

ENVIRONMENT NORMALIZATION FOR ROBUST SPEECH RECOGNITION USING DIRECT CEPSTRAL COMPARISON

Fu-Hua Liu, Richard M. Stern, Pedro J. Moreno, and Alejandro Acero

Department of Electrical and Computer Engineering
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

ABSTRACT

In this paper we describe and evaluate a series of new algorithms that compensate for the effects of unknown acoustical environments (or changes in environment) through the use of compensation vectors that are added to the cepstral representations of speech that is input to a speech recognition system. These compensation vectors are obtained from direct frame-by-frame comparisons of the cepstral representations of speech that is simultaneously recorded in the training environment and various testing environments, but the algorithms do not make use of such “stereo” speech data in analyzing speech from an unknown environment. In the proposed paper we will compare the improvement in recognition accuracy provided by the algorithms using common standard ARPA speech recognition corpora. For example, the normalization algorithm known as MFCDCN provided a 22% reduction in word error rate when compared to results obtained using cepstral mean normalization on the 1992 ARPA WSJ/CSR corpus, and a 56.6% reduction in error rate compared to baseline processing. A family of new algorithms, PDCN, which accomplish the environment normalization inside the decoder are described and evaluated in the same corpus. A substantial word error rate reduction, 66.8%, can be achieved by combining MFCDCN and PDCN in the system with cepstral mean normalization compared to baseline system.

1. INTRODUCTION

The need for speech recognition systems and spoken language systems to be robust with respect to their acoustical environment has become more widely appreciated in recent years (*e.g.* [1]). Many approaches have been considered in the development of robust speech recognition systems including techniques based on autoregressive analysis, the use of special distortion measures, the use of auditory models, and the use of microphone arrays, among many other approaches (as reviewed in [1,2]).

Over the past few years, CMU and other sites have developed a series of algorithms that reduce the effects of environmental variability on speech recognition accuracy [*e.g.* 2,3]. The CMU normalization algorithms are based on three different types of approaches. The first approach is that of *cepstral remapping based on a structural model of the acoustical degradation*. An example of this type of processing is the codeword-dependent cepstral normalization algorithm (CDCN), which assumes that the effects of environmental distortion can be characterized as unknown addi-

tive noise combined with unknown linear filtering [4]. The second approach to environmental normalization is that of *high-pass filtering of cepstral coefficients*, as exemplified by the various RASTA algorithms [5] and the practice of cepstral mean removal. The third approach, which is the focus of this paper, is based on *direct cepstral comparisons* of simultaneously-recorded data from different environments on a frame-by-frame basis. We describe some of the more useful cepstral-comparison algorithms in the next section.

2. ENVIRONMENTAL NORMALIZATION USING CEPSTRAL COMPARISON

Environment-normalization algorithms based on cepstral comparison all assume that differences between the training and testing environments can be characterized by an additive correction to the cepstral vectors that represent the speech. The compensation vectors are estimated empirically on the basis of direct frame-by-frame comparisons of the cepstral representations of speech that is simultaneously recorded in the training environment and various testing environments (“stereo data”). The individual algorithms differ in the way the compensation vectors are estimated from training data, and in the way in which the need for stereo data is circumvented when the recognition system analyzes speech from an unknown environment. This general approach has become much more popular with the availability of the ARPA Wall Street Journal corpus, which in its initial phase contained about 31,000 utterances of stereo data recorded in 16 different acoustical environments.

2.1. The SDCN and FCDCN algorithms

SDCN. The simplest compensation algorithm, *SNR-Dependent Cepstral Normalization* (SDCN) [2], applies an additive correction in the cepstral domain that depends exclusively on the instantaneous SNR of the signal. This compensation vector equals the average difference in cepstra between simultaneous stereo recordings of speech samples from both the training and testing environments at each SNR in the testing environment. At high SNRs, this compensation vector primarily compensates for the effects of unknown linear filtering, while at low SNRs the vector provides a form of noise subtraction. The SDCN algorithm is simple and effective, but it requires environment-specific training.

FCDCN. *Fixed codeword-dependent cepstral normalization* (FCDCN) [2] is similar to SDCN, but it provides a greater number of compensation vectors. At each SNR the observed cepstra in the testing environment are also clustered, based on a VQ codebook.

The FCDCN algorithm applies an additive correction that depends on both the instantaneous SNR of each frame of input speech, and the VQ codeword location to which the cepstral compensation vector is closest. FCDCN compensation provides greater recognition accuracy than SDCN, but it also requires environment-specific training.

