

# Dynamic Speech Parameterization for Text-Independent Phone Segmentation

Anaía S. Cherniz

María E. Torres

Hugo L. Rufiner

**Abstract**—In this work, a dynamic speech parameterization based on the continuous multiresolution divergence is used to modify a text-independent phone segmentation algorithm. This encoding is employed as input and also replaces an stage of the segmentation procedure responsible for the estimation of the intensity of changes in signal features. The segmentation performance of this representation has been compared with the original algorithm using as input a classical Melbank parameterization and speech representation based on the continuous multiresolution divergence. The results indicate that the modification here proposed increases the ability of the algorithm to perform the segmentation task. This suggests that continuous multiresolution divergence provides valuable information related to acoustic features that take into account phoneme transitions. Moreover, this parameterization gives enough information for its direct use without further processing.

## I. INTRODUCTION

Parameterization of speech is a fundamental step in multiple applications that allow to represent the signal with a few coefficients where the most important properties of speech are highlighted [5]. In this paper we present a new speech parameterization scheme, that takes advantage of nonlinear properties of the speech signal.

Information about the changes in the dynamics of speech signal can be hidden in different time-varying features. Tools based on entropy notions have been used to characterize the complexity degree of physiological signals. Their use has been extended over different time-scale distributions. The multiresolution entropy gives account of the temporal evolution of Shannon or Tsallis entropies computed over the wavelets coefficients. Continuous multiresolution entropy (CME) [6] has shown to be robust to additive white noise, in the detection of slight changes in the underlying nonlinear dynamics corresponding to physiological signals [2]. Recently, speech parameterization based on CME and continuous multiresolution divergence (CMD), using the Shannon entropy and the Kullback-Leibler distance respectively, have been used to perform text-independent phone segmentation, giving promising results [3].

This work was supported by Universidad Nacional de Entre Ríos (UNER), Universidad Nacional del Litoral (UNL) and CONICET, under Projects PID 6111 and 6106, PAE 37122 and PAE-PICT-2007-00052.

A. S. Cherniz is with Faculty of Engineering, UNER, Ruta 11 Km. 10, O. Verde, E. Ríos, Argentina. [acherniz@bioingenieria.edu.ar](mailto:acherniz@bioingenieria.edu.ar)  
M. E. Torres is with Faculty of Engineering, UNER and Centro de I+D en Señales, Sistemas e Inteligencia Computacional, Fac. de Ing. y Cs. Hídricas, UNL, Sta. Fe, Argentina and CONICET. [metorres@santafe-conicet.gov.ar](mailto:metorres@santafe-conicet.gov.ar)

H. L. Rufiner is with Faculty of Engineering, UNER and Centro de I+D en Señales, Sistemas e Inteligencia Computacional, Fac. de Ing. y Cs. Hídricas, UNL, Sta. Fe, Argentina and CONICET. [lrufiner@bioingenieria.edu.ar](mailto:lrufiner@bioingenieria.edu.ar)

In this paper we use a speech parameterization based on the CMD, computed with the Kullback-Leibler distance, to perform an automatic speech segmentation procedure using a modified version of the text-independent phone segmentation algorithm proposed by Esposito and Aversano [4].

The text-independent phone segmentation algorithm proposed in [4] is a novel method that accomplishes the phonetic segmentation based on the detection of spectral instability in multiple frequency bands. The algorithm works on an arbitrary number of time-varying features, obtained through a short-term analysis of the speech signal. In the case of text-independent segmentation methods, no prior knowledge of the linguistic content contained in the waveform is needed. These procedures can be useful to perform the segmentation when a phonetic transcription is unavailable or inaccurate, or in applications like speaker or language identification systems, concatenative speech synthesis, among others [1], [4]. The authors have been proved that this algorithm gives better performance than other methods of the same class.

The results of the segmentation obtained with the modification here introduced are compared with those obtained using the original algorithm with a classical Melbank encoding and the CMD-based parametrization.

## II. METHODS

In this section we introduce the main characteristics of the speech encodings based on divergence, the text-independent algorithm used and the experiments performed.

### A. Parameterization based on Continuous Multiresolution Divergence

The stages of speech signal parameterization based on CMD are outlined in Fig. 1.

