

# Data Based Filter Design for RASTA-like Channel Normalization in ASR

*Carlos Avendano, Sarel van Vuuren and Hynek Hermansky*

Oregon Graduate Institute of Science & Technology  
P.O. Box 91000, Portland, OR 97291 USA

## ABSTRACT

RASTA processing has proven to be a successful technique for channel normalization in automatic speech recognition (ASR). We present two approaches to the design of RASTA-like filters from training data. One consists of finding the solution to a constrained optimization problem on the feature time trajectories while the other uses Linear Discriminant Analysis (LDA). Whereas LDA is often applied to one or a few frames of the feature vectors we apply LDA to feature time trajectories. Both approaches result in similar filters which are consistent with the ad hoc designed RASTA filter.

## 1. Introduction

Relatively unstructured data-driven systems are the mainstream in today's ASR. These systems acquire their structure from the large amounts of training data and are susceptible to failure when used in conditions that assume a different structure than that acquired during the training.

It is our belief that more knowledge-constrained and structured designs will result in simpler and ultimately more reliable systems. However, if we are to hardwire any constraints into the system, it is crucial that these constraints are based on well tested, reliable and relevant knowledge.

Some reasonable constraints may be implied by properties of the human hearing process and we have been relatively successful when incorporating them into ASR [1, 2]. On the other hand, it is hard to deny the power of real speech data. Thus, we support using the speech data, as long as they are used in a way to provide permanent and reusable knowledge.

Along this line, we came to realize that since speech developed to optimally use the properties of human auditory perception, any relevant auditory knowledge may have its counterpart in the structure of the acoustic speech signal, and the constraints derived from the data may either correct or support the knowledge-based constraints.

### 1.1. RASTA Processing

For the past several years we have been working on the incorporation of temporal auditory masking into speech processing [2]. As an engineering simulation of this powerful auditory constraint we pro-

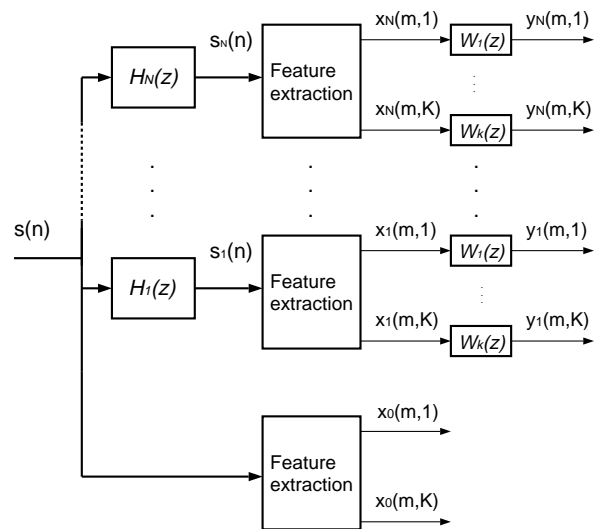


Figure 1: Problem setup block diagram

posed to filter out slow and fast changes in the trajectories of an critical logarithmic short-time spectrum of speech. The initial ad hoc form of the RASTA filters was optimized on a relatively small series of ASR experiments with noisy telephone digits. The experiments yielded a filter with a spectral zero at zero modulation frequency and a pass-band approximately from 1 to 12 Hz. The same filter was used for all frequency components.

Optimizations using ASR experiments are costly and there is no guarantee that the solutions obtained will not be specific to a given ASR problem. Any data-based optimization which would avoid using a specific ASR paradigm is desirable.

In the current work we use two optimization techniques to design a set of RASTA filters from realistic noisy data. The first is a constrained optimization technique which uses the natural time-varying structure of the continuous speech data. The second technique is based on a Linear Discriminant Analysis (LDA) [7] which uses phoneme-labeled speech data. The data used in both cases consisted of parallel speech recordings corrupted by different channels.

## 2. Filter Design by Constrained Optimization

In this section we find a set of filters for channel normalization by constructing a constrained optimization program. The design criteria is the distance minimization between the processed features when they are obtained from speech corrupted by different communication channels.