Figure 1 illustrates some typical compensation vectors obtained with the FCDCN algorithm, computed using the ARPA standard close-talking Sennheiser HMD-414 microphone and the unidirectional desktop PCC-160 microphone used as the testing environment. The vectors are computed at the extreme SNRs of 0 and 29 dB, as well as at 5 dB. These curves are obtained by calculating the cosine transform of the cepstral compensation vectors, so they provide an estimate of the effective spectral profile of the compensation vectors. The horizontal axis represents frequency, warped nonlinearly according to the mel scale. The maximum frequency corresponds to the Nyquist frequency, 8000 Hz. We note that the spectral profile of the compensation vector varies with SNR, and that especially for the intermediate SNRs the various VQ clusters require compensation vectors of different spectral shapes.

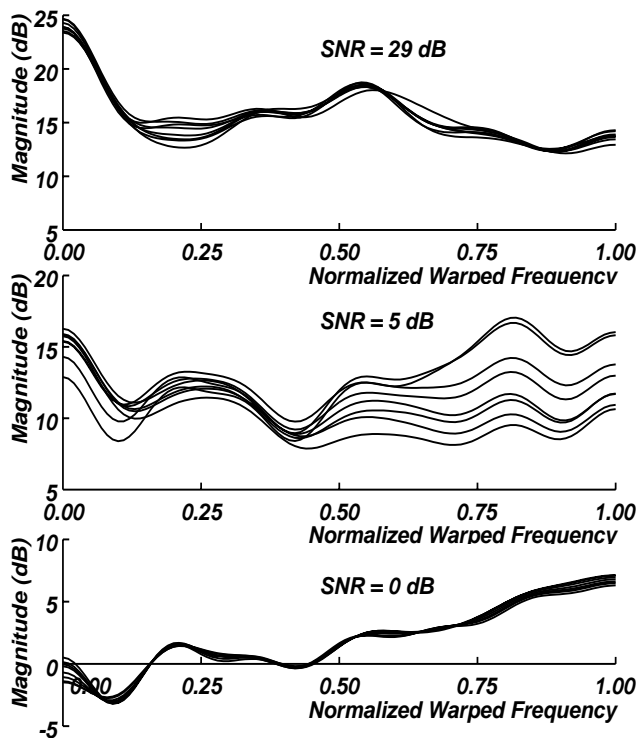


Figure 1: Comparison of compensation vectors using the FCDCN method with the PCC-160 unidirectional desktop microphone, at three different signal-to-noise ratios. The maximum SNR used by the FCDCN algorithm is 29 dB.

2.2. MFDCN and related algorithms

MFDCN. Multiple fixed codeword-dependent cepstral normalization (MFDCN) is a simple extension of the FCDCN algorithm [2] that has the advantage of not requiring that the identity of the testing environment be known *a priori*. In MFDCN, compensation vectors are precomputed in parallel for each of a set of testing environments, using the FCDCN training procedure. When a testing utterance from some unknown environment is input to

the recognition system, the system first determines which of the testing environments in the training data is most similar to the current testing environment. Compensation vectors for the chosen testing environment are applied to normalize the utterance according to the expression

$$\hat{\mathbf{x}}_i = \mathbf{z}_i + \mathbf{r}[k_i, l_i, e]$$

where k_i , l_i and e are the VQ codeword index, instantaneous frame

SNR, and the index of the chosen environment, respectively, and $\hat{\mathbf{x}}$, \mathbf{z} , and \mathbf{r} are the compensated (transformed) data, original data and compensation vectors, respectively.

Environment selection. We have made use of two schemes for environment selection. In the first procedure, referred to as *selection by compensation*, compensation vectors computed using each of the possible testing environments are applied successively to the incoming test utterance. The environment e is chosen that minimizes the average residual VQ distortion (D_e) over the entire utterance as follows,

$$D_e = \sum_i \|\mathbf{z}_i + \mathbf{r}[k_i, l_i, e] - \mathbf{c}[k_i]\|$$

where e is the index for the testing environment and $\mathbf{c}[k]$ is the k^{th} codeword.

In the second approach, referred to as *environment-specific VQ*, codebooks that are specific to each environment are generated from the original uncompensated speech. Environment selection is accomplished by vector quantizing the incoming test utterance using each environment-specific codebook in turn and choosing the (uncompensated) testing environment that is closest to the incoming speech in terms of VQ distortion. Using data from the 11/92 DARPA Wall Street Journal corpus, the selection-by-compensation method produces environment-selection errors 28.8% of the time for data from one of the 15 “secondary” environments and no selection errors for data obtained using the close-talking Sennheiser microphone used in the training data. The environment-specific VQ approach produces a 14.2% misjudgment rate for data using secondary microphones and 0.3% for Sennheiser mic data. Both methods produce similar speech recognition accuracy. The latter method is similar in spirit to the approach used by BBN [4], in which a classification is performed to select one of seven groups of acoustical environments for each incoming utterance.

