

# QUASI-CONTINUOUS LOCAL CODEBOOK FEATURES FOR MULTILINGUAL ACOUSTIC PHONETIC MODELLING

Frank Diehl and Asunción Moreno

Universitat Politècnica de Catalunya (UPC)  
Jordi Girona 1-3, 08034 Barcelona, Spain  
{frank, asuncion}@gps.tsc.upc.es

## ABSTRACT

In this article we present a method for defining the question set used for the induction of acoustic phonetic decision trees. The method is data driven resulting in an ordered feature space in contrast to the usual categorical one consisting of phonetic attribute values. Visualization of the feature space verifies that the derived characteristics are meaningful. We apply the features to a multilingual speech recognition task, showing that comparable results to the standard method, using question sets devised by human experts, can be derived.

## 1. INTRODUCTION

A central question in the design of an automatic speech recognition (ASR) system is the definition of proper acoustic phonetic entities. This question gets even more important when trying to share the acoustic feature space among different languages, trying to exploit acoustic similarities between languages for reducing the size of the overall parameter space.

The common state of the art approach for defining the acoustic models is the use of a phonetic decision tree. Such a tree constitutes a functional mapping from a feature to a model domain defining for all phonetic circumstances of the input domain a proper hidden Markov model (HMM) in the output domain.

One crucial topic of this mapping function is the definition of the input domain. A standard approach is the use of phonetic features assigned to the phonemes. This usually works quite good, but exhibits the problem that the designer is dependent on phonemical knowledge. In case of a multilingual system design this might be a severe problem. Not only knowledge for one but for all languages is needed. Additionally, this knowledge should be comparative. This is in strong contrast to the information found in phonetic dictionaries and textbooks assigned to one specific language. Information and transcriptions given there usually follow the principle of phonological contrast [1] using a broad transcription. This leads to the problem that equal phonetic features may be assigned to clearly distinguishable phonemes of different languages.

To cope with these problems, several methods were proposed to construct the features data driven. Ciprian et al. [2] present an approach based on a bottom-up clustering using bigram phone statistics. Beulen et al. [3] also propose a bottom-up approach but already incorporate acoustic information by making use of the HMMs seen in the database. Both methods were presented in a monolingual framework, but at least the second one should be applicable in the multilingual case too. As a third approach we mention the construction of phonetic broad classes by a confusion matrix, Byrne et al. [4], Žgank et al. [5].

This work was granted by the CICYT under contract TIC2002-04447-C02.

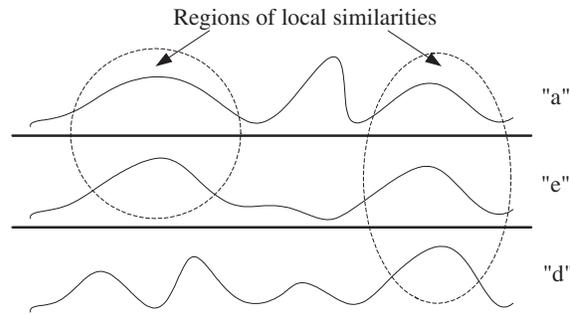


Fig. 1. Local codebook similarities.

In [6] local similarities between the probability density functions of HMMs are identified, and used for constituting phonetic features. For identifying the similarities, advantage is taken of the prototype character of the single mixture components of the codebooks of semicontinuous HMMs (SCHMM). In this work we extend this method of so called "local codebook features" (LCB) by transforming the high dimensional and low populated LCB feature space to a low dimensional subspace using a principal component analysis.

The paper is organized as follows. In Section 2 LCB features are presented, followed by Section 3 discussing the question generation. Section 4 gives a system overview followed by Section 5 with the test set up. Test results are presented in Section 6 and the conclusions are given in Section 7.

## 2. LOCAL CODEBOOK FEATURES

The basic idea behind LCB-features [7], [6] is that similar phonetic properties should cause similar shaped probability density functions (PDF) on a local scale. In Figure 1, this idea is depicted. It shows the assumed probability density functions of the phonemes 'a', 'e', and 'd', with two locally similar regions. The similarities might be caused by common properties of the phonemes as e.g. 'voiced' or 'open'.

Constructing features based on this idea means to identify such similarities in the PDFs of HMMs. In the case of continuous HMMs, this might be a quite hard task, but in the case of SCHMMs, it results in simple vector calculus. Assuming the PDFs being constructed by a set of prototype mixture components, the calculation can be reduced to the prototype weights.

