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

1	Speech Recognition System Based On Hidden Markov Model
2	Concerning the Moroccan Dialect DARIJA
3	Dr. A. EL GHAZI ¹ and C. $DAOUI^2$
4	¹ FST beni Mellal
5	Received: 9 July 2011 Accepted: 8 August 2011 Published: 23 August 2011
6	
7	Abstract In this work, we present a system for automatic speech recognition on the Moroccan dialect

- 8 In this work, we present a system for automatic speech recognition on the Moroccan dialect.
- We used the hidden Markov model to model the phonetic units corresponding to words taken
- ¹⁰ from the training base. The results obtained are very encouraging given the size of the
- ¹¹ training set and the number of people taken to the registration. To demonstrate the flexibility
- ¹² of the hidden Markov model we conducted a comparison of results obtained by the latter and
- ¹³ dynamic programming.

14

15 Index terms— Hidden Markov Model (HMM), MFCC, DTW, Acoustic vectors.

16 1 INTRODUCTION

he system of automatic speech recognition (ASR) can transcribe a voice message, extracting linguistic information from an audio signal. The system uses hidden Markov model [13] (Hidden Markov Model: HMM) to model words units and sentence of a language. In this work, the interest is to model the Moroccan dialect and implement a recognition system that converts a signal into a meaningful message that can be used thereafter. There is a several applications for a speech recognition system of Moroccan dialect. Most interesting are the man-machine dialogue, ie the passage of oral telephone calls; learning Moroccan dialect and systems helping people with disabilities [1]. The Moroccan dialect is a very important part of popular culture and covers almost the different regions.

The significance of ASR, several free systems have been developed, among the best known: HTK [11] and CMU Sphinx [2][3]. We used the last, it is based on Hidden Markov Model [3] and widely used in the field of speech recognition. In this context, our work focuses on the establishment of foundations for building a system of automatic speech recognition concerning the Moroccan dialect based on Sphinx4 [1].

In the following, we will outline the work done by starting with a theoretical approach to the hidden Markov model and dynamic programming (Section 2). Then, we present in brief a description of Moroccan Author ? ? ? ? ¥ : T raitement de l'information, Faculté des Sciences et Techniques PB 523, Béni Mellal, Maroc. Emails : hmadgm@yahoo.fr, daouic@yahoo.com, najlae_idrissi@yahoo.fr, fakfad@yahoo.fr, bbouikhlene@yahoo.fr. dialect (section 5). The comparison results obtained from the hidden Markov model and dynamic programming are given in Section 6. And it ends with a conclusion and outlook in Section 8.

³⁴ **2 II.**

35 3 THEORETICAL BASES a) Hidden Markov Model

The hidden Markov model is a stochastic system capable [19], after a learning phase, to estimate the likelihood of observation sequence was generated by this model. The case represents a set of acoustic vectors of a speech signal. The hidden Markov model can be seen as a set of discrete states and transition between these states, it can be defined by all of the following parameters: N : the number of model states $A = \{??, ???\} = P(??, ??, ???)$ (?? ??) : is a matrix of size N * N. It characterizes the transition matrix between states of the model. The transition probability to state j depends only on the state i :P(?? ?? = ??/?? ???1 = ??, ?? ???2 = ??, ?)=P(?? ?? = ??/?? ???1 = ??) (1) B = {?? ?? (?? ??) } = P (?? ?? ??), where j ? [1, N] is the set of emission probabilities of the observation ot when the system is in the state qj. The shape of this probability determines the type of HMM used. In this work, we use a continuous probability density [19] defined by the bellow relation: b(o, m, v) = ?(o, m, v) = 1?(2?) n |c| e ? 1 2 (o n ?m i) c ?1 (o n ?m i)? (2) Where: O: Observation trame C: covariance matrix(diagonal) C = 1 n?1 ? (o k ?m k)? * (o k ?m k)

48 4 September

Taking into account several pronunciations of a word requires the use of a multi-Gaussian probability density [21] that the resulting probability is given by:Bj(?? ??) = ? ?? ???? * ?? ??=1 ?? ?? (?? ??)(3)

k: number of Gaussian C ij: Gaussian weight of i in j B j (o t) :observation probability at time t for state j b) Dynamic programming Dynamic programming [18] is one of the algorithms used in speech recognition domain, the principle is to compare two speech signals based on the distance between two matrices corresponding to the coefficients of Mel [18] of the two signals. Calculates the Euclidean distance between two vectors is sound by using the relationship:??(??, ??) = ?? (?? ?? ?? ??) 2 ?? ??=1, ?? = ? ?? 1 ? ?? ?? ?? ?? ?? ?? ?? ?? ??

