



RECOGNITION OF ISOLATED
COMPLEX MONO- AND
BI-MANUAL 3D HAND GESTURES

Agnès Just ^a Olivier Bernier ^b

Sébastien Marcel ^a

IDIAP-RR 03-63

FEBRUARY 2004

TO APPEAR IN

Proceedings of the sixth International Conference on Automatic Face and
Gesture Recognition 2004

Dalle Molle Institute
for Perceptual Artificial
Intelligence • P.O.Box 592 •
Martigny • Valais • Switzerland

phone +41 - 27 - 721 77 11
fax +41 - 27 - 721 77 12
e-mail secretariat@idiap.ch
internet
<http://www.idiap.ch>

^a IDIAP, CP 592, 1920 Martigny, Switzerland

^b France Telecom R&D, Technopole Anticipa, 22307 Lannion, France

RECOGNITION OF ISOLATED COMPLEX MONO- AND BI-MANUAL 3D HAND GESTURES

Agnès Just

Olivier Bernier

Sébastien Marcel

FEBRUARY 2004

TO APPEAR IN

Proceedings of the sixth International Conference on Automatic Face and Gesture Recognition 2004

Abstract. In this paper, we address the problem of the recognition of isolated complex mono- and bi-manual hand gestures. In the proposed system, hand gestures are represented by the 3D trajectories of blobs. Blobs are obtained by tracking colored body parts in real-time using the EM algorithm. In most of the studies on hand gestures, only small vocabularies have been used. In this paper, we study the results obtained on a more complex database of mono- and bi-manual gestures. These results are obtained by using a state-of-the-art sequence processing algorithm, namely Hidden Markov Models (HMMs), implemented within the framework of an open source machine learning library.

1 Introduction

Nowadays, Human-Computer Interaction (HCI) is usually done using keyboards, mice or graphic boards. The use of hand gestures for HCI can help people to communicate with computers in a more intuitive way. The potential power of gestures has already been demonstrated in applications that use the hand gesture input to control a computer while giving a presentation for instance. Other possible applications of gesture recognition techniques include computer-controlled games, teleconferencing, robotics or the manipulation of objects by CAD designers.

In gestural HCI, the use of video cameras is more natural than any dedicated acquisition device (such as data-gloves for instance) but is also much more challenging. In video-based hand gesture recognition (HGR) it is necessary to distinguish two aspects of hand gestures: the *static* aspect and the *dynamic* aspect. The *static* aspect is characterized by a pose or configuration of the hand in an image. The *dynamic* aspect is defined either by the trajectory of the hand, or by the sequence of hand postures in a sequence of images. Furthermore, there is two sub-problems to address when dealing with dynamic hand gesture recognition: spotting and classification.

On one hand, spotting aims at identifying the beginning and/or the end of a gesture given a continuous stream of data. Usually, this stream of data is made by a random sequence of known gestures and non-gestures. On the other hand, given an isolated gesture sequence, classification outputs the class the gesture belongs to.

In this paper, we will focus on the classification of isolated hand gestures. First, we present an overview of related work on HGR. In section 3, we describe our approach to capture mono- and bi-manual 3D hand gestures, and we describe the database. Then, in section 4 we present experimental results using a baseline algorithm for HGR, namely Hidden Markov Models (HMMs). Finally, we discuss the results and conclude.

2 Related Work

Dynamic HGR is a **sequence processing** problem that can be accomplished by using various techniques.

Darell and Pentland in [3] used a vision-based approach to model both object and behavior. The object views were represented using sets of view models. This approach allowed them to learn their model by observation. The disadvantage of this method is that complex articulated objects have a very large range of appearances. Therefore, they used a representation based on interpolation of appearance from a relatively small number of views. The gesture classification is performed by stereotypical space-time patterns (i.e. the gestures) matched with stored gesture patterns using dynamic time warping (DTW). This system was tested on only two gestures. And images were focused on the hand. The experiment was also user-dependent since each of the seven users were involved in both the training and testing phases.