The idea can be explained with the aid of Fig. 1. A speech signal  $s(n)$  is corrupted by  $N$  different channels  $H_j(z)$ . After an appropriate feature extraction procedure we have a set of  $NK$  corrupted time trajectories  $x_j(m, k)$ , where  $k = 1, \dots, K$  is the feature index,  $m$  is the decimated time index and  $j=1, 2, \dots, N$ . Notice that we also show a set of trajectories of clean speech which will be used to set the constraints described in the next section.

Ideally we would like to apply a transformation  $W_k(z)$  on the time trajectories such that the outputs  $y_j(m, k)$  for the  $k$ th feature are as similar as possible to each other, that is, we seek channel independence. A trivial solution to this problem is to set  $W_k(z) = \bar{0}$  so the need to constrain the solution to some reasonable value is obvious. As will be seen later, the constraints will determine the behaviour of the filters at certain modulation frequencies.

### 2.1. Technique

Considering that the effect of the channel is multiplicative in the frequency domain we chose our transformation  $W_k(z)$  to be a linear finite impulse response (FIR) filter applied to  $x(m, k) = \log[a(m, k)]$ , where  $a(m, k)$  are critical band energy trajectories. Notice that in this case we assume one-dimensional (frequency specific) filters but the methods should still be applicable to a more general case. From now on we will call  $x(m, k)$  the critical bands omitting the term logarithmic.

The FIR filter for each critical band was derived in the following way:

Let

$$J_k = E \left\{ \sum_{j=1}^{N-1} \sum_{i=j+1}^N [y_j(m, k) - y_i(m, k)]^2 \right\} \quad (1)$$

be the objective function for the  $k$ th critical band.

From (1) we see that the objective functions are defined as the expected value (with respect to time) of the Euclidean distance between the outputs  $y_j(m, k)$  of the  $W_k(z)$  filter produced by  $x_j(m, k)$  for  $j = 1, \dots, N$ . These quadratic functions have a global minimum at  $W_k(z) = \bar{0}$ . To avoid the trivial solution in which all filter coefficients are set to zero we need to impose a set of constraints.

### 2.2. Experimental Design

To derive the filters, three parallel speech recordings were used. A sample of clean speech was taken from the TIMIT database (approximately 2 minutes). The other two samples were taken from the corresponding speech of the NTIMIT database (telephone channel) and TIMIT recorded through a cellular telephone channel. Auditory frequency band trajectories for the three recordings were computed by a

weighted sum of their short-term power spectrum as proposed in [1]. The logarithm of these trajectories was taken to produce  $x_0(m, k)$ ,  $x_1(m, k)$  and  $x_2(m, k)$  respectively.

The same design was applied to each critical band independently. For simplicity we will drop the frequency index  $k$  and it should be understood that the following procedure was applied  $K$  times.

Using (1) we find the objective function

$$J = \mathbf{w}^T R_{x_1, x_1} \mathbf{w} + \mathbf{w}^T R_{x_2, x_2} \mathbf{w} - \mathbf{w}^T R_{x_1, x_2}^v \mathbf{w} \quad (2)$$

In (2),  $\mathbf{w} = [w_0, w_1, \dots, w_{L-1}]^T$  is the vector of filter coefficients,  $R_{x_j, x_i}$  refers to the cross-correlation matrix between  $x_j(m)$  and  $x_i(m)$ , and  $R^v = R + R^T$ . In this case  $N = 2$  and we need to set two constraints ( $j=1, 2$ ), one for each of the outputs of the filter.

**Constraints** A reasonable constraint is to restrict the energy at the output of the filter to be a fraction of the energy of the input signal. While this constraint avoids the trivial solution, it imposes no restriction on the characteristics of the output signal and thus it will be less effective if we need to preserve relevant information about speech. Interestingly, by using this constraint we found that the filters had very similar responses (i.e. band-pass with strong dc suppression and narrow pass-band) as the so called delta cepstrum processing [4].