Then, calculate the minimum distance by traversing the element of the matrix obtained using the relation: ∂ ??" ∂ ??"(??, ??) = min ? ∂ ??" ∂ ??"(?? ? 1, ??) + ??(??, ??) ∂ ??" ∂ ??"(?? ? 1, ?? ? 1) + 2. ??(??, ??) ∂ ??" ∂ ??" ∂ ??"(??, ?? 1) + ??(??, ??)

60 The final distance is:?? = δ ??" δ ??"(??,??)

⁶¹ 5 ??+??

62 Where: I, J: Length acoustic arrays corresponding two signals.

63 6 III.

64 7 EXTRACTION OF ACOUSTIC PARAMETERS a) Pre-

65 treatment

The speech signals used were acquired using a microphone. The noise intra sentence was deleted manually using the tool wavsurfer. The digitized signals will be represented by a family (xn) n? ??1, k] where k is the total number of samples. After, the signal is sampled using the computer's sound card with a frequency Fs = 16kHz ie taking values follows a period 1/FS seconds.

70 8 i. Mel coefficients

Parameterization of speech signals is to extract the coefficients of Mel. This stage is based on the Mel scale to model the perception of speech in a manner similar to the human ear, linear up to 1000 Hz and logarithmically above [22]. The importance of the logarithmic scale appears when using a broad bench of values as it helps to space the small value and approach large values. The digitized signals must be further processed for use in the recognition phase. To do this the pre-emphasis is performed to meet the high frequencies:? ?? = 1 ? 0.97 * ?? ?? ?! (5)

Then the signal is segmented into frames each frame contains N sample of speech and includes almost 30ms of speech, to do this we use a sliding time window of size 256. The successive windows overlap by half of their size ie 128 points in common between two successive windows. In this work we used the Hamming window [23]:w(n) $= 0.54 + 0.46 * \cos(2? * n N?1)$ (6)

In the next step the signal spectrum is calculated, it can introduce the signal (time domain) in frequency domain using the fast Fourier transform FFT:??(??) = 1 ?? ? ??(??)?? ???? 2??(?? ??) ???1 ??=0 (7) $??????(\delta ??"\delta ??") = 2595 * \log (1 + \delta ??"\delta ??" 700) (8-1) ??(??, ??) = ? ??(??, ??) * ??????(??, ??) ??/2 ??=0$ (8-2)

The speech signal can be seen as the convolution in the time domain excitation signal g (n) and the vocal tract impulse response h (n):??(??) = δ ??" δ ??"(??) * ?(??)(9)

The application of the logarithm of the model on this equation gives:???? δ ??" δ ??"|??(??)| = ???? δ ??" δ ??"|??(??)|(10)

Finally, to obtain the coefficients of Mel applying the inverse Fourier transform defined by:FFT ?1 $\{X(i, n)\}$ = x(n) = 1 N? X(i, n)e jk2?(n N) N 2 ?1 k= N 2(11)

91 This gives a vector of coefficients on each Hamming window. The number of filter adopted in this

92 9 September

⁹³ To simulate the functioning of the human ear, we filter the signal through a bank of filters that each have ⁹⁴ a triangular response bandwidth. The filters are spaced so that their evolution is the Mel scale [22]. The ⁹⁵ approximate formula of the scale of Mel:

work is 12, it added the first and second derivatives of these coefficients, which gives in total 39 coefficients.

97 Figure 1 gives a summary of the extraction of Mel coefficients (MFCC).

⁹⁸ 10 Signal pre-emphasis and framing of Hamming

99 11 LEARNING

After the extraction phase, the speech signal is represented by a matrix N* 39 which N is the number of windows in the signal. The audio files used in the learning phase must be segmented into phonemes; each word corresponds to a sequence of phonemes. Each of these will be represented by a hidden Markov model with three states, each state is characterized by:

- -Vector averages for a state i, is given by:???? = 1 ?? ? ?? ?? ?? ?? ?? =1
- , n: number of vectors for each state O k : Observation vector number k.
- -Covariance matrix for state i:?????? = 1 ???1 ? (?? ?? ?? ?? ??)? * (?? ??=1 ?? ?? ?? ??) (12)

The calculation of the mean vector and covariance matrix is performed for each Gaussian. In this paper we use five Gaussian so there will be five vehicles and five averages covariance matrices for each state. The calculation of the probability of resulting observation for each state is realized by the relationship 3.

