

# A High-Dimensional Subband Speech Representation and SVM Framework for Robust Speech Recognition

Jibrán Yousafzai<sup>\*†</sup>, *Member, IEEE*    Zoran Cvetković<sup>†</sup>, *Senior Member, IEEE*  
 Peter Sollich<sup>‡</sup>    Matthew Ager<sup>‡</sup>

**Abstract**— This work proposes a novel support vector machine (SVM) based robust automatic speech recognition (ASR) front-end that operates on an ensemble of the subband components of high-dimensional acoustic waveforms. The key issues of selecting the appropriate SVM kernels for classification in frequency subbands and the combination of individual subband classifiers using ensemble methods are addressed. The proposed front-end is compared with state-of-the-art ASR front-ends in terms of robustness to additive noise and linear filtering. Experiments performed on the TIMIT phoneme classification task demonstrate the benefits of the proposed subband based SVM representation: it outperforms the standard cepstral front-end in the presence of noise and linear filtering for signal-to-noise ratio (SNR) below 12-dB. A combination of the proposed front-end with a conventional representation such as MFCC yields further improvements over the individual front-ends across the full range of noise levels.

**Index Terms**—Speech recognition, robustness, subbands, support vector machines.

## I. INTRODUCTION

AUTOMATIC speech recognition (ASR) systems suffer a severe performance degradation in the presence of environmental distortions, in particular additive and convolutive noise. Humans, on the other hand, exhibit a very robust behavior in recognizing speech even in extremely adverse conditions. The central premise behind the design of state-of-the-art ASR systems is that front-ends based on the non-linear compression of speech such as Mel-Frequency Cepstral Coefficients (MFCC) [1] and Perceptual Linear Prediction (PLP) coefficients [2], when combined with appropriate language and context modelling techniques, can bring the recognition performance of ASR close to that of humans. However, the effectiveness of context and language modelling depends critically on the accuracy with which the underlying sequence of elementary phonetic units is predicted [3], and this is where there are still significant performance gaps between humans and ASR systems. Humans recognize isolated speech units above the level of chance already at  $-18$ -dB SNR, and significantly above it at  $-9$ -dB SNR [4]. At such high noise levels, human speech recognition performance exceeds that of state-of-the-art ASR systems by over an order of magnitude. Even in quiet conditions, the machine phone error rates for nonsense syllables are significantly higher than human error rates [3, 5–7]. Although there are a number of factors preventing conventional ASR systems from reaching the human benchmark, several studies [7–12] have attributed the marked

difference between human and machine performance to the fundamental limitations of the ASR front-ends. These studies suggest that the large amount of redundancy in speech signals, which is removed in the process of extracting cepstral features such as Mel-Frequency Cepstral Coefficients (MFCC) [1] and Perceptual Linear Prediction (PLP) coefficients [2], is in fact needed to cope with environmental distortions. Among these studies, work on human speech perception [7, 9, 11, 12] has shown explicitly that the information reduction that takes place in the conventional front-ends leads to a severe degradation in human speech recognition performance; furthermore, in noisy environments a high degree of correlation was found between human and machine errors when recognizing speech with distortions introduced by typical ASR front-end processing. Over the years, techniques such as cepstral mean-and-variance normalization (CMVN) [13, 14], vector Taylor series (VTS) compensation [15] and the ETSI advanced front-end (AFE) [16] have been developed that aim to explicitly reduce the effects of noise on the short-term spectra, in order to make the ASR front-ends less sensitive to noise. However, the distortion of the cepstral features caused by additive noise and linear filtering critically depends on the speech signal, filter characteristics, noise type and noise level, in such a complex fashion that feature compensation or adaptation are very challenging and so far not sufficiently effective [14].

In our previous work we showed that using acoustic waveforms directly, without any compression or nonlinear transformation, can improve the robustness of ASR front-ends to additive noise [17]. In this paper, we propose an ASR front-end derived from the decomposition of speech into its frequency subbands, to achieve additional robustness to additive noise as well as linear filtering. This approach draws its motivation primarily from the experiments conducted by Fletcher [18], which suggest that the human decoding of linguistic messages is based on decisions within narrow frequency subbands that are processed quite independently of each other. This reasoning further implies that accurate recognition in any subband should result in accurate recognition overall, regardless of the errors in other subbands. While this theory has not been proved, and some studies on the subband correlation of speech signals [19, 20] have questioned its validity, there are additional technical reasons for considering classification in frequency subbands. First of all, decomposing speech into its frequency subbands can be beneficial since it allows a better exploitation of the fact that certain subbands may inherently provide better separation of some phoneme classes than others. Secondly, the effect of wideband noise in sufficiently narrow subbands can be

The authors are with the Division of Engineering<sup>†</sup> and the Department of Mathematics<sup>‡</sup> at King's College London (e-mail: {jibrán.yousafzai, zoran.cvetkovic, peter.sollich, matthew.ager}@kcl.ac.uk).

approximated as that of narrowband white noise, so that the compensation of features becomes approximately independent of the spectral characteristics of the additive noise and linear filtering. Moreover, appropriate ensemble methods for aggregation of the decisions in individual frequency subbands can facilitate selective de-emphasis of unreliable information, particularly in the presence of narrowband noise.

Previously, the subband approach has been used in [21–27] and resulted in marginal improvements in recognition performance over its full band counterpart. But the front-ends employed in these studies were subband-based variants of cepstral features or multi-resolution cepstral features. In contrast, our proposed front-end features are extracted from an ensemble of subband components of high-dimensional acoustic waveforms and thus retain more information about speech that is potentially relevant to discrimination of phonetic units than the corresponding cepstral representations. In addition to investigating the robustness of the proposed front-end to additive noise, we also assess its robustness to linear filtering due to room reverberation. This form of distortion causes temporal smearing of short-term spectra, which degrades the performance of ASR systems. This can be attributed primarily to the use of analysis windows for feature extraction in the conventional front-ends such as MFCC that are much shorter than typical room impulse responses. Furthermore, the distortion caused by linear filtering is correlated with the underlying speech signal. Hence, conventional methods for robust ASR that are tuned for recognition of data corrupted by additive noise only will not be effective in reverberant environments. Several speech dereverberation techniques that rely on multi-channel recordings of speech such as [28, 29] exist in the literature. However, these considerations extend beyond the scope of this paper and instead, standard single channel feature compensation methods for additive noise and linear filtering such as VTS and CMVN compensation are used throughout this paper.

Robustness of the proposed front-end to additive noise and linear filtering is demonstrated by its comparison with the MFCC front-end on a phoneme classification task; this task remains important in comparing different methods and representations [21, 30–39]. The improvements achieved in the classification task can be expected to extend to continuous speech recognition [40, 41], where SVMs have been employed in hybrid frameworks [41, 42] with hidden Markov models (HMMs) as well as in frame-based architectures using the token passing algorithm [43]. Our results demonstrate the benefits of the subband classification in terms of robustness to additive noise and linear filtering. The subband-waveform classifiers outperform even MFCC classifiers that are trained and tested under matched conditions for signal-to-noise ratios below 6-dB. Furthermore, in classifying noisy reverberant speech, the subband classifier outperforms the MFCC classifier compensated using VTS for all signal-to-noise ratios (SNRs) lying below a crossover point between 12-dB and 6-dB. Finally, their convex combination yields further performance improvements over both individual classifiers.

The paper is organized as follows: the proposed subband classification approach is described in Section II. Experimental

results that demonstrate its robustness to additive noise and linear filtering are presented in Section III. Finally, Section IV draws some conclusions and suggests possible future work towards application of the proposed front-end in continuous speech recognition tasks.

## II. SUBBAND CLASSIFICATION USING SUPPORT VECTOR MACHINES

Support vector machines (SVMs) are receiving increasing attention as a tool for speech recognition applications due to their good generalization properties [17, 35, 41, 42, 44–46]. Here we use them in conjunction with the proposed subband-based representation, with the aim of improving the robustness of the standard cepstral front-end to noise and filtering. To this end we construct a fixed-length representation that could potentially be used as the front-end for a continuous speech recognition systems based on *e.g.* hidden Markov models (HMMs) [41–43], as highlighted in Section IV. Dealing with variable phoneme length has been addressed by means of generative kernels such as Fisher kernels [45, 47] and dynamic time-warping kernels [48], but lies beyond the scope of this paper.