Finite state machine (FSM) was the first technique applied to sequence processing. It was applied to gestures by Davis and Shah [4]. Each gesture was decomposed into four distinct phases. Since the four phases occurred in fixed order, a FSM was used to guide the flow and to recognize seven gestures. Experiments were using a close-up on the hand.

Hong et al. [5] proposed another approach based on FSM, that used 2D positions of the centers of the user's head and hands as features. Their system permitted to recognize in real-time four mono-manual gestures.

But the most important technique, widely used for dynamic HGR, is Hidden Markov Models. This approach is inspired by the success of the application of HMMs both in speech recognition and in hand-written character recognition fields [8]. Starner and Pentland in [10] used 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. They used networks of HMMs to recognize a sequence of gestures taken from the American Sign Language. Training was performed by labeling each sign with the corresponding video stream. They used language modeling to segment the different signs. The Viterbi decoding algorithm was both used with and without a strong grammar based on the known form of the sentences. With a lexicon of forty words, they obtained 91,9% of accuracy in the test phase. Unfortunately, these results are almost impossible to reproduce.

More recently, Marcel et al. [6] have proposed Input/Output Hidden Markov Models (IOHMMs). An IOHMM is based on a non-homogeneous Markov chain where emission and transition probabilities depend on the input. HMMs are in fact based on homogeneous Markov chains since the dynamic of the system is deter-

mined only by the transition probabilities which are time independent. Their system was able to recognize four types of gestures with 98,2% of accuracy. Furthermore, this database is publicly available from the Internet.

In most of the studies on hand gestures, small vocabulary has been used. In the next section, we describe a more important database for the recognition of mono- and bi-manual 3D hand gestures.

3 Mono- and Bi-manual 3D Hand Gestures

Hand gestures involve not only one hand but also the two hands. Furthermore, gestures occur in a 3D space and not in a 2D image plane. In the proposed system, hand gestures are represented by the 3D trajectories of blobs. Blobs are obtained by tracking colored body parts in real-time using the EM algorithm. This approach is similar to the statistical region approach for person tracking of [11], for gesture recognition of [7], or for deformable models of [9].

3.1 Tracking Blobs in 3D

A detailed description of the 3D blob tracking algorithm can be found in [2]. This algorithm tracks head and hands in near real-time (12Hz) using two cameras (Figure 1).



Figure 1: Top: left and right captured images. Middle: left and right images with projected ellipsoids. Down left: ellipsoids projection on the frontal plane. Down right: ellipsoids projection on the side plane (the cameras are on the left side).

The algorithm is based on simple preprocessing followed by the use of a **statistical model** linking the observations (resulting from the preprocessing stage) to the parameters: the position of the hands and the head. Preprocessing consists of background subtraction followed by specific colors detection, using a simple color lookup table. The **statistical model** is composed of four ellipsoids, one for each hand, one for the head and one for the torso. Each one is projected on both camera planes as an ellipse. A Gaussian probability density function with the same center and size is associated with each ellipse. The parameters of the model (positions and orientations of the ellipsoids) are adapted to the pixels detected by the preprocessing stage. This adaptation simultaneously takes into account the detected pixels from the two cameras, and is based on the maximum likelihood principle. The EM algorithm is used to obtain the maximum of the likelihood.

3.2 Gesture Database

The database used in this paper has been obtained using the tracking method described above. The database consists of 16 gestures (Table 1) carried out by 20 different persons. Most of gestures are mono-manual and some are bi-manual (*fly*, *swim* and *clap*).



Figure 2: From top-left to bottom-right, a frame-by-frame decomposition of a “swim” gesture from the point of view of the right camera.

The use of gloves with distinct colors permits to avoid occlusion problems that occur with bi-manual gestures. The person performing the gesture wears gloves of different colors and a sweat-shirt of a specific color different from the skin color and different from the glove colors in order to help the segmentation of hands, head and torso.

For each person and each gesture, there are 5 sessions and 10 shots per session. All the gestures start and end in the same rest position (the hands lying along the thighs). The temporal segmentation was manually accomplished after a recording session. For each gesture, a trajectory for each blob has been generated. Finally, the database is composed of 1000 trajectories per gesture.

