

# DISTORTION-CLASS WEIGHTED ACOUSTIC MODELING FOR ROBUST SPEECH RECOGNITION UNDER GSM RPE-LTP CODING

Juan M. Huerta and Richard M. Stern

Carnegie Mellon University  
Department of Electrical and Computer Engineering  
Email: {juan, rms}@speech.cs.cmu.edu

## ABSTRACT

We present a method to reduce the effect of GSM RPE-LTP coding by combining two sets of acoustic models during recognition, one set trained on distorted speech and the other trained on clean speech. During recognition, the posterior probabilities of an utterance are calculated as a sum of the posteriors of the individual models, weighted according to the distortion class each state in the model represents. We analyze the origin of spectral distortion to the long-term residual introduced by the RPE-LTP coding and discuss how this distortion varies according to phonetic class. For the Research Management corpus, the proposed method reduces the degradation in frame-by-frame phonetic recognition accuracy introduced by coding by more than 20 percent.

## 1. INTRODUCTION

Speech coding reduces the accuracy of speech recognition systems [4]. As speech recognition application in cellular and mobile environments becomes ubiquitous, robust recognition in these conditions becomes crucial. Even though the speech coding introduces one of the several factors that contribute to the degradation in recognition accuracy, it is necessary to understand and compensate for this degradation in order to achieve a system's full potential. We focus on the acoustic distortion introduced by the full-rate GSM coding [1]. This distortion can be traced to the quantization of the log-area ratio (LAR) and to the quantization of the downsampling performed in the RPE-LTP process. The distortion introduced in the residual signal affects recognition to a larger extent than the quantization of the LAR coefficients undergo [2].

In Section 2 of this paper, we discuss the origin of the distortion in the RPE-LTP codec. We observe that based on the "predictability" of the short-term residual signal, the RPE-LTP will be able to minimize the error in the quantized long-term residual. This predictability is later shown to be related to phonetic characteristics of the signal. In Section 3 we show that the relative spectral distortion introduced in the quantized long-term residual tends to be concentrated around three levels, and that the amount of relative spectral distortion can be loosely related to the relative degradation introduced in the recognition accuracy by GSM coding. In Section 4 we

separate the set of phonemes into cluster their relative spectral distortion. In Section 5 we describe a method to weight two sets of acoustic models the distortion categories introduced in Section 6. We describe the results of recognition experiments techniques in Section 6.

## 2. THE RPE-LTP CODEC AS A SOURCE OF ACOUSTIC DEGRADATION

The full-rate GSM codec decomposes the speech signal into a set of LAR coefficients and a short-term residual. The LAR coefficients are quantized and transmitted to the decoder while the short-term residual is segmented into subframes and coded by an LTP code. In this section we discuss the degradation that exist in the RPE-LTP code. Figures 1 and 2, a simplified version of the RPE-LTP codec and a simplified version of a RPE-LTP decoder respectively.

Figure 1 depicts an ideal codec able to reconstruct a signal that is identical to the original signal. Even though this codec would not achieve a full reconstruction of the input sequence it becomes crucial. Even though the speech coding introduces one of the several factors that contribute to the degradation in recognition accuracy, it is necessary to understand and compensate for this degradation in order to achieve a system's full potential. We focus on the acoustic distortion introduced by the full-rate GSM coding [1]. This distortion can be traced to the quantization of the log-area ratio (LAR) and to the quantization of the downsampling performed in the RPE-LTP process. The distortion introduced in the residual signal affects recognition to a larger extent than the quantization of the LAR coefficients undergo [2].

The ideal RPE-LTP codec works as follows: the short-term residual (the short-term residual) and the long-term residual (the long-term residual) are reconstructed and compared to a predicted sequence (the short-term residual) and the long-term residual (the long-term residual) produced by the predictor and block (the LTP section). The difference between these two signals is what the predictor is unable to predict, which we will refer to as the "innovation sequence". In the absence of acoustic degradation, our ideal codec will generate a reconstructed time sequence that is identical to the original signal. This reconstructed sequence will be identical to the original signal by definition. In reality, the RPE-LTP coder, does not transmit the necessary information for the decoder to reconstruct the signal exactly. The information that will result in the reconstructed signal (called the reconstructed signal) will approximate the original signal to the extent possible. The energy of the innovation signal which is on how good the predictor module (RPE)

"follow" or predict the next subframe of the time observed in a frame. We observe sequence based on previous reconstructed subframes