A more reasonable constraint was found by not allowing the distance between the filter outputs  $y_j(m, k)$  and the original uncorrupted speech features  $x_0(m, k)$  (see bottom of Fig. 1) to be large. This similarity constraint can be more or less restrictive depending on the amount of error allowed, and the resulting filters will have different characteristics depending on this factor. Notice that the availability of clean speech is only necessary for this particular constraint and is not a general requirement of the technique.

The constraints proposed above can be written as:

$$E\{[\hat{x}_0(m) - \hat{y}_j(m)]^2\} < c_j \quad (3)$$

In these constraints we removed the dc component from the original clean speech features and from the output of the filter to obtain  $\hat{x}_0(m)$  and  $\hat{y}_j(m)$  respectively. This normalization is needed in order to make a fair comparison of the clean and corrupted features (since adding or removing a constant in the logarithmic domain corresponds to modifying the power of the signals in the linear domain).

Writing (3) in matrix notation we get:

$$S_{\hat{x}_0} + \mathbf{w}^T R_{x_j, x_j} \mathbf{w} - 2\mathbf{w}^T \mathbf{p}_{\hat{x}_0, x_j} - \mathbf{w}^T \hat{\mathbf{x}}_j \hat{\mathbf{x}}_j^T \mathbf{w} < c_j \quad (4)$$

Here  $S_{\hat{x}_0}$  is the power of  $\hat{x}_0(m)$ ,  $\hat{\mathbf{x}}_j$  is a vector with all elements equal to the mean of  $x_j(m)$ , and  $\mathbf{p}_{\hat{x}_0, x_j}$  is the cross-correlation vector between  $\hat{x}_0(m)$  and  $x_j(m)$ .

The constants  $c_j$  were initially chosen using a -3dB criteria on the power of the inputs to the filter and were varied systematically until

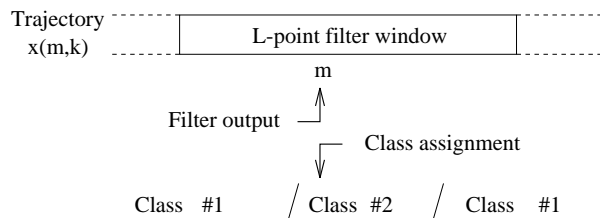
no feasible solution could be found. The last feasible point was chosen to be the new Constraint-Optimized (COP) filter. The optimization program described by (2) and (4) is non-linear with non-linear constraints and was solved using sequential quadratic programming [5].

### 3. Filter Design Using Linear Discriminant Analysis

In this section we find a set of filters using LDA <sup>1</sup>. Our approach is different from previous works [3, 8] in that we apply LDA to feature *time trajectories* rather than just to feature frames. We will call this technique DISCRIMINANT OPTIMIZED (DISCO) filtering.

#### 3.1. Technique

We start with speech which is labeled into different classes. We wish to design a filter that will maximize discriminability between these classes. The technique is best explained with the aid of Fig. 2. The figure shows an FIR filter acting on the  $k$ 'th time trajectory of feature  $x(m)$  (which in this case is obtained by concatenation of  $x_0(m, k), x_1(m, k)$  and  $x_2(m, k)$ ).



**Figure 2:** DISCO filtering technique. Illustrated is an FIR filter centered at time  $m$  and acting on the  $k$ 'th time trajectory of feature  $x$ .

Referring to the figure, view the features on which the filter is acting at time  $m$  as a vector  $x_m$ . For a filter of length  $L$ , this vector lies in an  $L$ -dimensional space. For every analysis step  $m$ , the output of the filter is a projection of the  $L$ -dimensional input vector onto a one-dimensional output space. Assigning a class to each vector  $x_m$ , we can use LDA to find the principal directions for which class discriminability is maximum at the output. The first principal component of this analysis is the  $L$ -point FIR filter for which linear discriminability is maximized. A separate LDA is carried out for each critical band time trajectory, thus yielding a separate FIR filter for each band.