Gesture trajectories correspond to 3D coordinates of center of the head, of the two hands and of the torso. They are produced with the natural hand (left hand for left-handed and right hand for right-handed persons). For the left-handed persons, trajectories have been mirrored. Figure 2 shows an example of the swim gesture sequence from the point of view of the right camera.

Furthermore, for each person and each session, a “Vinci” sequence has been recorded (Figure 3). This sequence gives the maximum arm spread.

Figure 4 and 5 present in a three dimensional space¹ the coordinates of the center of each blob (head, torso and hands) for a “swim” gesture sequence and a “hello” gesture sequence respectively.

4 Experimental Results

In this section, we present baseline results obtained using HMMs on the proposed mono- and bi-manual database. The open source machine learning library used for all experiments is Torch <http://www.torch.ch>.

¹the z axis is the vertical axis of the person.

Name	Description	R M/B
Stop/yes	Raised hand on the head level and facing palm	M
No/wipe	Idem with movements from right to left	R M
Raise hand	Raised hand higher than the head	M
Hello/wave	Idem with movements from right to left	R M
Left	Hand on the hip level, movements to the left	R M
Right	Hand on the hip level, movements to the right	R M
Up	Hand on the hip level, movements to the up	R M
Down	Hand on the hip level, movements to the down	R M
Front	Hand on the hip level, movements to the front	R M
Back	Hand on the hip level, movements to the back	R M
Swim	Swimming mimic gesture	R B
Fly	Flying mimic gesture	R B
Clap	On the torso level, clap the hands	R B
Point left	On the torso level, point to the left	M
Point front	On the torso level, point to the front	M
Point right	On the torso level, point to the right	M

Table 1: Description of the 16 gestures. A hand gesture could involve one hand (**Mono-manual**) or both hands (**Bi-manual**). The gesture could be also a **R**epetitive movement such as *clap*.



Figure 3: Example of images of the “Vinci” sequence from the point of view of the left camera (on the left) and from the point of view of the right camera (on the right).

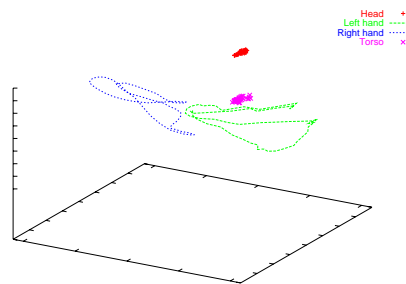


Figure 4: 3D coordinates of the center of each blob (head, torso, left hand and right hand) for a “swim” gesture.

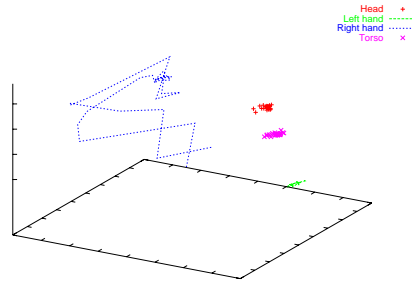


Figure 5: 3D coordinates of the center of each blob (head, torso, left hand and right hand) for a “hello” gesture.

4.1 Hidden Markov Models

A Hidden Markov Model (HMM) [8] is a statistical machine learning algorithm which models sequences of data. It consists of a set of N states, called hidden states because non-observable. It also contains transition probabilities between these states and emission probabilities from the states to model the observations. The data sequence is thus factorized over time by a series of hidden states and emission from these states.

Let q_t be the state, y_t be the output (observation) at time t . The emission probability $P(y_t | q_t = i), \forall i = 1 \dots N$ depends only on the current state q_t . The transition probability between states $P(q_t = i | q_{t-1} = j), \forall i, j = 1 \dots N$ depends only on the previous state.

In order to efficiently use HMMs, it is necessary to impose a topology to the state graph. This topology limits the number of free parameters and allows to inject in the model an *a priori* knowledge on the nature of the data. The left-right topology (Figure 6), has been used in our experiments.



Figure 6: Left-right topology

Then, the training of a HMM can be carried out using the *Expectation-Maximization* (EM) algorithm [1]. For classification problems, a distinct HMM has to be trained for each of the considered classes. Finally, in order to discriminate gesture sequences during the test, a naive Bayes classifier could be used, assuming equal prior probabilities for each gesture.