Interpolated FCDCN. The MFDCN algorithm described above applies compensation from the single environment in the training set that is believed to have acoustical characteristics that most closely resemble those of the testing environment. In some cases, however, the testing environment does not closely resemble any single environment in the training set. In that case, interpolating the compensation vectors of several environments may be more helpful than using compensation vectors from a single (incorrect) environment.

For these reasons, the Interpolated Fixed Codeword Dependent Cepstral Normalization algorithm (IFDCN) estimates compensation vectors for new environments by linear interpolation of several of the compensation vectors that had been precomputed for environments in the training database:

$$\hat{\mathbf{r}}[k, l] = \sum_{e=1}^E f_e \cdot \mathbf{r}[k, l, e]$$

where $\hat{\mathbf{r}}[k, l]$, $\mathbf{r}[k, l, e]$, and f_e are the estimated compensation vectors, the environment-specific compensation vector for the e^{th} environment, and the weighting factor for the e^{th} environment, respectively.

The weighting factors for each environment are also based on residual VQ distortion:

$$f_e = \frac{p(e|\bar{\mathbf{Z}})}{\sum_{i=1}^E p(i|\bar{\mathbf{Z}})} = \frac{\exp\{D_e/(2\sigma^2)\}}{\sum_{i=1}^E \exp\{D_i/(2\sigma^2)\}}$$

where σ is the codebook standard deviation using clean speech and $\bar{\mathbf{Z}}$ is the testing utterance. With the present training and testing data we have generally used a value of 3 for E .

2.3. Phone-Dependent Cepstral Normalization (PDCN) and related algorithms

All of the compensation techniques described above compensate for mismatches between training and testing conditions of the cepstral vectors that are input to the classifier. This approach is the easiest to implement, as the environmental normalization algorithms are external to the recognition system. Nevertheless, we believe that further improvements in recognition accuracy can be obtained by exploiting more directly the speech knowledge contained in the acoustic models.

In this section, we describe a new family of algorithms, referred to as phone-dependent normalization procedures, which compensate for environmental variation based on the presumed phoneme identity of individual acoustical segments during the search process. This approach has the advantage that information from the acoustic-phonetic and language models as well as the constraints arising from the search process can be brought to bear in determining the most effective form of environmental compensation.

PDCN. In the current implementation of phone-dependent cepstral normalization (PDCN), we develop compensation vectors that are specific to individual phonetical events, using a base phone set of 51 phonemes, including silence but excluding other types of non-lexical events.

Labelled phonetic segments for training PDCN compensation are produced by running the decoder in supervised mode using the correct transcription of the incoming speech. For each phoneme, compensation vectors are derived by averaging the difference between cepstral coefficients obtained from the training environment and a given target environment, using the same stereo pairs of training sentences that were used for MFCDCN. This approach is similar to SDCN except that different compensation vectors are calculated according to phonetic identity rather than instantaneous frame SNR values. The compensation vectors in PDCN are described as follows,

$$\bar{c}[p] = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} (\bar{x}_t - \bar{z}_t) \delta(f_t - p)}{\sum_{i=1}^N \sum_{t=1}^{T_i} \delta(f_t - p)}$$

where f_t is the phoneme at time t , and T_i is the length of the i^{th} utterance out of N sentences from each of training environments in stereo databases.

The SPHINX-II system uses the senone [6,8], a generalized state-based probability density function, as the basic unit to compute the likelihood from acoustical models. The probability of observing senone s at time t in the cepstral vector $\bar{\mathbf{z}}(t)$ of incoming speech can be expressed as

$$Prob_s(t) = \sum_{i=1}^N w_i N_{\bar{z}(t)}(\mu_i, \sigma_i)$$

[IS THIS A PDF OR A PROBABILITY?] where i stands for the index of the best N Gaussian mixtures of senone s at time t , and μ_i , σ_i and w_i are the corresponding mean, variance, and probabilistic weight for the i^{th} mixture in senone s .

