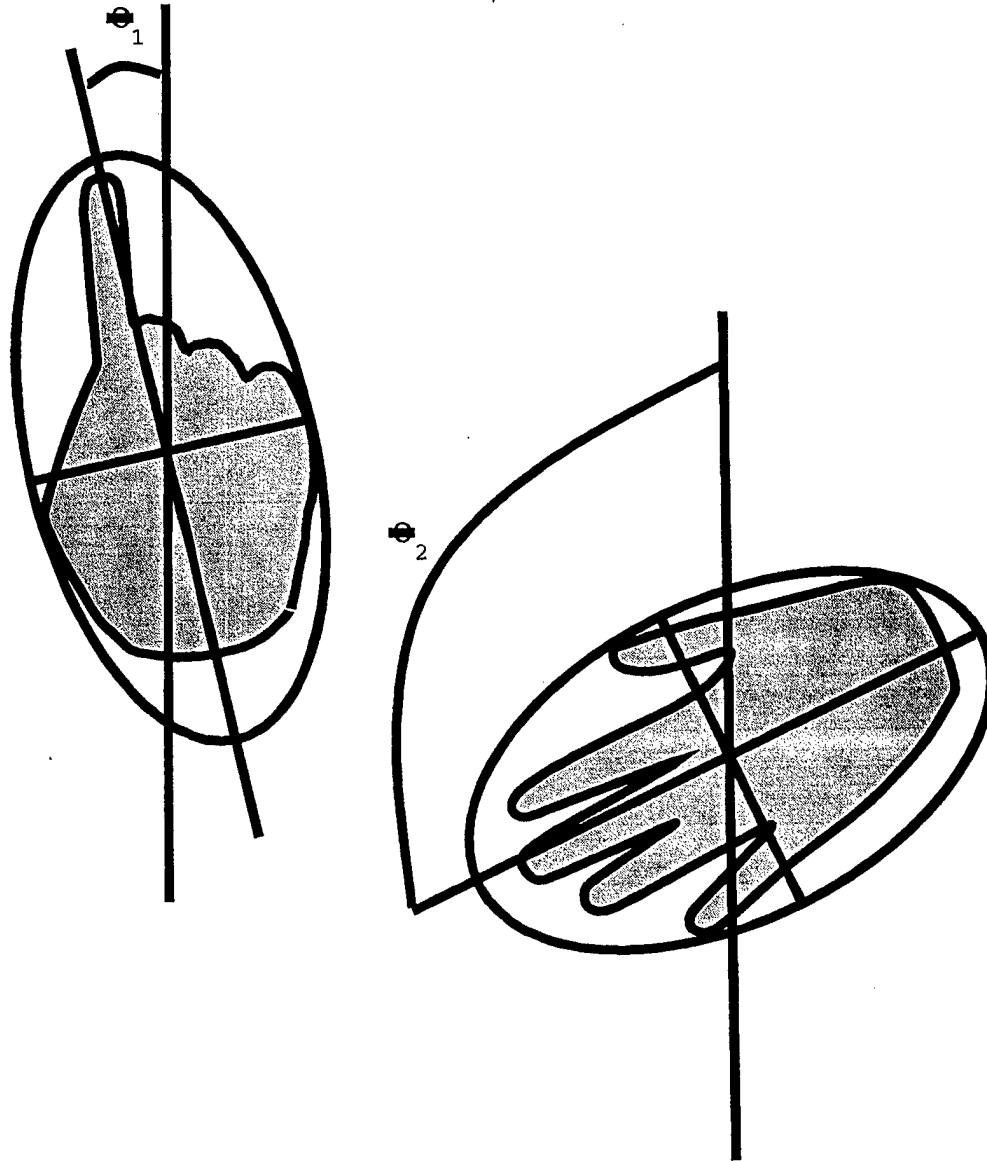


Figure 4-4: Bounding ellipses generated by tracking code.

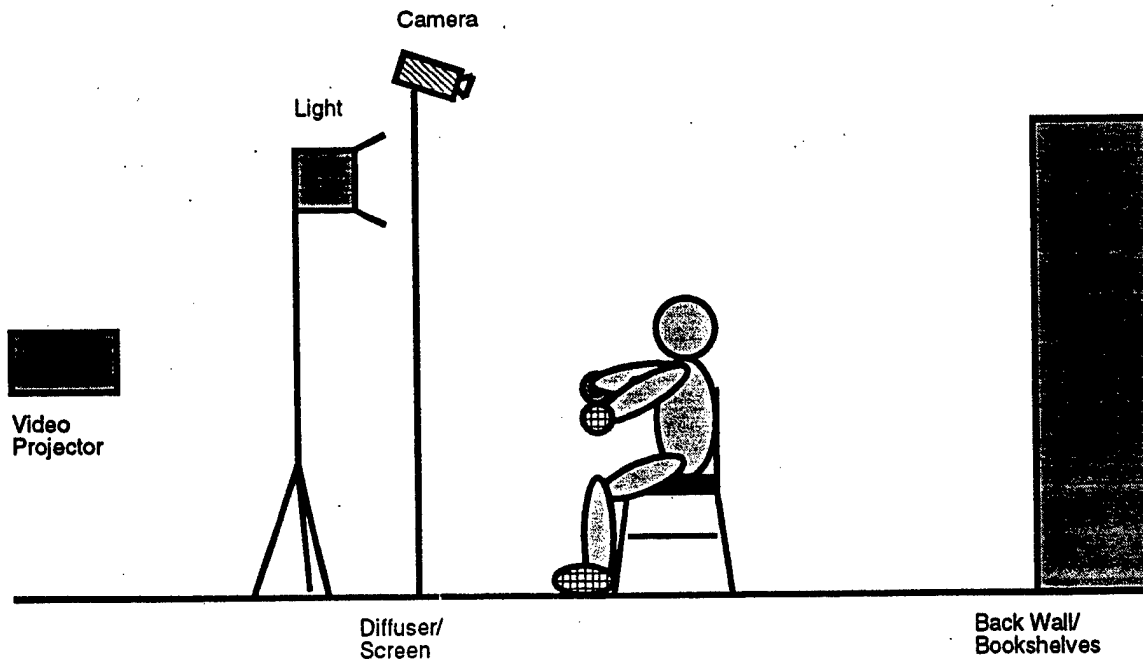


Figure 4-5: Tracking environment.

a by-product of the growing step and is stored as the seed for the next frame. Given the resultant bitmap and centroid, second moment analysis is performed as described earlier.

4.4 Selecting an HMM Topology

Previously, a 3 state model was used to represent the gesture "I." While this seems sufficient to model a simple gesture, such a model will not do as well in general. For more complex signs such as "table" which involves repetitive motion (onset; up, down, up motion of the right hand patting the left forearm; and offset), more states are necessary. In fact, a particular sign requires a different number of states if the gesturer is allowed to abbreviate the sign. However, through use of skip transitions, where the model has a certain probability of skipping part of the modeled sign, abbreviated signs can be accommodated (see Figure 4-7).

Skip transitions can also be used to avoid the task of determining the topology of each sign. For example, the model in Figure 4-7 can be trained to accommodate a three state sign like "I" through use of the skip states as well as a five state sign like "table," where the skip transitions would have low weights. Thus, once the minimum and maximum number

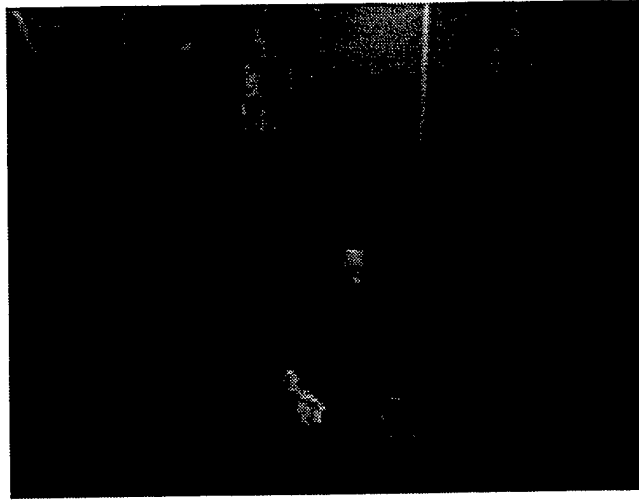


Figure 4-6: View from the tracking camera.

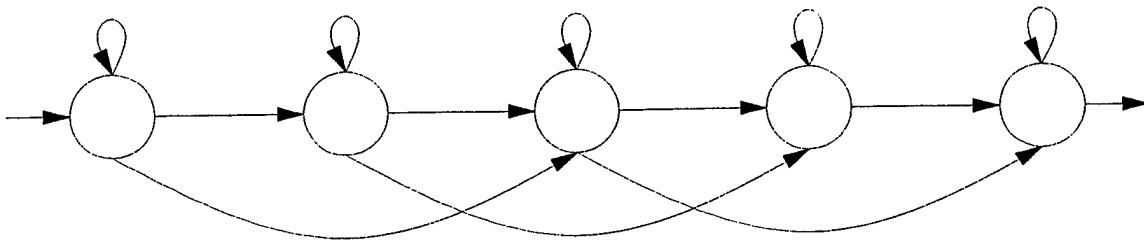


Figure 4-7: A more generalized topology

of states required is determined, a common topology may be used for all signs without too much loss of power (an added benefit is ease of coding).

4.5 Training an HMM network

When using HMM's to recognize strings of data, such as continuous speech, cursive handwriting, or ASL sentences, several methods can be brought to bear for training and recognition. Generally, individual models are concatenated together to model the larger language structures. Continuous speech recognition efforts have particularly advanced this field.

In speech, the fundamental unit (at least for these purposes) is the phoneme. Models for the individual phonemes can be trained separately, but this training is of limited utility if the goal is to recognize continuous sentences. In particular, the "co-articulation" effects of several phonemes spoken together may cause phonemes taken from continuous speech to differ significantly from phonemes spoken in isolation. An initial solution to this problem is to train on phonemes that have been manually segmented from continuous speech. When first addressing a task, manual segmentation is often worthwhile, especially when there is a small amount of training data and statistics on the field are not generally available. However, manual segmentation often has error and ambiguity. Furthermore, since the goal is recognition and not segmentation, it does not matter if the final system uses co-articulation to help recognize the phonemes. In fact, co-articulation is beneficial in that it provides context for recognizing the words of the sentence. To take advantage of this, two forms of context dependent training are used, embedded and context training.

Embedded training addresses the issue of segmentation. While initial training of the models might rely on manual segmentation or dividing the evidence evenly among the models, embedded training trains the models *in situ* and allows these boundaries to shift through a probabilistic entry into the initial states of each model [43].

Context training uses the co-occurrence of two or more fundamental units to allow recognition of blocks of units, which have more evidence than single units alone. In speech recognition, two and three phoneme blocks (biphones and triphones) are generally used. Note that recognition of these blocks might violate the first-order Markov assumption that was made earlier. However, by unrolling these blocks in the Viterbi process, the first-order

assumption can be preserved.

A final use of context in speech recognition is on the word level. Statistical grammars relating the probability of the co-occurrence of two or more words can be used to weight the recognition process. Grammars that associate two words are called bigrams, whereas grammars that associate three words are called trigrams.

In ASL, the fundamental unit is the sign. Since most signs represent whole words, context training occurs at the sentence level. In fact, three sign contexts (trisines) might represent most of a sentence. With finger spelling however, a more direct relation of sign to phoneme, word to word, and sentence to sentence can be established. This may present a challenge in higher level language modeling in the future, but a simple solution is simply to treat finger spelling as if it was on a word level. Statistical grammars may have a place in helping to merge these two levels of signing. With exclusive word signing, grammars are still useful. Generally, context training occurs at the model level using the training data provided. However, grammars are trained separately and can be trained solely on potential word orders. For example, in speech, a grammar may be trained on the articles from the Wall Street Journal from the past three years, while speech is only available for a small fraction of these articles. Thus, grammars might provide additional constraints on the data and simplify recognition.

A common mistake in neural net and HMM training is to provide too little training data. How many examples are enough? This, of course, is highly dependent on the task. In general, the data should be representative of the task domain. Training samples should include as much variance as is possible and reasonable for the task. Another issue is the role of context. While the scope of this thesis includes forty signs, trisine contexts may expand this number to 64,000 three sign combinations. In practice, the constraints of the language curtail this exponentiation. In fact, given the sentence structure used for the task proposed, only 2580 distinct trisines and 364 bisines are possible. The largest class of trisines, "pronoun verb noun," has 1080 members. The largest class of bisines, "verb noun," has 180 members. Note that while context is a powerful tool, it is not necessary for each context to be present in the training data for that context to appear as a result from recognition. Training can be "pooled" from small contexts or even the individual units to gain evidence for a new context during recognition. Of course, the best situation is for all

contexts to be seen before recognition, but this is not always possible when a recognition task includes models with a low probability of occurrence.

Weighing these factors against the time and expense of collecting data, 500 sentences were determined sufficient for an attempt at the task. Separating the database into 400 training and 100 test sentences gives an average of between 20 to 80 training examples of each sign (depending on class) with a good likelihood that each business occurs at least once.

Chapter 5

Experimentation

The results of the handtracking discussed in the last section can be seen in Figure 5-1. Tracking markers are overlaid on the camera images. The center of the arrow tracks the right hand while the incomplete diamond tracks the left. The length and width of the indicators note the length of the major and minor axes of the bounding ellipses while the angle of the indicators show the angle of the principle axes. Occasionally tracking would be lost (generating error values of 0) due to lighting effects, but recovery was fast enough (within a frame) so that this was not a problem. The 5 frame/sec rate was maintained within a tolerance of a few milliseconds. However, frames where tracking of a hand was lost were deleted. Thus, a constant data rate was not guaranteed.