### A. Support Vector Machines

A binary SVM classifier estimates a decision surface that jointly maximizes the margin between the two classes and minimizes the misclassification error on the training set. For a given training set  $(\mathbf{x}_1, \dots, \mathbf{x}_p)$  with corresponding class labels  $(y_1, \dots, y_p)$ ,  $y_i \in \{+1, -1\}$ , an SVM classifies a test point  $\mathbf{x}$  by computing a score function,

$$h(\mathbf{x}) = \sum_{i=1}^p \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) + b \quad (1)$$

where  $\alpha_i$  is the Lagrange multiplier corresponding to the  $i^{\text{th}}$  training sample,  $\mathbf{x}_i$ ,  $b$  is the classifier bias – these parameters are optimized during training – and  $K$  is a kernel function. The class label of  $\mathbf{x}$  is then predicted as  $\text{sgn}(h(\mathbf{x}))$ . While the simplest kernel  $K(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \mathbf{x}, \tilde{\mathbf{x}} \rangle$  produces linear decision boundaries, in most real classification tasks the data is not linearly separable. Nonlinear kernel functions implicitly map data points to a high-dimensional feature space where the data could potentially be linearly separable. Kernel design is therefore effectively equivalent to feature-space selection, and using an appropriate kernel for a given classification task is crucial. Commonly used is the polynomial kernel,  $K_p(\mathbf{x}, \tilde{\mathbf{x}}) = (1 + \langle \mathbf{x}, \tilde{\mathbf{x}} \rangle)^\Theta$ , where the polynomial order  $\Theta$  in  $K_p$  is a hyper-parameter that is tuned to a particular classification problem. More sophisticated kernels can be obtained by various combinations of basic SVM kernels. Here we use a polynomial kernel for classification with cepstral features (MFCC) whereas classification with acoustic waveforms in frequency subbands is performed using a custom-designed kernel described in the following subsection.

For multiclass problems, binary SVMs are combined via error-correcting output codes (ECOC) methods [49, 50]. In this work, for an  $M$ -class problem we train  $N = M(M - 1)/2$

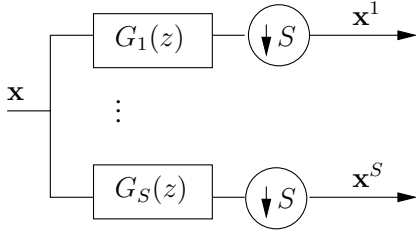


Fig. 1: Decomposition of phonemes into its subband components using an  $S$ -channel cosine modulated filter bank.

binary pairwise classifiers, primarily to lower the computational complexity. The training scheme can be captured in a coding matrix  $w_{mn} \in \{0, 1, -1\}$ , i.e. classifier  $n$  is trained only on data from the two classes  $m$  for which  $w_{mn} \neq 0$ , with  $\text{sgn}(w_{mn})$  as the class label. One then predicts for test input  $\mathbf{x}$  the class that minimizes the loss  $\sum_{n=1}^N \chi(w_{mn} f_n(\mathbf{x}))$  where  $f_n(\mathbf{x})$  is the output of the  $n^{\text{th}}$  binary classifier and  $\chi$  is a loss function. We experimented with a variety of loss functions, including hinge, Hamming, exponential and linear. The hinge loss function  $\chi(z) = \max(1 - z, 0)$  performed best and is therefore used throughout.

### B. Kernels for Subband Classification

The acoustic waveform features used for classification in frequency subbands are obtained from fixed-length,  $D$  samples long, acoustic waveform segments that will be denoted by  $\mathbf{x}$ . These features will be studied in comparison with MFCC features obtained from the same speech segments  $\mathbf{x}$ . To obtain the subband features, each waveform segment  $\mathbf{x}$  is first decomposed into its subband components,  $\mathbf{x}^s, s = 1, \dots, S$ , by means of an  $S$ -channel maximally-decimated perfect reconstruction cosine modulated filter bank (CMFB) [51], as shown in Figure 1. Several other subband decompositions such as discrete wavelet transform, wavelet packet decomposition and discrete cosine transform also achieved comparable, albeit somewhat inferior performance. A summary of the classification results obtained with different subband decompositions in quiet conditions is presented in Section III-B. The CMFB consists of a set of orthonormal analysis filters

$$g_s[k] = \frac{1}{\sqrt{S}} g[k] \cos\left(\frac{2s-1}{4S} (2k-S-1)\pi\right), \quad (2)$$

where  $s = 1, \dots, S$ ,  $k = 1, \dots, 2S$ , and  $g[k] = \sqrt{2} \sin(\pi(k-0.5)/2S)$ , is a low-pass prototype filter. Such a filter bank implements an orthogonal transform, hence the collection of the subband components is a representation of the original waveform in a different coordinate system [51]. A maximally-decimated filter bank was chosen primarily because the sub-sampling operation avoids introducing additional redundancies and thus limits the overall computational burden. However, we believe that redundant expansions of speech signals obtained using over-sampled filter banks could be advantageous to effectively account for the shift invariance of speech.

For classification in frequency subbands, an SVM kernel is constructed by partly following steps from our previous

work [17], which attempted to capture known invariances or express explicitly the waveform properties which are known to correlate with phoneme identity. First, an even kernel is constructed from a baseline polynomial kernel  $K_p$  to account for the sign-invariance of human speech perception as

$$K_e(\mathbf{x}^s, \mathbf{x}_i^s) = K'_p(\mathbf{x}^s, \mathbf{x}_i^s) + K'_p(\mathbf{x}^s, -\mathbf{x}_i^s) \quad (3)$$

where  $K'_p$  is a modified polynomial kernel given by

$$K'_p(\mathbf{x}^s, \mathbf{x}_i^s) = K_p\left(\frac{\mathbf{x}^s}{\|\mathbf{x}^s\|}, \frac{\mathbf{x}_i^s}{\|\mathbf{x}_i^s\|}\right) = \left(1 + \left\langle \frac{\mathbf{x}^s}{\|\mathbf{x}^s\|}, \frac{\mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right\rangle\right)^\Theta. \quad (4)$$

The kernel  $K'_p$ , which acts on normalized input vectors, will be used as a baseline kernel for the acoustic waveforms.

Next, the temporal dynamics of speech are explicitly taken into account by means of features that capture the evolution of energy in individual subbands. To obtain these features, waveform synthesized from each decimated subband component  $\mathbf{x}^s$  is first divided into  $T$  frames which are denoted by  $\mathbf{x}^{t,s}, t = 1, \dots, T$ . Then a vector of the frame energies is formed as  $\boldsymbol{\omega}^s = [\log \|\mathbf{x}^{1,s}\|^2, \dots, \log \|\mathbf{x}^{T,s}\|^2]$ . Finally, time differences [52] of  $\boldsymbol{\omega}^s$  are evaluated to form the dynamic subband feature vector  $\boldsymbol{\Omega}^s$  as  $\boldsymbol{\Omega}^s = [\boldsymbol{\omega}^s \ \Delta \boldsymbol{\omega}^s \ \Delta^2 \boldsymbol{\omega}^s]$ . This dynamic subband feature vector  $\boldsymbol{\Omega}^s$  is then combined with the corresponding acoustic waveform subband component  $\mathbf{x}^s$  in evaluating a kernel  $K_\Omega$  given by

$$K_\Omega(\mathbf{x}^s, \mathbf{x}_i^s, \boldsymbol{\Omega}^s, \boldsymbol{\Omega}_i^s) = K_e(\mathbf{x}^s, \mathbf{x}_i^s) K_p(\boldsymbol{\Omega}^s, \boldsymbol{\Omega}_i^s), \quad (5)$$

where  $\boldsymbol{\Omega}_i^s$  is the dynamic subband feature vector corresponding to the  $s^{\text{th}}$  subband component  $\mathbf{x}_i^s$  of the  $i$ -th training point  $\mathbf{x}_i$ .

On the other hand, the standard polynomial kernel  $K_p$  is used for classification with the cepstral representations, since they consist of analogous subband features in the compressed domain, while feature standardization by cepstral mean-and-variance normalization (CMVN) [13] ensures that feature vectors typically already have unit norm.

### C. Ensemble Methods

For each binary classification problem, decomposing an acoustic waveform into its subband components produces an ensemble of  $S$  classifiers. The decision of the subband classifiers  $s = 1, \dots, S$  in the ensemble, given by

$$f^s(\mathbf{x}^s, \boldsymbol{\Omega}^s) = \sum_i \alpha_i^s y_i K_\Omega(\mathbf{x}^s, \mathbf{x}_i^s, \boldsymbol{\Omega}^s, \boldsymbol{\Omega}_i^s) + b^s, \quad (6)$$

are then aggregated using ensemble methods to obtain the binary classification decision for a test waveform  $\mathbf{x}$ . Here  $\alpha_i^s$  and  $b^s$  are respectively the Lagrange multiplier corresponding to  $\mathbf{x}_i^s$  and the bias of the  $s^{\text{th}}$  subband binary classifier.