As before, compensated cepstral vectors are formed by adding the compensation vector to incoming cepstra, $\hat{x}_p(t) = \bar{z} + \bar{c}[p]$, on a frame-by-frame basis. This is a simple process in the present implementation because each senone corresponds to only one distinctive base phoneme. As a result, senone probabilities can be calculated directly in terms of *compensated* incoming speech vectors, by assuming the phonetic identity that corresponds to a given senone. Using this approach, the senone probability with PDCN is re-written as

$$Prob_s(t) = \sum_{j=1}^N w_j N_{\hat{x}_p(t)}(\mu_j, \sigma_j)$$

where j is the index of the best N Gaussian mixtures for senone s at time t with respect to the PDCN normalized cepstral vector $\hat{x}_p(t)$.

Compensation vectors are calculated by the decoder during the process of searching for the optimal sequence of states in the HMM, and scores used to evaluate hypotheses are calculated using the compensated cepstral vectors. The increase in computation incurred by PDCN is very minor and arises primarily from an increase in the number of vector quantization operations performed on the 51 alternatives for each cepstral vector. [I DON'T UNDERSTAND THIS SENTENCE.]

SNR-Dependent PDCN. The performance of PDCN can be further improved by further partitioning the compensation vectors in terms of SNR (as is done with SDCN and FCDCN). The estimation of compensation vectors for SPDCN can then be expressed as

$$\bar{c}[p, l] = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} (\bar{x}_t - \bar{z}_t) \delta(f_t - p) \delta(s_t - l)}{\sum_{i=1}^N \sum_{t=1}^{T_i} \delta(f_t - p) \delta(s_t - l)}$$

where s_t is the instantaneous frame SNR of $\bar{z}(t)$. [CHECK NOTATION CONSISTENCIES.] We chose a range of 30 dB of SNR in our current implementation.

Interpolated PDCN (IPDCN). PDCN, like SDCN and FCDCN assumes the existence of a database of utterances recorded in stereo in the training and testing environments. As in the case of MFCDCN, the PDCN algorithm can be extended to cases where the testing environment is unknown by developing ensembles of

PDCN compensation vectors for a variety of testing environments, and applying to incoming utterance either the set of compensation vectors from the “closest” environment used to train the algorithm (MPDCN), or an interpolation of compensation vectors from several of the closest environments (IPDCN). In the current implementation of IPDCN, we use environment-specific VQ for environment selection to obtain the 3 closest environments with the best 4 Gaussian mixtures contributing to the interpolation weights.

3. EXPERIMENTAL RESULTS

The MFCDCN, IFDCDCN, PDCN and related algorithms were evaluated using the SPHINX-II recognition system [6] in the context of the ARPA 5000-word closed-vocabulary task consisting of dictation of sentences from the Wall Street Journal. The system is trained on WSJ0 training corpus and has 7000 senones. The testing corpus consists of utterances from a set of “secondary” microphones including desktop microphones, stand-mounted microphones and telephone handsets and speakerphones. We also compared recognition accuracy for the same system using two types of cepstral high-pass filtering: the RASTA filter [5] as implemented in the SRI ARPA system [7], and cepstral mean normalization (CMN).

Table 1 compares recognition accuracy obtained using the various processing schemes along with the corresponding reduction of word error rate with respect to the baseline (no processing). The system was trained on the standard Sennheiser closetalking HMD-414 microphone (CLSTLK), and tested using either the CLSTLK mic or one of several secondary microphones (OTHER). The word error rate of 38.5% obtained by testing with alternate microphones (compared to 8.1% with the CLSTLK mic) demonstrates the effect of the mismatch of training and testing environments. The error rate with alternate mics is reduced by 44.4% using high-pass filtering algorithms, and by 56.6% using MFCDCN and IFDCDCN. Further improvement is obtained when MFCDCN or IFDCDCN is combined with a high-pass filtering technique like CMN. The PDCD can generate 59.2% word error rate reduction for alternate mics with CMN. This result indicates that PDCN may suffer from mistakes that the decoder make in terms of environment compensation. Reduce the effect of decoder’s errors, PDCN combined with MFCDCN can yield further word error rate reduction, 66.8%.

PROCESSING METHOD	CLSTL K mic	% Dec.	OTHER mics	% Dec.
Baseline	8.1	–	38.5	–
RASTA	9.0	-11.1	28.0	27.3
CMN	7.6	6.2	21.4	44.4
MFCDCN	8.1	0	16.7	56.6
IFDCDCN	8.4	-3.7	16.7	56.6

Table 1: Percentage of word errors and corresponding error rate reduction for different processing schemes on the test corpus for the ARPA 11/92 5000-word, closed-vocabulary task using sentences from the Wall Street Journal.