4.2 Preprocessing the Database

As a first step, a normalization has been performed on all gesture trajectories. We suppose that each gesture occurs in a cube centered on the torso and of vertex size the maximum spread given by the “Vinci” sequence. This cube is then normalized to reduce the vertex to one. Finally, the range of x , y and z coordinates varies between -0.5 and 0.5 . The 3D coordinates of the head and torso are almost stationary. Thus, we keep the normalized 3D trajectories of both hands only. This lead to an input feature vector of size 6.

4.3 Results

In order to find the optimal hyper-parameters (number of states, number of Gaussians per state) for the HMM, the following experimental protocol has been used. For experimental purposes, the database has been split into three subsets: the training set, the validation set and the test set.

	Training set	Evaluation set	Test set
minimum number of frames	12	6	10
average number of frames	25	24	28
maximum number of frames	64	71	89

Table 2: Minimum, average and maximum number of frames for the different subsets

The training set and the validation set contain 5 subjects each. The test set contains 10 subjects. For each subject, all recordings from all shots have been used. Table 2 provides the minimum/average and maximum number of frames per sequence for each subset of the database.

	stop	no/wipe	raise	hello/wave	left	right	up	down	front	back	swim	fly	clap	point left	point front	point right
stop	64.8	29	0.6	1	0.6	1.6	10.8	5.2	2.8	10.6	0.2	0.2	0.4	0.8	2.6	1.8
no/wipe	27.6	63.2	2.4	3.8	0	1.8	7.4	1.4	4.4	11	0	0	0	0.4	0	1.2
raise	1	0	48.2	26.2	0	0	0	0	0	0	0	0	0	3.6	6.8	4.6
hello/wave	4.4	1.6	48.6	68.4	0	0	0	0	0	0	0	0	0	0	0	2.2
left	0	0	0	0	68.8	5.6	4.8	13.8	3.8	3.2	0	0	0	4.8	1.8	0.6
right	0	0	0	0	7.6	73	4.6	4.8	1.2	2	0	0	0	0	0	3.2
up	0.4	0.2	0	0	8	7.8	34.4	23.8	7.6	11.4	0	0	0	0	0.2	0.6
down	0	1.6	0	0	5.4	2	2.4	22.8	2.2	1.2	0	0	0	0	0	0
front	0.6	0.4	0	0	1.6	0.8	2.6	1.2	58.8	0.2	0	0	0	0	2	6.6
back	0.6	1	0	0	2.8	1	26.2	18.2	7.6	53.6	0	0	0	0	1	0.6
swim	0.2	0	0	0.6	0	3	0	0	0	0.2	97.8	2	0.6	0	0	0.2
fly	0	0	0	0	0	0	0	0	0	0	0	97.6	0	0	0	0
clap	0.2	0	0	0	2	0	2.2	2.6	0	0.4	2	0.2	99	0	0	0
point left	0	0	0.2	0	0.6	0	1.8	3.2	5.6	2.4	0	0	0	46.6	10.8	6.2
point front	0	0.2	0	0	2.6	2.4	0.2	2.2	3.8	1.8	0	0	0	25.8	66.6	1.8
point right	0.2	2.8	0	0	0	1	2.6	0.8	2.2	2	0	0	0	18	8.2	70.4

Table 3: Confusion matrix for a HMM on the test set (rows: desired, columns: obtained) in %. Bold types correspond to the well-classified gestures.

Several HMMs have been trained on the training set for different values of the hyper-parameters. All these models were tested on the validation set. The best HMM was obtained for 15 states and 1 Gaussian per state. Then, a HMM with 15 states and 1 Gaussian per state has been re-trained using both the training subset and the validation subset. After training, this model has been applied on the test set. Table 3 shows the obtained results. The average error rate on the 16 classes is equal to 35.38%.