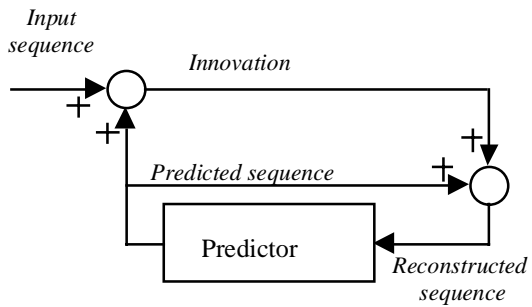


Figure 1. Simplified block diagram of an ideal RPE

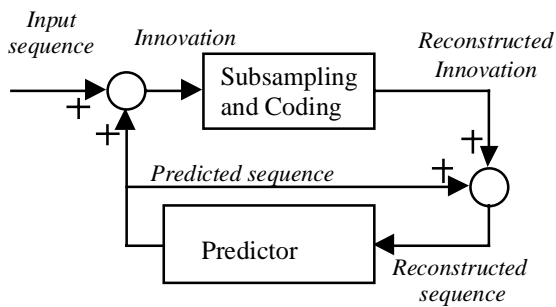


Figure 2. Simplified block diagram of an actual RPE

### 3. THE RPE-LTP INDUCED SPECTRAL DISTORTION

#### 3.1 Relative spectral distortion introduced by the RPE-LTP coder

We use the relative log spectral distance to measure the dissimilarity between the reconstructed and the original innovation sequences. Let  $E(\omega)$  be the power spectrum of the innovation sequence and  $E_R(\omega)$  be the power spectrum of the quantized innovation sequence. We integrate the absolute value of the difference between the log power spectra and normalize it by the log of the power spectrum of the innovation signal.

$$L_1 = \int_0^\pi |\log(E(\omega) - \log(E_R(\omega)))| d\omega / \int_0^\pi |\log(E(\omega))| d\omega$$

#### 3.2 Distribution of the relative log spectral distortion

We computed the relative log spectral distortion introduced by the RPE-LTP coder to the RM corpus. Figure 3 is a histogram that shows the log frequency of observing various levels of relative log spectral distortion. The horizontal axis is the amount of relative

distortion introduced by the RPE-LTP coder. We observe that the frames suffer only a small amount of distortion, and high distortion, separated by breakpoint 33 and 67 percent. We also observe that most of the time the LTP section of the coder do a reasonably good job of predicting the residual signal.

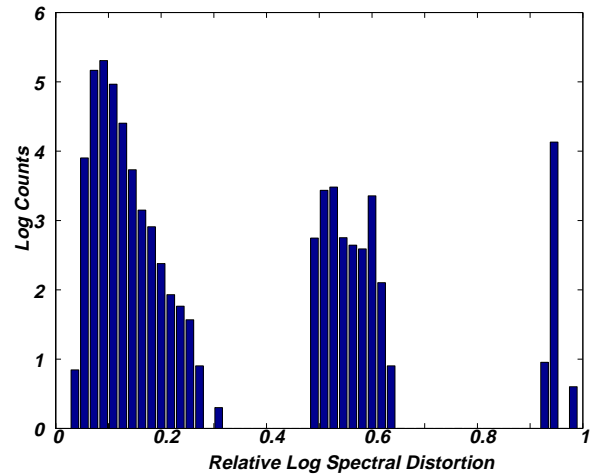


Figure 3. Histogram of the log distribution of relative log spectral distortion introduced to a large sample

#### 3.3 Impact of relative log spectral distortion on phonetic recognition

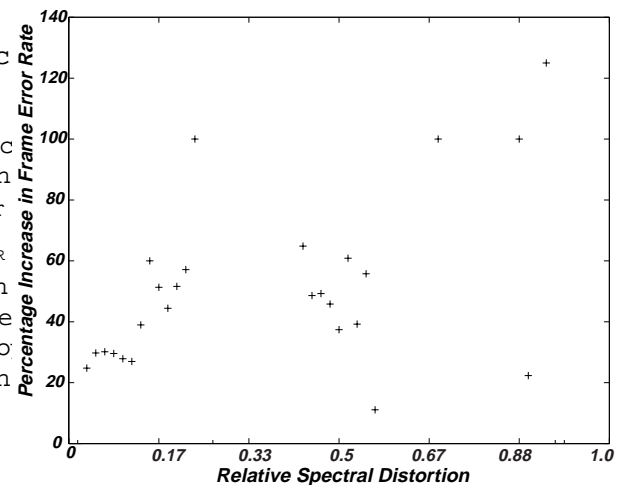


Figure 4. Relative degradation in frame-based phonetic accuracy as a function of the relative spectral distortion introduced by the RPE-LTP coder.