Of the 500 sentences collected, six had to be thrown out due to subject error or outlier signs. Each sign ranged from approximately 1 to 3 seconds in length. No intentional pauses were given between signs within a sentence, but the sentences themselves were distinct.

Table 5.1: Percentage of words correctly recognized

	<i>test on training</i>	<i>test on independent test set</i>
grammar	99.4%	97.0%
no grammar	95.9% (93.5% accuracy) (D=13, S=88, I=60, N=2470)	95.0% (90.7% accuracy) (D=3, S=22, I=21, N=495)

To avoid the boot strapping segmentation process, the evidence in the sentences was evenly distributed between the words. Initial estimates for the means and variances of the output probabilities were provided by iteratively using Viterbi alignment on the training

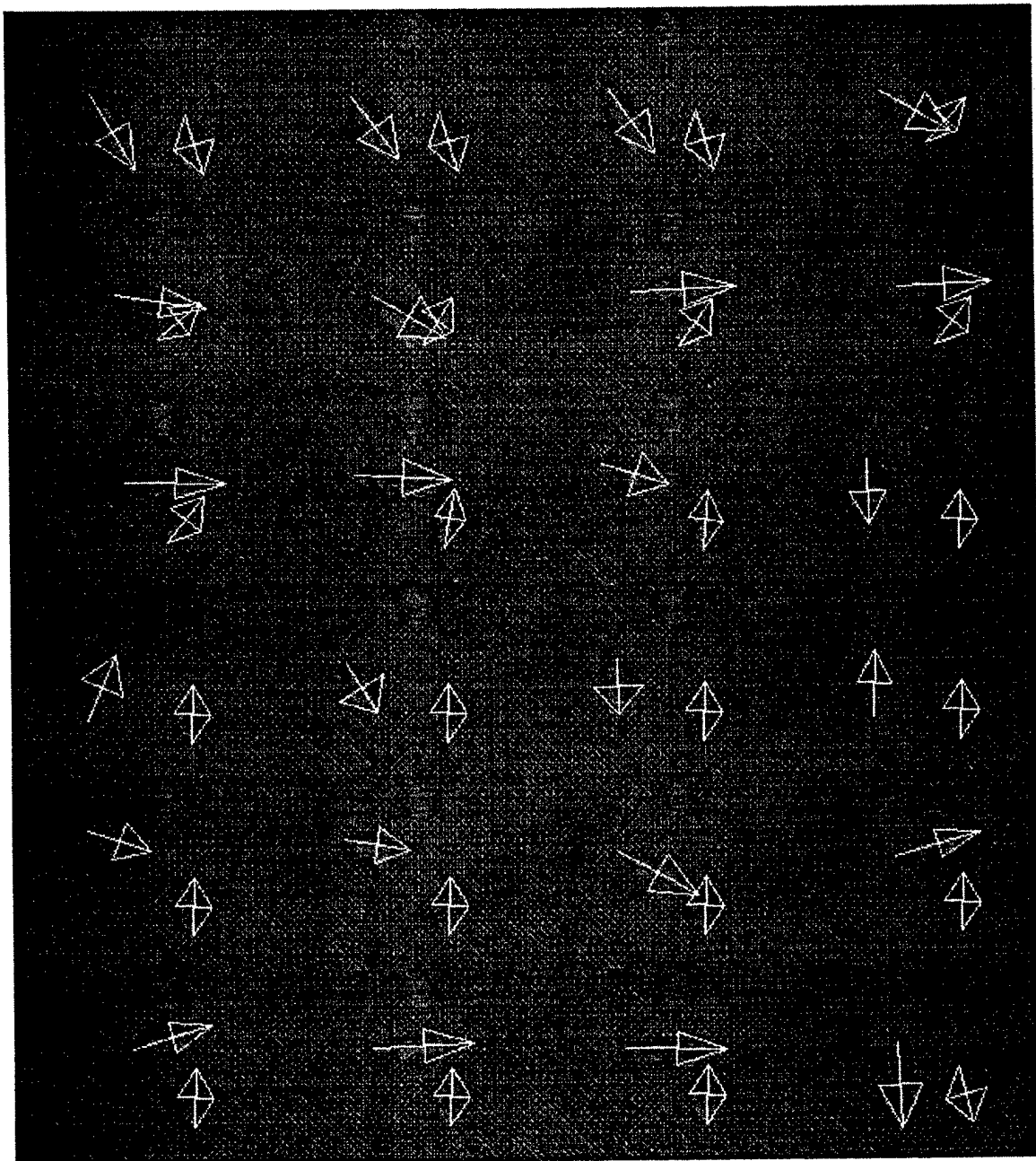


Figure 5-1: Hand tracking of the second half of a sentence "paper yellow we."

data and then recomputing the means and variances by pooling the vectors in each segment. Entropic's Hidden Markov Model ToolKit (HTK) was used as a basis for this step and all other HMM modeling and training tasks. The results from the initial alignment program are fed into a Baum-Welch re-estimator, whose estimates are, in turn, refined in embedded training which ignores any initial segmentation. For recognition, HTK's Viterbi recognizer was used both with and without a strong grammar based on the known form of the sentences. Gesture recognition occurs at a rate five times faster than real time. Contexts were not used, since a similar effect could be achieved with the strong grammar given this data set.

Word recognition results are shown in Table 5.1. When testing on training, all 494 sentences were used for both the test and train sets. For the fair test, the sentences were divided into a set of 395 training sentences and a set of 99 independent test sentences. The 99 test sentences were not used for any portion of the training. Given the strong grammar (pronoun, verb, noun, adjective, pronoun), insertion and deletion errors were not possible since the number and class of words allowed is known. Thus, all errors are substitutions when the grammar is used (and accuracy equals percent correct). However, without the grammar the recognizer is allowed to match the observation vectors with any number of the 40 vocabulary words in any order. Thus, deletion (D), insertion (I), and substitution (S) errors are possible. The absolute number of errors of each type are listed in Table 5.1 as well. The accuracy measure is calculated by subtracting the number of insertion errors from the number of correct labels and dividing by the total number of signs. Note that, since all errors are accounted against the accuracy rate, it is possible to get large negative accuracies (and corresponding error rates of over 100%). Most insertion errors occurred at signs with repetitive motion.

Chapter 6

Analysis and Discussion

While these results are far from being sufficient to claim a “working system” for full ASL recognition, they do show that this approach is promising. The high recognition rate on the training data indicates that the HMM topologies are sound and that the models are converging. Even so, the remaining 6.5% error rate (error rates will be based on accuracy measures) on the “no grammar” case indicates that some fine tuning on the feature set and model is in order. The 3.0% error rate on the independent test set shows that the models are generalizing well. However, a close look at the text produced by the recognition process shows some of the limitations of the feature set. Since the raw positions of the hands were used, the system was trained to expect certain gestures in certain locations. When this varied due to subject seating position or arm placement, the system could become confused. A simple fix to this problem would be to use position deltas in the feature vector instead.

Examining the errors made when no grammar was used shows the importance of finger position information. Signs like “pack,” “car,” and “gray” have very similar motions. In fact, the main difference between “pack” and “car” is that the fingers are pointed down for the former and clenched in the latter. Since this information was not available to the model, confusion could occur. While recovering general and/or specific finger position may be difficult in real time in the current testing area, simple palm orientation could be used for discrimination instead. In the current system, a simple implementation would be to paint the back of the gloves a different color than the palm.

A more interesting problem in the no grammar results was that signs with repetitive or long gestures were often inserted twice for each actual occurrence. In fact, insertions caused

almost as many errors as substitutions. Thus, a sign "shoes" might be recognized as "shoes shoes," which is a viable hypothesis without a language model. However, both problems can be addressed using context training or a statistical or rule-based grammar.

Using context modeling as described before may improve recognition accuracy. While the rule-based grammar explicitly constrained the word order, statistical context modeling would have a similar effect while leaving open the possibility of different sentence structures. In addition, bisine and trisine contexts would help fine-tune the training on the phrase level. However, trisine modeling would not support the tying of the beginning pronoun to the ending pronoun as the grammar does. If task oriented or domain centered sentences were used instead of randomly generated sentences, context modeling and a statistical grammar would improve performance considerably. For example, the random sentence construction allowed "they like pill gray they" which would have a low probability of occurrence in everyday conversation. As such, context modeling would tend to suppress this sentence in recognition unless strong evidence was given for it. In speech and handwriting, a factor of 2 and factor of 4 cut in error rate can be expected for application of contexts and grammars respectively [36]. While ASL does not have the same hierarchy of components as speech and handwriting (letters to words to sentences), the factor of 3 decrease in error when the grammar was applied hints at similar performance increases at the sentence level.

While extending this recognition system to the full 6000 word ASL lexicon would present unforeseen problems, some basic improvements could be made to begin adapting the system to the task:

- Use deltas instead of absolute positions. An alternative is to determine some feature on the subject from which the positions can be measured (for example, the centroid of the subject).
- Add finger and palm tracking information. Exact position information may not be necessary. Some simple starting features may be how many fingers are visible along the contour of the hand and whether the palm is facing up or down.
- Collect appropriate domain or task oriented data and perform context modeling.

These improvements do not address the subject independence issue. Just as in speech, making a system which can understand different subjects with their own variations of the

language involves collecting data from many subjects. Until such a system is tried, it is hard to estimate the number of subjects and the amount of data that would comprise a suitable training database. In general, the training database should span the space of required input, with many examples of each context. Independent recognition often places new requirements on the feature set as well. While the modifications mentioned above may be sufficient initially, the development process is highly empirical.

So far, finger spelling has been ignored. However, incorporating finger spelling into the recognition system is a very interesting problem. Of course, changing the feature vector to address finger information is vital to the problem, but adjusting the context modeling is also of importance. With finger spelling, a closer parallel can be made to speech recognition. Here, trisine context is at a lower level than grammar modeling and will have more of an effect. A point of inquiry would be switching between the different modes of communication. Can trisine context be used across finger spelling and signing? Is it beneficial to switch to a separate mode for finger spelling recognition? Can natural language techniques be applied, and if so, can they also be used to address the spatial positioning issues in ASL? The answers to these questions may be key in creating an unconstrained sign language recognition system.

Chapter 7

Summary and Future Work

An unencumbered way of recognizing a subset of American Sign Language (ASL) through the use of a video camera has been shown. Using hidden Markov models, low error rates were achieved on both the training data (.6%) and an independent test set (3%). Instead of invoking complex models of the hands, simple position and bounding ellipse tracking at 5 frames/sec were used to generate the requisite feature vectors. These feature vectors were then converted to text by the recognition process at five times faster than real time.

The recognition system presented does not purport to be a solution to machine ASL recognition. Issues such as finger spelling and the spatial positioning aspects of ASL were ignored. However, immediate extensions of the system toward this goal have been presented. In fact, a constrained, person dependent, complete lexicon system may be possible using the principles described here. Unfortunately, collecting a suitable training database is prohibitive for the 6000 word task.

An immediate future goal is to train the system on a native ASL signer. Now that a stable training process has been created, collecting this data will be simpler, and the experimental results will provide an initial benchmark for the utility of the system.

Research into removing the constraint of using colored gloves for tracking is underway. While a primitive system has already been created, tracking the hands in front of the face and across the variable background of clothing has proven difficult with the current low resolution, single camera environment. Faster hardware, better lighting and calibration, and more sophisticated algorithms are being applied to create a system which can provide not only the current level of tracking, but also an improved, more stable set of features for

recognition.

While the current system's components each run in parallel in real time, a direct connection between the two modules has not been established. With this addition, immediate sentence translation should be possible. Since recognition time increases with the log of the size of the lexicon, considerable expansion should be possible while still providing immediacy.

Currently the system ignores semantic context given by facial expressions while signing. By adding expression tracking techniques demonstrated by Essa and Darrell [11], this information might be recovered in the current system. An active camera to provide higher resolution "focus of attention" images might also be added to the current apparatus. This would help alleviate the constraint of constant position when signing as well as provide better data for the tracking software.

Bibliography

- [1] Adaptive Optics Associates. *Multi-Trax User Manual*. Adaptive Optics Associates, Cambridge, MA, 1993.
- [2] L. E. Baum. An inequality and associated maximization technique in statistical estimation of probabilistic functions of markov processes. *Inequalities*, 3:1-8, 1972.
- [3] L. E. Baum, T. Petrie, G. Soules, and N. Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *Ann. Math. Stat.*, 41:164-171, 1970.
- [4] Richard A. Bolt and Edward Herranz. Two-handed gesture in multi-modal natural dialog. In *Proceedings of UIST '92, Fifth Annual Symposium on User Interface Software and Technology*, Monterey, CA, 1992.
- [5] L. Campbell. Recognizing classical ballet setps using phase space constraints. Master's thesis, Massachusetts Institute of Technology, 1994.
- [6] J. Cassell, C. Pelachaud, N. Badler, M. Steedman, B. Achorn, T. Becket, B. Douville, S. Prevost, and M. Stone. Animated conversation: Rule-based generation of facial expression, gesture, and spoken intonation for multiple conversational agents. In *Computer Graphics (SIGGRAPH '94 Proceedings)*, pages 413-420, July 1994.
- [7] C. Charayaphan and A.E. Marble. Image processing system for interpreting motion in American Sign Language. *Journal of Biomedical Engineering*, 14:419-425, September 1992.
- [8] T.J. Darrell and A.P. Pentland. Space-time gestures. *IEEE Conf. on Computer Vision and Pattern Rec.*, pages 335-340, 1993.