Learning the model tends to maximize the logarithm of the probability of observation called the likelihood, to do this we use the Baum-Welch algorithm [15], whose steps are: 1-Initializing the model -Creation of HMM for each state -Initialization of the initial probability vector ? with a higher probability for the first state and non-zero for the other two remaining states.

-Initialization of the transition matrix with probabilities respecting any transitions that the sum is equal to 1 and the model is a left-right (upper diagonal). 2 -Maximization: In this step, each iteration updates the model parameters and calculate again the likelihood. The updating of the model parameters is done via the following relations: With?? ???? = ? ?? ?? (?? ,??) ?? ??=1 ? ?? (??) ?? ??=1 , 1 ? ?? ? ?? , 1 ? ?? ? ?? (?

¹²¹ Cjk is the weight of the Gaussian k relative to the state j and the coefficients ? and ? are calculated by the ¹²² Forward-Backward algorithm [15].

123 V.

124 **12 RECOGNITION**

The principle of recognition can be explained as the calculation of the probability P(W / O): the probability that a sequence of words W is the signal S and to determine the word sequence that maximizes this probability.

that a sequence of words W is the signal S and to determine the word sequence that maximizes this probabili According to Bayes formula the probability P(W/S) can be written:

¹²⁸ 13 P(W/S)=P(w).P(S/W)/P(S) (2)

With: -P (W): Prior probability of word sequence W:(Sample language). -P (S / W): Probability of signal S, given the sequence of words W (Acoustic Model).

-P (S): probability of the acoustic signal S (independent of W).

The figure 2 shows the various stages of recognition, as a first step the signal is treated to extract acoustic

vectors, based on these vectors the acoustic model is built from the HMM of phonemes learned on the training rate corpus . The succession of phonemes HMMS form the words models.

135 14 PRESENTATION OF MOROCCAN DIALECT

The Moroccan dialect called Darija is the popular language broadcast in almost all regions of the country. This dialect is a communication tool widely used and is different from one region to another. The dialect Darija contains almost Arabic words in addition to a regional component, the difference between classical Arabic and dialect Darija is at the pronunciation.

140 15 EXPERIMENTAL RESULTS

¹⁴¹ 16 a) Learning base

The learning base used in our system contains 2500 pronunciation, the characteristics of the training set are illustrated in the following table:

¹⁴⁴ 17 Duration of the training set

¹⁴⁵ 18 Number of pronunciations 1h40min of pronunciations. 2500

¹⁴⁶ recorded pronunciation independently and in different situ-

147 ations

148 Table1 : Characteristics of the learning base

The construction of the training set was made by taking the pronunciation of Arabic numerals 0 to 9 in the Moroccan dialect, Table **??** shows the formation of the learning base. The test database contains 300 different pronunciations including noisy audio files. The recognition quality is measured by calculating the rate of recognition given by equation (3

153 19 September

The efficiency of dynamic programming appears on the audio files not noisy. The disadvantage is that the execution time increases proportionally with the length of the file, which influence the time of recognition. In comparison with dynamic programming, hidden Markov model can model a word by a sequence of phonemes and sentence by a sequence of word models, which makes this process more effective and more appropriate to be implemented in systems Recognition advanced.

¹⁵⁹ **20 VIII.**

160 21 CONCLUSION

This work enables the establishment of a voice recognition system of the Moroccan dialect. This article can give an idea about the phonetics used for the recognition of the language. In comparison with dynamic programming, the results obtained by the hidden Markov model are very satisfactory despite the limited number of speakers and size of the database. This shows the importance of stochastic and probabilistic modeling in the field of

165 recognition.

Based on what has been achieved in this work, we'll build a system of passing oral phone call on the Moroccan
 dialect integrated into mobile phones, helping people with disabilities and people who do not dial telephone numbers.