In order to analyze the relation that exists between relative recognition performance and the amount of relative log spectral distortion introduced by the RPE-LTP block, we performed two phonetic recognition experiments: one training and testing using

non-GSM coded) speech data and a second experiment of the normalized histograms of the using speech that underwent GSM coding. We show below. These five classes exhibit the frame accuracy for each amount of relative spectral distortion in each of the three distortion both for when GSM was present and when it was absent in Section 3.2. Classes 1 and 3 was absent. Figure 4 depicts how much the frame-based phonetic recognition error rate increased due to the introduction of the GSM coding as a function of relative spectral distortion. We observe that most generally concentrated in the low speaking, the phonemes that produced a greater amount of distortion due to GSM coding also suffered greater amounts of frame error rate.

#### 4. RELATING RELATIVE LOG SPECTRAL DISTORSION PATTERNS TO PHONETIC CLASSES

##### 4.1 Clustering Phonetic-classes using the relative log-spectral distortion

We grouped the 52 phonetic units used by our system into phonetic clusters by incrementally the closest histograms of the counts of the spectral distortion for the frames associated with each phonetic unit in the corpus. The distance between a pair of histograms was calculated using normalized histograms. The clustering yielded five classes, as shown in the table below.

Class	Phonemes in class	Frame Acc.	Frame Acc., GSM	% Degradation
1	P D K D	46.76%	44.10%	5%
2	I X B B D D D D F M N J H N G V W Y D X Z H A X S H	68.28%	58.34%	31.34%
3	D G K P T T D C H	70.88%	65.94%	16.96%
4	F H H S Z T H T S	72.39%	57.63%	53.43%
5	I Y O W O Y U W L R A A A E A H A X R E H A O I H E R A W A Y U H E Y	74.37%	67.90%	25.24%

Table 1. Phonetic classes generated by automatic clustering the phonemes distortion histograms

Classes 1, 3, and 4 encompass the majority of the consonants, Class 4 being mostly fricatives. Class 5 includes the remaining consonants and a couple of vowels, while Class 2 encompasses the vowels, diphthongs and semivowels. Silence units were omitted from this analysis. We can see that even though the classes do not exactly correspond to phonetic classes, they achieve a reasonable partition between broad classes of sounds. The pattern of distributions of relative log-spectral distortion introduced by the RPE-LTP is strongly related to phonetic properties of the signal.

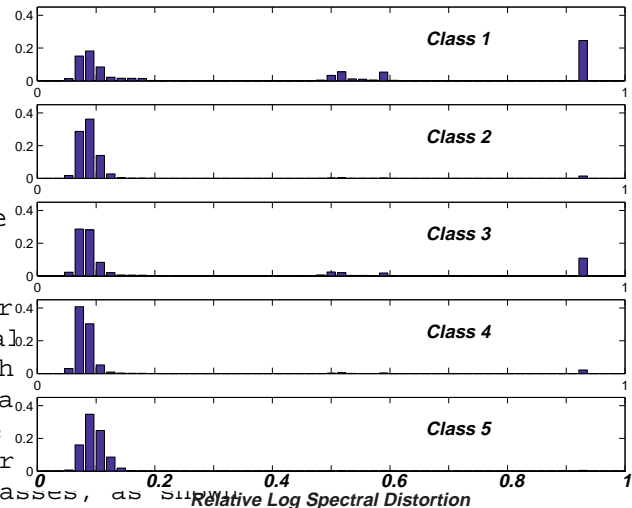


Figure 5 Histograms of the relative spectral distortion for five phonetic classes of Table 1.

#### 5. ACOUSTIC MODEL WEIGHTING

Acoustic modeling for HMM-based speech recognition commonly makes use of mixtures of Gaussian distribution representing a set of tied states. The probability that an observed vector has been in a certain state is thus expressed by

$$b_j(o_t) = \sum_{k=1}^{M_j} c_{j,k} N(o_t, \mu_{j,k}, C_{j,k})$$

The term  $c_{j,k}$  expresses the prior probability of the  $k$ th Gaussian component of the  $j$ th HMM. For a given state  $j$ , the sum of  $c_{j,k}$  over all  $k$  is equal to 1.