Given a discrete signal  $\vec{s} = \{s[k]\}$  of length  $K$ , we first obtain the discretized version of the continuous wavelet transform (CWT). This is the decomposition  $\{d[j, k]\} \in \mathbb{R}^{J \times K}$ , with scales  $j = 1, \dots, J$ ,  $J \in \mathbb{Z}$ , and time sample  $k \in \mathbb{Z}$ . For the sake of notational simplicity, for a fixed scale  $j$  the CWT coefficient's temporal evolution will be named as  $\{d_j[k]\}$  in what follows, with  $d_j[k] = d[j, k]$ .

We perform now a windowing over each  $d_j[k]$ , using a set  $W^j(m, L, \Delta) = \{d_j[k], k = l + m\Delta, l = 1, \dots, L\}$  of  $m = 0, 1, \dots, M$  rectangular sliding windows, which determine the frames of analysis. The selection of  $L$  and  $\Delta$  is accomplished in agreement with the windowing performed in order to obtain the Melbank parameterization of the speech signal used in [4].

In order to obtain the divergence along each scale  $j$ , we estimate first the probability distribution over each frame of analysis by means of the regular histogram. Accordingly, we perform an equipartition which provides a subset  $I_n^j$  of  $N$  disjoint subintervals. We denote with  $p_m^j(I_n^j)$  the probability that a given  $d_j[k] \in W^j(m, L, \Delta)$  belongs to the interval  $I_n^j$ . Therefore, for each window  $W^j(m, L, \Delta)$  a set  $P^j[m]$  of  $N$  probabilities  $p_m^j(I_n^j)$  is obtained. Observe that here  $m$  represents the time-evolution at the considered scale  $j$ .

Having in mind  $P^j[m]$ , we now consider a second set of  $N$  probabilities,  $R^j[m] = \{r_m^j(I_n^j), n = 1, \dots, N\}$ , corresponding to the next window  $W^j(m+1, L, \Delta)$ . Using  $P^j[m]$  and  $R^j[m]$  we can compute the Kullback-Leibler divergence over each set of consecutive windows:

$$\mathcal{D}_d[j, m] = \sum_{n=1}^N p_m^j(I_n^j) \ln \left( \frac{p_m^j(I_n^j)}{r_m^j(I_n^j)} \right). \quad (1)$$

Observe that here  $\mathcal{D}_d$  stands for  $\mathcal{D}_d(P, R)$ . The probability reference has been skipped in order to make the notation more readable. At each fixed scale  $j$  and for each fixed  $m$ , the divergence value corresponding to the wavelet coefficients on the window  $W^j(m, L, \Delta)$  is computed. This procedure, when accomplished for all the scales, gives the matrix  $\{\mathcal{D}_d[j, m], j = 1, \dots, J, m = 0, \dots, M\}$ , denoted as **CMD**, where  $CMD(j, m) = \mathcal{D}_d[j, m]$ , corresponds to the continuous multiresolution divergence.

Finally, the principal component analysis (PCA) method is used to extract the time-varying features that will compose the CMD-based parameterization. The matrix of principal components is computed as:  $\mathbf{Y} = \mathbf{Q}^T \mathbf{CMD}^*$ , where  $\mathbf{Q}$  is the eigenvector matrix of  $\sigma_{\mathbf{CMD}} = \mathbf{U}\mathbf{U}^T$ , and  $\mathbf{U} = \mathbf{CMD}^*$  is the statistical normalized matrix associated to **CMD**.

The CMD-based parameterization is obtained using the first eight rows of  $\mathbf{Y}$ , associated with the eight larger values of  $\mathbf{\Lambda}$  (the diagonal matrix of eigenvalues):  $\tilde{y}_i = \{y_i[m]\}$ , with  $i = 1, \dots, \mathcal{J}$  ( $\mathcal{J} = 8$  in this case). The elements  $y_i[m]$  of the components  $\tilde{y}_i$  are now the new features that represent the frame  $m$ .

The chosen  $\mathcal{J} = 8$  components for the new parameterization is in agreement with the amount of features of the Melbank encoding scheme used to perform the experiments. Moreover, from the point of view of PCA, the first eight components have more than 95% of the total variability of the divergence of the wavelet transform of the speech signal.

For more details on CMD computation, see [3].