- [9] Trevor Darrell and Alex Pentland. Attention-driven expression and gesture analysis in an interactive environment. MIT Media Lab Perceptual Computing Group Technical Report No. 312, Massachusetts Institute of Technology, 1994.
- [10] B. Dorner. Hand shape identification and tracking for sign language interpretation. In *IJCAI Workshop on Looking at People*, 1993.
- [11] Irfan Essa, Trevor Darrell, and Alex Pentland. Tracking facial motion. Technical Report 272, MIT Media Lab Vision and Modeling Group, 20 Ames St, Cambridge MA, 1994. To appear in IEEE Workshop on Nonrigid and articulated Motion, Austin TX, Nov 94.
- [12] W. T. Freeman and M. Roth. Orientation histograms for hand gesture recognition. Technical Report 94-03, Mitsubishi Electric Research Labs., 201 Broadway, Cambridge, MA 02139, 1994.
- [13] Yang He and Amlan Kundu. Planar shape classification using hidden markov models. In *Proc. 1991 IEEE Conf. on Computer Vision and Pattern Rec.*, pages 10-15. IEEE Press, 1991.
- [14] David Hogg. Model-based vision: a program to see a walking person. *Image and Vision Computing*, 1(1):5-20, Feb 1983.
- [15] Berthold Horn. *Robot Vision*. MIT Press, New York, 1986.
- [16] X. D. Huang, Y. Ariki, and M. A. Jack. *Hidden Markov Models for Speech Recognition*. Edinburgh University Press, Edinburgh, 1990.
- [17] Tom Humphries, Carol Padden, and Terrence J. O'Rourke. *A Basic Course in American Sign Language*. T. J. Publishers, Inc., Silver Spring, MD, 1980.
- [18] B. H. Juang. Maximum likelihood estimation for mixture multivariate observations of markov chains. *AT&T Technical Journal*, 64:1235-1249, 1985.
- [19] Myron W. Krueger. *Artificial Reality II*. Addison-Wesley Publishing Company, Reading, Massachusetts, 1991.

- [20] F. Kubala, A. Anastasakos, J. Makhoul, L. Nguyen, R. Schwartz, and G. Zavaliagos. Comparative experiments on large vocabulary speech recognition. In *ICASSP 94*, 1994.
- [21] A. Kundu, Y. He, and P. Bahl. Handwritten word recognition: a hidden markov model based approach. volume 22, pages 283–297, 1989.
- [22] L. R. Liporace. Maximum likelihood estimation for multivariate observations of markov sources. *IEEE Trans. Information Theory*, IT-28:729–734, 1982.
- [23] Pattie Maes, Trevor Darrell, Bruce Blumberg, and Alex Pentland. The ALIVE system: full-body interaction with animated autonomous agents. MIT Media Lab Perceptual Computing Group Technical Report No. 257, Massachusetts Institute of Technology, 1994.
- [24] H. Morita, S. Hashimoto, and S. Ohteru. A computer music system that follows a human conductor. *Computer*, 24(7):44–53, July 1991.
- [25] Eadweard Muybridge. *Human and Animal Locomotion*, volume 1-2. Dover Publications, Inc., Mineola, N.Y., 1979.
- [26] R. Nag, K. H. Wong, and F. Fallside. Script recognition using hidden markov models. In *ICASSP 86*, 1986.
- [27] S. Niyogi and E. Adelson. Analyzing gait with spatiotemporal surfaces. Technical Report 290, MIT Media Lab Vision and Modeling Group, 20 Ames St, Cambridge MA, 1994. To appear in IEEE Workshop on Nonrigid and articulated Motion, Austin TX, Nov 94.
- [28] D. H. Parish, G. Sperling, and M.S. Landy. Intelligent temporal subsampling of american sign language using event boundaries. *Journal of Experimental Psychology: Human Perception and Performance*, 16(2):282–294, 1990.
- [29] H. Poizner, U. Bellugi, and V. Lutes-Driscoll. Perception of american sign language in dynamic point-light displays. volume 7, pages 430–440, 1981.
- [30] L. R. Rabiner and B. H. Juang. An introduction to hidden markov models. *IEEE ASSP Magazine*, pages 4–16, January 1996.

- [31] J. M. Rehg and T. Kanade. DigitEyes: vision-based human hand tracking. School of Computer Science Technical Report CMU-CS-93-220, Carnegie Mellon University, December 1993.
- [32] K. Rohr. Towards model-based recognition of human movements in image sequences. *CVGIP: Image Understanding*, 59(1):94-115, Jan 1994.
- [33] Jennifer Schlenzig, Edd Hunter, and Ramesh Jain. Recursive identification of gesture inputers using hidden markov models. In *Proc. of the Second Annual Conference on Applications of Computer Vision*, pages 187-194, 1994.
- [34] Jeffrey Mark Siskind and Quaid Morris. A maximum-likelihood approach to visual event perception. manuscript, 1995.
- [35] G. Sperling, M. Landy, Y. Cohen, and M. Pavel. Intelligible encoding of ASL image sequences at extremely low information rates. *Computer Vision, Graphics, and Image Processing*, 31:335-391, 1985.
- [36] Thad Starner, John Makhoul, Richard Schwartz, and George Chou. On-line cursive handwriting recognition using speech recognition methods. In *ICASSP 94*, 1994.
- [37] W. C. Stokoe, D. C. Casterline, and C. G. Groneberg. *A Dictionary of American Sign Language on Linguistic Principles*. Linstok Press, London, 1976.
- [38] Y. Suenaga, K. Mase, M. Fukumoto, and Y. Watanabe. Human reader: an advanced man-machine interface based on human images and speech. *Systems and Computers in Japan*, 24(2):88-101, 1993.
- [39] T. Takahashi and F. Kishino. Hand gesture coding based on experiments using a hand gesture interface device. *SIGCHI Bulletin*, 23(2):67-73, 1991.
- [40] Shinichi Tamura and Shingo Kawasaki. Recognition of sign language motion images. volume 21, pages 343-353, 1988.
- [41] V. C. Tartter and K. C. Knowlton. Perceiving sign language from an array of 27 moving spots. volume 289, pages 676-678, 1981.

- [42] J. Yamato, J. Ohya, and K. Ishii. Recognizing human action in time-sequential images using hidden markov model. In *Proc. 1992 IEEE Conf. on Computer Vision and Pattern Rec.*, pages 379–385. IEEE Press, 1992.
- [43] S. J. Young. *HTK: Hidden Markov Model Toolkit V1.5*. Cambridge University Engineering Department Speech Group and Entropic Research Laboratories Inc., Washington DC, December 1993.

Visual Recognition of American Sign Language Using Hidden Markov Models

Thad Starner and Alex Pentland
Perceptual Computing Section, The Media Laboratory,
Massachusetts Institute of Technology
Room E15-383, 20 Ames Street, Cambridge MA 02139, USA
E-mail: thad@media.mit.edu, sandy@media.mit.edu

Abstract

Hidden Markov models (HMM's) have been used prominently and successfully in speech recognition and, more recently, in handwriting recognition. Consequently, they seem ideal for visual recognition of complex, structured hand gestures such as are found in sign language. We describe an HMM-based system for recognizing sentence level American Sign Language (ASL) which attains a word accuracy of 99.2% without explicitly modeling the fingers.

1 Introduction

There has been a resurging interest in recognizing human hand gestures. While there are many interesting domains, one of the most structured sets of gestures are those belonging to any of the several sign languages. In sign language, each gesture already has assigned meaning, and strong rules of context and grammar may be applied to make recognition tractable.

To date, most work on sign language recognition has employed expensive wired "datagloves" which the user must wear [19]. In addition, these systems have mostly concentrated on finger signing, where the user spells each word with hand signs corresponding to the letters of the alphabet [3]. However, most signing does not involve finger spelling but, instead, gestures which represent whole words. This allows signed conversations to proceed at about the pace of spoken conversation.

In this paper we describe an extensible system which uses a single color camera to track hands in real time and interprets American Sign Language (ASL) using Hidden Markov Models (HMM). The hand tracking stage of the system does not attempt to produce a fine-grain description of hand shape; studies have shown that such detailed information may not be necessary for humans to interpret sign language [12, 16]. Instead, the tracking process produces only a coarse description of hand shape, orientation, and trajectory. Currently we require that the user wear inexpensive colored gloves to facilitate the hand tracking frame rate and stability. This shape, orientation, and trajectory information is then input to a HMM for recognition of the signed words.

Hidden Markov models have intrinsic properties which make them very attractive for sign language recognition. Explicit segmentation on the word level is not necessary for either training or recognition [18]. Language and context models can be applied on several different levels, and much related development of this technology has already been done by the speech recognition community. Consequently, sign language

recognition seems an ideal machine vision application of HMM technology, offering the benefits of problem scalability, well defined meanings, a pre-determined language model, a large base of users, and immediate applications for a recognizer.

American Sign Language (ASL) is the language of choice for most deaf people in the United States. ASL uses its own grammar instead of borrowing from English. This grammar allows more flexibility in word placement and sometimes uses redundancy for emphasis. Another variant, English Sign Language has more in common with spoken English but is not in widespread use in America. ASL consists of approximately 6000 gestures of common words with finger spelling used to communicate obscure words or proper nouns.

Conversants in ASL may describe a person, place, or thing and then point to a place in space to temporarily store that object for later reference [16]. For the purposes of this experiment, this aspect of ASL will be ignored. Furthermore, in ASL the eyebrows are raised for a question, held normal for a statement, and furrowed for a directive. While there has been work in recognizing facial gestures [4], facial features will not be used to aid recognition in the task addressed.

While the scope of this work is not to create a person independent, full lexicon system for recognizing ASL, a desired attribute of the system is extensibility towards this goal. Another goal is to allow the creation of a real-time system by guaranteeing each separate component (tracking, analysis, and recognition) runs in real-time. This demonstrates the possibility of a commercial product in the future, allows easier experimentation, and simplifies archiving of test data. "Continuous" sign language recognition of full sentences is desired to demonstrate the feasibility of recognizing complicated series of gestures. Of course, a low error rate is also a high priority.

Table 1: ASL Vocabulary Used

<i>part of speech</i>	<i>vocabulary</i>
pronoun	I you he we you(pl) they
verb	want like lose dontwant dontlike love pack hit loan
noun	box car book table paper pants bicycle bottle can wristwatch umbrella coat pencil shoes food magazine fish mouse pill bowl
adjective	red brown black gray yellow

In this recognition system, sentences of the form

“personal pronoun, verb, noun, adjective, (the same) personal pronoun” are to be recognized. This sentence structure emphasizes the need for a distinct grammar for ASL recognition and allows a large variety of meaningful sentences to be randomly generated using words from each class. Table 1 shows the words chosen for each class. Six personal pronouns, nine verbs, twenty nouns, and five adjectives are included making the total lexicon number forty words. The words were chosen by paging through Humphries *et al.* [8] and selecting those which would provide coherent sentences when generating random sentences. At first a naive eye was used to avoid ambiguities in the selected signs, but this was shortly subsumed by the coherency constraint.

2 Machine Sign Language Recognition

Attempts at machine sign language recognition have begun to appear in the literature over the past five years. However, these systems have generally concentrated on isolated signs and small training and test sets. Tamura and Kawasaki demonstrated an early image processing system which could recognize 20 Japanese signs based on matching cheremes [20]. Charayaphan and Marble [2] demonstrated a feature set that could distinguish between the 31 isolated ASL signs in their training set (which also acted as the test set). Takahashi and Kishino in [19] discuss a Dataglove-based system that could recognize 34 of the 46 Japanese kana alphabet gestures (user dependent) using a joint angle and hand orientation coding technique. The test user made each of the 46 gestures 10 times to provide data for principle component and cluster analysis. A separate test set was created from five iterations of the alphabet by the user, with each gesture well separated in time.

3 Previous Use of Hidden Markov Models in Gesture Recognition

While the continuous speech recognition community adopted HMM's many years ago, these techniques are just now entering the vision community. Most early work was limited to handwriting recognition [10, 11]. More recently, He and Kundu [5] report using continuous density HMM's to classify planar shapes. Another early effort by Yamato *et al.* [21] uses discrete HMM's to successfully recognize image sequences of six different tennis strokes among three subjects. This experiment is significant because it used a 25x25 pixel quantized subsampled camera image as a feature vector. Even with such low-level information, the model could learn the set of motions to perform respectable recognition rates. Schlenzig *et al.* [15] also use hidden Markov models for visual gesture recognition. The gestures are limited to “hello,” “good-bye,” and “rotate.” The authors report “intuitively” defining the HMM associated with each gesture and imply that the normal Baum-Welch re-estimation method was not implemented. However, this study shows the continuous gesture recognition capabilities of HMM's by recognizing gesture sequences.

4 Hidden Markov Modeling

While a substantial body of literature exists on HMM technology [1, 7, 13, 22], this section briefly outlines a traditional discussion of the algorithms. After outlining the fundamental theory in training and testing of a discrete HMM, this result is then generalized to the continuous density case used in the experiments. For broader discussion of the topic, [7, 17] are recommended.

A time domain process demonstrates a Markov property if the conditional probability density of the current event, given all present and past events, depends only on the j th most recent events. If the current event depends solely on the most recent past event, then the process is a first order Markov process. While the order of words in American Sign Language is not truly a first order Markov process, it is a useful assumption when considering the positions and orientations of the hands of the signer through time.

The initial topology for an HMM can be determined by estimating how many different states are involved in specifying a sign. Fine tuning this topology can be performed empirically. While different topologies can be specified for each sign, a four state HMM with skip transitions was determined to be sufficient for this task (Figure 1).

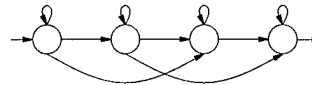


Figure 1: The four state HMM used for recognition.