We can consider the amount of distortion that a phonetic class undergoes while evaluating a probability using several models that represent different distortion regions. We can express this by introducing a function  $f$  that weights the probability of the  $k$ th model depending on the distortion of the observed frame, and indicate the model representing a given phonetic class. This function also depends on the prior probability of the  $k$ th model. This function can also be expressed as a weighted version  $\sum_{k=1}^{M_j} c_{j,k} p_{j,k}$ . The resulting probability becomes

$$b_{j,d}(o_t) = \sum_{p=1}^2 \sum_{k=1}^M f(c_p[j,k], d_t, j) N_p(o_t, \mu_{jk}, C_{jk})$$

The function can also be dependent on the distortion class the model represents. This weight more the clean models for states that model phonemes that suffer small average GSM distortion. Alternatively, one can make the function  $f$  depend on knowledge of the instantaneous relative distortion of each frame if information is available. This function should give more weight to the distorted models when the relative distortion is greater.

## 6. SPEECH RECOGNITION EXPERIMENTS

Phonetic recognition experiments were performed using the Resource Management corpus and the SPHINX system. Our basic acoustic models consisted of 2500 tied states, each modeled by a mixture of two Gaussians. One set of models was trained on clean speech while the other set of models was trained with speech that underwent GSM coding. We obtained our reference transcription by performing automatic forced alignment on the 600 RPE-10P utterances using their pronunciations. We compared with the phone accuracy and the frame accuracy. Frame accuracy is calculated considering insertions and deletions of phones in the hypothesis. We evaluated the frame accuracy by making a frame-by-frame comparison of the output of the phonetic recognizer with the acoustic models that had been trained on clean and GSM speech, using clean and GSM models. We also show the results when models are trained using both GSM-coded speech and clean speech (multi-style trained models).

Testing data	Training data	Phone Accuracy	Frame Accuracy
Clean	Clean	63.3%	69.86%
GSM	Clean	55.7%	61.20%
GSM	GSM	60.1%	64.50%
GSM	GSM+Clean Multistyle	60.1%	64.64%

Table 2. Baseline frame accuracy and phoneme recognition accuracy in Resource Management

From Table 2 we can see that the effect of having GSM codec distortion during training and recognition reduces frame accuracy by 5.35 points. Recognition accuracy can be improved somewhat by using multi-style trained models.

### 6.1 Experiments using weighted acoustic models

We performed recognition experiments using weighted acoustic models by considering different mixing factors for two acoustic models, trained GSM-coded and clean speech. Table 3 shows results weighting both models equally. We

also made these weights dependent of the phone using three automatically-clustered phonetic models differently, either of these two weights provide better results than using multi-class models. The 3-class optimized weights reduce degradation in frame error rate introduced by coding by 27%, from 5.36% to 3.87%. The frame weighting improves recognition accuracy approximately equally for all three phonetic

Weighted modeling	Phone Accuracy	Frame Accuracy
Equal weights	61.2%	65.77%
3-class optimized weights	61.2%	65.99%

Table 3. Frame accuracy obtained when using equally weighted models and when using phonetic-class based weights

## 7. CONCLUSIONS

In this paper we examined the distortion in the RPE-10P GSM codec to the short-term parameters related to the relative distortion produced by the quantized long-term parameters. We observed that this distortion influences the performance and that it can be related to the properties of the speech signal. We introduced a mixture of acoustic models that had been trained on clean and distorted speech by taking into consideration the models of distortion suffered by the phonetic

## 8. REFERENCES

- [1] European Telecommunication Standards Institute, "Phase 1 of the European digital telecommunications system (Phase 1) digital speech processing functions (GSM 06.01)", 1994.
- [2] Huerta, J. M. and Stern R. M., "Speech Recognition with GSM Codec Parameters" Proc. ICSLP-98, 1998.
- [3] Kroon, P., Deprettere, E. F., Sluyter, R. F., "Excitation - A Novel Approach to Effective and Efficient Multi-pulse Coding of Speech", IEEE Trans. on Speech and Signal Processing, 34:1054-1063, 1986.
- [4] Lilly, B. T. and Palani, K. K. "Effect of Speech Codec Distortion on Speech Recognition Performance", Proc. ICSLP-96, 1996.