Fig. 1. Scheme of the stages of the parameterization based on CMD.

### B. Text-independent speech segmentation algorithm

The speech segmentation algorithm, proposed by Esposito and Aversano [4], performs the segmentation task working on the time-varying features obtained through the short-term analysis of the speech signal. This is regulated by three operational parameters:  $\alpha$ ,  $\beta$  and  $\gamma$ , which are chosen in

agreement with the parameterization used as input. Figure 2 depicts its stages.

Parameter  $\alpha$  is used to estimate the intensity of an abrupt change through the jump function computation:

$$\mathcal{F}_i^\alpha[m] = \left| \sum_{\mu=m-\alpha}^{m-1} \frac{y_i[\mu]}{\alpha} - \sum_{\mu=m+1}^{m+\alpha} \frac{y_i[\mu]}{\alpha} \right|. \quad (2)$$

A relative thresholding procedure, with parameter  $\beta$ , identifies the frame  $m^*$  where a possible transition from one phoneme to another is localized. The relative height  $\eta$  is computed as:

$$\eta = \min \left\{ \mathcal{F}_i^\alpha[m^*] - \mathcal{F}_i^\alpha[u], \mathcal{F}_i^\alpha[m^*] - \mathcal{F}_i^\alpha[v] \right\}, \quad (3)$$

where  $\mathcal{F}_i^\alpha[u]$  and  $\mathcal{F}_i^\alpha[v]$  are two valleys of function  $\mathcal{F}_i^\alpha$  in the interval  $[u, v] \subset [\alpha, M - \alpha]$  and  $m^* \in [u, v]$

The frame  $m^*$ , corresponding to a peak of equation (2), is considered as a possible phone transition and stored in the binary matrix  $\mathbf{T} = \{T(i, m)\}$  (equal to 1 if a valid transition has been detected for the time-sequence  $i$  at the frame  $m$  and 0 otherwise) when  $\eta$  exceeds the threshold  $\beta$ .

It has been observed that sharp transitions do not occur simultaneously for each component of the speech features, even though they take place in a close time interval. A fitting procedure takes care of combining the different transition events detected in the neighboring of frame  $m^*$  into a unique indication of phone boundary. Parameter  $\gamma$  determines the width of the interval  $V = [m, m+1, \dots, m+\gamma-1]$  where this barycenter is individuated, using the function:

$$\mathcal{G}[c] = \sum_{\mu=c}^{c+\gamma-1} \sum_{i=1}^{\mathcal{J}} T(i, \mu) |\mu - c|, \quad c \in V. \quad (4)$$

The possible barycenter of interval  $V$  is the frame  $\tilde{c}$  where:

$$\mathcal{G}[\tilde{c}] = \min_{c \in V} \mathcal{G}[c], \quad m \leq \tilde{c} \leq m + \gamma - 1. \quad (5)$$

For each frame  $m$ , the value  $\tilde{\mathcal{G}}[m]$  indicates how many barycenters  $\tilde{c}$  have been found on it. This leads to a new function where the peaks correspond to the indication of a possible phone boundary.

In [4] various standard multi-band representations of speech signals have been tested, using different number of parameters. The authors have proved that the Melbank encoding with 8 coefficients, provides the best results. The Melbank parameterization is a standard short-term processing of speech signal, which is based on a bank-of-filters model [5].

