

Introduction

- Information on the type of distortion corrupting a signal can be used to inform the choice of appropriate enhancement algorithms.
- Most existing methods focused on detecting a single and specific type of distortion in a signal.
- In [1], we proposed a method to classify four major types of distortion in vowels directly from MFCCs extracted from speech signals.
- Limitations of [1]:
 - MFCCs encode not only distortion in signals, but also other variability (speaker, articulation and disorder).
 - Distortion classification decision is made by majority vote over all frames, and the computation time increases with increasing signal length.
- In this paper, distortion in variable duration recordings is modeled with a fixed-length, low-dimensional vector.

Distortion Modeling

- Channel variability** can be produced artificially by corrupting the clean recording by different types and levels of distortion.
- Method:**
 - Fitting a Gaussian mixture model (GMM) to the features of a recording.
 - Assuming that the GMM mean supervector of the r^{th} recording from the s^{th} speaker can be decomposed as:

$$\mathbf{M}_{s,r} = \mathbf{m} + \mathbf{V}\mathbf{y}_s + \mathbf{U}\mathbf{x}_{s,r} + \mathbf{D}\mathbf{z}_s. \quad (1)$$

- Definitions:**
 - \mathbf{m} is speaker- and channel-independent supervector,
 - \mathbf{V} is a rectangular matrix of low rank with high speaker variability
 - \mathbf{y}_s is the speaker factor
 - \mathbf{U} is a rectangular matrix of low rank with high channel variability
 - $\mathbf{x}_{s,r}$ is the channel factor containing channel related information
 - \mathbf{D} is a diagonal matrix describing any remaining speaker variability
 - \mathbf{z}_s is the speaker-specific residual factor
 - The factors $\mathbf{x}_{s,r}$, \mathbf{y}_s and \mathbf{z}_s are assumed to be independent of each other and have a standard normal prior distribution.

- Estimating the matrices \mathbf{V} , \mathbf{U} , \mathbf{D} , and the vectors $\mathbf{x}_{s,r}$, \mathbf{y}_s and \mathbf{z}_s [2]:**
 - Train \mathbf{V} , assuming that \mathbf{U} and \mathbf{D} are zero.
 - Estimate \mathbf{U} given the estimate of \mathbf{V} and assuming that \mathbf{D} is zero.
 - Estimate the residual matrix \mathbf{D} given the estimates of \mathbf{V} and \mathbf{U} .
 - $\mathbf{x}_{s,r}$, \mathbf{y}_s and \mathbf{z}_s are then calculated given the estimates of \mathbf{V} , \mathbf{U} and \mathbf{D} .

Channel Factor and Subspace Estimation

- The channel factor $\mathbf{x}_{s,r} \sim N(\boldsymbol{\mu}_{s,r}, \boldsymbol{\Lambda}_{s,r})$ and the channel subspace \mathbf{U} are estimated by applying an EM algorithm [2].

- In the E-step**, using a random initialization of \mathbf{U} , the posterior distribution of the channel factor is calculated as:

$$\boldsymbol{\mu}_{s,r} = E[\mathbf{x}_{s,r}] = (\mathbf{I} + \mathbf{U}^T \boldsymbol{\Sigma}^{-1} \mathbf{N}_s \mathbf{U})^{-1} \mathbf{U}^T \boldsymbol{\Sigma}^{-1} \mathbf{f}_{s,r} \quad (2)$$

$$\boldsymbol{\Lambda}_{s,r} = E[\mathbf{x}_{s,r} \mathbf{x}_{s,r}^T] = \boldsymbol{\mu}_{s,r} \boldsymbol{\mu}_{s,r}^T + (\mathbf{I} + \mathbf{U}^T \boldsymbol{\Sigma}^{-1} \mathbf{N}_s \mathbf{U})^{-1}. \quad (3)$$

- In the M-step**, the channel subspace is updated by solving the equations:

$$\mathbf{U}_i \boldsymbol{\Theta}_c = \boldsymbol{\Psi}_i. \quad (4)$$

- Definitions:**

- $\boldsymbol{\Sigma}$ is a block-diagonal matrix entries form the covariance matrix of the c^{th} mixture of the UBM,
- $N_{s,r,c} = \sum_{l=1}^L \gamma_{c,l}$ and $\mathbf{f}_{s,r,c} = \sum_{l=1}^L \gamma_{c,l} [\boldsymbol{\rho}_l - (\mathbf{m}_c + \mathbf{V}_c \mathbf{y}_s)]$ are the zero- and first order statistics for each speaker s , recording r and mixture component c .
- $\boldsymbol{\rho}_l$ is the acoustic features of the l^{th} frame
- \mathbf{I} is an identity matrix,
- \mathbf{N}_s is a block-diagonal matrix which its entries are $(\sum_r N_{s,r,c}) \mathbf{I}$
- $\mathbf{f}_{s,r}$ is a vector constructed by concatenation of $\mathbf{f}_{s,r,c}$
- $\gamma_{c,l}$ is the posterior probability of the c^{th} mixture generating $\boldsymbol{\rho}_l$,
- \mathbf{m}_c and \mathbf{V}_c are, respectively, the subvector of \mathbf{m} and the submatrix of \mathbf{V} of mixture component c .
- $\boldsymbol{\Theta}_c = \sum_s \sum_r N_{s,r,c} \boldsymbol{\Lambda}_{s,r}$ $c = 1, \dots, C$
- $\boldsymbol{\Psi}_i$ is the i^{th} row of $\boldsymbol{\Psi} = \sum_s \sum_r \mathbf{f}_{s,r,c} \boldsymbol{\mu}_{s,r}^T$