1) *Uniform Aggregation*: Here the decisions of the subband classifiers in the ensemble are assigned uniform weights. Majority voting is the simplest uniform aggregation scheme commonly used in machine learning. In our context it is equivalent to forming a meta-level score function as

$$h(\mathbf{x}) = \sum_{s=1}^S \text{sgn}(f^s(\mathbf{x}^s, \boldsymbol{\Omega}^s)), \quad (7)$$



then predicting the class label as  $y = \text{sgn}(h(\mathbf{x}))$ . In addition to this conventional majority voting scheme, which maps the scores in individual subbands to the corresponding class labels ( $\pm 1$ ), we also considered various smooth squashing functions, *e.g.* sigmoidal, as alternatives to the  $\text{sgn}$  function in (7), and obtained similar results. To gain some intuition about the potential of ensemble methods such as the majority voting in improving classification performance, consider the ideal case when the errors of the individual subband classifiers in the ensemble are independent with error probability  $p < 1/2$ . Under these conditions, a simple combinatorial argument shows that the error probability  $p_e$  of the majority voting scheme is given by

$$p_e = \sum_{s=\lceil S/2 \rceil}^S \binom{S}{s} p^s (1-p)^{S-s} . \quad (8)$$

where the largest contribution to the overall error is due to the term with  $s = \lceil S/2 \rceil$ . For large ensemble cardinality  $S$ , this error probability can be bounded as:

$$p_e < p^{\lceil S/2 \rceil} (1-p)^{S-\lceil S/2 \rceil} \sum_{s=\lceil S/2 \rceil}^S \binom{S}{s} \approx \frac{1}{2} (4p(1-p))^{S/2} . \quad (9)$$

Therefore, in ideal conditions, the ensemble error decreases exponentially in  $S$  even with this simple aggregation scheme [53, 54]. However, it has been shown that there exists a correlation between the subband components of speech and the resulting speech recognition errors in individual frequency subbands [19, 20]. As a result, the majority voting scheme may not yield significant improvements in classification performance, particularly at low SNRs. Uniform aggregation schemes further suffer from a major drawback: they do not exploit the differences in the relative importance of individual subbands in discriminating among specific pairs of phonemes. To remedy this, we use stacked generalization [55] as discussed next, to explicitly learn weighting functions specific to each pair of phonemes for non-uniform aggregation of the outputs of base-level SVMs.

2) *Stacked Generalization*: Our practical implementation of stacked generalization [55] consists of a hierarchical two-layer SVM architecture, where the outputs of subband base-level SVMs are aggregated by a meta-level linear SVM. The decision function of the meta-level SVM classifier is of the form

$$h(\mathbf{x}) = \langle \mathbf{f}(\mathbf{x}), \mathbf{w} \rangle + v = \sum_s w^s f^s(\mathbf{x}^s, \Omega^s) + v , \quad (10)$$

where  $\mathbf{f}(\mathbf{x}) = [f^1(\mathbf{x}^1, \Omega^1), \dots, f^S(\mathbf{x}^S, \Omega^S)]$  is the base-level SVM score vector of the test waveform  $\mathbf{x}$ ,  $v$  is the classifier bias, and  $\mathbf{w} = [w^1, \dots, w^S]$  is the weight vector of the meta-level classifier. Note each of the binary classifiers will have its own weight vector, determined from an independent development/validation set  $\{\tilde{\mathbf{x}}_j, \tilde{y}_j\}$ . Each weight vector can, therefore, be expressed as

$$\mathbf{w} = \sum_j \beta_j \tilde{y}_j \mathbf{f}(\tilde{\mathbf{x}}_j), \quad (11)$$

where  $\mathbf{f}(\tilde{\mathbf{x}}_j) = [f^1(\tilde{\mathbf{x}}_j^1, \tilde{\Omega}_j^1), \dots, f^S(\tilde{\mathbf{x}}_j^S, \tilde{\Omega}_j^S)]$  is the base-level SVM score vector of the training waveform  $\tilde{\mathbf{x}}_j$ , and  $\beta_j$  and  $\tilde{y}_j$  are the Lagrange multiplier and class label corresponding to  $\mathbf{f}(\tilde{\mathbf{x}}_j)$ , respectively. While a base-level SVM assigns a weight to each supporting feature vector, stacked generalization effectively assigns an additional weight  $w^s$  to each subband based on the performance of the corresponding base-level subband classifier. ECOC methods are then used to combine the resulting meta-level binary classifiers for multiclass classification.

An obvious advantage of the subband approach for ASR is that the effect of environmental distortions in sufficiently narrow subbands can be approximated as similar to that of narrow-band white noise. This, in turn, facilitates the compensation of features independently of the spectral characteristics of the additive and convolutive noise sources. In a preceding paper [17], we proposed an ASR front-end based on the full-band acoustic waveform representation of speech, where a spectral shape adaptation of the features was performed in order to account for the varying strength of contamination of the frequency components due to the presence of colored noise. In this work, compensation of the features is performed using standard approaches such as cepstral mean-and-variance normalization (CMVN) and vector Taylor series (VTS), which do not require any prior knowledge of the additive and convolutive noise sources. Furthermore, we found that the stacked generalization also depends on the level of noise contaminating its training data. The weight vectors of the stacked classifiers can then be tuned for classification in a particular environment by introducing similar distortion to its training data. In scenarios where a performance gain over a wide range of SNRs is desired, a multi-style training approach can be employed that offers a reasonable compromise between various test conditions. For instance, a meta-level classifier can be trained using the score feature vectors of noisy data or the score feature vectors of a mixture of clean and noisy data.

Note that since the dimension,  $S$ , of the score feature vectors that form the input to the stacked subband classifier is very small compared to the typical MFCC or waveform feature vectors, only a very limited amount of data is required to learn optimal weights of the meta-level classifiers. As such, stacked generalization offers flexibility and some coarse frequency selectivity for the individual binary classification problems, and can be particularly useful in de-emphasizing information from unreliable subbands. The experiments presented in this paper show that the subband approach attains major gains in classification performance over its full-band counterpart [17] as well as over state-of-the-art front-ends such as MFCC.

### III. EXPERIMENTAL RESULTS

#### A. Experimental Setup

Experiments are performed on the ‘si’ (diverse) and ‘sx’ (compact) sentences of the TIMIT database [56]. The training set consists of 3696 sentences from 168 different speakers. For testing we use the core test set which consists of 192 sentences from 24 different speakers not included in the training set. The development set consists of 1152 sentences uttered by 96 male

and 48 female speakers not included in either the training or the core test set, with speakers from 8 different dialect regions. In training the meta-level subband classifiers, we use a small subset, randomly selecting an eighth of the data points in the complete TIMIT development set. The glottal stops /q/ are removed from the class labels and certain allophones are grouped into their corresponding phoneme classes using the standard Kai-Fu Lee clustering [57], resulting in a total of  $M = 48$  phoneme classes and  $N = M(M - 1)/2 = 1128$  classifiers. Among these classes, there are 7 groups for which the contribution of within-group confusions is not counted toward the multiclass error, again following standard practice [35, 57]. Initially, we experimented with different values of the hyperparameters for the binary SVM classifiers but decided to use fixed values for all classifiers as parameter optimization had a large computational overhead but only a small impact on the multiclass classification error: the degree of  $K_p$  is set to  $\Theta = 6$  and the penalty parameter (for slack variables in the SVM training algorithm) to  $C = 1$ .

To test the classification performance in noise, each TIMIT test sentence is normalized to unit energy per sample and then a noise sequence is added to the entire sentence to set the sentence-level SNR. Hence for a given sentence-level SNR, signal-to-noise ratio at the level of individual phonemes will vary widely. Both artificial noise (white, pink) and recordings of real noise (speech-babble) from the NOISEX-92 database are used in our experiments. White noise was selected due to its attractive theoretical interpretation as probing in an isotropic manner the separation of phoneme classes in different representation domains. Pink noise was chosen because  $1/f$ -like noise patterns are found in music melodies, fan and cockpit noises, in nature etc. [58–60]. In order to further test the classification performance in the presence of linear filtering, noisy TIMIT sentences are convolved with an impulse response with reverberation time  $T_{60} = 0.2\text{sec}$ . This impulse response is one that was measured using an Earthworks QTC1 microphone in the ICSI conference room [61] populated with people; its magnitude response  $R(e^{j\omega})$  is shown in Figure 2, where we also show the spectrum of an impulse response corresponding to a different speaker position in the same room,  $R'(e^{j\omega})$ . While the substantial difference between these filters is evident from their spectra and spectral colorations (defined as a ratio of the geometric mean to the arithmetic mean of spectral magnitude),  $R'(e^{j\omega})$  can be viewed as an approximation of the effect of  $R(e^{j\omega})$  on the speech spectrum and is used in some of our experiments for training of the cepstral and meta-level subband classifiers in order to reduce the mismatch between training and test data.