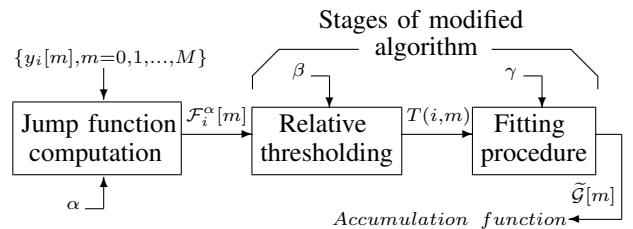


Fig. 2. Scheme of the segmentation algorithm described in Sec. II-B

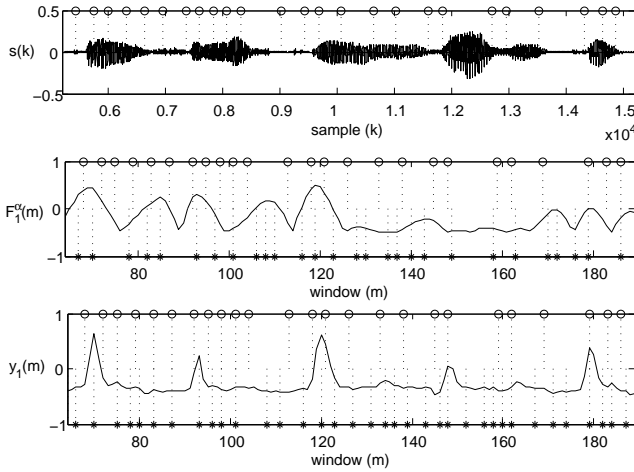


Fig. 3. (a) Speech signal with Albayzin labeling (upper dotted lines with circle markers). (b) Jump function computation  $\mathcal{F}_1^\alpha[m]$  corresponding to the Melbank feature  $y_1[m]$  of signal (a); upper dotted lines with circle markers indicate the database labeling and inferior dotted lines with star markers the segmentation obtained using the Melbank representation as input for the original algorithm. (c) Feature  $y_1[m]$  obtained using the CMD-based encoding of the signal displayed in (a); upper dotted lines with circle markers indicate the database labeling and inferior dotted lines with star markers the segmentation obtained using the modified algorithm. Operational parameters used for the automatic segmentation procedures are:  $\alpha = 6$ ,  $\beta = 0.05$  and  $\gamma = 3$ .

### C. Phone segmentation experiments

The speech signal encoding based on CMD (Sec. II-A) has been used as input for the text independent speech segmentation algorithm described in Sec. II-B, with good results [3]. In the present work, this parameterization is used not only as input for the algorithm but also to replace its first stage. In Fig. 2 are showed what stages are used with this new approach, that will be referred as the modified algorithm, to differentiate it from the original algorithm described in Sec. II-B.

As an example, we show in Fig. 3 the time evolution of one feature of Melbank speech parameterization and CMD-based encoding. In Fig. 3(a) a part of a labeled speech signal is shown. Fig. 3(b) shows the time evolution of the jump function,  $\mathcal{F}_1^\alpha$ , computed over the Melbank feature  $\bar{y}_1$ . The segmentation obtained with this parameterization is indicated with the inferior dotted lines with star markers. Fig. 3(c) depict the feature  $\bar{y}_1$  obtained through the CMD-based encoding described in Sec. II-A. The inferior dotted lines with star markers showed in 3(c) correspond to the segmentation obtained using the modified algorithm with the parameterization based on CMD. Upper dotted lines with circle indicate the database reference labeling in the three figures.

It can be observed from Fig. 3 that the tool based on CMD detect phone boundaries that are ignored by jump function, for example, the third inferior dotted line indication in 3(c)). Also the modified algorithm detects some points that are missing when jump computation is used, such as the sixth an seventh inferior dotted line indication of 3(c). These findings motivate the use of the CMD based parameterization as a part

of the algorithm.

The segmentation performance of the modification here proposed was compared with the original algorithm (II-B) using as inputs a Melbank parameterization and the CMD-based encoding described in II-A.

### D. Signals and Database

A subset of the Albayzin speech corpus, consisting of 600 sentences, 200 words vocabulary, was used [7]. Speech utterances had 3.55 sec. mean phrase duration, and they were spoken by 6 males and 6 females from the central area of Spain (average age 31.8 years). The labeled speech files, consisting in a hand-corrected HMM forced alignment segmentation, recorded the position of the phone boundaries expressed in milliseconds and were used as the reference segmentation. Each phrase in the corpus has been normalized in mean, pre-emphasized and Hamming windowed in segments of 20 ms length, shifted 10 ms. An 8 coefficients encoding were used for the three evaluated conditions.