The Proposed Method

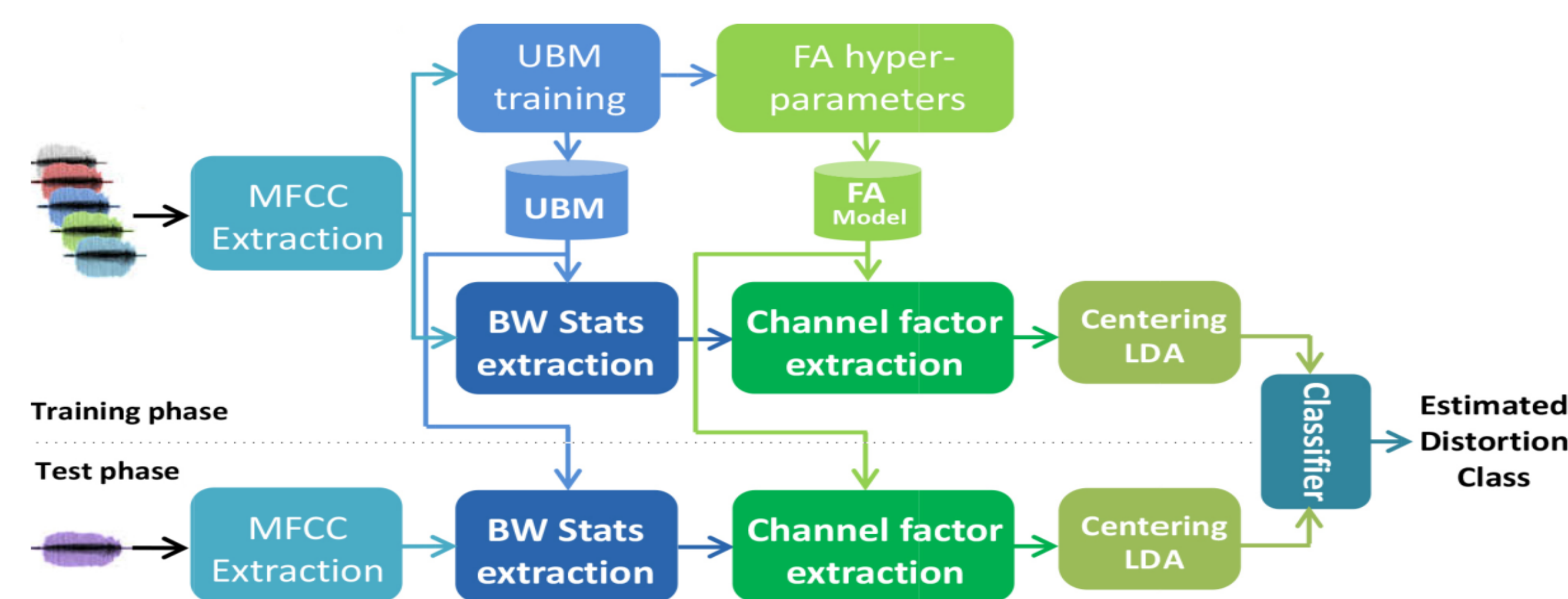


Fig. 1: The block diagram of the proposed distortion classification system.

Experimental Setup

- Database:**
 - Parkinson's voice database (sustained vowels, 750 telephone recordings).
- Distortion Classes:**
 - Additive noise (white Gaussian, babble, office ambient noises)
 - Reverberation (8 different real room impulse responses)
 - Peak clipping (clipping level: 0.3, 0.4, 0.5, 0.6)
 - Coding (6.3 kbps, 9.6 kbps and 16 kbps CELP codecs)
- Acoustic features:**
 - 39 dimensional vector (12 MFCCs + frame energy + Δ + $\Delta\Delta$)
- Distortion Modeling:**
 - GMM with 256 mixtures
 - Speaker factor dim.: 0
 - Channel factor dim.: 210
- Classifiers:**
 - SVM with RBF kernel
 - PLDA