PROCESSING METHOD	CLSTL K mic	% Dec.	OTHER mics	% Dec.
CMN+MFCDCN	8.1	0	14.5	62.3
CMN+IFDCDCN	8.4	-3.7	14.8	61.6
CMN+PDCN	–	–	15.7	59.2
CMN+MFCDCN+N+PDCN	–	–	12.8	66.8

Table 1: Percentage of word errors and corresponding error rate reduction for different processing schemes on the test corpus for the ARPA 11/92 5000-word, closed-vocabulary task using sentences from the Wall Street Journal.

METHOD	CLSTLK mic	% Dec.	OTHER mics	% Dec.
Baseline	8.1	–	38.5	–
CMN+MFCDCN	8.1	0	16.1	58.2
+PDCN	8.1	0	14.8	61.6
CMN+IFDCDCN	8.4	-3.7	14.8	61.6
+IPDCN	8.4	-3.7	13.5	64.9

Table 2: Recognition accuracy obtained for the same task as in Table 1, but with the testing environments excluded from the corpus used to develop the compensation vectors.

Table 2 summarizes results obtained with the MFCDCN and Interpolated FCDCN algorithms when the actual testing environment was excluded from the set of data used to develop the compensation vectors. We also compare PDCN and IPDCN for the same data. Comparing the results of Table 2 to those of Table 1, it is seen that the removal of the correct environment from the training data causes recognition accuracy using MFCDCN to degrade slightly. PDCN then can provide 8% error rate reduction over just MFCDCN. Similarly conclusion can be draw for Interpolated FCDCN and IPDCN, demonstrating the ability of IFDCDCN and IPDCN to provide compensation, even for environments that are not part of the training data.

4. SUMMARY

MFCDCN and IFDCDCN accomplishes environment normalization in the front-end stage with a major advantage of maintain the same configuration of the recognition system. On the other hand, a family of PDCN algorithms that work inside the recognizer provide a method to make use of useful information from the search for environment normalization. When evaluated in the context of experimental results using the WSJ/CSR task, they can yield a substantial error reduction compared with baseline without any processing on alternate mic data. The word error reductions are 56.6%, 62.3, 59.2% and 66.8% for MFCDCN, MFCDCN with CMN, PDCN with CMN, and PDCN with MFCDCN and CMN,

respectively. For the experiments in which the testing environment doesn't resemble any environment in the training set, compensation vectors from IFCDCN provide about 10% word error rate reduction compared with MFCDCN. Note that though PDCN is evaluated in SPHINX-II with a semi-continuous density HMM framework, it can easily be implemented in systems with a discrete-density HMM or continuous-density HMM framework.

ACKNOWLEDGMENTS

This research was sponsored by the Department of the Navy, Naval Research Laboratory, under Grant No. N00014-93-2005. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government. The authors would like thank Mei-Yuh Hwang, Xue-Dong Huang, Yoshiaki Ohshima, Raj Reddy and the rest of the CMU speech group for their contributions to this work.

REFERENCES

1. Juang, B.H., "Speech Recognition in Adverse Environments", *Computer Speech and Language* 5:275-294, 1991.
2. Acero, A., *Acoustical and Environmental Robustness in Automatic Speech Recognition*, Kluwer Academic Publishers, Boston, MA, 1993
3. Liu, F.H., Acero, A., and Stern, R.M., "Efficient Joint Compensation of Speech for the Effects of Additive Noise and Linear Filtering", *ICASSP-92*, pages 865-868, March 1992.
4. Schwartz, R., Anastasakos, T., Kubala, F., Makhoul, J., Nguyen, L., and Zavalagkos, G., "Comparative Experiments on Large Vocabulary Speech Recognition", *Proc. ARPA Human Language Technology Workshop*, March, 1993
5. Hermansky, H., Morgan, N., and Hirsch, H., "Recognition of Speech In Additive and Convolutional Noise Based on RASTA Processing", *ICASSP-93*, pages II-83~86, April, 1993.
6. Huang, X., Alleva, F., Hwang, M., and Rosenfeld, R., "An Overview of the SPHINX-II Speech Recognition System", *Proc. ARPA Human Language Technology Workshop*, March 1993.
7. Murveit, H., Butzberger, M., and Weitraub, M., "Reduced Channel Dependence for Speech Recognition", *Proc. DARPA Spoken and Natural Language*, Feb., 1992
8. Hwang, M.Y. and Huang, X.D.: "Shared-Distribution Hidden Markov Models for Speech Recognition", *IEEE Trans. on Speech and Audio Processing*, Vol 1, No. 10, 1993