There are three key problems in HMM use. These are the evaluation, estimation, and the decoding problems. The evaluation problem is that given an observation sequence and a model, what is the probability that the observed sequence was generated by the model ($Pr(\mathbf{O}|\lambda)$) (notational style from [7])? If this can be evaluated for all competing models for an observation sequence, then the model with the highest probability can be chosen for recognition.

$Pr(\mathbf{O}|\lambda)$ can be calculated several ways. The naive way is to sum the probability over all the possible state sequences in a model for the observation sequence:

$$Pr(\mathbf{O}|\lambda) = \sum_{\text{all } S} \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t}(O_t)$$

However, this method is exponential in time, so the more efficient forward-backward algorithm is used in practice. The following algorithm defines the forward variable α and uses it to generate $Pr(\mathbf{O}|\lambda)$ (π are the initial state probabilities, a are the state transition probabilities, and b are the output probabilities).

- $\alpha_1(i) = \pi_i b_i(O_1)$, for all states i (if $i \in S_I$, $\pi_i = \frac{1}{n_I}$; otherwise $\pi_i = 0$)
- Calculating $\alpha_t(j)$ along the time axis, for $t = 2, \dots, T$, and all states j , compute

$$\alpha_t(j) = \left[\sum_i \alpha_{t-1}(i) a_{ij} \right] b_j(O_t)$$

- Final probability is given by

$$Pr(\mathbf{O}|\lambda) = \sum_{i \in S_F} \alpha_T(i)$$

The first step initializes the forward variable with the initial probability for all states, while the second step inductively steps the forward variable through time. The final step gives the desired result $Pr(\mathbf{O}|\lambda)$, and it can be shown by constructing a lattice of states and transitions through time that the computation is only order $O(N^2T)$. The backward algorithm, using a process similar to the above, can also be used to compute $Pr(\mathbf{O}|\lambda)$ and defines the convenience variable β .

The estimation problem concerns how to adjust λ to maximize $Pr(\mathbf{O}|\lambda)$ given an observation sequence \mathbf{O} . Given an initial model, which can have flat probabilities, the forward-backward algorithm allows us to evaluate this probability. All that remains is to find a method to improve the initial model. Unfortunately, an analytical solution is not known, but an iterative technique can be employed.

Using the actual evidence from the training data, a new estimate for the respective output probability can be assigned:

$$\bar{b}_j(k) = \frac{\sum_{t \in O_j = v_k} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$

where $\gamma_t(i)$ is defined as the posterior probability of being in state i at time t given the observation sequence and the model. Similarly, the evidence can be used to develop a new estimate of the probability of a state transition (\bar{a}_{ij}) and initial state probabilities ($\bar{\pi}_i$).

Thus all the components of model (λ) can be re-estimated. Since either the forward or backward algorithm can be used to evaluate $Pr(\mathbf{O}|\lambda)$ versus the previous estimation, the above technique can be used iteratively to converge the model to some limit. While the technique described only handles a single observation sequence, it is easy to extend to a set of observation sequences. A more formal discussion can be found in [1, 7, 22].

While the estimation and evaluation processes described above are sufficient for the development of an HMM system, the Viterbi algorithm provides a quick means of evaluating a set of HMM's in practice as well as providing a solution for the decoding problem. In decoding, the goal is to recover the state sequence given an observation sequence. The Viterbi algorithm can be viewed as a special form of the forward-backward algorithm where only the maximum path at each time step is taken instead of all paths. This optimization reduces computational load and additionally allows the recovery of the most likely state sequence. The steps to the Viterbi are

- Initialization. For all states i , $\delta_1(i) = \pi_i b_i(O_1)$; $\psi_1(i) = 0$
- Recursion. From $t = 2$ to T and for all states j , $\delta_t(j) = \text{Max}_i[\delta_{t-1}(i) a_{ij}] b_j(O_t)$; $\psi_t(j) = \text{argmax}_i[\delta_{t-1}(i) a_{ij}]$
- Termination. $P = \text{Max}_{s \in S_F} [\delta_T(s)]$; $s_T = \text{argmax}_{s \in S_F} [\delta_T(s)]$

- Recovering the state sequence. From $t = T - 1$ to 1, $s_t = \psi_{t+1}(s_{t+1})$

In many HMM system implementations, the Viterbi algorithm is used for evaluation at recognition time. Note that since Viterbi only guarantees the maximum of $Pr(\mathbf{O}, S|\lambda)$ over all state sequences S (as a result of the first order Markov assumption) instead of the *sum* over all possible state sequences, the resultant scores are only an approximation. However, [13] shows that this is often sufficient.

So far the discussion has assumed some sort of quantization of feature vectors into classes. However, instead of using vector quantization, the actual probability densities for the features may be used. Baum-Welch, Viterbi, and the forward-backward algorithms can be modified to handle a variety of characteristic densities [9]. In this context, however, the densities will be assumed to be Gaussian. Specifically,

$$b_j(O_t) = \frac{1}{\sqrt{(2\pi)^n |\sigma_j|}} e^{\frac{1}{2}(O_t - \mu_j)' \sigma_j^{-1} (O_t - \mu_j)}$$

Initial estimations of μ and σ may be calculated by dividing the evidence evenly among the states of the model and calculating the mean and variance in the normal way. Whereas flat densities were used for the initialization step before, the evidence is used here. Now all that is needed is a way to provide new estimates for the output probability. We wish to weight the influence of a particular observation for each state based on the likelihood of that observation occurring in that state. Adapting the solution from the discrete case yields

$$\bar{\mu}_j = \frac{\sum_{t=1}^T \gamma_t(j) O_t}{\sum_{t=1}^T \gamma_t(j)}$$

and

$$\bar{\sigma}_j = \frac{\sum_{t=1}^T \gamma_t(j) (O_t - \bar{\mu}_j)(O_t - \bar{\mu}_j)^t}{\sum_{t=1}^T \gamma_t(j)}$$

For convenience, μ_j is used to calculate $\bar{\sigma}_j$ instead of the re-estimated $\bar{\mu}_j$. While this is not strictly proper, the values are approximately equal in contiguous iterations [7] and seem not to make an empirical difference [22]. Since only one stream of data is being used and only one mixture (Gaussian density) is being assumed, the algorithms above can proceed normally, incorporating these changes for the continuous density case.

5 Recovering Hands from Video

Previous systems have shown that, given some constraints, relatively detailed models of the hand can be recovered from video images [3, 14]. However, many of these constraints conflict with recognizing ASL in a natural context, either by requiring simple, unchanging backgrounds (unlike clothing), not allowing occlusion, requiring carefully labelled gloves, or being difficult to run in real time.

Since real-time recognition is a goal in this project, several compromises were made. The subject wears distinctly colored gloves on each hand (a yellow glove

for the right hand and an orange glove for the left) and sits in a chair before the camera. Figure 2 shows the view from the camera's perspective and gives an impression of the quality of video that is used. Color NTSC composite video is captured and analyzed at a constant 5 frames per second at 320 by 243 pixel resolution on a Silicon Graphics Indigo 2 with Galileo video board. To find each hand initially, the algorithm scans the image until it finds a pixel of the appropriate color. Given this pixel as a seed, the region is grown by checking the eight nearest neighbors for the appropriate color. Each pixel checked is considered to be part of the hand. This, in effect, performs a simple morphological dilation upon the resultant image that helps to prevent edge and lighting aberrations. The centroid is calculated as a by-product of the growing step and is stored as the seed for the next frame. Given the resultant bitmap and centroid, second moment analysis is performed as described earlier.

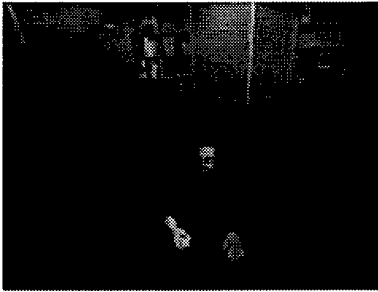


Figure 2: View from the tracking camera.

6 Feature Extraction

Previous experience has shown that it is often best to start simple and evolve a feature set [18]. Since finger spelling is not allowed and there are few ambiguities in the test vocabulary based on individual finger motion, a relatively coarse tracking system may be used. Based on previous work, it was assumed that a system could be designed to separate the hands from the rest of the scene. Traditional vision algorithms could then be applied to the binarized result. Aside from the position of the hands, some concept of the shape of the hand and the angle of the hand relative to horizontal seemed necessary. Thus, an eight element feature vector consisting of each hand's x and y position, angle of axis of least inertia, and eccentricity of bounding ellipse was chosen. The eccentricity of the bounded ellipse was found by determining the ratio of the square roots of the eigenvalues that correspond to the matrix

$$\begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix}$$

where a , b , and c are defined as

$$a = \int \int_{I'} (x')^2 dx' dy'$$

$$b = \int \int_{I'} x' y' dx' dy'$$

$$c = \int \int_{I'} (y')^2 dx' dy'$$

(x' and y' are the x and y coordinates normalized to the centroid)

The axis of least inertia is then determined by the major axis of the bounding ellipse, which corresponds to the primary eigenvector of the matrix [6]. Note that this leaves a 180 degree ambiguity in the angle of the ellipses. To address this problem, the angles were only allowed to range from -90 to +90 degrees.

7 Training an HMM network

When using HMM's to recognize strings of data, such as continuous speech, cursive handwriting, or ASL sentences, several methods can be used to bring context to bear in training and recognition. A simple context modeling method is embedded training. While initial training of the models might rely on manual segmentation or, in this case, evenly dividing the evidence among the models, embedded training trains the models *in situ* and allows boundaries to shift through a probabilistic entry into the initial states of each model [22].

Generally, a sign can be affected by both the sign in front of it and the sign behind it. In speech, this is called "co-articulation." While this can confuse systems trying to recognize isolated signs, the context information can be used to aid recognition. For example, if two signs are often seen together, recognizing the two signs as one group may be beneficial.

A final use of context is on the word level. Statistical grammars relating the probability of the co-occurrence of two or more words can be used to weight the recognition process. Grammars that associate two words are called bigrams, whereas grammars that associate three words are called trigrams. Rule-based grammars can also be used to aid recognition.

8 Experimentation

The handtracking system as described earlier worked well. Occasionally tracking would be lost (generating error values of 0) due to lighting effects, but recovery was fast enough (within a frame) so that this was not a problem. A 5 frame/sec rate was maintained within a tolerance of a few milliseconds. However, frames were deleted where tracking of one or both hands was lost. Thus, a constant data rate was not guaranteed.

Of the 500 sentences collected, six had to be thrown out due to subject error or outlier signs. In general, each sign ranged from approximately 1 to 3 seconds in length. No intentional pauses were placed between signs within a sentence, but the sentences themselves were distinct.

Initial estimates for the means and variances of the output probabilities were provided by iteratively using Viterbi alignment on the training data (after initially dividing the evidence equally among the words in the sentence) and then recomputing the means and variances by pooling the vectors in each segment. Entropic's Hidden Markov Model ToolKit (HTK) is used as a basis for this step and all other HMM modeling and training tasks. The results from

Table 2: Word accuracy

	<i>on training</i>	<i>on indep. test set</i>
grammar	99.5%	99.2%
no gram.	92.0% (97% corr.) (D=9, S=67, I=121, N=2470)	91.3% (97% corr.) (D=1, S=16, I=26, N=495)

the initial alignment program are fed into a Baum-Welch re-estimator, whose estimates are, in turn, refined in embedded training which ignores any initial segmentation. For recognition, HTK's Viterbi recognizer is used both with and without a strong grammar based on the known form of the sentences. The actual recognition step runs at a rate five times faster than real time. Contexts are not used, since a similar effect could be achieved with the strong grammar given this data set.

Word recognition results are shown in Table 2. When testing on training, all 494 sentences were used for both the test and train sets. For the fair test, the sentences were divided into a set of 395 training sentences and a set of 99 independent test sentences. The 99 test sentences were not used for any portion of the training. Given the strong grammar (pronoun, verb, noun, adjective, pronoun), insertion and deletion errors were not possible since the number and class of words allowed is known. Thus, all errors are substitutions when the grammar is used (and accuracy is equivalent to percent correct). However, without the grammar, the recognizer is allowed to match the observation vectors with any number of the 40 vocabulary words in any order. Thus, deletion (D), insertion (I), and substitution (S) errors are possible. The absolute number of errors of each type are listed in Table 2. The accuracy measure is calculated by subtracting the number of insertion errors from the number of correct labels and dividing by the total number of signs. Note that, since all errors are accounted against the accuracy rate, it is possible to get large negative accuracies (and corresponding error rates of over 100%). Most insertion errors occurred at signs with repetitive motion.

9 Analysis and Discussion

While these results are far from being sufficient to claim a "working system" for ASL recognition, they do show that this approach is promising. The high recognition rate on the training data indicates that the HMM topologies are sound and that the models are converging. Even so, the remaining 6.5% error rate on the "no grammar" case (error rates will be based on accuracy measures) indicates that some fine tuning on the feature set and model is in order. The 0.8% error rate on the independent test set shows that the models are generalizing well. However, a close look at the text produced by the recognition process shows some of the limitations of the feature set. Since the raw positions of the hands were used, the system was trained to expect certain gestures in certain locations. When this varied due to subject seating position or arm placement, the system could

become confused. A simple fix to this problem would be to use position deltas in the feature vector instead.

Examining the errors made when no grammar was used shows the importance of finger position information. Signs like "pack," "car," and "gray" have very similar motions. In fact, the main difference between "pack" and "car" is that the fingers are pointed down for the former and clenched in the latter. Since this information was not available in the model, confusion could occur. While recovering general and/or specific finger position may be difficult in real time in the current testing area, simple palm orientation could be used for discrimination. In the current system, a simple implementation would be to paint the back of the gloves a different color than the palm.