### E. Indexes of segmentation performance evaluation

The percentage of correctly detected phone boundaries ( $PC$ ) and the percentage of erroneously inserted points ( $PI$ ) were computed in order to evaluate the obtained segmentation. The  $PC$  index relates the number of correctly detected boundaries,  $B_C$ , with the overall number of phone boundaries contained in the database,  $B_T$ , using a tolerance of  $\pm 20$  ms. The  $PI$  index relates the number of phone boundaries erroneously detected  $B_I = B_D - B_C$  ( $B_D$  is the whole number of segmentation points detected by the algorithm), with the total number of frames  $F_T$  in the signal:

$$PC = 100 \left( \frac{B_C}{B_T} \right) \quad PI = 100 \left( \frac{B_I}{F_T} \right). \quad (6)$$

Another indexes derived from the theory of signal detectability, similar to those defined above but mathematically more accurate, have been assessed: the false alarm rate  $P_{fa}$  and the missed detection rate  $P_{md}$ . They are defined as:

$$P_{fa} = \frac{B_I}{F_T - B_T} 100 \quad P_{md} = \frac{B_T - B_C}{B_T} 100. \quad (7)$$

With the values of  $P_{fa}$  and  $P_{md}$  we have constructed receiver operating characteristic (ROC) curves for each of the encoding schemes evaluated.

## III. RESULTS AND DISCUSSION

We now present and discuss the results of the phone segmentation task obtained with the speech encoding based on CMD (Sec. II-A) into the algorithm of Sec. II-B. We also compare the modification proposed in this paper with the original algorithm using both, Melbank and CMD-based parameterizations.

It can be observed for Table I, that for almost all the set of operational parameters, the modified algorithm gives better results when  $\gamma \leq 3$ . The proposed method increases the correctly detected bounds decreasing also the erroneously inserted points. This can be appreciated by comparing similar  $PC$  and  $PI$  indexes without having in mind the operational

TABLE I

PERCENTAGE OF CORRECTLY DETECTED PHONE BOUNDARIES ( $PC$ ) AND PERCENTAGE OF ERRONEOUSLY INSERTED POINTS ( $PI$ ) OBTAINED FOR THE MODIFIED ALGORITHM AND THE ORIGINAL ALGORITHM USING AS INPUTS THE MELBANK AND CMD-BASED PARAMETERIZATION. DIFFERENT OPERATIONAL PARAMETERS WERE TESTED. PARAMETER  $\alpha = 6$  WHEN ORIGINAL ALGORITHM WAS USED.

Parameters		Modified algorithm		Original algorithm			
$\gamma$	$\beta$	$PC$	$PI$	Melbank		CMD	
2	0.01	99.54	22.88	97.33	17.50	93.07	16.21
	0.05	97.81	14.10	95.44	16.10	91.27	9.26
	0.1	95.06	11.45	92.06	11.06	87.13	7.71
3	0.01	93.91	18.44	93.21	15.22	94.10	15.57
	0.05	94.62	11.22	90.44	12.01	91.47	9.83
	0.1	92.69	9.13	87.40	6.85	88.15	8.64
4	0.01	95.29	15.52	86.92	13.84	95.44	14.98
	0.05	94.17	9.25	83.55	9.21	93.05	9.98
	0.1	91.47	7.38	79.44	5.59	89.43	8.91

parameters used. For example, the modified algorithm gives  $PC=97.81\%$  and  $PI=14.10\%$  ( $\gamma = 2$ ,  $\beta = 0.05$ ) and the original algorithm using Melbank parameterization gives  $PC=97.33\%$  and  $PI=17.50\%$  ( $\gamma = 2$ ,  $\beta = 0.01$ ). Also, the modified algorithm give  $PC=94.62\%$  and  $PI=11.22\%$  ( $\gamma = 3$ ,  $\beta = 0.05$ ) and the original algorithm using now the CMD-based encoding give  $PC=94.10\%$  and  $PI=15.57\%$  ( $\gamma = 3$ ,  $\beta = 0.01$ ). Notice that there is no need to set parameter  $\alpha$  in the modified algorithm case. The optimal  $PC$  and  $PI$  indexes are 100% and 0% respectively.

In Fig. 4 we show the ROC curves, constructed with the  $P_{fa}$  and  $P_{md}$  indexes in order to compare the performance of the modified algorithm using two values for parameter  $\gamma$ . We also compare this approach with the original algorithm using the two encoding schemes considered in this work. High values of  $P_{fa}$  occur because many over-segmentations have been obtained. High values of  $P_{md}$  indicate that the algorithm has under-segmented the signal. As can be seen from the figure, when  $P_{fa}$  is low,  $P_{md}$  is high and viceversa. The closer the curve is to the bottom and to the left axes, the more accurate is the detection. We can see that the modified algorithm gives better results when  $\gamma = 3$ . It can be seen that this approach allows to dramatically reduce the missed detection rate (less than 20%) with low increments in the false alarm rate.

