

Noise-Robust Speech Features Based on Cepstral Time Coefficients

Ja-Zang Yeh, and Chia-Ping Chen

Presented at Rocling 2009

Outline

- Introduction
- Cepstral Time Coefficients
- Experiments
- Conclusion

Outline

- **Introduction**
- Cepstral Time Coefficients
- Experiments
- Conclusion

Introduction

- A front-end of a speech recognition system may consist of several stages.
 - In the early stage : spectral subtraction and Wiener filter
 - In the middle stage : pre-emphasis and hamming window
 - In the post-processing stage: normalization, temporal information integration.

Introduction(Cont.)

- We investigate novel features based on simple transformation post-processing methods.
 - Insert a window of static cepstral vectors in a matrix and then apply the *discrete cosine transform*(DCT)
 - Coefficients after DCT is called *cepstral time coefficients*
 - Resultant matrix is called *cepstral time matrix* (CTM)
 - Further apply normalization and routines for delta and acceleration feature extraction
 - Combined with the static MFCC features

Outline

- Introduction
- **Cepstral Time Coefficients**
- Experiments
- Conclusion

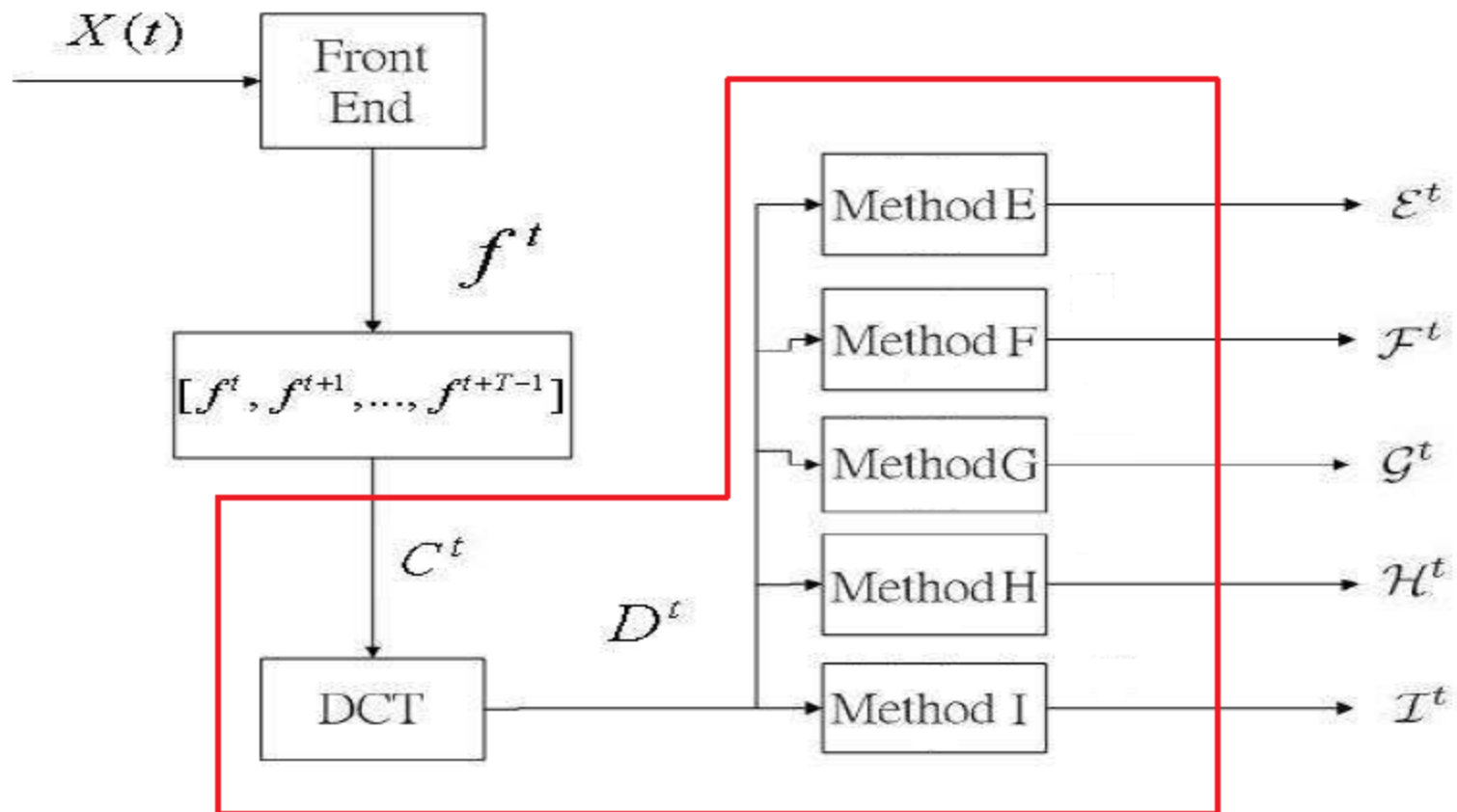


Figure 1: The block diagram of the proposed feature transformation methods.

Cepstral Time Coefficients

- We first insert a fixed number of adjacent feature vectors in a matrix

$$C^t \triangleq \begin{pmatrix} C_{11}^t & \cdots & C_{1T}^t \\ \vdots & \ddots & \vdots \\ C_{K1}^t & \cdots & C_{KT}^t \end{pmatrix} \triangleq \begin{bmatrix} f^t & f^{t+1} & \cdots & f^{t+T-1} \end{bmatrix}$$

K : the feature vector dimension

f^t : feature vector of frame t