Acoustic waveforms segments  $\mathbf{x}$  are extracted from the TIMIT sentences by applying a 100ms rectangular window at the centre of each phoneme and are then decomposed into subband components  $\{\mathbf{x}^s\}_{s=1}^S$  using a cosine-modulated filter bank (see (2)). We conducted experiments to examine the effect of the number of filter bank channels  $S$  on classification accuracy. Generally, decomposition of speech into wider subbands does not effectively capture the frequency-specific dynamics of speech and thus results in relatively poor performance. On the other hand, decomposition of speech into

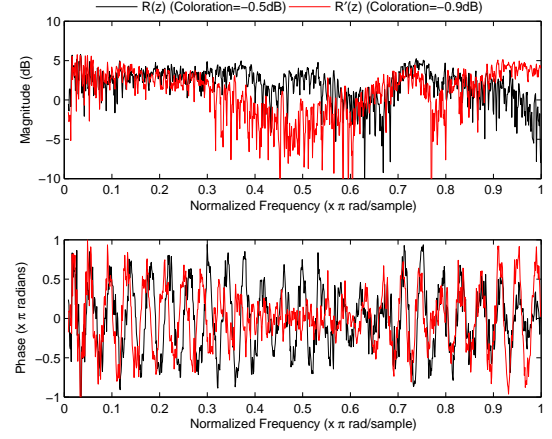


Fig. 2: Frequency response of the ICSI conference room filters with spectral coloration -0.5-dB and -0.9-dB. Here, spectral coloration is defined as the ratio of the geometric mean to the arithmetic mean of spectral magnitudes.  $R(z)$  is used to add reverb to the test data whereas  $R'(z)$ , a proxy filter recorded at a different location in the same room, is used for the training of cepstral and meta-level subband classifiers.

sufficiently narrow subbands improves classification performance as demonstrated in [22], but at the cost of an increase in overall computational complexity. In order to avoid the latter, the number of filter bank channels is limited to  $S = 16$  for all results presented in this paper. The dynamic subband feature vector,  $\Omega^s$ , is computed by extracting  $T = 10$  equal-length (25ms with an overlap of 10ms) frames around the centre of each phoneme, yielding a vector of dimension 30. These feature vectors are further standardized within each sentence of TIMIT for the evaluation of kernel  $K_\Omega$ . Note that the training of base-level SVM subband classifiers is always performed with clean data. The development subset is then used for training of the meta-level subband classifiers as learning the optimal weights requires only a very limited amount of data. Several scenarios are considered for training of the meta-level classifiers:

1. **Anechoic clean training** - training the meta-level SVM classifier with the base-level SVM score vectors obtained from anechoic clean data.
2. **Anechoic multi-style training** - training the meta-level SVM classifier with the base-level SVM score vectors of anechoic data containing a mixture of clean waveforms and waveforms corrupted by white noise at 0-dB SNR,
3. **Reverberant multi-style training** - training the meta-level SVM classifier with the base-level SVM score vectors of reverberant data containing a mixture of clean waveforms and waveforms corrupted by white noise at 0-dB SNR. Two particular cases of this scenario are considered. (a) The development data for training as well as the test data are convolved with the same filter  $R(e^{j\omega})$ . This case provides a lower bound on the classification error by assuming that exact knowledge of the convolutive distortion is available. (b) The development data for training is convolved with  $R'(e^{j\omega})$  whereas the test data is convolved with  $R(e^{j\omega})$ . This probes the effects of

a mismatch of the linear filter used for convolution with the training and test data. Since the exact properties of the linear filter corrupting the test data are usually difficult to determine, this scenario is more practical and its performance is expected to lie within the bracket formed by the two scenarios mentioned above, *i.e.* anechoic training, and reverberant training and testing using the same filter.

4. **Matched training** - training and testing with the meta-level classifier in conditions of identical noise level and type. Results for this scenario are shown only in the presence of additive noise.

To obtain the cepstral (MFCC) representation, each sentence is converted into a sequence of 13 dimensional feature vectors. These feature vectors are further combined with their first and second order time derivatives to form a sequence of 39 dimensional feature vectors. Then,  $T = 10$  frames (with frame duration of 25ms and a frame rate of 100 frames/sec) closest to the center of a phoneme are concatenated to give a representation in  $\mathbb{R}^{390}$ . Noise compensation of the MFCC features is performed via the vector Taylor series (VTS) method, which has been extensively used in the recent literature and can be considered as state-of-the-art. VTS estimates the distribution of noisy speech given the distribution of clean speech, a segment of noisy speech, and the Taylor series expansion that relates the noisy speech features to the clean ones, and then uses it to predict the unobserved clean cepstral feature vectors. In our experiments, a Gaussian mixture model (GMM) with 64 mixture components was used to learn the distribution of the Mel-log spectra of clean training data. In addition to VTS, cepstral mean-and-variance normalization (CMVN) [13, 14] is performed to standardize the cepstral features, fixing their range of variation for both training and test data. CMVN computes the mean and variance of the feature vectors across a sentence and standardizes the features so that each has zero mean and a fixed variance. The following training-test scenarios are considered for classification with the cepstral front-end:

1. **Anechoic training with VTS** - training of the SVM classifiers is performed with anechoic clean speech and the test data is compensated via VTS.
2. **Reverberant training with VTS** - training of the SVM classifiers is performed with reverberant clean speech with feature compensation of the test data via VTS. Again, two particular cases in this scenario are considered. (a) The clean training data and the noisy test data are convolved with the same linear filter,  $R(e^{j\omega})$ . (b) The data used for training of the SVM classifiers as well as learning of the distribution of log-spectra in VTS feature compensation is convolved with  $R'(e^{j\omega})$  while the test data is convolved with  $R(e^{j\omega})$ .
3. **Matched training** - In this scenario, the training and testing conditions are identical. Again, this is an impractical target; nevertheless, we present the results (only in the presence of additive noise) as a reference, since this setup is considered to give the optimal achievable performance with cepstral features [14, 62, 63].

Next, we present the results of TIMIT phoneme classification with the setup detailed above.

### B. Results: Robustness to Additive Noise

First we compare various frequency decompositions and ensemble methods for subband classification. A summary of their respective classification errors in quiet condition is presented in Table I. Stacked generalization yields significantly better results than majority voting; it consistently achieves over 10% improvement over majority voting for all subband decompositions considered here. Among these decompositions, classification with the 16-channel cosine-modulated filter bank achieves the largest improvement of 5.5% over the composite acoustic waveforms [17] and is therefore selected for further experiments.

TABLE I: Errors obtained with different subband decompositions [51] (listed in the left column) and aggregation schemes for subband classification in quiet condition.

Subband Analysis	ERROR [%]	
	Maj. Voting	Stack. Gen.
Level-4 wavelet decomposition	43.7	<b>31.8</b>
Level-4 wavelet packet decomposition	45.1	<b>33.1</b>
DCT (16 uniform-width bands)	44	<b>32.6</b>
16-channel CMFB	42.4	<b>31.2</b>
Composite Waveform [17]	<b>36.7</b>	

Let us now consider classification of phonemes in the presence of additive noise; robustness of the proposed method to both additive noise and linear filtering will be discussed in Section III-C. In Figure 3, we compare the classification in frequency subbands using ensemble methods with composite acoustic waveform classification (as reported in [17]) in the presence of white and pink noise. The dashed curves correspond to subband classification using ensemble methods, in particular, uniform combination (majority voting) and stacked generalization with different training scenarios for meta-level classifiers (see Section III-A): quiet (scenario 1), multi-style (scenario 2), and matched (scenario 4). The results show that stacked generalization achieves significantly better performance than uniform aggregation. The majority voting scheme even performs poorly in comparison with the composite acoustic waveforms across all SNRs. On the other hand, even the stacked subband classifier trained only in quiet conditions (scenario 1) improves over the composite waveform classifier in low noise conditions. However, its performance then degrades relatively quickly in high noise because its meta-level binary classifiers are trained to assign weights to different subbands that are tuned for classification in quiet. To improve the robustness to additive noise, the meta-level classifiers can be trained in a multi-style manner (scenario 2). Figure 4 shows the resulting weights (mean  $\pm$  standard deviation across  $N = 1128$  binary classifiers) for the  $S = 16$  subbands. Relatively high weights are assigned to the low frequency subband components. This is reasonable as these subbands hold a substantial portion of speech energy and can provide reliable discriminatory information in the presence of wideband noise. The large amount of variation in the assigned weights as indicated by the error bars is consistent with the