A more interesting problem with the no grammar results was that signs with repetitive or long gestures were often inserted twice for each actual occurrence. In fact, insertions caused more errors than substitutions. Thus, a sign "shoes" might be recognized as "shoes shoes," which is a viable hypothesis without a language model. However, both problems can be addressed using context training or a statistical or rule-based grammar.

Using context modeling as described before may improve recognition accuracy. While the rule-based grammar explicitly constrained the word order, statistical context modeling would have a similar effect while leaving open the possibility of different sentence structures. In addition, bisine (two sign) and trisine (three sign) contexts would help fine-tune the training on the phrase level. However, trisine modeling would not support the tying of the beginning pronoun to the ending pronoun as does the grammar. If task oriented or domain centered sentences were used instead of randomly generated sentences, context modeling and a statistical grammar would improve performance considerably. For example, the random sentence construction allowed "they like pill gray they" which would have a low probability of occurrence in everyday conversation. As such, context modeling would tend to suppress this sentence in recognition unless strong evidence was given for it.

While extending this recognition system to the full 6000 word ASL lexicon would present many problems, some basic improvements could be made in order to begin adapting the system to the task:

- Use deltas instead of absolute positions. An alternative is to determine some feature on the subject from which the positions can be measured (for example, the centroid of the subject).
- Add finger and palm tracking information. This may be as simple as how many fingers are visible along the contour of the hand and whether the palm is facing up or down.
- Collect appropriate domain or task oriented data and perform context modeling.

These improvements do not address the subject independence issue. Just as in speech, making a system which can understand different subjects with their own variations of the language involves collecting data from many subjects. Until such a system is tried, it is hard to estimate the number of subjects and the amount of data that would comprise

a suitable training database. Independent recognition often places new requirements on the feature set as well. While the modifications mentioned above may be sufficient initially, the development process is highly empirical.

So far, finger spelling has been ignored. However, incorporating finger spelling into the recognition system is a very interesting problem. Of course, changing the feature vector to address finger information is vital to the problem, but adjusting the context modeling is also of importance. With finger spelling, a closer parallel can be made to speech recognition. Here, triseme context is at a lower level than grammar modeling and will have more of an effect. A point of inquiry would be switching between the different modes of communication. Can triseme context be used across finger spelling and signing? Is it beneficial to switch to a separate mode for finger spelling recognition? Can natural language techniques be applied, and if so, can they also be used to address the spatial positioning issues in ASL? The answers to these questions may be key in creating an unconstrained sign language recognition system.

10 Conclusion

We have shown an unencumbered way of recognizing American Sign Language (ASL) through the use of a video camera. Through use of hidden Markov models low error rates were achieved on both the training set and an independent test set without invoking complex models of the hands. With a larger training set and context modeling, lower error rates are expected and generalization to a freer, person-independent ASL recognition system should be obtainable.

References

- [1] L. Baum. An inequality and associated maximization technique in statistical estimation of probabilistic functions of markov processes. *Inequalities*, 3:1-8, 1972.
- [2] C. Charayaphan and A. Marble. Image processing system for interpreting motion in American Sign Language. *Journal of Biomedical Engineering*, 14:419-425, Sept. 1992.
- [3] B. Dörner. Hand shape identification and tracking for sign language interpretation. In *IJCAI Workshop on Looking at People*, 1993.
- [4] I. Essa, T. Darrell, and A. Pentland. Tracking facial motion. IEEE Workshop on Nonrigid and articulated Motion, Austin TX, Nov 94.
- [5] Y. He and A. Kundu. Planar shape classification using hidden markov models. In *Proc. 1991 IEEE Conf. on Comp. Vision and Pat. Rec.*, p. 10-15. IEEE Press, 1991.
- [6] B. Horn. *Robot Vision*. MIT Press, NY, 1986.
- [7] X. Huang, Y. Ariki, and M. Jack. *Hidden Markov Models for Speech Recognition*. Edinburgh Univ. Press, Edinburgh, 1990.
- [8] T. Humphries, C. Padden, and T. O'Rourke. *A Basic Course in American Sign Language*. T. J. Publ., Inc., Silver Spring, MD, 1980.
- [9] B. Juang. Maximum likelihood estimation for mixture multivariate observations of markov chains. *AT&T Technical Journal*, 64:1235-1249, 1985.
- [10] A. Kundu, Y. He, and P. Bahl. Handwritten word recognition: a hidden markov model based approach. 22: 283-297, 1989.
- [11] R. Nag, K. Wong, and F. Fallside. Script recognition using hidden markov models. In *ICASSP 86*, 1986.
- [12] H. Poizner, U. Bellugi, and V. Lutes-Driscoll. Perception of american sign language in dynamic point-light displays. 7: 430-440, 1981.
- [13] L. Rabiner and B. Juang. An introduction to hidden markov models. *IEEE ASSP Magazine*, p. 4-16, Jan. 1996.
- [14] J. Rehg and T. Kanade. DigitEyes: vision-based human hand tracking. School of Computer Science Technical Report CMU-CS-93-220, Carnegie Mellon Univ., Dec. 1993.
- [15] J. Schlenzig, E. Hunter, and R. Jain. Recursive identification of gesture inputers using hidden markov models. In *Proc. of the Second Annual Conf. on Appl. of Comp. Vision*, p. 187-194, 1994.
- [16] G. Sperling, M. Landy, Y. Cohen, and M. Pavel. Intelligible encoding of ASL image sequences at extremely low information rates. *Comp. Vision, Graphics, and Image Proc.*, 31:335-391, 1985.
- [17] T. Starner. Visual Recognition of American Sign Language Using Hidden Markov Models. Master's thesis, MIT Media Laboratory, Feb. 1995.
- [18] T. Starner, J. Makhoul, R. Schwartz, and G. Chou. On-line cursive handwriting recognition using speech recognition methods. In *ICASSP*, 1994.
- [19] T. Takahashi and F. Kishino. Hand gesture coding based on experiments using a hand gesture interface device. *SIGCHI Bulletin*, 23(2):67-73, 1991.
- [20] S. Tamura and S. Kawasaki. Recognition of sign language motion images. volume 21, p. 343-353, 1988.
- [21] J. Yamato, J. Ohya, and K. Ishii. Recognizing human action in time-sequential images using hidden markov model. In *Proc. 1992 IEEE Conf. on Comp. Vision and Pat. Rec.*, p. 379-385. IEEE Press, 1992.
- [22] S. Young. *HTK: Hidden Markov Model Toolkit V1.5*. Cambridge Univ. Eng. Dept. Speech Group and Entropic Research Lab. Inc., Washington DC, Dec. 1993.

BABYL OPTIONS:

Version: 5

Labels:

Note: This is the header of an rmail file.

Note: If you are seeing it in rmail,

Note: it means the file has no messages in it.

Human Powered Wearable Computing

Thad Starner

Perceptual Computing Section, The Media Laboratory,
Massachusetts Institute of Technology
Room E15-383, 20 Ames Street, Cambridge MA 02139, USA
E-mail: thad@media.mit.edu

Abstract

Batteries add size and weight to present day notebook and wearable computers. This paper explores the possibility of removing this impediment by using the operator's natural body functions to generate power for his computer. Power generation through leg motion is presented in depth as well as a survey of less practical methods such as generation by breath or blood pressure, body heat, and finger and limb motion.

1 Introduction

Many of today's portable devices such as notebook computers, PDA's, and radios require extra battery packs or frequent rechargings for extended use. However, the human body itself is a tremendous source of energy, using between 70,000 and 1,400,000 calories per hour depending on the activity (see Table 1). In fact, trained athletes can expend close to 9.5 million calories/hr for short bursts [15], and the minimum amount of energy expended by a human is

$$(70,000 \text{ calories/hr}) \left(\frac{4.19 \text{ J}}{\text{calorie}} \right) \left(\frac{1 \text{ hr}}{3600 \text{ sec}} \right) = 81 \text{ W}$$

If only a small fraction of this power could be harnessed conveniently and unobtrusively, batteries *per se* could be eliminated. However, difficulties arise from the acquisition, regulation, and distribution of such power.

Recent technology makes these tasks easier. Computers are now small enough to disappear into the user's clothing. With such small devices, the main power consumers, namely the CPU and storage, could be located near the implemented power source. Interface devices, such as keyboards, displays, and speakers have limitations on their placement on the body. Such devices may generate their own power, share in a power distribution system with the main generator (wired or wireless), or use extremely long lasting batteries. Thus, depending on the user interface desired, wires may not be needed in powering a wearable computer.

In the following sections, power generation from breath, body heat, blood transport, arm motion, typing, and walking are discussed. While some of these ideas are fanciful, each has its own peculiar benefits and may be applied to other domains such as medical systems, general consumer electronics, and user interface sensors. More attention is given to typing and walking since these processes currently seem promising sources of power.

Table 1: Human energy expenditures for selected activities (note that the traditional dietetic Calorie is actually 1,000 calories) [15].

activity	Kcal/hr
Sleeping	70
Lying quietly	80
Sitting	100
Standing at ease	110
Conversation	110
Eating meal	110
Strolling	140
Driving car	140
Playing violin or piano	140
Housekeeping	150
Carpentry	230
Hiking, 4 mph	350
Swimming	500
Mountain climbing	600
Long distance run	900
Sprinting	1400

2 Breath

An average person of 68 kg (150 lbs.) has an approximate air intake of 30 liters a minute [15]. However, available breath pressure is only 2% above atmospheric pressure [7, 17]. Also, air intake has limitations in generating power due to human physiological constraints [17]. Thus, the available energy is

$$W = p\Delta V =$$

$$\left(\frac{0.02 \times 10^5 \text{ kg}}{\text{msec}^2} \right) \left(\frac{30 \text{ l}}{\text{min}} \right) \left(\frac{1 \text{ min}}{60 \text{ sec}} \right) \left(\frac{1 \text{ m}^3}{1000 \text{ l}} \right) = 1.0 \text{ W}$$

. During sleep the breath rate, and therefore the available power, may drop in half, while increased activity may increase the breath rate. Forcing an elevated breath pressure with an aircraft-style pressure mask can increase the effective power by a factor of 2.5, but it causes significant stress on the subject [9].

Interfering with breathing involves breath masks which encumber the user. For some professionals such as military aircraft pilots, astronauts, or handlers of hazardous materials, such masks are already in place. However, the efficiency of a turbine and generator combination is only about 40% [10], and any attempt to tap this energy source would provide additional load on the user. Thus, the 0.4W of recoverable energy has to be weighed against the other more convenient methods discussed in the following sections.

Another way to generate power from breathing is to fasten a tight band around the chest of the user. From empirical measurements, there is a 2.5 cm change in chest circumference when breathing normally, and up to a 5 cm change while breathing deeply. A large amount of force can be maintained over this interval. Assuming 10 breaths per minute and an ambitious 100 N force applied over the 0.05 m distance, the maximum power that can be generated is

$$(100N)(0.05m)\left(\frac{10\text{breaths}}{\text{min}}\right)\left(\frac{1\text{min}}{60\text{sec}}\right) = 0.83W$$

A ratchet and flywheel attached to an elastic band around the chest might be used to recover this energy. Since the mechanical to electrical efficiency of generators is very good, most of the losses will come from friction within the mechanism itself. Thus, with careful design, most of this energy may be recoverable. However, 0.83 W is still a relatively small amount of power for the inconvenience.

3 Body Heat

Since the human body's prime means of eliminating waste energy is heat, a natural idea is to try to harness this process. However, Carnot efficiency puts an upper limit on how well this waste heat can be recovered. Assuming body temperature and a relatively low room temperature (68° F), the Carnot efficiency is

$$\frac{T_H - T_L}{T_H} = \frac{(310K - 293K)}{310K} = 5.5\%$$

In a warmer environment (80° F) the Carnot efficiency is

$$\frac{T_H - T_L}{T_H} = \frac{(310K - 300K)}{310K} = 3.2\%$$

Referring to the value in Table 1 for sitting, a total of 116 W of power is available. Using a Carnot engine to model the recoverable energy yields 3.7-6.4 W. For more extreme temperature differences, higher efficiencies may be achieved; however, robbing the user of heat in adverse environmental temperatures is not practical.

The above efficiencies assume that all of the heat radiated by the human body is captured and perfectly transformed into power. However, such a system would encapsulate the user in something similar to a wetsuit. Knowing that this is unacceptable in most situations, localized heat engines can be used. Unfortunately, the low temperature necessary for high efficiency causes the body to restrict blood flow to the limbs [9]. When the skin surface encounters cold air, a rapid constriction of the blood vessels in the skin allows the skin temperature to approach the temperature of the air so that heat exchange is reduced. This, in turn, allows the heat to be maintained in the torso. Thus, any heat exchanger should be located on the torso or head. However, even under the best of conditions (basal, non-sweating), evaporative heat loss accounts for 25% of the total heat loss. Water diffusing through the skin, sweat glands in the palms and soles keeping the skin pliable, and

expulsion of water-saturated air from the lungs accounts for this "insensible perspiration." [9] Thus, even assuming all heat exchange is performed by the "core" areas, the maximum energy available, without trying to reclaim heat expended by the latent heat of vaporization, is 2.8-4.8 W.