C^t : the matrix whose column vectors are the T consecutive feature vectors starting from frame t

Cepstral Time matrix

- The cepstral time matrix at frame t , D^t , is related to C^t by the discrete-cosine transform(DCT).

$$D_{i:}^t = DCT(C_{i:}^t)$$

$D_{i:}^t$: the i -th **row** of matrix D^t

D_{in}^t : the n -th cepstral time coefficient (CTC) of channel i at frame t

- Our matrix index starts from 1 instead of 0

$$D_{in}^t = \sum_{\tau=1}^T C_{i\tau}^t \cos\left(\frac{(2\tau-1)(n-1)\pi}{2T}\right)$$

Method E

- Dividing the first column of D^t by the number of frames , while leaving other columns unchanged.

$$\begin{cases} E_{:1}^t = D_{:1}^t / T \\ E_{:n}^t = D_{:n}^t, \quad n \neq 1 \end{cases}$$

- Then we apply the delta and acceleration feature extraction steps.

$$\begin{cases} \cup \\ E_{:2}^t = E_{:2}^t - E_{:1}^t \\ \cup \\ E_{:3}^t = E_{:3}^t - 2E_{:2}^t + E_{:1}^t \end{cases}$$

Method E(Cont.)

- We add the $\overset{\cup}{E}_{:2}^t$ and $\overset{\cup}{E}_{:3}^t$ to the static MFCCs, resulting in a feature vector of

$$\mathcal{E}^t = \begin{bmatrix} C_{:1}^t \\ \overset{\cup}{E}_{:2}^t \\ \overset{\cup}{E}_{:3}^t \end{bmatrix}$$

Method F

- Normalize the feature values in the first column to the range of $[-1, 1]$. Let F^t be defined by

$$\begin{cases} F_{:1}^t = D_{:1}^t / N^t \\ F_{:n}^t = D_{:n}^t, \quad n \neq 1 \end{cases}$$

N^t : the maximum magnitude in the first column.

i.e. $N^t = \max_d |D_{d1}^t|$

- The remaining operations are similar to Method E.

$$\begin{cases} \overset{\cup}{F}_{:2}^t = F_{:2}^t - F_{:1}^t \\ \overset{\cup}{F}_{:3}^t = F_{:3}^t - 2F_{:2}^t + F_{:1}^t \end{cases}$$

Method F(Cont.)