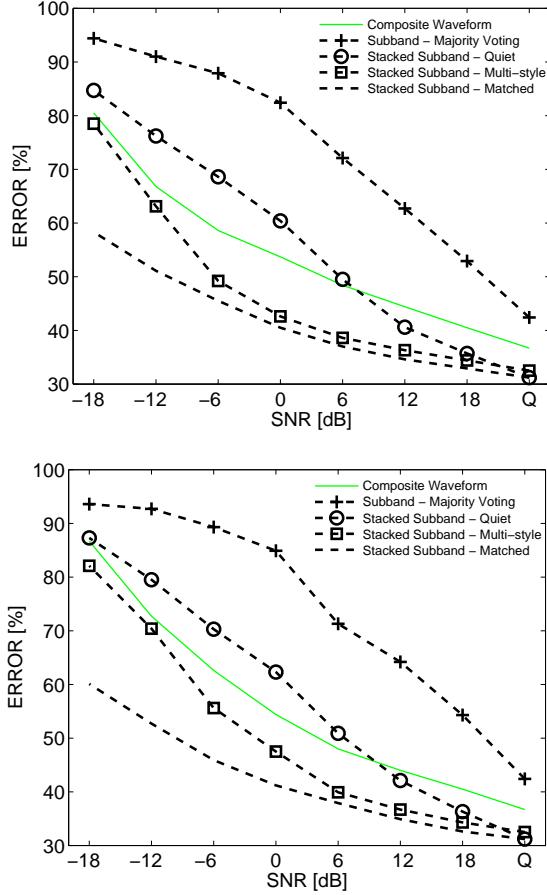


Fig. 3: Ensemble methods for aggregation of subband classifiers and their comparison with composite acoustic waveform classifiers (as reported in [17]) in the presence of white noise (top) and pink noise (bottom). The curves correspond to uniform combination (majority voting) and stacked generalization with different training scenarios for the meta-level classifiers. The multi-style stacked subband classifier is trained only with the small development subset (one eighth randomly selected score vectors from the development set) consisting of clean and white-noise (0-dB SNR) corrupted anechoic data. The classifiers are then tested on data corrupted with white noise (matched noise type) and pink noise (mismatched).

variation of speech data encountered by the  $N = 1128$  binary phoneme classifiers. In terms of the resulting performance (Figure 3), the multi-style subband classifiers consistently improves over the composite waveform classifier as well as the stacked subband classifier trained in quiet condition. Overall, across the range of SNRs considered, it achieves average improvements of 6.8% and 5.9% over the composite waveform classifier in the presence of white (matched noise type) and pink (mismatched noise type) noise, respectively. As expected, the stacked subband classifier trained in matched conditions, finally, outperforms the other classifiers in all noise conditions.

Next, we compare the performance of the multi-style subband classifier with the VTS-compensated MFCC classifier and the composite acoustic waveform classifier [17] in the presence of additive white and pink noise. These results along with classification with the stacked subband classifier and MFCC classifier, both in matched training-test conditions,

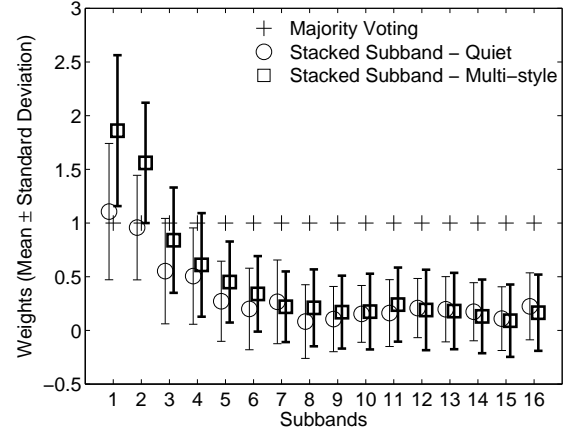


Fig. 4: Weights (mean  $\pm$  standard deviation across  $N = 1128$  binary classifiers) assigned to  $S = 16$  subbands by the multi-style meta-level classifiers and by the meta-level classifiers trained in quiet conditions.

are presented in Figure 5. The results show that the stacked subband classifier exhibits better classification performance than the VTS-compensated MFCC classifier for SNR below 12-dB whereas the performance crossover between MFCC and composite acoustic waveform classifiers occurs later, between 6-dB and 0-dB SNR. The stacked subband classifier achieves average improvements of 8.7% and 4.5% over the MFCC classifier in the presence of white and pink noise, respectively. Moreover, and quite remarkably, the stacked subband classifier also significantly improves over the MFCC classifier trained and tested in matched conditions for SNRs below a crossover point between 6-dB and 0-dB SNR, even though its meta-level classifiers are trained only using clean data and data corrupted by white noise at 0-dB SNR and the number of data points used to learn the optimal weights amounts only to a small fraction of the data set used for training of the MFCC classifier in matched conditions. In particular, an average improvement of 6.5% in the phoneme classification error is achieved by the multi-style subband classifier over the matched MFCC classifier for SNRs below 6-dB in the presence of white noise.

In [64] we showed that the MFCC classifiers suffer performance degradation in case of a mismatch of the noise type between training and test data. On the other hand, the stacked subband classifier degrades gracefully in a mismatched environment as shown in Figure 5. This can be attributed to the decomposition of acoustic waveforms into frequency subbands where the effect of wideband colored noise on each binary subband classifier can be approximated as that of narrow-band white noise. In comparison to the result reported in [33], where a 77.8% error was obtained at 0-dB SNR in pink noise using a second-order regularized least squares algorithm (RLS2) trained using MFCC feature vectors with variable length encoding, our proposed method achieves a 30% improvement in similar conditions, using only a fixed length representation. Figure 5 also shows a comparison of the stacked subband classifier with the MFCC classifier when both are trained and tested in matched conditions. The matched-condition subband classifier significantly outperforms

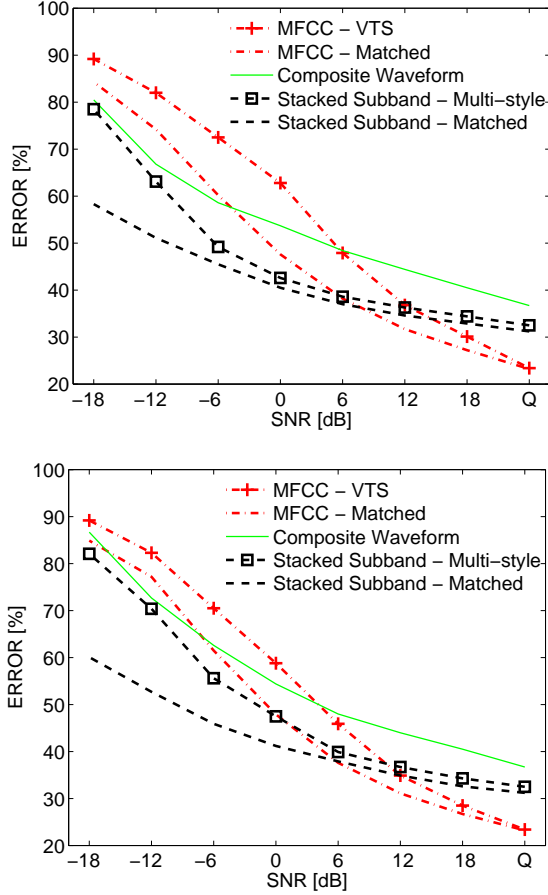


Fig. 5: SVM classification in the subbands of acoustic waveforms and its comparison with MFCC and composite acoustic waveform classifiers in the presence of white noise (top) and pink noise (bottom). The multi-style stacked subband classifier is trained only with a small subset of the development data (one eighth randomly selected score vectors from the development set) consisting of clean and white-noise (0-dB SNR) corrupted data. In the matched training case, noise levels as well as noise types of training and test data are identical for both MFCC and stacked subband classifiers.

the matched MFCC classifier for SNRs below 6-dB. Around 13% average improvement is achieved by the subband classifier over the MFCC classifier for SNRs below 6-dB, in the presence of both white and pink noise. This suggests that the high-dimensional subband representation obtained from acoustic waveforms provides a better separation of phoneme classes than cepstral representations.

### C. Results: Robustness to Linear Filtering

We now consider classification in the presence of additive noise as well as linear filtering. First, Figure 6 presents results of the ensemble subband classification using stacked generalization with multiple training-test scenarios (see Section III-A) in the presence of white and pink noise. To reiterate, three different scenarios are considered for training of the multi-style stacked subband classifier: one involves training the meta-level classifiers with the base-level SVM score vectors of the development subset consisting of clean and white-noise (0-dB SNR)

corrupted anechoic data, one involves training with the score vectors of the same development data convolved with  $R'(e^{j\omega})$  (mismatched reverberant conditions) while the last involves training in matched reverberant conditions, *i.e.* training with the same development subset convolved with  $R(e^{j\omega})$ . These classifiers, which we refer to as *anechoic* and *reverberant* multi-style subband classifiers (see Section III-A), are then tested on data corrupted by white, pink or speech-babble noise, and convolved with  $R(e^{j\omega})$ . Similar to our findings in the previous section, the results in Figure 6 show that the anechoic multi-style subband classifier consistently improves over the stacked subband classifier trained only in quiet condition. Moreover, the reverberant multi-style subband classifiers (both matched and mismatched) further reduce the mismatch with the test data and hence improve the performance further. For instance, in the presence of pink noise and linear filtering, the subband classifiers trained in mismatched and matched reverberant conditions attain average improvements of 6% and 8.5% across all SNRs over the anechoic multi-style subband classifier, respectively. Note that an accurate measurement of the linear filter corrupting the test data may be difficult to obtain in practical scenarios. Nonetheless, classification results in matched reverberant condition are presented as a lower bound on the error. On the other hand, the mismatched reverberant case can be considered as a more practical solution to the problem and its performance lies as expected between anechoic training and matched reverberant training.

