

# A BLOCK COSINE TRANSFORM AND ITS APPLICATION IN SPEECH RECOGNITION

Jingdong Chen †\*, Kuldip K. Paliwal † and Satoshi Nakamura \*

† School of Microelectronic Engineering, Griffith University  
Brisbane, QLD 4111, Australia

\* ATR Spoken Language Translation Research Labs  
Kyoto, 619-0288, Japan

E-mail: {jingdong.chen, nakamura}@slt.atr.co.jp, k.paliwal@me.gu.edu.au

## ABSTRACT

Noise robust speech recognition has become an important area of research in recent years. The fact that human listeners can recognize speech in the presence of strong noise inspires researchers to imitate some aspects of human auditory perception in automatic speech recognition. This has led to sub-band based speech recognition in which the full-band speech is split into several sub-bands and where each sub-band is processed separately. The resulting multi-band features can be combined in various ways for carrying out speech recognition task. Reported results have shown the superiority of this technique for speech recognition in strong noise conditions. In this paper, we will briefly review the multi-band feature extraction. We will then propose a block discrete cosine transform (BDCT) with its kernel transformation matrix being derived from the decomposition of the kernel of the discrete cosine transform (DCT). We show that the BDCT approximates the DCT in keeping information in decorrelating a sequence. When the BDCT is applied to the mel frequency filter bank energies (FBEs) to replace the DCT to convert them to cepstral coefficients, a new kind of MFCCs is yielded. We call these new features Block discrete cosine transform based MFCCs (BMFCCs) and show that a sub-band processing idea is implicit in the BMFCCs since the BDCT automatically divides the mel frequency FBEs into two sub-bands. We will report various speech recognition results using the BMFCCs as well as the comparison with the multi-band MFCCs and full-band MFCCs to elaborate the properties of the BMFCCs.

## 1. INTRODUCTION

Significant advances have been made in recent years in the area of automatic speech recognition. It is now possible to use a speech recognition successfully in a controlled environment. However, the performance of a speech recognizer suffers dramatic degradation when there is a mismatch between training and testing environments [1-3]. There are many factors that contribute to this mismatch. The main factor, that causes the mismatch, is the presence of ambient background noise in the speech signal. Maintaining good recognition accuracy in noisy conditions has become one of the challenging areas of research currently.

The fact that human listeners can achieve and secure very high recognition accuracy even in the conditions in which the signal to noise ratio becomes extremely low inspires researchers to mimic some aspects of human auditory perception in automatic speech recognition. Psychoacoustic evidence shows that human beings process speech on a narrow band basis. An intuitive way to imitate the auditory system is to split the full-band speech into several sub-bands and represent each sub-band individually. This has led to a technique called sub-band based speech recognition.

One straightforward way to achieve sub-band representation

is to divide the full-band speech signal into several sub-bands and convert each sub-band spectrum into several cepstral features. These sub-band features are then concatenated together as a single feature vector and used for speech recognition [4-7,12-14]. Results reported in the literature have shown the advantages of these multi-band features for noisy speech recognition. However, the performance for clean speech is often poorer as the features from different sub-bands may be correlated. Furthermore, the number of sub-bands and the boundaries for each sub-band are empirical values which have to be manually adjusted to gain good performance for a given recognition task.

An alternative is to model the sub-band features independently and to combine the likelihood score at some segmental level [8-10]. Such a combination may enable a more flexible way to manipulate the sub-band features to permit further enhancement in performance. An unsolved problem with this approach is how to determine the weighting function to guarantee at least a sub-optimal combination of sub-band features.

In this paper, we will firstly review the extraction of multi-band MFCCs. We will then propose a block discrete cosine transform (BDCT) with its kernel transformation matrix being derived from the decomposition of the kernel of the discrete cosine transform (DCT). We show that the proposed new transform behaves similarly with the DCT in keeping information in decorrelating a sequence.

When the BDCT is applied to the representation of the mel frequency cepstrum in replacing the DCT, a new type of MFCCs is obtained. We call these new MFCCs Block discrete cosine transform based MFCCs (BMFCC). It is found that the BDCT automatically divide the power spectrum into two sub-bands, hence a sub-band processing idea is implicit in the new cepstral features.

Various speech recognition experiments are carried out to test the properties of the BDCT and the BMFCCs. We will report some results as well the comparison with the multi-band and full-band MFCCs to elaborate the properties of these new features.

## 2. MULTI-BAND MFCC

Let  $E = [e_1, e_2, \dots, e_N]$  denote a sequence of log filter bank energies, where  $N$  is the number of filters in the filter bank, then the full-band cepstrum is computed from a DCT,

$$X_F = [C]E \quad (1)$$

Suppose we divide the speech signal into  $M$  sub-bands. In order to compute cepstral coefficients for the  $m$ th sub-band, we process each sub-band signal with a filter bank having  $N_m$  filters. Thus these filter bank energies are given as





Feature set	Clean speech	30dB	20dB	15dB	10dB
FB MFCC	88.0	88.0	87.6	87.1	84.7
MB MFCC	87.8	87.8	87.7	87.3	84.6
BMFCC	89.5	89.5	89.5	89.4	87.8

Table 3. Speech recognition accuracy (%) in car noise condition ( Car-Volvo-340 120 km/h, 4<sup>th</sup> gear).