For the sake of simplicity, we assume during derivation of the method, that the beforehand trained incontextual HMMs have only one state and refer to only one codebook. In the following, this al-

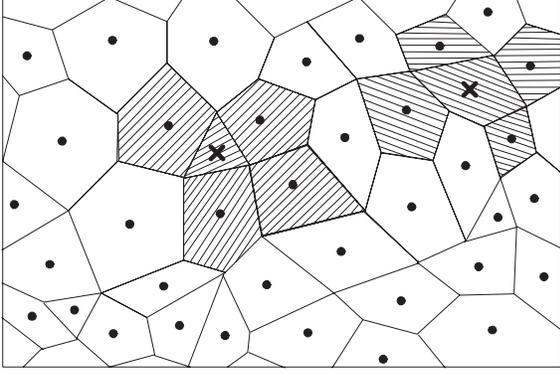


Fig. 2. Locality and neighborhood.

lows to leave out the state and codebook indices. Hence, the PDF of one SCHMM is given as

$$G_{mix}(i) = \sum_{l=1}^L c_{il} \cdot G(\mu_l, \cdot) \quad i \in \{1, \dots, I\}, \quad (1)$$

with  $L$  the codebook size and  $i$  the HMM index.  $G(\mu_l, \cdot)$  names the  $l^{th}$  mixture component with mean vector  $\mu_l$ . For the local search, we define a locality just by a mixture component, and therefore by every index  $\tilde{l}$  out of the  $L$  mixture components. I.e. there are  $L$  localities. We also define the neighborhood of  $\tilde{l}$  as the  $\tilde{L}$  mixture components closest to the locality  $\tilde{l}$ .

Figure 2 depicts the concept for a two dimensional case. It shows a part of a codebook with class regions and the mean values (black dots). Two localities are accentuated by marking the corresponding means by crosses instead of dots. The neighborhood is set to  $\tilde{L} = 5$ , and the equivalent regions with the five closest mean values to the localities are drawn hatched.

For a formal derivation of the method, we define the distance  $d(\tilde{l}, l)$  between the mixture components by Equation (2)

$$d(\tilde{l}, l) = \langle \mu_{\tilde{l}}, \mu_l \rangle \quad l, \tilde{l} \in \{1, \dots, L\}. \quad (2)$$

Identifying for each locality  $\tilde{l}$  the  $\tilde{L}$  closest neighbors is done by evaluating Equation 2 for all  $l$  and  $\tilde{l}$ . As result we get for each locality  $\tilde{l}$  an index set  $S_{\tilde{l}}$  naming the indices of the  $\tilde{L}$  closest mixture components of locality  $\tilde{l}$ . Using  $\min^{(n)}$  to signify the " $n^{th}$  smallest value of" we can express  $S_{\tilde{l}}$  as

$$S_{\tilde{l}} = \left\{ l \mid \arg \min_{1 \leq l \leq L}^{(n)} d(\tilde{l}, l), \quad n \in \{1, \dots, \tilde{L}\} \right\}, \quad (3)$$

where  $\tilde{l} \in \{1, \dots, L\}$ .

Applying these index sets to the PDF for each HMM  $i$ , it results into local weights vectors  $\tilde{c}_{\tilde{l}i}$  comprising of each HMM  $i$  the weights of the codebook mixture components in the neighborhood of the locality  $\tilde{l}$ . Hence,  $\tilde{c}_{\tilde{l}i}$  represents the local shape of the underlying PDF. Although this is not completely correct, the remaining mixture components also contribute to the probability mass at the locality  $\tilde{l}$ , it is a reasonable simplification. Grouping all local weights vectors of a locality  $\tilde{l}$  together, we get the  $\tilde{L} \times I$  matrix  $\tilde{C}_{\tilde{l}}$

$$\tilde{C}_{\tilde{l}} = [\tilde{c}_{\tilde{l}1}, \dots, \tilde{c}_{\tilde{l}I}]. \quad (4)$$

In this vector formulation, the term "locally similar" for the PDFs of different HMMs is equivalent to vectors  $\tilde{c}_{\tilde{l}i}$  being close together in the vector space spanned by them. Searching for local similarities is therefore reduced to the search of suitable clusters in the