Figure 7 compares the classification performances of the subband and VTS-compensated MFCC classifiers trained under three different scenarios (see Section III-A) in the presence of linear filtering, and pink and speech-babble noise. The first, anechoic training, represents an agnostic case that does not rely on any information at all regarding the source of the convolutive noise  $R(e^{j\omega})$ . The reverberant mismatch case uses a proxy reverberation filter  $R'(e^{j\omega})$  in order to reduce the mismatch of the training and the reverberant test environments up to a certain degree, whereas the reverberant matched case employs accurate knowledge of the reverberation filter  $R(e^{j\omega})$  in the training of the MFCC classifiers and the meta-level subband classifiers. These training scenarios are respectively represented by squares, stars and circles in Figure 7. The results show that the comparisons of the stacked subband classifiers and MFCC classifiers under the different training regimes exhibit similar trends. Generally speaking, the MFCC classifier outperforms the corresponding subband classifier in quiet and low noise conditions, while the latter yields significant improvements in high noise conditions, with a crossover at moderate noise levels. For example, the anechoic subband classifiers yields better classification performance than the anechoic MFCC classifier for SNRs below a crossover point between 12-dB and 6-dB. Quantitatively similar conclusions apply to the comparative performances of the MFCC and subband classifiers in the reverberant training scenarios. Under the three different training regimes and two different noise types, the subband classifiers attain an average improvement of 8.2% over the MFCC classifiers across all SNRs below 12-dB. Note that in the reverberant training scenarios, the MFCC classifier is trained with the complete TIMIT reverberant train-



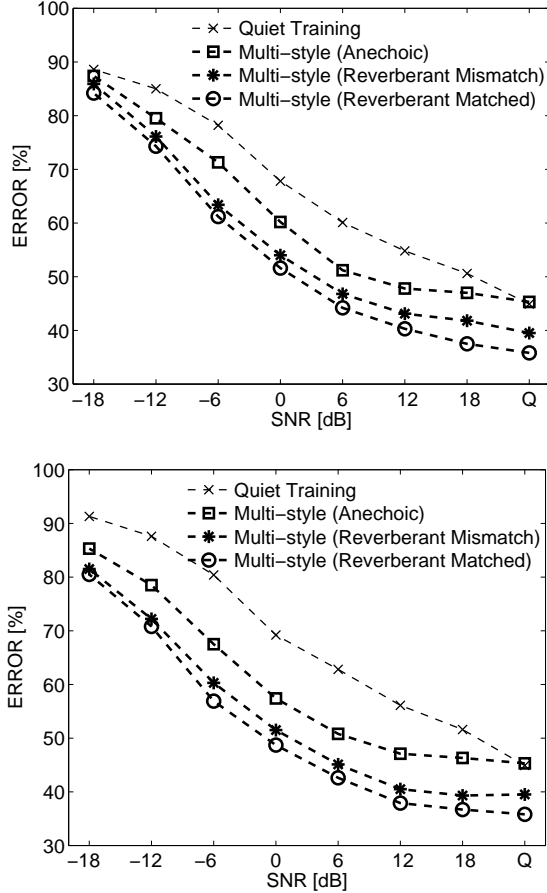


Fig. 6: Classification in frequency subbands using ensemble methods in the presence of the linear filter  $R(e^{j\omega})$  with white noise (top) and pink noise (bottom). The curves correspond to stacked generalization with different training scenarios for the meta-level subband classifier.

ing set. On the other hand, the meta-level subband classifier is trained using the reverberant development subset with a number of data points less than 4% of that in the TIMIT training set. Moreover, the dimension of the feature vectors that form the input to the meta-level classifiers is almost 24 times smaller than that of the MFCC feature vectors. As such, the subband approach offers more flexibility in terms of training and adaptation of the classifiers to a new environment.

Since an obvious performance crossover between the subband and MFCC classifiers exists at moderate SNRs, we also consider a convex combination of the scores of the SVM classifiers with a combination parameter  $\lambda$  as discussed in [17]. Here  $\lambda = 0$  corresponds to the MFCC classification whereas  $\lambda = 1$  corresponds to the subband classification. The combination approach was also motivated by the differences in the confusion matrices of the two classifiers (not shown here). This suggests that the errors of the subband and MFCC classifiers may be independent up to a certain degree and therefore a combination of the two may yield better performance than either of classifiers individually. Two different values of the combination parameter  $\lambda$  are considered. First, the value of  $\lambda$  is set to  $1/2$  which corresponds to the arithmetic mean of the MFCC and subband SVM classifier scores. In

the second case, we set the combination parameter  $\lambda$  to a function  $\lambda_{\text{emp}}(\sigma^2)$  which approximates the optimal combination parameter values for an independent development set. This approximated function was determined empirically in our experiments with composite waveforms [17] and is given by  $\lambda_{\text{emp}}(\sigma^2) = \eta + \zeta/[1 + (\sigma_0^2/\sigma^2)]$ , with  $\eta = 0.2$ ,  $\zeta = 0.5$  and  $\sigma_0^2 = 0.03$ . Note that  $\lambda_{\text{emp}}(\sigma^2)$  requires an estimate of the noise variance ( $\sigma^2$ ) which we obtained for our data using the decision-directed estimation algorithm [65, 66].

Figure 8 compares the classification performance of the subband and MFCC classifiers with their convex combination in the presence of speech-babble noise and filtering with  $R(e^{j\omega})$ , under the anechoic and reverberant mismatched training regimes. The combined classification with  $\lambda_{\text{emp}}$  consistently outperforms either of the individual classifiers across all SNRs. For instance, under the anechoic training of the classifiers, the combined classification with  $\lambda_{\text{emp}}$  attains a 5.3% and 7.2% average improvement over the subband and MFCC classifiers respectively, across all SNRs considered. Even the combined classification via a simple averaging of the subband and MFCC classifiers by setting  $\lambda = 1/2$  provides a reasonable compromise between classification performance achieved within the two representation domains. While the performance of the combined classifier with  $\lambda = 1/2$  degrades only slightly (approximately 2%) for SNRs above a cross over point between 18-dB and 12-dB, it achieves relatively far greater improvements in high noise. e.g. under the anechoic training regime, the combined classifier with  $\lambda = 1/2$  attains a 13% and 4.2% improvement over the MFCC and subband classifiers at 0-dB SNR, respectively. Quantitatively similar conclusions apply in the reverberant mismatched training scenario as shown in Figure 8.

#### IV. CONCLUSIONS

In this paper we studied an SVM front-end for robust speech recognition that operates in frequency subbands of high-dimensional acoustic waveforms. We addressed the issues of kernel design for subband components of acoustic waveforms and the aggregation of the individual subband classifiers using ensemble methods. The experiments demonstrated that the subband classifiers outperform the cepstral classifiers in the presence of noise and linear filtering for SNRs below 12-dB. While the subband classifiers do not perform as well as the MFCC classifiers in low noise conditions, major gains across all noise levels can be attained by a convex combination of both classifier types [17].