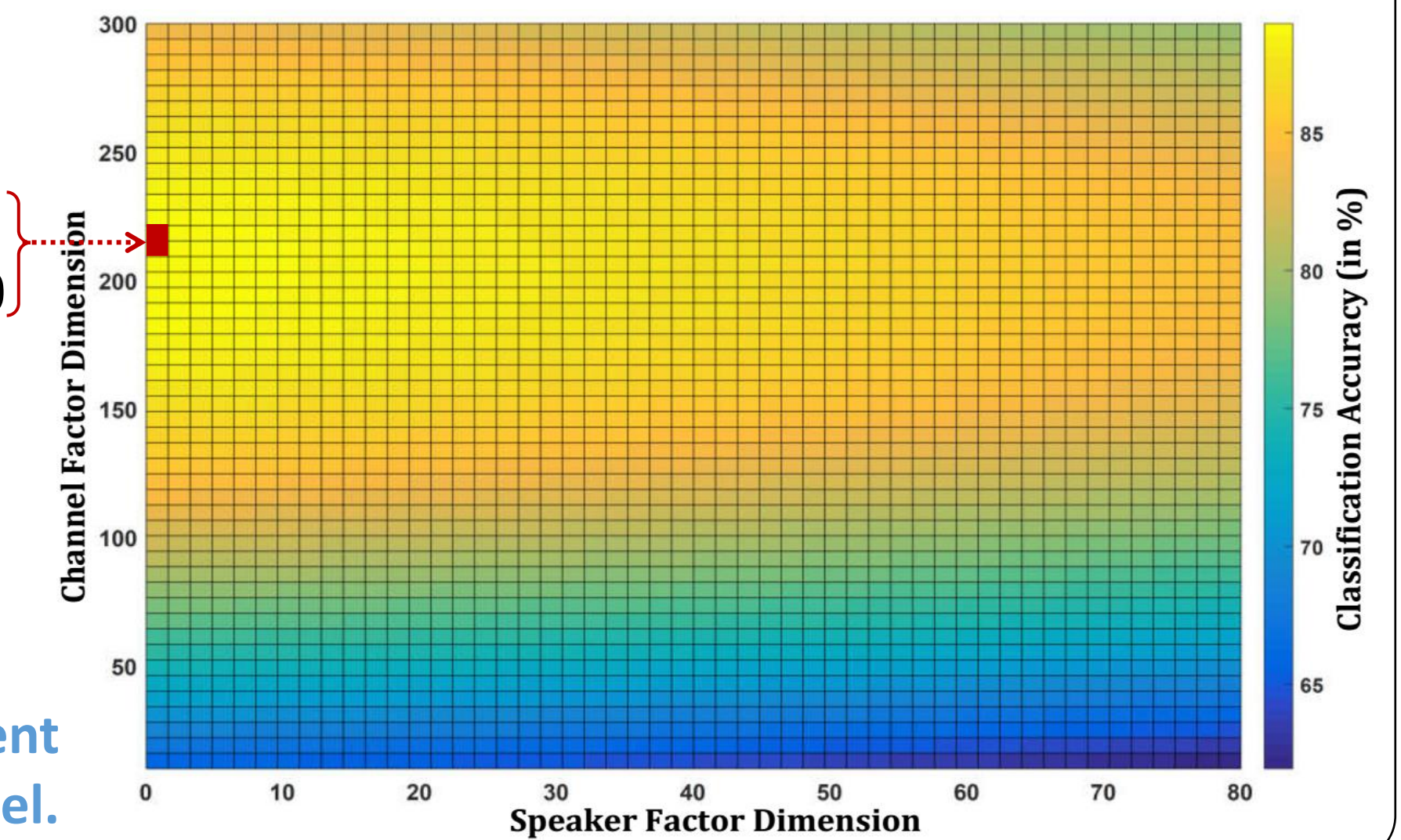


Fig. 2: Performance of different configuration of the FA model.

Results

Table 1: Comparison of [1] and the proposed method before and after pre-processing channel vectors using LDA. Results are in the form of mean \pm STD.

System	Clean	Noisy	Reverb.	Clipped	Coded	Overall
Baseline	55 \pm 11	97 \pm 4	77 \pm 4	82 \pm 7	85 \pm 9	79 \pm 3
PLDA	100 \pm 0	0 \pm 0	0 \pm 0	0 \pm 0	0 \pm 0	20 \pm 0
PLDA + LDA	77 \pm 4	98 \pm 2	86 \pm 4	82 \pm 2	93 \pm 3	87 \pm 1
SVM	28 \pm 18	33 \pm 5	31 \pm 16	35 \pm 14	68 \pm 12	39 \pm 4
SVM + LDA	78 \pm 3	97 \pm 2	87 \pm 4	85 \pm 2	93 \pm 3	88 \pm 1

Conclusions

- Distortion in variable duration signals is modeled by a fixed-length, low-dimensional vector which is more suitable for classification algorithms.
- Channel vectors are more robust to small changes in signal characteristics than MFCCs, they are more suitable for distortion classification in pathological voices.

References

- A. H. Poorjam, J. R. Jensen, M. A. Little, and M. G. Christensen, "Dominant distortion classification for pre-processing of vowels in remote biomedical voice analysis," in INTERSPEECH, 2017, pp. 289–293.
- P. Kenny, P. Ouellet, N. Dehak, V. Gupta, and P. Dumouchel, "A Study of inter-speaker variability in speaker verification," IEEE Trans. Audio. Speech. Lang. Processing, vol. 16, no. 5, pp. 980–988, 2008.