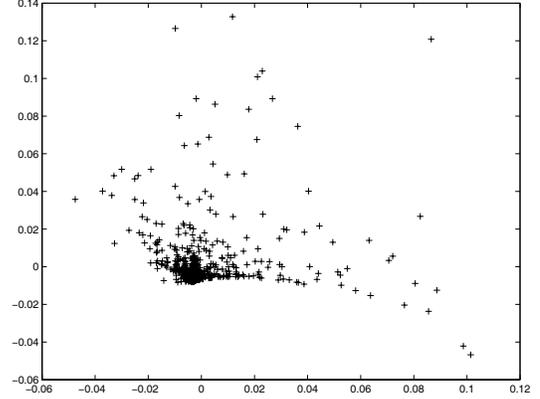


Fig. 3. Sammon map of local weights vectors.

space  $\text{span}\{\tilde{c}_{\tilde{l}1}, \dots, \tilde{c}_{\tilde{l}I}\}$ . For identify meaningful clusters we apply K-means clustering to the samples  $\tilde{C}_{\tilde{l}}$ .

In Figure 3 we visualize the elements falling in one cluster by a nonlinear multidimensional scaling [8] trying to preserve the relative distances between elements. The figure shows the resulting Sammon Map of a typical  $\tilde{C}_{\tilde{l}}$  sample consisting out of 732 vectors. The mapping is performed from dimension  $\tilde{L} = 6$  to the plane using the euclidian distance. The relative mean mapping error of the vector distances is given by less than 7%, indicating the map being a good representation of the original sample.

In Figure 3 three regions can be identified. There is a highly crowded region around the origin. These samples correspond to HMMs without significant probability mass within the locality under investigation. It follows a less crowded broad stripe clearly showing some structure. Here some meaningful cluster can be defined. Finally, there is an individual entry in the right upper corner. It is caused by an HMM showing a quite different structure at the current locality than the others.

We proceed by constructing for each HMM a LCB feature vector stating for each cluster whether the corresponding vectors  $\tilde{c}_{\tilde{l}i}$  falls in it. This results in binary LCB feature vectors  $f_i$  indicating with 0 and 1 whether HMM  $i$  contribute to a cluster or not. Putting together the vectors  $f_i$  leads to the LCB feature matrix  $F$ ,

$$F = [f_1, f_2, \dots, f_I]. \quad (5)$$

Matrix  $F$  consists out of  $I$  columns, one for each HMM. The number of rows is given by the number of localities used for defining the LCB features, multiplied by the number of clusters identified per locality. This number can be quite high. Assuming 256 localities, according to a codebook with 256 mixture components, and  $K = 5$  cluster, we get  $256 \cdot 5 = 1280$  rows. That is, the feature space constituted by matrix  $F$  is of dimension 1280, whereas only 30–50 measurements, the typical number of phonemes, are available.

Reducing the dimensionality of the feature space is done by a principal component analysis (PCA) of matrix  $F$ . The correlation matrix  $FF'$  is build and the eigenvalue problem

$$FF'U = UD \quad (6)$$

is solved. The resulting eigenvector and eigenvalue matrices  $U$  and  $D$  define the final, quasi-continuous feature space  $\tilde{F}$  by the transformation

$$\tilde{F} = \tilde{U}'F. \quad (7)$$

Matrix  $\tilde{U}$  is a reduced version of the eigenvector matrix  $U$ . It consists out of the eigenvectors with the biggest eigenvalues, contributing strongest to the new feature space.

The resulting feature space is given by  $\text{span}(\tilde{F}')$  with dimension  $\text{rank}(\tilde{F})$ , typically in the range 10 – 15. The features defined by matrix  $\tilde{F}$  are given by its columns consisting of continuous values, defining for each phoneme a point in the feature space. In reality these features are not continuous although the number of possible feature vectors is quite high. We derive them by the transformation of the binary matrix  $F$  to  $\tilde{F}$ . Therefore, the number of different features is at the most given by all possible permutations of a vector  $f_i$ . With an assumed dimension of 1280 we get  $2^{1280}$  permutations. This leads to an very fine grid in the final feature space  $\text{span}(\tilde{F}')$  and the name "quasi-continuous LCB features".