Figure 1:

666-666-666-666

Figure 2: Speech

1 Figure 3: Fig 1:

168

 $^{^1 \}odot$ 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XV Version I 3

Speech Recognition System Based on Hidden Markov Model Concerning the Moroccan Dialect DARIJA ?? = number of words recognized size of the test The results obtained are shown in Table 4. Test database Results 300 different pronunciations T = 91%introducing more noisy audio files Tab.4 : Results for the recognition system of the dialect Darija The comparison of the results was made on noisy audio data. Table 5 illustrates the results obtained. HMMDTW Execution Time Very Plus fast then 10s for a big wav files $91\% \ 60\%$ Recognition rate Tab.5 : Results of comparison between the HMM and DTW

[Note: © 2011 Global Journals Inc. (US) Global Journal of Computer Science and Technology Volume XI Issue XV Version I 4 2011]

Figure 4:

21 CONCLUSION

- [Galley et al. ()] , P Galley , B Grand , & S Rossier . 2006. (reconnaissance vocale Sphinx-4" EIA de Fribourg
 mai)
- [Amour et al. ()], M Amour, A Bouhjar, & F Boukhris, Ircam. 2008. (initiation à la langue Amazigh" 2004.
 10. RAP: Thèse. Benjamin LECOUTEUX)
- 173 [Cornijeol and Miclet ()] Apprentissage 14. Artificielle-méthode et concept, A Cornijeol, L Miclet . 1988.
- 174 [Cornijeol and Miclet ()] Apprentissage artificielleméthode et concept, L Cornijeol, Miclet . 1988.
- 175 [Resch ()] Automatic Speech Recognition with HTK, B Resch . 2003.
- [Chan et al. (2005)] Building Speech Applications Using Sphinx and Related Resources, A Chan, Evandro Gouvêa
 & Rita, Singh. http://docpp.sourceforge.net August 2005.
- 178 [Pellegrini ()] Durée suivi la voix parlée garce au modèle caché, T Pellegrini , Raphael . 1989.
- [Fang ()] From Dynamic Time Warping (DTW) to Hidden Markov Model, Chunsheng Fang . 2009. HMM)
 University of Cincinnati
- 181 [Robiner ()] Fundamentales of speech recognition, Juang Robiner . 1993.
- 182 [Reweis ()] Hidden Markov-Modele-Sam, Reweis . 1980.
- [Dr and Drygajlo] Introduction aux statistiques gaussiennes et à la reconnaissance statistique de formes, . A Dr
 , Drygajlo . Ecole Polytechnique Fédérale de Lausanne
- [Semet G and Treffot ()] 'La reconnaissance de la parole avec les MFCC'. G Semet & G , Treffot . TIPE juin
 2002.
- [Sigurdsson et al.] 'Mel Frequency Cepstral Coefficients: An Evaluation of Robustness of MP3Encoded Music'. S
 Sigurdsson, Kaare Brandt Petersen, Tue Lehn-Schiøler. Informatics and Mathematical Modelling Technical
 University of Denmark Richard Petersens Plads -Building 321 DK-2800, (Kgs. Lyngby -Denmark)
- [Al Ani] Modèles de Markov Cachés (Hidden Markov Models (HMMs)), T Al Ani . Paris / LIRIS. (Laboratoire
 19. A2SI-ESIEE)
- IJamoussi ()] 'Méthodes statistiques pour la compréhension automatique de la parole'. S Jamoussi . Ecole
 doctorale IAEM Lorraine, 2004.
- [Divejver and Killer ()] 'Pattern recognition'. J Divejver , Killer . Pattern Recognition: a statistical approach,
 1982. Prentice Hall.
- [Ali Sadiqui Noureddine Chenfour ()] 'Reconnaissance de la parole arabe basé sur CMU Sphinx'. Ali Sadiqui &
 Noureddine Chenfour . Séria Informatica 1 2010. VIII fasc.
- [Semet Gaetan et al. ()] Reconnaissance de la parole avec les coefficients MFCC' TIPE jiun, & Semet Gaetan ,
 'Treffo , Grégory . 2002.
- [Pellegrini and Duée ()] 'Suivi de la voix parlée grâce aux modèles de Markov Caché'. T Pellegrini , R Duée .
 lieu : IRCAM 1 place Igor Stravinsky 75004 PARIS jiun, 2003.
- 202 [Gonzales and Thomson ()] Syntactic pattern recognition, R Gonzales, M Thomson . 1986.
- 203 [Satori and Harti] 'Système de la reconnaissance de la reconnaissance automatique de la parole'. H Satori , & M
- Harti . Faculté des Sciences, B.P. 1796, Dhar Mehraz Fès, (Maroc)