Fig. 2 also gives an idea of how we assign a class to the input vector  $x_m$ . In the current work, we assume a non-causal FIR filtering. Thus, we align the center of the input vector  $x_m$  with the underlying label. This is not the only way to do the class assignment; e.g. had we decided for a causal FIR filter we would have aligned the class label with the first element of the  $x_m$ . The next input vector is  $x_{m+1}$ , i.e. it is formed by shifting the  $x_m$  by one analysis step.

In our work the vector space for the LDA is constructed from seg-

<sup>1</sup> As far as we know LDA was first used in speech by Melvyn Hunt [3].

ments of time trajectories of a single speech feature over a relatively long (typically at least a syllable) span of time. This is different from other LDA-based approaches (eg. [8]) that apply LDA to a feature vector or to a relatively short block of feature vectors.

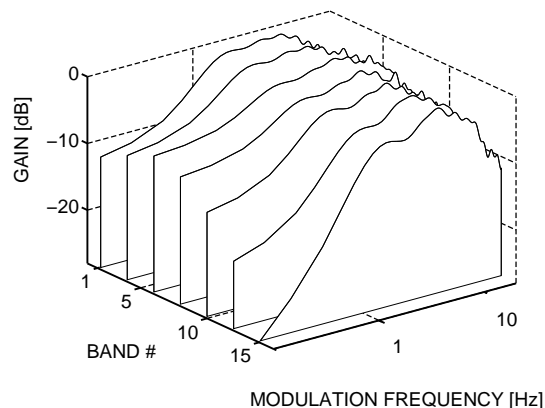
#### 3.2. Scope

In this paper we find a filter for each feature trajectory  $x(m, k)$  separately but it should be understood that our technique is much more general:

- Our technique in principle extends to multivariate filtering or non-linear discriminants [7].
- It could allow for more than just the principal component of the LDA to be used.
- It allows for causal or non-causal filter design.
- It is not limited to any particular domain (e.g. since the aim was to remove linear distortions, the domain was the logarithmic critical band power spectrum while if the aim was to remove additive noise the domain might have been different).
- It applies to different types and sets of classes. In this paper we use broad TIMIT phonetic classes<sup>2</sup>.

### 4. Results

In this section we present the results obtained by the techniques described above. If no distinction is made the reader should assume that the results discussed apply to both approaches.



**Figure 3:** DISCO filter bank

Experimental designs were obtained for every critical band (in this case we used 15 bands) using both approaches and using the same data. For different bands we found no substantial difference in the frequency responses of the filters except for bands where speech had little or no energy and were dominated by noise. We illustrate this fact by showing the magnitude frequency responses for the DISCO

<sup>2</sup> We used the following classes: closure, back vowel, mid vowel, front vowel, fricative, stop, flap and affricate, retroflex, and nasal, but did not use silence and non-speech.

**Table 2:** Search effort and recognition results for the word conditioned (WC) and the time conditioned (TC) tree search methods on the NAB'94 H1 development data.

LM type	search method	average number of					DEL-INS-SUB = TOTAL	WER [%]
		states	arcs	trees	word ends	LM recomb.		
bigram $PP_{bi} = 205.4$	WC	5022	1416	20	86	86	185-198-879 = 1262	17.1
		9335	2593	28	156	156	181-197-860 = 1238	16.8
		26493	7222	45	540	540	179-195-846 = 1220	16.5
		52786	13968	60	560	560	179-193-842 = 1214	16.4
	TC	6175	1781	27	117	1539	177-217-884 = 1278	17.3
		11280	3101	31	197	2976	180-199-857 = 1236	16.7
		30019	7998	37	520	8761	179-195-845 = 1219	16.5
		60692	15877	43	1075	18081	179-193-840 = 1212	16.4
trigram $PP_{tri} = 135.5$	WC	4714	1392	29	80	80	128-211-757 = 1096	14.8
		18734	5430	70	294	294	120-206-721 = 1047	14.2
		48940	13877	112	702	702	120-204-717 = 1041	14.1
		86772	23688	145	1223	1223	121-200-709 = 1030	13.9
	TC	8573	2459	28	136	1477	119-216-735 = 1070	14.5
		15103	4194	31	234	2806	118-205-720 = 1043	14.1
		24914	6780	34	388	4858	118-201-717 = 1036	14.0
		71757	18953	42	1202	14293	118-199-715 = 1032	14.0

is between 27 and 50, which seems to correlate with the average word duration.