### 3. QUESTION GENERATION

For model definition we use a binary decision tree [7] splitting nodes according to a binary question and an entropy based impurity measure. In the classical approach, a binary question is composed out of phonetic attribute values assigned to the SAMPA representations of the phonemes we use in our system. The attribute values are taken from corresponding IPA descriptions [1] of the phonemes. Examples are given in Table 1.

A binary question is composed out of one up to two phonetic at-

**Table 1.** Examples of phonetic attributes.

	p	b	i	e
Attri. 1	consonant	consonant	vowel	vowel
Attri. 2	obstruent	obstruent	front	front
Attri. 3	plosive	plosive	close	close-mid
Attri. 4	-	-	short	short

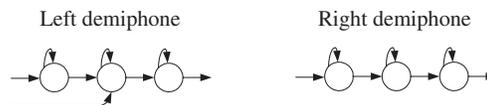
tribute values associated with the context of a model. According to Table 1, a question may be: "Is the context of the model plosive or close?". This is done for all possible compound questions of an attribute and for all attributes.

In case of LCB features, we started by training incontextual models (3 states) for each language. This is followed by extracting the features as described in Section 2. The features are based on common multilingual mel-cepstrum coefficients (MFCC) codebooks. The neighborhood is chosen to  $\tilde{L} = 6$  and, for each locality K-means clustering is applied leading to 14 cluster. Together with a codebook size of 256 we get  $256 \cdot 14 = 3584$  cluster. A phoneme is assigned to a cluster if any of its states matches it. With 47 phonemes for German, 44 for English and 31 for Spanish the intermediate matrix  $F$  has dimension  $3584 \times 122$ . PCA leads to the final feature matrices  $\tilde{F}$  of dimension  $15 \times 122$ . A clipping of matrix  $\tilde{F}$  for some German phonemes is presented in Table 2. These features are directly overtaken for constructing question by

**Table 2.** LCB features.

Phoneme					
2:	2.177	-0.786	-3.199	2.566	...
6	1.680	0.447	2.749	-2.175	...
9	-2.066	-1.169	2.562	1.833	...
a:	1.609	-2.587	-1.795	1.721	...

the decision tree. A question corresponding to Table 2 is: "Is the third component of the model's feature less than  $-1.795$ ?" Quasi-continuous LCB features constitute a fundamental change



**Fig. 4.** Demiphone topology.

in the process of model definition by a decision tree. Instead of the usual categorical variables, as plosive, voiced, etc. we now have to handle ordered variables in the definition space of the tree.

### 4. SYSTEM OVERVIEW

The system we use works with SCHMMs. Every 10ms twelve mel-cepstrum coefficients (MFCC) (and the energy) using cepstral mean subtraction are computed. First and second order differential MFCCs plus the differential energy are employed. For each sub-feature, a codebook is constructed consisting of 256 and 32 (delta energy) gaussian mixtures, respectively. Common multilingual codebooks are used.

Acoustic phonetic modelling is done by demiphones, Mariño et al. [9]. They can be thought of as triphones which are cut in the middle giving a left and a right demiphone, see Figure 4.

The concept of the demiphones also influences the process of the LCB feature construction. Instead of monophones we construct the LCB features on left and right incontextual demiphones leading to two independent features spaces, one for the left and one for the right demiphones.

### 5. TEST SET UP

Training and testing the systems are performed using the SpeechDat-II fixed telephone databases. Three languages are used: Spanish (S), English (E), and German (G). For each language, a 1000 speaker training and a 400 speaker test part is extracted. For training, we use phonetical rich sentences. For testing, phonetical rich words mixed with application words are used. The test suit consists of isolated words avoiding the need of a language model. Detailed statistics on the corpora are given in Table 3.

	Phrases	Females	Males	Grammar size
S-Training	7994	500	500	-
S-Test	2644	200	200	1438
E-Training	8089	500	500	-
E-Test	2554	200	200	1254
G-Training	7540	500	500	-
G-Test	2700	200	200	1314

**Table 3.** Training and test data statistics.

State tying without a-priori distinction between different base phones is done. That is, each tree performs the model definition over the whole phoneme set.

### 6. TEST RESULTS

The actual performance tests for the quasi-continuous LCB features are preceded by an investigation of the feature space's topology. That is, we try to examine the neighborhood relationships within  $\text{span}(\tilde{F})$ . For reasonable features it is expected to find a partitioning in vowels, fricatives, etc. Accessing the question is done by building a Sammon map [8] for the left feature space