From these experiment results, we can make following observations:

1. The BMFCCs yield the best performance in clean speech and high SNR conditions. This demonstrates the superiority of the BDCT to the DCT in cepstral analysis for speech recognition.
2. All sub-band based features are more robust than the full-band MFCCs to three types of noise investigated.
3. BMFCCs are more robust to the machine gun noise and car noise than the multi-band MFCCs. While its robustness to speech noise is slightly poorer than that of the multi-band MFCCs. The reason for this is not clear. We will perform further experiments with various speech databases and with more kinds of real noise before we draw a conclusion.

## 7. CONCLUSION

In this paper, we proposed a block DCT and applied it to the mel frequency cepstral analysis in speech recognition. We showed that a sub-band processing idea is implicit in the new features (BMFCCs). Experiment results based on a continuous density isolated speech recognizer revealed that the new MFCCs yield better recognition accuracy than full-band MFCCs in noise as well as in clean speech conditions.

BMFCCs were also compared with the multi-band MFCCs in terms of their recognition performance. The results showed that the BMFCCs were able to yield better performance in clean speech and various noisy speech environments. Thus, BMFCCs are more robust than the multi-band MFCCs under various noise conditions.

In this paper, experiments are only performed for small vocabulary isolated speech recognition tasks. Work is in progress to test the BMFCCs and various sub-band based front-ends for large vocabulary continuous speech recognition in clean as well as noise conditions.

## REFERENCE

- [1] D. S. Pallett, *et al*, "1996 Preliminary Broadcast News Benchmark", Proceedings of the 1997 DARPA Speech Recognition Workshop, International Conference Center Chantilly, Virginia, February 2-5, 1997
- [2] D. S. Pallett, *et al*, "1997 Broadcast News Benchmark Test Results: English and Non-English", Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop, Lansdowne Conference Resort, Lansdowne, Virginia, February 8-11, 1998.
- [3] D. S. Pallett, *et al* "1998 Broadcast News Benchmark Test Results: English and Non-English Word Error Rate Performance Measures", Proceedings of the DARPA Broadcast News Workshop, Hilton at Washington Dulles Airport Herndon, Virginia, February 28-March 3, 1999.
- [4] P. McCourt, S. Vaseghi and N. Harte, "Multi-Resolution Cepstral Features for Phoneme Recognition Across Speech Sub-bands," ICASSP'98, Seattle, USA, PP.557-560.
- [5] S. Okawa, E. Bocchieri and A. Potamianos, "Multi-Band Speech Recognition in Noisy Environments," ICASSP'98, Seattle, USA, PP.641-644.
- [6] A. Chen, S. Vaseghi and P. McCourt, "Transformation of Full-Band to Sub-band HMMS for Speech Recognition in Noisy Car Environments," ARSU'99, Keystone, Colorado, December 1999, Vol. 1, PP. 31-34.
- [7] B. Doherty, S. Vaseghi and P. McCourt, "Linear Transformations in Sub-band Groups for Speech Recognition," EUROSPEECH'99, Budapest, Hungary, September, 1999, PP. 1359-1366.
- [8] H. Bourlard and S. Dupont, "A new ASR Approach Based on Independent Processing and Recombination of Partial Frequency Bands," ICSLP'96, Philadelphia, October 1996.
- [9] H. Bourlard and S. Dupont, "Subband-based Speech Recognition," ICASSP'97, PP. 1251-1254.
- [10] S. Okawa, T. Nakajima and K. Shiria, "A Recombination Strategy for Multi-band Speech Recognition Based on Mutual Information Criterion," EUROSPEECH'99, PP. 603-606.
- [11] K. K. Paliwal, "Decorrelated and Liftered Filter-Bank Energies for Robust Speech Recognition," EUROSPEECH'99, PP. 85-88.
- [12] R. Chengalvarayan, "Hierarchical Subband Linear Prediction Cepstral (HSLPC) Features for HMM-Based Speech Recognition," ICASSP'99, PP. 409-412.
- [13] K. Yoshida K. Takagi and K. Ozeki, "Speaker Identification Using Subband HMMS," EUROSPEECH'99, PP. 1019-1022.
- [14] S. Rao and W. A. Pearlman, "Analysis of Linear Prediction, Coding, and Spectral Estimation from Subbands," IEEE Trans. on Information Theory, Vol. 42, No. 4, 1996, PP. 1160-1178.
- [15] A. D. Poularikas, "The Transforms and Applications Handbook," IEEE Press, 1995.
- [16] K. R. Rao and R. Yip, "Discrete Cosine Transform: Algorithm, Advantages, Application," Academic Press, Inc, 1990.
- [17] A. Varga, *et al*, "The Moise-92 Study on the Effect of Additive Noise on Automatic Speech Recognition," DRA Speech Research Unit, St. Andrew's Rd., Malvern, Worcestershire, WR14 3PS UK.