This work primarily focused on comparison of different representations in terms of the robustness they provide. To this end, experiments were conducted on the TIMIT phoneme classification task. However, the results reported in this paper also have implications for the construction of ASR systems. In future work, we plan to investigate extensions to the proposed technique to recognition of continuous speech. One straightforward approach would be to pre-process the speech signals using the combination of subband and cepstral SVM classifiers. Using error-correcting output codes, feature vectors of class probabilities/scores can be derived for progressive

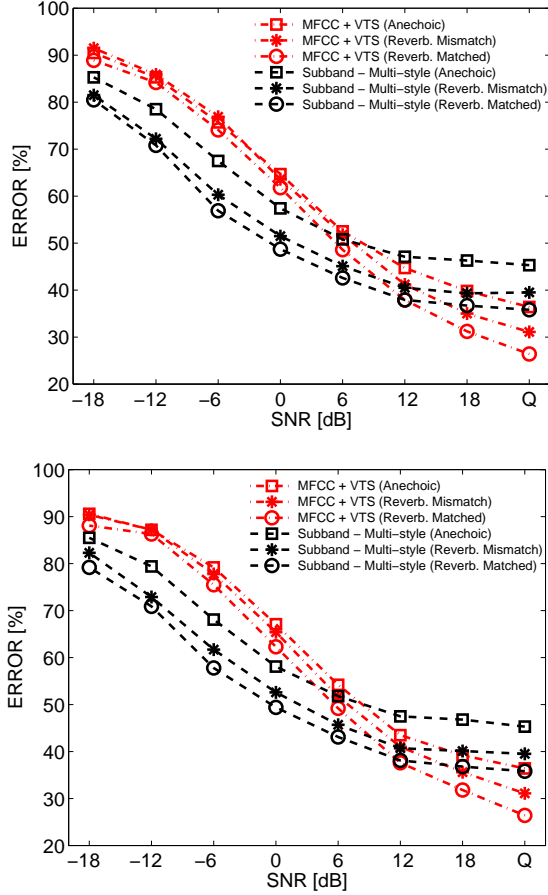


Fig. 7: Classification with the subband and VTS-compensated MFCC classifiers trained under three different scenarios: anechoic training (squares), reverberant mismatched training (stars) and reverberant matched training (circles). Classification results for test data contaminated with pink noise (top) and speech-babble noise (bottom), and linear filter  $R(e^{j\omega})$  are shown.

frames of speech. An HMM can then be trained with these feature vectors for recognition of continuous speech. Alternatively, the proposed technique can also be integrated with other approaches such as the hybrid phone-based HMM-SVM architecture [41, 42] and the token-passing algorithm [43] for continuous speech recognition. In the former, a baseline HMM system would be required to perform a first pass through the test data, generating for each utterance a set of possible segmentations into phonemes. The best segmentations can then be re-scored by the combined SVM classifier to predict the final phoneme sequence. This approach has provided improvements in recognition performance over HMM baselines on both small and large vocabulary recognition tasks, even though the SVM classifiers were constructed solely from the cepstral representations [41, 42]. However, this HMM-SVM hybrid solution can also limit the efficiency of SVMs due to possible errors in the segmentation stage. In [43], a recognizer based solely on SVMs was proposed. In particular, it employs SVMs to classify each frame of speech as belonging to a particular part of a phoneme and then determines the chain of recognized phonemes and words using the token-passing algorithm. These

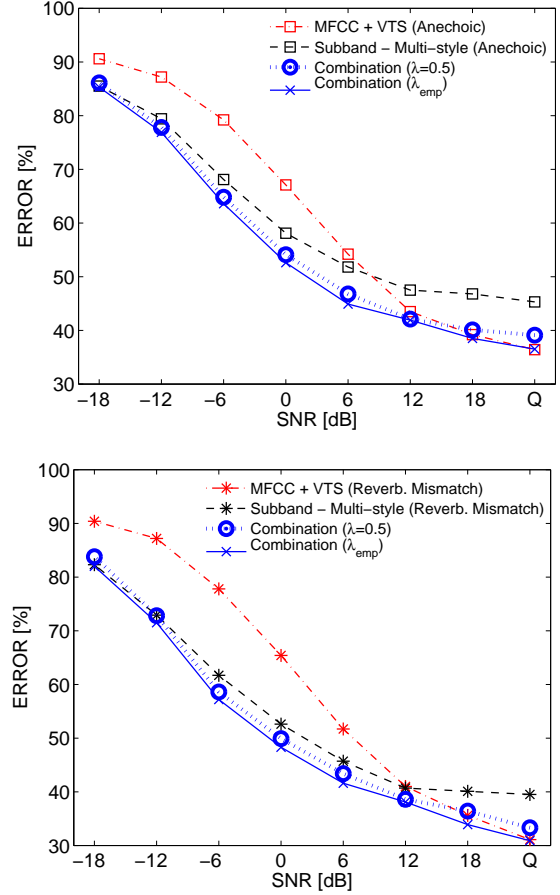


Fig. 8: Comparison of the classification performance of the subband and MFCC classifiers with their convex combination in the presence of speech-babble noise and filtering with  $R(e^{j\omega})$ , under anechoic training (top) and reverberant mismatched training regimes (bottom). Results are shown for two different settings of the combination parameter,  $\lambda$ .

extensions will be the subject of a future study.

#### ACKNOWLEDGMENT

The authors would like to thank Jont Allen, Bishnu Atal, and Andreas Buja for encouragement and inspiration.

#### REFERENCES

- [1] S. B. Davis and P. Mermelstein, "Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences," *IEEE Trans. ASSP*, vol. 28, pp. 357–366, 1980.
- [2] H. Hermansky, "Perceptual Linear Predictive (PLP) Analysis of Speech," *J. Acoust. Soc. Amer.*, vol. 87, no. 4, pp. 1738–1752, April 1990.
- [3] R. Lippmann, "Speech Recognition by Machines and Humans," *Speech Comm.*, vol. 22, no. 1, pp. 1–15, 1997.
- [4] G. Miller and P. Nicely, "An Analysis of Perceptual Confusions among some English Consonants," *J. Acoust. Soc. Amer.*, vol. 27, no. 2, pp. 338–352, 1955.
- [5] J. B. Allen, "How do humans process and recognize speech?," *IEEE Trans. Speech & Audio Proc.*, vol. 2, no. 4, pp. 567–577, 1994.
- [6] J. Sroka and L. Braida, "Human and Machine Consonant Recognition," *Speech Comm.*, vol. 45, no. 4, pp. 401–423, 2005.