Placing a heat engine on the torso is inconvenient with clothing. In order to recover body heat efficiently, some sort of "wetsuit" is necessary. However, covering the entire torso in a wetsuit apparatus is uncomfortable, hard to remove in many situations, and has physiological limitations for extended use. On the other hand, the neck offers a good location for a tight seal, access to major centers of blood flow, and ease of removal by the user. The neck is approximately 1/15 of the surface area of this "core" region. Thus, a maximum of 0.20-0.32 W can be recovered conveniently by such a neck brace. The head may also be a convenient heat source for some applications where protective hoods are already in place. The surface area of the head is about 3 times that of the neck and can provide 0.60-0.96 W of power given optimal conversion. Even so, trying to increase the efficiency of heat loss from the body is detrimental to the user's health over extended time.

4 Blood Pressure

While powering electronics from blood pressure may seem impractical, in actuality the numbers are quite surprising. Assuming an average blood pressure of 100 mm of Hg (normal desired blood pressure is 120/80 above atmospheric pressure), a rest heart rate of 60 beats per minute, and a heart stroke volume of 70 ml passing through the aorta [4]:

$$(100\text{mmHg})\left(\frac{1.013 \times 10^5 \text{ kg/msec}^2}{760\text{mmHg}}\right) = \frac{1.33 \times 10^4 \text{ kg}}{\text{msec}^2}$$

$$\left(\frac{1.33 \times 10^4 \text{ kg}}{\text{msec}^2}\right)\left(\frac{60\text{beats}}{\text{min}}\right)\left(\frac{1\text{min}}{60\text{sec}}\right)\left(\frac{.07\text{l}}{\text{beat}}\right)\left(\frac{1\text{m}^3}{1000\text{l}}\right) = 0.93W$$

While this rate can easily double when running, harnessing this power is difficult. Adding a turbine to the system would increase the load on the heart, perhaps dangerously so. However, even if 2% of this power is harnessed, low power microprocessors and sensors could run. Thus, self-powering medical sensors and prostheses could be created.

5 Notebook computer power

Current notebook computers offer a unique method of generating power. Simply opening the computer may supply power. However, this one action needs to provide power for the entire session (otherwise, users would be forced to flap their computers open and closed during their session). From some simple testing, the maximum force that is reasonable for a user to exert in opening his machine is

$$(20\text{lbs})\left(\frac{2.2\text{kg}}{\text{lb}}\right)\left(\frac{9.8\text{m}}{\text{sec}^2}\right) = 430N$$

Assuming a maximum of 0.5m of motion in opening the computer,

$$(430N)(0.5m) = 215J$$

Thus, for a 10 minute use, a maximum of

$$\frac{215J}{600sec} = 0.36W$$

would be available. At an hour's usage, the rate drops to 0.06 W. For many current applications, these power rates are inadequate.

6 Finger motion

Keyboards will continue to be a major interface for computers into the next decade [22]. As such, typing may provide a useful source of energy. On a one-handed chording keyboard (HandyKey's Twiddler), 130 grams is necessary to depress a key the required 1 mm to register. Thus,

$$\left(\frac{.13kg}{keypress}\right)\left(\frac{9.8m}{sec^2}\right)(0.001m) = \frac{1.3mJ}{keypress}$$

is necessary to type. Assuming a moderately skilled typist (40 wpm), and taking into account multiple keypress combinations, an average of

$$\left(\frac{1.3mJ}{keypress}\right)\left(\frac{5.3keypresses}{sec}\right) = 6.9mW$$

of power is generated. A fast QWERTY typist (90 wpm) presses about 7.5 keys per second. A typical keyboard requires about 40-50 grams to depress a key by the 0.5 cm necessary to register a keystroke (measured on a DEC PC 433 DX LP). Thus, a QWERTY typist may generate

$$\left(\frac{.05kg}{keypress}\right)\left(\frac{9.8m}{sec^2}\right)(0.005m)\left(\frac{7.5keypresses}{sec}\right) = 19mW$$

of power. Unfortunately, neither method provides enough continuous power to sustain a wearable computer, especially considering that the user will not be continuously typing on the keyboard. However, there may be enough energy in each keystroke for each key to "announce" its character to a nearby receiver [11]. For example, the keyboard may have a permanent magnet in its base. Each key would then have an embedded coil that would generate a current when the key was moved. Another possibility is to use piezoelectric material which bends at each keystroke to generate power (around 12% efficiency). Thus, a wearable, wireless keyboard may be possible.

7 Upper limb motion

Comparing the amount of energy expended playing a violin or housekeeping versus the amount spent standing at ease using Table 1, up to 30 Kcal/hr or

$$\frac{30Kcal}{hr}\left(\frac{4.19J}{calorie}\right)\left(\frac{1hr}{3600sec}\right) = 35W$$

of power is put into moving the upper limbs. Using the results from [3] on a 58.8 kg man, the lower arm plus hand masses 1.4 kg, the upper arm 1.8 kg, and the whole arm 3.2 kg. The distance through which the center of mass of the lower arm moves for a full bicep curl is 0.335 m, while raising the arm fully over the head moves its center of mass about 0.725 m. Empirically, bicep curls can be performed at a maximum rate of 2 curls/sec, and lifting the arms above

the head at 1.3 lifts/sec. Thus, the maximum power generated by bicep curls is

$$(1.8kg)\left(\frac{9.8m}{sec^2}\right)(0.335m)\left(\frac{2curls}{sec}\right)(2arms) = 24W$$

while the maximum power generated by arm lifts is

$$(3.2kg)\left(\frac{9.8m}{sec^2}\right)(0.725m)\left(\frac{1.3lifts}{sec}\right)(2arms) = 60W$$

Obviously, housekeeping and violin playing do not involve as much strenuous exercise as these experiments. However, these calculations do show that there is plenty of energy to be recovered from an active user. The task at hand, though, is to recover energy without burdening the user. A much more reasonable number, even for a user in an enthusiastic gestural conversation, is to divide the bicep curl power by a factor of 8. Thus, the user might make one arm gesture every two seconds. This, then, creates a total of 3 W of power and is no longer in the domain of recoverable energy without encumbering the user. By double loading the user and mounting a pulley system on the belt, 1.5 W might be recovered (assuming 50% efficiency from loss due to friction and the small parts involved), but the system would be extremely inconvenient.

However, a less encumbering system would involve mounting pulley systems in the elbows of a jacket. The take-up reel of the pulley system could be spring-loaded so as to counter-balance the weight of the user's arm. Thus, the system would generate power from the change in potential energy of the arm on the downstroke and not require additional energy by the user on the upstroke. The energy generation system, the CPU, and the interface devices could be incorporated into the jacket. Thus, the user would simply don his jacket to use his computer. However, any pulley or piston generation system would involve many inconvenient moving parts and the addition of significant mass to the user. A more innovative solution would be to use piezoelectric materials at the joints which would generate charge from the movement of the user. Thus, no moving parts *per se* would be involved, and the jacket would not be significantly heavier than a normal jacket. In order to take a closer look at such a solution, some more background is necessary.

7.1 Piezoelectric materials

Piezoelectric materials create electrical charge when mechanically stressed. Among the natural materials with this property are quartz, human skin, and human bone. Table 2 shows properties of common piezoelectric materials: polyvinylidene fluoride (PVDF) and lead zirconate titanate (PZT). Several advanced treatments of piezoelectricity and data sheets are included in the references for convenience [2, 5, 6, 14, 19].

The coupling constant shown in Table 2 is the efficiency with which a material converts mechanical to electrical energy. The subscripts on some of the constants indicate the direction or mode of the mechanical/electrical interactions (see 1 from [18]). The 31 mode, commonly used in industry, applies an electrical charge in the 3 axis resulting in bending in the

Table 2: Piezoelectric characteristics of PVDF and PZT (adapted from [6, 2, 18]).

Property	Units	PVDF	PZT
Density	$\frac{10^3 kg}{m^3}$	1.78	7.6
Relative permittivity	$\frac{\epsilon}{\epsilon_0}$	12	1700
Elastic modulus	$\frac{10^{10} N}{m}$	0.3	4.9
Piezoelectric constant	$\frac{10^{-12} C}{N}$	$d_{31}=20$ $d_{33}=30$	$d_{31}=180$ $d_{33}=360$
Coupling constant	$\frac{CV}{Nm}$.11	$k_{31}=0.35$ $k_{33}=0.69$

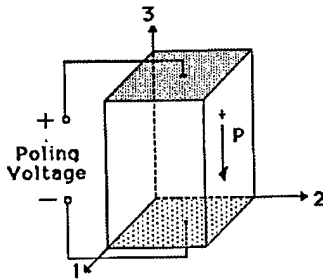


Figure 1: Definition of axes for piezoelectric materials. Note that the electrodes are mounted on the 3 axis [18].

1 axis (d_{31}). Conversely, mechanical bending in the 1 axis will produce an electrical charge in the 3 axis.

The most efficient energy conversion, as indicated by the coupling constants in Table 2, comes from compressing PZT (d_{33}). Even so, the amount of effective power that could be transferred this way is minimal since compression follows the formula

$$\Delta H = \frac{FH}{AY}$$

where F is force, H is the unloaded height, A is the area over which the force is applied, and Y is the elastic modulus. The elastic modulus for PZT is $4.9 \times 10^{10} N/m^2$. Thus, it would take an incredible force to compress the material a small amount. Since energy is defined as force through distance, the effective energy generated through human-powered compression of PZT would be vanishingly small, even with perfect conversion.

On the other hand, bending a piece of piezoelectric material along its 31 mode is much easier. Because it is brittle, PZT does not have much range of motion in this direction. Maximum surface strain for this material is 5×10^{-4} . Surface strain can be defined as

$$S = \frac{xt}{L_c^2}$$

where x is the deflection, t is the thickness of the beam, and L_c is the cantilever length. Thus, the maximum strain before failure for a relatively long beam (20 cm) of a piezoceramic thin sheet (0.002

cm) is

$$x = \frac{(S)(L_c^2)}{t} = \frac{(5 \times 10^{-4})(0.2^2 m)}{0.00002 m} = 1 cm$$

Thus, PZT is unsuitable for jacket design.

PVDF, on the other hand, is very flexible. In addition, it is easy to handle and shape, exhibits good stability over time, and does not depolarize when subjected to very high alternating fields. The cost, however, is that PVDF's coupling constant is significantly lower than PZT. Because of edge effects, shaping PVDF reduces the effective coupling. Furthermore, the material's efficiency can be degraded by operating climate and the number of plies used.

With specifications to be discussed in the next section, a $116 cm^2$ 40 ply triangular plate with a center metal shim deflected 5 cm by 68 kg 3 times every 5 seconds results in the generation of 1.5 W of power according to an industry expert [12]. Thus, out of a total of $(68 kg)(9.8 m/sec^2)(0.05 m)(0.6 hits/sec) = 20 W$ of mechanical power input, 7.5% of this power is converted to electricity. Applying this model to a power generating jacket, a total of 0.23 W could be produced through conversational arm motion by placing PVDF in a power generating jacket.

8 Walking

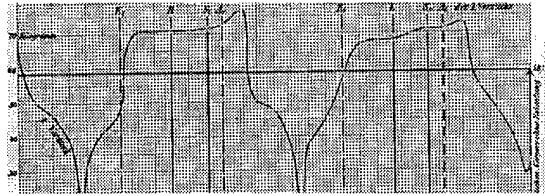


Figure 2: Effective force perpendicular to the ground of a foot of a 58.7 kg man while walking [3].

Using the legs is one of the most power consuming activities the human body performs. In fact, a 68 kg man walking at 3.5 mph or 2 steps/sec produces 280 kcal/hr or 324 W of power [15]. Comparing this to a standing or strolling rate implies that up to half this power is being used for moving the legs. While walking, the traveller puts up to 30% more than his body weight on the balls of his feet (Figure 2). However, calculating the minimum amount of power that can be generated (using the fall of the heel) reveals that

$$(68 kg)(9.8 m/sec^2) = 670 N/step$$

$$(670 N/step)(0.050 m)(2 steps/sec) = 67 W$$

of power is available. This is a promising result given the previous analyses. In addition, some of this power could be stored, providing a constant power supply even when the user is not walking.

8.1 Piezoelectric shoe inserts

Pursuing this further, consider using PVDF shoe inserts for recovering some of this power. There are many advantages to this tactic. First, a 40 ply pile would be only $(28 \mu m)(40) = 1.1 mm$ thick (without electrodes). In addition, the natural shape change of

the shoe when walking provides the necessary deflection for generating power from the piezoelectric pile. PVDF is easy to cut into an appropriate shape and is very durable [6, 2]. Thus, the inserts could be easily put into shoes without serious redesign of the shoe or moving parts.

Detailing the numbers of the previous section, a small women's shoe has a footprint of approximately 116cm^2 . In order to get 5 cm of deflection with a 40 ply pile, 68 kg of effective weight is needed. However, assuming maximum deflection at the end of a step (balls of the feet), the user needs to only weigh 52 kg (115 lbs.). While the numbers given in the last section were for a 15.2 cm by 15.2 cm triangular 40 ply pile, the value can be used to approximate the amount of power an appropriately shaped piezoelectric insert could produce. Thus, scaling the previous 1.5 W at .6 deflections/sec,

$$(1.5W) \left(\frac{2\text{steps/sec}}{.6\text{steps/sec}} \right) = 5W$$

of electrical power could be generated at a brisk walk by a 52 kg user.

8.2 Rotary generator conversion

Through the use of a cam and piston or ratchet and flywheel mechanism, the motion of the heel might be converted to electrical energy through more traditional rotary generators. Generally, the efficiency for a normal generator is very good. However, the added mechanical friction of the stroke to rotary converter reduces this efficiency. A normal car engine, which contains all of these mechanisms plus suffers from inefficient fuel combustion, attains 25% efficiency. However, more realistically, 50% efficiency is a good approximation. Thus, conservatively, 17-34 W might be recovered from a "mechanical" generator.