So far, we have not considered the computational effort for the language model recombination. This computational effort is much greater in the time-conditioned search, where, for each word triple (or word pair) under consideration, a separate optimization over the unknown word boundary has to be carried out. In contrast, in the word conditioned method, the optimization over the word boundaries is already taken into account when starting up the trees. Typically in the time conditioned search, the effort for the language model recombination is higher by a factor of approximately 10 for both bigram and trigram language models. All in all, the experiments indicate that both search methods are suitable for large vocabulary recognition and are comparable in terms of efficiency.

**Acknowledgment.** This work was partly carried out within the project MAIS and supported by the Commission of the European Community (MLAP LRE-63-036).

## REFERENCES

1. H. Aust, M. Oerder, F. Seide, V. Steinbiss: "Experience with the Philips Automatic Train Timetable Information System", Proc. Second IEEE Workshop on Voice Technology for Telecommunications Applications, pp. 67-72, Kyoto, September 1994.
2. C. Dugast, R. Kneser, X. Aubert, S. Ortmanns, K. Beulen, H. Ney: "Continuous Speech Recognition Tests and Results for the NAB'94 Corpus", Proc. ARPA Spoken Language Technology Workshop, Austin, TX, pp. 156-161, January 1995.
3. R. Haeb-Umbach, H. Ney: "Improvements in Time-Synchronous Beam Search for 10,000-Word Continuous Speech Recognition", IEEE Trans. on Speech and Audio Processing, Vol. 2, pp. 353-356, April 1994.
4. F. Kubala: "Design of the 1994 CSR Benchmark Tests", Proc. ARPA Spoken Language Technology Workshop, Austin, TX, pp. 41-46, January 1995.
5. H. Ney, D. Mergel, A. Noll, A. Paeseler: "Data Driven Organization of the Dynamic Programming Beam Search for Continuous Speech Recognition", IEEE Trans. on Signal Processing, Vol. SP-40, No. 2, pp. 272-281, Feb. 1992.
6. H. Ney: Search Strategies for Large-Vocabulary Continuous-Speech Recognition. NATO Advanced Studies Institute, Bunion, Spain, June-July 1993, pp. 210-225, in A.J. Rubio Ayuso, J.M. Lopez Soler (eds.): "Speech Recognition and Coding - New Advances and Trends", Springer, Berlin, 1995.
7. J.J. Odell, V. Valtchev, P.C. Woodland, S.J. Young: "A One-Pass Decoder Design for Large Vocabulary Recognition", Proc. ARPA Spoken Language Technology Workshop, Plainsboro, NJ, pp. 405-410, March 1994.
8. M. Oerder, H. Ney: "Word Graphs: An Efficient Interface Between Continuous Speech Recognition and Language Understanding", Proc. Int. Conf. on Acoustics, Speech and Signal Processing, Minneapolis, MN, Vol. II, pp. 119-122, April 1993.
9. S. Ortmanns, H. Ney, A. Eiden: "Language-Model Look-Ahead for Large Vocabulary Speech Recognition", Proc. Int. Conf. on Spoken Language Processing, Philadelphia, PA, October 1996.
10. H. Sakoe: "Two-Level DP Matching - A Dynamic Programming-Based Pattern Matching Algorithm for Connected Word Recognition", IEEE Trans. on Acoustics, Speech and Signal Processing, Vol. ASSP-27, pp. 588-595, December 1979.
11. V. Steinbiss, B.-H. Tran, H. Ney: "Improvements in Beam Search", Proc. Int. Conf. on Spoken Language Processing, Yokohama, pp. 2143-2146, September 1994.