From the results, we observe that bi-manual gestures are very well classified. Few mistakes happen between “swim” and “clap” gestures. If we now have a look to the mono-manual gestures, we notice a few things. First, there is a misclassification between “stop” and “no/wipe” gestures, and between “raise” and “hello/wave” gestures. If we refer to table 1, the only difference between these four gestures is an oscillatory movement of the hand from the left to the right. Thus, HMMs have problems to model this oscillatory movement. This certainly comes from the 3D blob tracking system which is near real-time (12Hz) and thus is not fast enough to capture this kind of movement.

Let us consider the positioning gesture category (“left”, “right”, “up”, “down”, “front” and “back” gestures). We see a block around the diagonal of the matrix. It shows first that HMMs differentiate quite accurately this category of gestures from the others. It shows also that they have difficulties to provide the correct class within this category. The discriminant aspect of these gestures is the dynamic of the hand (Table 1). Then due to the acquisition frequency of the system, HMMs have difficulties at modeling the dynamic which would help to discriminate accurately the gestures from this category.

Finally, if we consider pointing gestures, we see that this action is quite well recognized, but with some mis-classifications in the direction (left, right, front).

5 Conclusions

In this paper, we addressed the problem of the recognition of isolated complex mono- and bi-manual hand gestures. Hand gestures were represented by the 3D trajectories of blobs obtained by tracking colored body parts in real-time using the EM algorithm.

We provide recognition results obtained on a complex database of 16 mono- and bi-manual gestures by a state-of-the-art sequence processing algorithm, namely Hidden Markov Models (HMMs), implemented within the framework of an open source machine learning library. The obtained results are encouraging. Bi-manual gestures are very well classified, and mono-manual gestures are fairly classified.

In the future, we will investigate new features in order to increase the performance of HMMs. We plan also to compare other sequence processing algorithms such as IOHMMs on this database.

Acknowledgments

This research has been carried out in the framework of the GHOST project, funded by France Telecom R&D (project number 021B276). This work was also funded by the Swiss National Science Foundation through the National Center of Competence in Research (NCCR) on "Interactive Multimodal Information Management (IM2)". The authors wish to thank J. Guerin and B. Rolland (from FTR&D Recherche et Développement) for recording and annotating the gesture database.

References

- [1] A.P.Dempster, N.M. Laird, and D.B. Rubin. Maximum-likelihood from incomplete data via the em algorithm. *Journal of Royal Statistical Society*, 39:1–38, 1977.
- [2] O. Bernier and D. Collobert. Head and hands 3d tracking in real time by the em algorithm. In *Proceeding of the IEEE ICCV Workshop on Recognition, Analysis, and Tracking of Faces and Gestures in Real-Time Systems (RATFG-RTS'01)*, 2001.
- [3] T. Darrell and A. Pentland. Space-time gestures. In *Proc. of the Conference on Computer Vision and Pattern Recognition*, pages 335–340, 1993.
- [4] J. Davis and M. Shah. Recognizing hand gestures. In *Proc. of European Conference on Computer Vision*, volume 1, pages 331–340, 1994.
- [5] P. Hong, M. Turk, and T.S. Huang. Gesture modeling and recognition using finite state machines. In *Proc. of the fourth International Conference on Automatic Face and Gesture Recognition*, 2000.
- [6] S. Marcel, O. Bernier, J.E. Viallet, and D. Collobert. Hand gesture recognition using input-ouput hidden markov models. In *Proc. of the FG'2000 Conference on Automatic Face and Gesture Recognition*, 2000.
- [7] V.I. Pavlovic, R. Sharma, and T.S. Huang. Visual interpretation of hand gestures for human-computer interaction: A review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19:677–695, 1997.
- [8] L. R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. In *Proceedings of the IEEE*, volume 77, 1989.
- [9] S. Sclaroff and J. Isidoro. Active blobs. In *Proc. of the International Conference on Computer Vision*, 1998.
- [10] T. E. Starner and A. Pentland. Visual recognition of american sign language using hidden markov models. In *Int. Workshop on Automatic Face- and Gesture-Recognition*, 1995.

- [11] C.R. Wren, A. Azarbayejani, T. Darrell, and A. Pentland. Pfinder: Real-time tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19:780–785, 1997.