How can this energy be recovered without creating additional load on the user? If the energy is generated on the downstroke of the heel to the floor, the user may feel as if he is walking up a hill. However, a spring system, where the energy is generated on the upstroke may be more beneficial. Previous experiments have shown that runners can improve their speeds by placing a specially designed spring in their sneakers [20]. This spring returns some of the energy from the downstroke of the heel to the upstroke. Normally this energy is lost to friction and the inelasticity of the runner's muscles and tendons (humans, unlike kangaroos, become less efficient the faster they run [15]). Thus, if such a spring could be designed for normal walking, and a ratchet and flywheel system coupled to the upstroke of the spring, energy could be generated while improving the walking efficiency of the user (Figure 3).

In practice, however, such a system may prove troublesome. In running, the downward force of the heel is greater than in walking. Thus, the increased resistance caused by a spring may not be as noticeable in the downstroke. Furthermore, not all of the energy from the compression of the spring can be recovered in the upstroke, not only for physical reasons, but also for the user interface. If the system is advertised as adding a "spring to your step," then the user should feel some return force from the spring,

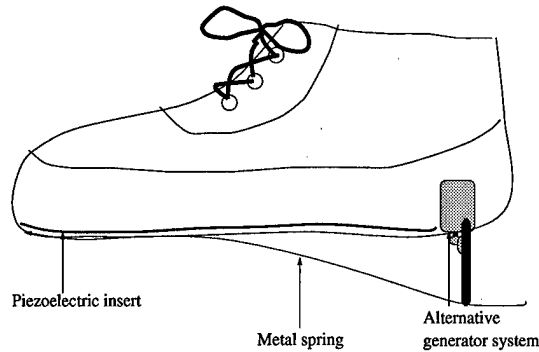


Figure 3: Simple diagram showing two shoe generation systems: 1) piezoelectric film insert or 2) metal spring with coupled generator system.

which puts a limitation on the power generation system. Assuming that 50% of the energy stored in the spring will be recovered by the power generation system, a maximum of 8-17W could be generated.

The issue still remains of how such a system would feel to the user. In order to address this issue, the author is prototyping apparatus which will simulate the effects of the piezoelectric insert and the mechanical spring systems.

8.3 Air resistance

A final suggested method of generating power is to harness air drag while the user is walking. At a 6mph run, only 3% of the expended energy is performed against air resistance [15]. While this is an impressive number, W, little of this energy could be harnessed without severely encumbering the user. At more reasonable walking speeds, the available power declines sharply. Thus, it seems pointless to pursue a hard-to-recover energy source which can only yield 3% of the user's total energy when leg motion may consume over 50% of the total energy during the same activity.

9 Storage considerations

Every power generation system proposed, with the possible exception of heat conversion, would require some power storage device for periods between power generation cycles. Thus, some attention is necessary to the efficiency of storage.

Electrical storage may be preferable due to its prevalence and miniaturization. First however, the power must be converted to a useable form. For the piezoelectric method, a step down transformer and regulator would be needed. Current strategies for converting 100V AC to 5V DC attain over 90% efficiency [13]. Care would be needed to properly match the high impedance of the piezo generator, and, due to the esoteric nature of the problem, the actual efficiency may be lower. For the other generation methods, power regulators would be needed as well. Here, a typical efficiency value would be 93%.

The most direct solution to electrical storage is to charge capacitors that can be drained for power during times when no generation is occurring. However, simply charging the capacitor results in the loss of half the available power [8]. Unfortunately, a purely

capacitive solution for the problem is also restricted by size. A small (under 1 cubic inch) 5 V supercapacitor can handle about 3 Farads. Thus, only

$$E = (.5)CV^2 = (.5)(3F)(5V)^2 = 37.5J$$

of energy could be stored. Correspondingly, for non-generative cycles of a minute,

$$\frac{38J}{60sec} = .62W$$

could be provided from a fully charged capacitor. This is acceptable to provide a energy reservoir for breathing, blood pressure, and body heat. However, over an hour this rate drops to 0.01 W. In order to provide even 1 W of power over this time interval, 100 such capacitors would be necessary. Thus, except for domains when the particular body action is continuously performed, capacitive storage is not suitable for upper limb motion, walking, or typing. In such cases, rechargeable batteries may be necessary.

Mechanical energy storage may be more attractive for some of the generation mechanisms described above. For example, with walking, flywheels, pneumatic pumps, and clock springs may prove more fruitful in storing power. However, the possibilities are numerous and coverage of the field is beyond the scope of this paper.

10 Power requirements for computing

A recent trend in computing is for more capability to be packed into smaller spaces with less power consumption. At first this trend was pushed by laptop computers. With the advent of pen computing and PDA's components have become even smaller and more manageable. Now it is possible to make a computer which can be worn and runs constantly [22]. For example, the author's wearable computer requires 5 W of power to run all components continuously (head mounted display, hard disk, 16 MHz 80286 CPU, memory, serial/parallel/PCMCIA ports, etc.) A standard off-the-shelf gel cell battery can provide this power for 8 hours. However, such a battery has a volume of 450 cubic centimeters. Better battery technology is available, and the author's computer is relatively inefficient on power consumption. For example, a PIC16C71 processor from Microchip Technology Inc. requires only 18 mW at 4 MHz and 0.1 mW at 32 KHz [16]. With memory instead of rotary disk storage, some driver circuitry, and a Private Eye head mounted display from Reflection Technology, a functional wearable computer (without communications) could be made with a power consumption of .5W.

While computing, display, communications, and storage technology may become efficient enough to require unobtrusive power supplies, the desire for the fastest cpu speeds and highest bandwidth possible will offset the trend. In addition, dependence on power cells requires the user to "plug in" occasionally. This is impossible in some military and professional contexts. If body motion is used, it may be significantly more convenient to shift weight from one foot to another, for example, than to search for an electrical outlet.

11 Conclusion

Each of the generation methods has its own strengths and weaknesses depending on the application. However, power generation through walking seems best suited for general purpose computing. The user can easily generate power when needed, and, in many cases, the user's every day walking may be sufficient. A surprising amount of power (8-17 W) may be recovered while walking at a brisk pace, possibly without stressing the user. If less power is needed, piezoelectric inserts may be used reducing the mechanical complexity of the generation system. However, issues of energy storage and human factors still have to be resolved. Thus, the natural next step is to prototype a generator.

12 Acknowledgements

Thanks to Michael Hawley and Neil Gershenfeld for suggesting this project. Thanks also to AMP, PSI, David A. Ross, Gerald Maguire, Russ Hoffman, Paul Picot and a host of folks on the wearables list for suggestions, references, initial calculations, and corrections.

References

- [1] R. Alexander and G. Goldspink (eds.). *Mechanics and Energetics of Animal Locomotion*. Chapman and Hall Ltd., London, 1977.
- [2] Piezo Film Vibration Sensors: New Techniques for Modal Analysis and data sheet. AMP Inc., Valley Forge, PA, USA.
- [3] W. Braune and O. Fischer. *The Human Gait*. (orig. published 1895-1904.) Springer-Verlag, Berlin 1987.
- [4] E. Braunwald (ed.). *Heart Disease: A Textbook of Cardiovascular Medicine*. W B Saunders Company, Philadelphia, 1980.
- [5] W. Cady. *Piezoelectricity*. Dover Publishers, Inc., New York, 1964.
- [6] J. Fraden. *AIP Handbook of Modern Sensors*. American Institute of Physics, New York, 1993.
- [7] Correspondance with Donna Wren, M.D.
- [8] N. Gershenfeld. *The Physics of Information Technology (preliminary draft)*. Cambridge, MA, USA, May 1995.
- [9] J. Gillies (ed.). *A Textbook of Aviation Physiology*. Pergamon Press, Oxford, 1965.
- [10] D.Halliday, R. Resnick, and K. Krane. *Physics 4th edition volumes 1 and 2 extended*. Wiley & Sons Inc., New York, 1992.
- [11] Idea suggested by Michael Hawley, MIT Media Laboratory.
- [12] Correspondance with Don Halvorsen, AMP Inc.
- [13] Correspondance with Russell Hoffman, Fore Systems Inc.
- [14] R. Holland and E. EerNisse. *Design of Resonant Piezoelectric Devices*. Research Monograph no. 56. MIT Press, Cambridge, MA 1969.

- [15] D. Morton. *Human Locomotion and Body Form*. The Williams & Wilkins Co., Baltimore, 1952.
- [16] Preliminary data sheet on PIC 16C71. Microchip Tech Inc., 1992.
- [17] Correspondance with Paul Picot, Roberts Research Institute.
- [18] Piezoelectric Motor/Actuator Kit, Intro to Piezoelectricity. Piezo System, Inc., Cambridge, MA.
- [19] N. Rogacheva. *The Theory of Piezoelectric Shells and Plates*. CRC Press, Boca Raton, 1994.
- [20] Scientific American Frontier (with Alan Alda) Produced by Chedd Angier Productions Watertown, MA
- [21] R. Shephard. *Human physiological work capacity*. Cambridge University Press, Cambridge, U.K., 1978.
- [22] T. Starner Wearable Computing Vision and Modeling Technical Report #318, MIT Media Laboratory, Cambridge, MA, 1993.
- [23] R. Stein, K. Pearson, R. Smith, and J. Redford (eds.). *Control of Posture and Locomotion*. Plenum Press, New York, 1973. very low level animal muscle stuff
- [24] J. Williams, H. Metcalfe, F. Trinklein, and R. Lefler. *Modern Physics*. Holt, Rinehart, and Winston Inc., New York, 1968.

Wearable Computing (and other strange ideas)

Thad Starner

Perceptual Computing Section, The Media Laboratory,
Massachusetts Institute of Technology
Room E15-383, 20 Ames Street, Cambridge MA 02139, USA
E-mail: thad@media.mit.edu

1 Introduction

Invention is not a fast process. Many scientists and inventors speak of sudden insights or a particular late night which led them to great progress. What they don't tell are the tales of working long hours to bring a project to fruition, having the patience to wait until the environment is right to start such a project, or even developing a work environment that encourages innovation. The projects that are reviewed here did not happen overnight, nor did they happen in isolation but with the help of many colleagues and sources of inspiration. Some of the projects are currently underway and have only crude prototypes or concept demonstrations. Others have many years of history behind them. However, each of these projects has led me closer to effective wearable computing.

2 Wearable Computing

What happens when man and machine become one? This is a question that often brings horror stories to mind. However, in real life "cyborgs" are commonplace. Pacemakers have increased longevity for many years. Bionic limbs and internal prostheses improve the quality of life for many people around the world. However, these devices are often considered necessary, in that they augment people with disabilities to approximate the normal human experience. What happens if we go further and try to augment humans with computers to help them achieve more than normal human existence? The goal of this work is to explore using wearable computing to enable the user, whether or not he is disabled.

2.1 Hardware

My current wearable computer consists of Reflection Technology's *PrivateEye*TM, a small display that sits in front of my eye with an adjustable focus from 10" to infinity; HandyKey's *Twiddler*TM, a one-handed chording keyboard and mouse which I am optimizing to improve my typing speed (currently 50 words per minute); a CPU box consisting of a PC class processor, 85M hard drive, and many communications interfaces; and a cellular telephone with modem to connect to the Internet (see Figure 1). This system resides in a shoulder bag which I carry with me almost continuously. A belt pack can also be used to carry the computer. I have been using this particular wearable computer since the summer of 1993. While the current hardware is inconvenient and obtrusive, the technology is advancing so rapidly that soon wearable computers may become undetectable upon casual inspection. We are looking into upgrading the current system to a new, more powerful, and more convenient wearable computer.

Simple applications of this basic interface include 24 hour access to stock market trading systems, courtroom assistants for trial lawyers, police accident reporting, computer system administration monitoring, medical reporting and precedent searching, personalized newspaper readers, news reporting assistants, portable World Wide Web browsers, and many others. However, I believe that the most interesting

applications involve augmented memory, augmented reality, and intellectual collectives.



Figure 1: Computer glasses interface to the wearable computer.

2.2 Remembrance Agent

A general rule of thumb is that computers are used 95% of the time to do simple word processing. In my experience, this number is higher for wearable computing. With constant access to a display and keyboard, the user will type everything, from "Todo" lists, to casual ideas, to his Master's thesis. In fact, portions of this submission have been written on my wearable computer. However, word processing uses only 1% of the available CPU time. This fact opens the door to a whole new field, augmented memory. An example of augmented memory is my Remembrance Agent (RA). As a user works, the computer uses the extra CPU time to look for documents with content similar to what the user is currently typing. The file names or salient lines from these correlated documents can then be displayed at the bottom of the user's word processor in order of similarity (Figure 2). Given a fast reference system (which has been demonstrated), such a display can be updated continually as the user types. Thus, the user is automatically presented with timely information. Organization is improved because the user always has several suggestions for file locations for new information. In addition, during professional talks or technical conversations the user has background material from his own work readily available for reference. Such a system is invaluable for students to keep their education fresh and easily accessible. For example, I now insist on an electronic copy of each text used in the classes I attend. This in itself leads to a new product, value-added texts. The annotations and corrections one student makes over the course of a term may be very valuable to another student.

Another use of the Remembrance Agent is for 24 hour access to other users' "minds." A query can be placed to a particular user's RA which the RA answers by sending back the file which closest matches the query. Thus, the user is not bothered with trivial interruptions for information. Taking this to an extreme, large portions of knowledge can be transferred between users via use of their RA's. On average, engineers in the United States change jobs every 4 years, causing a serious training problem for many companies. As one engineer moves out of a job, he can transfer that portion of his "mind" which applies to the task to his replacement. Thus, while the new

engineer may not have all the active knowledge of his predecessor, relevant information automatically appears in a timely fashion as the new engineer learns his tasks. Finally, given the intimate nature of the Remembrance Agent interface and the large amount of information at its disposal (I have typed over 2 million characters on my wearable computer), the RA's knowledge may be used to help prioritize incoming information for the user.