- We add $\overset{\cup}{F}_{:2}^t$ and $\overset{\cup}{F}_{:3}^t$ to the static MFCCs, resulting in a feature vector of

$$\mathcal{F}^t = \begin{bmatrix} C_{:1}^t \\ \overset{\cup}{F}_{:2}^t \\ \overset{\cup}{F}_{:3}^t \end{bmatrix}$$

Method G

- In Method G, we add the first and second columns of D^t to the static MFCC vector,

$$\mathcal{G}^t = \begin{bmatrix} C_{:1}^t \\ D_{:1}^t \\ D_{:2}^t \end{bmatrix}$$

Method H

- In Method H, we add the second and third columns of D^t to the static MFCC vector,

$$\mathcal{H}^t = \begin{bmatrix} C_{:1}^t \\ D_{:2}^t \\ D_{:3}^t \end{bmatrix}$$

Method I

- In Method I, we simply use the zeroth, first, and second cepstral time coefficients,

$$\mathcal{I}^t = \begin{bmatrix} D_{:1}^t \\ D_{:2}^t \\ D_{:3}^t \end{bmatrix}$$

Outline

- Introduction
- Cepstral Time Coefficients
- **Experiments**
- Conclusion

Evaluation corpus

- Evaluated by Aurora 3 corpus.
- Aurora 3 consisting of digit-string utterances in Danish, German, Finnish and Spanish.
- It provides a platform for fair comparison between systems of different front-ends.

Back end

- We use HTK as the back end to run experiments.
- Implemented by HMM, and use word-level model.
- Each word model consists of 16 emitting states, and each consists of 3 Gaussian components.
- The silence model consists of 3 state, and each state consists of 6 Gaussian components.

Experiments set

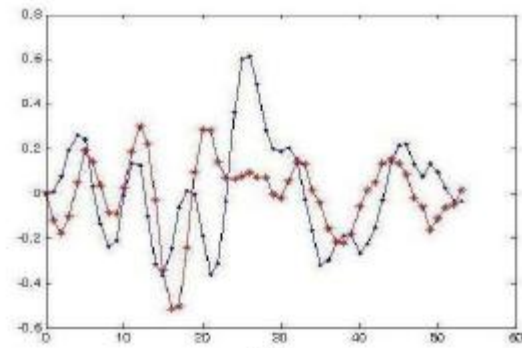
- We first evaluate and decide to use $T = 15$.
- For the static features we use 12 MFCC features and the log energy, making $K = 13$.
- Our baseline, it simply uses the MFCC, delta and acceleration features.

Relative improvement

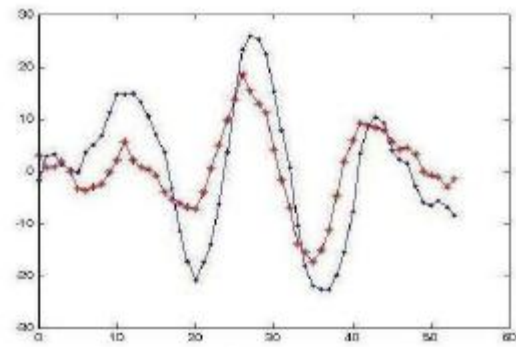
	German	Spanish	Finnish	Danish
Method E	-12.4	16.2	16.5	16.3
Method F	-10.5	22.4	10.8	16.3
Method G	-58.1	-29.0	-42.9	-19.2
Method H	7.5	26.6	25.4	23.2
Method I	-10.8	19.8	8.5	13.1

Discussion

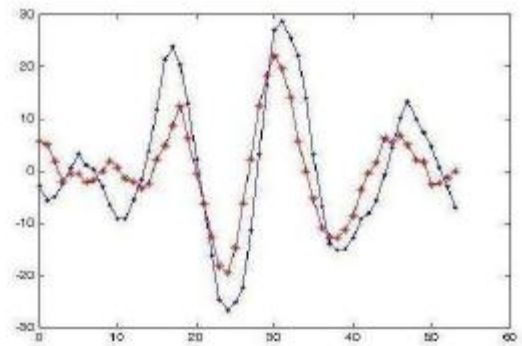
- Method E and F are similar, and they have similar performance level.
- Method G and H concludes that zeroth CTC is detrimental of recognition.
- We try scheme of normalizing and dividing.



(B)



(F)



(H)

Outline

- Introduction
- Cepstral Time Coefficients
- Experiments
- **Conclusion**

Conclusion

- In this paper, we use five difference feature sets based on the cepstral time coefficients.
- The combination of raw MFCC and the second and the third columns of CTM yields the best.