These results indicate that the CMD-based parameterization provides information related with abrupt changes in the speech signal, improving the phone segmentation algorithm performance. This approach is particularly suitable for accurate segmentation requirements, with low over-segmentation, and reduces the amount of parameters to be estimated in the segmentation algorithm.

#### IV. CONCLUSIONS

In this work we have proposed a modification to the algorithm introduced in [4], using the continuous multireso-

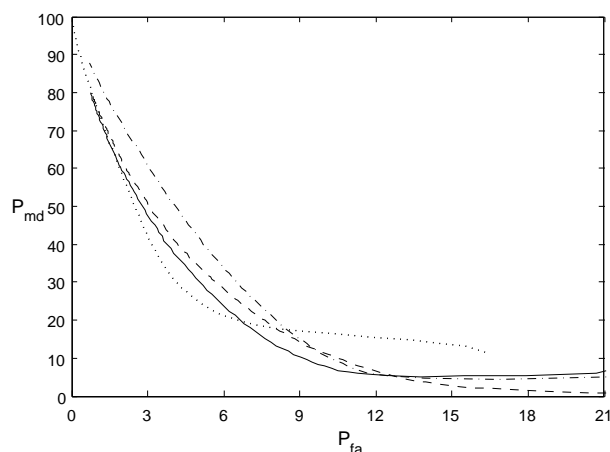


Fig. 4. ROC curves for the phone segmentation algorithm using the modified algorithm with  $\gamma = 3$  (solid line) and  $\gamma = 2$  (dashed line), and the original algorithm using Melbank (dotted line) and CMD-based parameterization (dash-dot line), with operational parameters  $\alpha = 6$  and  $\gamma = 3$ . Parameter  $\beta$  varies from 0 to 1 with increments of 0.01.

lution divergence, based on Kullback-Leibler distance. This encoding scheme has been used as input for the phone segmentation algorithm, replacing its first stage, corresponding to the jump function computation. An interesting side effect of this modification is the reduction in the number of tunable parameters of the algorithm. The performance of this approach has been compared with the original algorithm using as input the classical Melbank representation and a CMD-based parameterization. The results indicate that the modification proposed increases the algorithm ability to perform the segmentation task. The number of boundaries are correctly detected, with low increments in the amount of erroneously inserted points.

This demonstrates that dynamic measures based on CMD provide valuable information related to acoustic features that take into account transitions from one phoneme to another. Furthermore, this information needs no the additional computation of the jump function to perform the segmentation.

#### REFERENCES

- [1] G. Almpandis and C. Kotropoulos. Phoneme segmentation using the generalized gamma distribution and small sample bayesian information criterion. *Speech Comm.*, vol. 50, 2008, pp 38–55.
- [2] M. Añino, M. Torres, and G. Schlotthauer. Slight parameter changes detection in biological models: A multiresolution approach. *Phys. A*, vol. 324, 2003, pp 3–4.
- [3] A. Cherniz, M. Torres, H. Rufiner, and A. Esposito. Multiresolution analysis applied to text-independent phone segmentation. *J. Phys. Conf. Ser.*, vol. 90, 2007.
- [4] A. Esposito and G. Aversano. Text independent methods for speech segmentation. In G. Chollet et al., editor, *Nonlinear Speech Modeling And Applications: Advanced Lectures and Revised Selected Papers*, Springer, Berlin, Germany, 2005, pp 261–290.
- [5] L. Rabiner and B. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, Englewood Cliffs, NJ; 1993.
- [6] M. Torres, L. Gamero, P. Flandrin, and P. Abry. On a multiresolution entropy measure. In *SPIE'97, Wavelet Applications in Signal and Image Processing V*, Washington, USA, vol. 3169, 1997, pp 400–407.
- [7] J. Diaz Verdejo and A. Peinado and A. Rubio and E. Segarra and N. Prieto and F. Casacuberta. Albayzin: a task-oriented Spanish speech corpus. In *Proc. of the First Int. Conf. on Language Resources and Evaluation*, Granada, vol. I, 1998, pp 497–502.