Figure 2: The Remembrance Agent interface.

2.3 Augmented Reality

Augmented reality refers to the combination the real and the virtual to assist the user in his environment. Applications include telemedicine, architecture, construction, devices for the disabled, and many others. Several large augmented reality systems already exist (for example, JIVE mentioned later), but a wearable computer with a small camera and digitizer opens a whole new set of applications.

2.3.1 Finger Tracking

One of the simplest applications of this camera-based wearable computer is finger tracking. Many pen computer users appreciate the pen interface for its drawing capability. However, given that a computer can visually track the user's finger, there is no need to use a pen (see Figure 3). Such an interface allows the user to replace normal computer pointing devices, such as the mouse, with his finger. The user can control the operating system in this manner or digitize an image and virtually annotate it.

2.3.2 Face Recognition

Several years ago, I demonstrated the Photobook face database system. Photobook has the ability to search a database of 8000 faces in approximately 1 second on the equivalent of a high-end 80486 system and return the top 40 closest matches for a given face (Figures 4 and 5). The system is surprisingly tolerant of lighting changes and facial hair differences. In addition, the user can specify a combination of faces for the search. This allows a "mugshot" application where a crime victim may be able to search mugbooks much quicker than ever before. An experienced user can find a particular person in the 8000 face database within a few mouse clicks.

With the addition of face-finding software, this system is being adapted for use in wearable computing. The goal for the system is to overlay names on faces as the user moves about the world. Markets include the police, reporters, politicians, the visually disabled (with an audio interface), and those with bad memories for faces (the author being included in this last group).

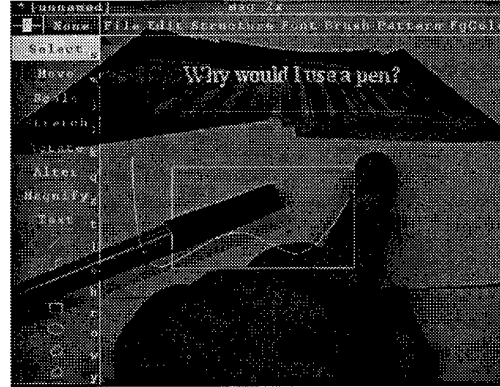


Figure 3: A prototype of a vision-based finger tracking system. As commercial computing becomes more powerful, the orange thimble will not be necessary.



Figure 4: A randomized initial screen from Photobook.

2.3.3 Visual Filter

Steve Mann and I have recently demonstrated a digital visual filter. The basic concept is to process video images digitally in real time to assist the user in everyday tasks. For example, users with low vision find that enhancing the edges in an image helps in face recognition. Another application is to map around "blind spots" in the visually disabled. Figure 6 shows yet another variation, digitally magnifying the image through use of a virtual fish eye lens for help in reading. While current wearable computers do not have the processing power to do these manipulations in real time, the video image can be transferred to a base station computer which does the transformation and resends the image to the user. Thus experimentation can be performed until wearable computers become powerful enough to manipulate video locally.

2.3.4 Navigation

The Global Position System (GPS) allows private users to find their position anywhere on the globe to within 100 meters. Naturally, hooking a GPS system to a wearable computer and mapping software allows the user to track himself while exploring a city. However, the resolution is not fine enough in many situations. By using optical flow (comparing consecutive images to determine the direction of motion) not only can the movement of a user's head be tracked, but warnings can be given of approaching objects for the visually disabled. By implementing a local beacons or a dead-reckoning system in the

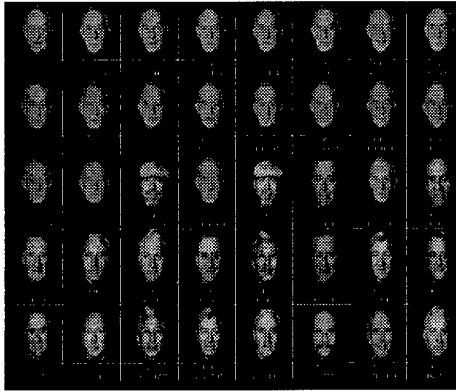


Figure 5: The results of a Photobook search on a face from the first screen.

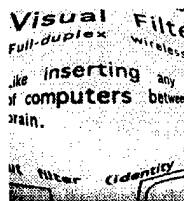


Figure 6: A digital video texture-mapped fisheye "lens" that can be used to help those with low vision.

workplace, much more advanced applications can be developed. Examples include virtual museum tour guides, automatic wiring and gas line view overlays in buildings and on streets, and a new computing environment I like to call the "reality" metaphor. The reality metaphor replaces the typical computer desktop metaphor by overlaying files onto real world objects. Thus, a filing cabinet may have a searchable index overlaid on it. Telephones may have virtual phone directories attached. Virtual 3D "post-it" notes and movies may be applied to objects. Recent electronic mail messages may be rendered on a co-worker's door (or the co-worker!) to remind the user of the last communication with that person. Again, such a system would help provide context-based information in a timely fashion.

2.3.5 Repair instruction

Since a manufacturer can control the markings on the inside of his product, why not instrument the product to allow a wearable camera system to track the object in the user's visual field? By simply putting three distinctive marks at known distances from each other, a wearable camera with known focal length can recover the 3D location of the plane defined by these three marks. By extrapolation from an on-line technical manual, the rest of the object's 3D location can be derived. Thus, when a repair technician walks up to a broken machine, the machine can transmit its diagnostics to the technician's wearable. The wearable automatically determines the problem, locates the 3D position of the object, and overlays specific 3D real-time step-by-step guidelines on the object for the technician to follow. A prototype of such a system is currently under way.

2.4 Intellectual Collectives

What happens when two users have constant access to each other's mind. This is the subject of intellectual collectives. With constant access to a keyboard, a wearable user tends to type every salient thought.

If these thoughts are instantaneously and continuously shown to another user, a tighter collaboration is formed than has ever been experienced. Some interesting questions include does this help in training, what are the privacy issues, how can such a system be made unobtrusive to the user's normal work habits, and does such a collaboration increase the efficiency of a small work group. With a pair of new wearable computers, I hope to begin experimenting with these questions with my undergraduate assistant.

2.5 Human-powered wearable computing

Batteries add size and weight to present-day notebook and wearable computers. Presently, I am working on a paper which explores the possibility of removing this impediment by using the operator's natural body functions to generate power for his computer. Several methods seem feasible, but the most practical generates power through walking. Assuming a spring, ratchet, and small generator system with 50% conversion efficiency in the shoes, 5.5 W (ignoring storage losses) can be continuously supplied if the user walks 10 minutes out of every hour. A more marketable system involves using piezoelectric PVDF film shoe inserts which can generate 4.5 W while the user walks at a brisk pace. Powering wearable computing in such a manner reduces the upkeep of the system and makes it more attractive for military field uses.

3 Previous work

Since being accepted to MIT in 1987, I have worked on many projects. These include an expert system that decorates bathrooms; ThingPaint, a way to apply 3D topographic paint on 3D objects using a standard computer 2D painting interface (Figures 7 and 8), MusicWorld, a virtual environment that uses Kalman filters to synchronize user motion, sound, and graphics (see Figure 9); on-line cursive handwriting recognition (which has arguably the best error rates for its task-Table1); JIVE, a wireless 3D virtual environment which uses only a camera and a microphone for its sensing (see Figure10); and a computer vision based American Sign Language recognizer (40 word lexicon, 5 word full sentence, 99.2% accuracy on an independent test set, see Figure 11 for an example of what the computer "sees"). Some of these projects will be hooked into the wearable computing project. For example, wearable computers can help the user interact with office-based shared virtual environments like JIVE by providing a keyboard interface and private data space when necessary. Sign language recognition may someday be possible using a wearable computer. Finally, Kalman filtering is necessary when trying to align the virtual and real in wearable augmented reality systems. Thus, much of my previous interdisciplinary experience will be called upon in the wearable computing project.

Knowing that wearable computing is a relatively new field, I have founded the wearable computing mailing list. Actively seeking out researchers, I am attempting to provide an atmosphere of cooperation and assistance. In this vein, several volunteers and I are starting the wearable web page which will contain technical information, pointers, and news related to the development of wearable computing. It is my hope that through these actions the benefits of wearable computing will reach the public sooner.

References

- [1] T. Starner. Human-Powered Wearable Computing. *Perceptual Computing TR#328, MIT Media Lab*

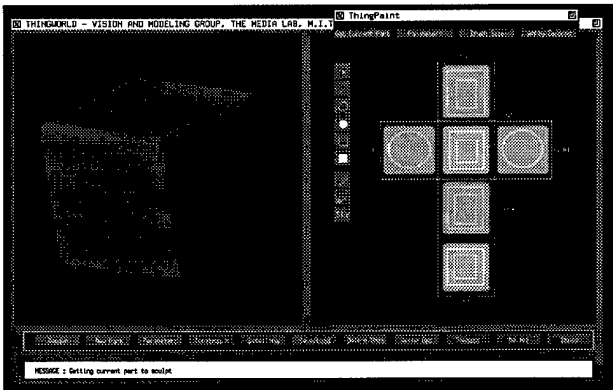


Figure 7: Drawing with topographic paint on the six orthogonal views of an object.

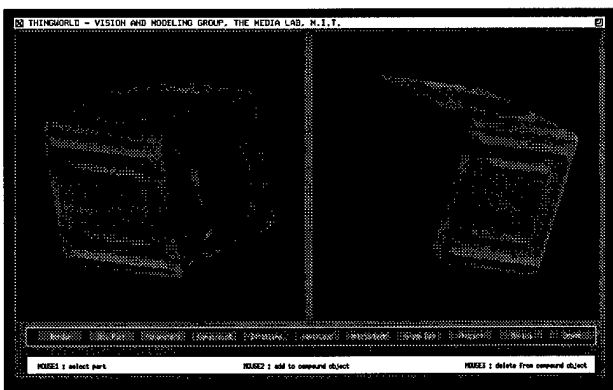


Figure 8: The results of applying the topographic paint to the 3D cube.

- [2] T. Starner and A. Pentland. Visual Recognition of American Sign Language Using Hidden Markov Models. To appear *IWAFGR'95*, Zurich, Switzerland.
- [3] K. Russell, T. Starner, and A. Pentland. Unencumbered Virtual Environments. To appear *IJCAI '95 Entertainment Workshope*, Montreal, Canada.
- [4] T. Starner. Wearable Computing. *Perceptual Computing TR#318, MIT Media Lab*
- [5] T. Starner, J. Makhoul, R. Schwartz, G. Chou. On-line Cursive Handwriting Recognition Using Speech Methods *ICASSP'94* Adelaide, Australia, Vol V, pp. 125-128.
- [6] A. Pentland, T. Starner, N. Etcoff, A. Masiou, O. Oliyide, and M. Turk. Experiments with Eigenfaces *IJCAI'93 Looking at People Workshop*, Chamberry, France 1993.
- [7] A. Pentland, I. Essa, M. Friedmann, B. Horowitz, S. Sclaroff, T. Starner. The Thingworld Modeling System book chapter in *Algorithms and Parallel VLSI Architectures Vol B*. E. Depretere and A. van der Veen (eds.), Elsevier 1991
- [8] M. Friedmann, T. Starner, and A. Pentland Synchronization in Virtual Realities *Presence*, 1(1), 1992.

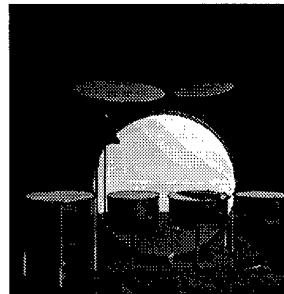


Figure 9: MusicWorld: a demonstration environment for synchronized sound and graphics.

subject	Subst.	Delet.	Insert.	Total
aim	2.7%	0.4%	1.4%	4.5%
dsf	3.6%	0.4%	1.2%	5.2%
rgb	3.3%	0.5%	1.7%	5.5%
shs	1.5%	0.1%	0.5%	2.1%
slb	2.9%	0.1%	1.3%	4.3%
wcd	2.1%	0.4%	0.5%	3.0%
ave.	2.8%	0.3%	1.1%	4.1%

Table 1: Wall Street Journal 25,595 word, writer dependent task word errors for six writers. Total word error is 4.1%.

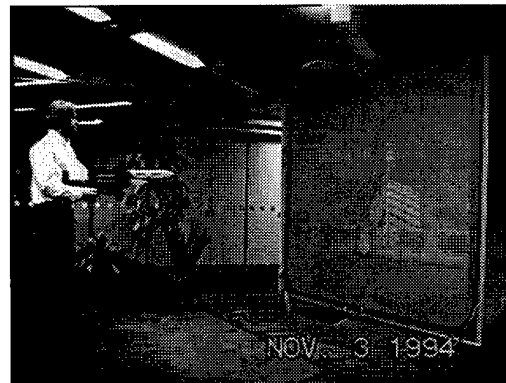


Figure 10: SURVIVE: a game implementation of JIVE using wireless 3D sensing from a camera and microphone.

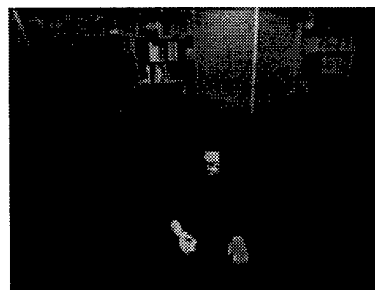


Figure 11: Real-time sign language recognition using only one video camera.