- [7] B. Meyer, M. Wächter, T. Brand, and B. Kollmeier, "Phoneme Confusions in Human and Automatic Speech Recognition," *Proc. INTERSPEECH*, pp. 2740–2743, 2007.
- [8] B.S. Atal, "Automatic Speech Recognition: a Communication Perspective," *Proc. ICASSP*, pp. 457–460, 1999.
- [9] S. D. Peters, P. Stubble, and J. Valin, "On the Limits of Speech Recognition in Noise," *Proc. ICASSP*, pp. 365–368, 1999.
- [10] H. Bourlard, H. Hermansky, and N. Morgan, "Towards Increasing Speech Recognition Error Rates," *Speech Comm.*, vol. 18, no. 3, pp. 205–231, 1996.
- [11] K. K. Paliwal and L. D. Alsteris, "On the Usefulness of STFT Phase Spectrum in Human Listening Tests," *Speech Comm.*, vol. 45, no. 2, pp. 153–170, 2005.
- [12] L. D. Alsteris and K. K. Paliwal, "Further Intelligibility Results from Human Listening Tests using the Short-Time Phase Spectrum," *Speech Comm.*, vol. 48, no. 6, pp. 727–736, 2006.
- [13] O. Viikki and K. Laurila, "Cepstral Domain Segmental Feature Vector Normalization for Noise Robust Speech Recognition," *Speech Comm.*, vol. 25, pp. 133–147, 1998.
- [14] C. Chen and J. Bilmes, "MVA Processing of Speech Features," *IEEE Trans. ASLP*, vol. 15, no. 1, pp. 257–270, 2007.
- [15] P. J. Moreno, B. Raj, and R. M. Stern, "A Vector Taylor Series Approach for Environment-Independent Speech Recognition," *Proc. ICASSP*, pp. 733–736, 1996.
- [16] ETSI standard doc., "Speech processing, Transmission and Quality aspects (STQ): Advanced front-end feature extraction," *ETSI ES 202 050*, 2002.
- [17] J. Yousafzai, Z. Cvetković, P. Sollich, and B. Yu, "Combined Features and Kernel Design for Noise Robust Phoneme Classification Using Support Vector Machines," *To appear in the IEEE Trans. ASLP*, 2011.
- [18] H. Fletcher, *Speech and Hearing in Communication*, Van Nostrand, New York, 1953.
- [19] J. McAuley, J. Ming, D. Stewart, and P. Hanna, "Subband correlation and robust speech recognition," *IEEE Trans. on Speech and Audio Proc.*, vol. 13, no. 5, pp. 956 – 964, 2005.
- [20] J. Ming, P. Jancovic, and F.J. Smith, "Robust speech recognition using probabilistic union models," *IEEE Trans. on Speech and Audio Proc.*, vol. 10, no. 6, pp. 403 – 414, Sept. 2002.
- [21] P. McCourt, N. Harte, and S. Vaseghi, "Discriminative Multi-resolution Sub-band and Segmental Phonetic Model Combination," *IET Electronics Letters*, vol. 36, no. 3, pp. 270 –271, 2000.
- [22] S. Thomas, S. Ganapathy, and H. Hermansky, "Recognition Of Reverberant Speech Using Frequency Domain Linear Prediction," *IEEE Signal Process. Letters*, vol. 15, pp. 681–684, 2008.
- [23] S. Tibrewala and H. Hermansky, "Subband Based Recognition Of Noisy Speech," *Proc. ICASSP*, pp. 1255–1258, 1997.
- [24] P. McCourt, S. Vaseghi, and N. Harte, "Multi-Resolution Cepstral Features for Phoneme Recognition across Speech Sub-Bands," *Proc. ICASSP*, pp. 557–560, 1998.
- [25] H. Bourlard and S. Dupont, "Subband-based Speech Recognition," *Proc. ICASSP*, pp. 1251–1254, 1997.
- [26] S. Okawa, E. Bocchieri, and A. Potamianos, "Multi-band Speech Recognition in Noisy Environments," *Proc. ICASSP*, pp. 641–644, 1998.
- [27] C. Cerisara, J.-P. Haton, J.-F. Mari, and D. Fohr, "A Re-combination Model for Multi-band Speech Recognition," *ICASSP*, pp. 717 –720 vol.2, 1998.
- [28] J. Flanagan, J. Johnston, R. Zahn, and G. Elko, "Computer-Steered Microphone Arrays for Sound Transduction in Large Rooms," *J. Acoust. Soc. Amer.*, vol. 78, no. 11, pp. 1508–1518, 1985.
- [29] M. Wu and D. Wang, "A Two-Stage Algorithm for One-Microphone Reverberant Speech Enhancement," *IEEE Trans. ASLP*, vol. 14, pp. 774–784, 2006.
- [30] H. Chang and J. Glass, "Hierarchical Large-Margin Gaussian Mixture Models for Phonetic Classification," *Proc. ASRU*, pp. 272–275, 2007.
- [31] D. Yu, L. Deng, and A. Acero, "Hidden Conditional Random Fields with Distribution Constraints for Phone Classification," *Proc. INTERSPEECH*, pp. 676–679, 2009.
- [32] F. Sha and L. K. Saul, "Large Margin Gaussian Mixture Modeling for Phonetic Classification and Recognition," *Proc. ICASSP*, pp. 265–268, 2006.
- [33] R. Rifkin, K. Schutte, M. Saad, J. Bouvrie, and J. Glass, "Noise Robust Phonetic Classification with Linear Regularized Least Squares and Second-Order Features," *Proc. ICASSP*, pp. 881–884, 2007.
- [34] A. Halberstadt and J. Glass, "Heterogeneous Acoustic Measurements for Phonetic Classification," *Proc. EuroSpeech*, pp. 401–404, 1997.
- [35] P. Clarkson and P. J. Moreno, "On the Use of Support Vector Machines for Phonetic Classification," *Proc. ICASSP*, pp. 585–588, 1999.
- [36] A. Gunawardana, M. Mahajan, A. Acero, and J. C. Platt, "Hidden Conditional Random Fields for Phone Classification," *Proc. INTERSPEECH*, pp. 1117–1120, 2005.
- [37] V. Pitsikalis and P. Maragos, "Analysis and Classification of Speech Signals by Generalized Fractal Dimension Features," *Speech Comm.*, vol. 51, pp. 1206–1223, 2009.
- [38] S. Dusan, "On the Relevance of Some Spectral and Temporal Patterns for Vowel Classification," *Speech Comm.*, vol. 49, pp. 71–82, 2007.
- [39] K. M. Indrebo, R. J. Povinelli, and M. T. Johnson, "Sub-banded Reconstructed Phase Spaces for Speech Recognition," *Speech Comm.*, vol. 48, no. 7, pp. 760–774, 2006.
- [40] A. Halberstadt and J. Glass, "Heterogeneous Measurements and Multiple Classifiers for Speech Recognition," *Proc. ICSLP*, pp. 995–998, 1998.
- [41] A. Ganapathiraju, J. E. Hamaker, and J. Picone, "Applications of Support Vector Machines to Speech Recognition," *IEEE Trans. Signal Proc.*, vol. 52, no. 8, pp. 2348–2355, 2004.
- [42] S. E. Krüger, M. Schaffner, M. Katz, E. Andelic, and A. Wendemuth, "Speech Recognition with Support Vector Machines in a Hybrid System," *Proc. INTERSPEECH*, pp. 993–996, 2005.
- [43] J. Padrell-Sendra, D. Martín-Iglesias, and F. Díaz-de Marí, "Support Vector Machines for Continuous Speech Recognition," *Proc. EUSIPCO*, 2006.
- [44] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, New York, 1995.
- [45] J. Louradour, K. Daoudi, and F. Bach, "Feature Space Mahalanobis Sequence Kernels: Application to SVM Speaker Verification," *IEEE Trans. ASLP*, vol. 15, no. 8, pp. 2465–2475, 2007.
- [46] A. Sloin and D. Burshtein, "Support Vector Machine Training for Improved Hidden Markov Modeling," *IEEE Trans. Signal Proc.*, vol. 56, no. 1, pp. 172–188, 2008.
- [47] T. Jaakkola and D. Haussler, "Exploiting Generative Models in Discriminative Classifiers," in *Adv. Neural Inf. Process. Syst.*, 1999, vol. 11, pp. 487–493.
- [48] R. Solera-Urena, D. Martín-Iglesias, A. Gallardo-Antolín, C. Peláez-Moreno, and F. Díaz-de Marí, "Robust ASR using Support Vector Machines," *Speech Comm.*, vol. 49, no. 4, pp. 253–267, 2007.
- [49] T. Dietterich and G. Bakiri, "Solving Multiclass Learning Problems via Error-Correcting Output Codes," *J. Artif. Intell. Res.*, vol. 2, pp. 263–286, 1995.
- [50] R. Rifkin and A. Klautau, "In Defense of One-Vs-All Classification," *J. Mach. Learn. Res.*, vol. 5, pp. 101–141, 2004.
- [51] Martin Vetterli and Jelena Kovacevic, *Wavelets and Subband Coding*, Prentice-Hall, Englewood Cliffs, NJ, 1995.
- [52] S. Furui, "Speaker-Independent Isolated Word Recognition using Dynamic Features of Speech Spectrum," *IEEE Trans. ASSP*, vol. 34, no. 1, pp. 52–59, 1986.
- [53] T. Dietterich, "Ensemble Methods in Machine Learning,"



*Lecture Notes in Computer Science: Multiple Classifier Systems*, pp. 1–15, 2000.

- [54] L. Hansen and P. Salamon, “Neural Network Ensembles,” *IEEE Trans. PAMI*, vol. 12, no. 10, pp. 993–1001, 1990.
- [55] D. Wolpert, “Stacked Generalization,” *Neural Networks*, vol. 5, no. 2, pp. 241–259, 1992.
- [56] J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallet, and N. Dahlgren, “TIMIT Acoustic-Phonetic Continuous Speech Corpus,” *Linguistic Data Consortium*, 1993.
- [57] K. F. Lee and H. W. Hon, “Speaker-Independent Phone Recognition Using Hidden Markov Models,” *IEEE Trans. ASSP*, vol. 37, no. 11, pp. 1641–1648, 1989.
- [58] R. F. Voss and J. Clarke, “1/f Noise in Music: Music from 1/f Noise,” *J. Acoust. Soc. Amer.*, vol. 63, no. 1, pp. 258–263, 1978.
- [59] B. J. West and M. Shlesinger, “The Noise in Natural Phenomena,” *American Scientist*, vol. 78, no. 1, pp. 40–45, 1990.
- [60] P. Grigolini, G. Aquino, M. Bologna, M. Lukovic, and B. J. West, “A Theory of 1/f Noise in Human Cognition,” *Physica A: Stat. Mech. and its Appl.*, vol. 388, no. 19, pp. 4192–4204, 2009.
- [61] “The ICSI Meeting Recorder Project - Room Responses,” Online Web Resource.
- [62] M. Holmberg, D. Gelbart, and W. Hemmert, “Automatic Speech Recognition with an Adaptation Model Motivated by Auditory Processing,” *IEEE Trans. ASLP*, vol. 14, no. 1, pp. 43–49, 2006.
- [63] R. Lippmann and E. A. Martin, “Multi-Style Training for Robust Isolated-Word Speech Recognition,” *Proc. ICASSP*, pp. 705–708, 1987.
- [64] J. Yousafzai, Z. Cvetković, and P. Sollich, “Towards Robust Phoneme Classification with Hybrid Features,” *Proc. ISIT*, pp. 1643–1647, 2010.
- [65] Y. Ephraim and D. Malah, “Speech Enhancement Using a Minimum Mean-Square Error Short-time Spectral Amplitude Estimator,” *IEEE Trans. ASSP*, vol. ASSP-32, pp. 1109–1121, 1984.
- [66] Y. Ephraim and D. Malah, “Speech Enhancement Using a Minimum Mean-Square Log-Spectral Amplitude Estimator,” *IEEE Trans. ASSP*, vol. ASSP-33, pp. 443–445, 1985.