# TEXT-DEPENDENT SPEAKER VERIFICATION UNDER NOISY CONDITIONS USING PARALLEL MODEL COMBINATION

*Lit Ping Wong and Martin Russell*

School of Electronic and Electrical Engineering
The University of Birmingham, Edgbaston, Birmingham B15 2TT, United Kingdom
L.P.Wong, M.J.Russell@bham.ac.uk

## ABSTRACT

In real speaker verification applications, additive or convolutive noise creates a mismatch between training and recognition environments, degrading performance. Parallel Model Combination (PMC) is used successfully to improve the noise robustness of Hidden Markov Model (HMM) based speech recognisers [5]. This paper presents the results of applying PMC to compensate for additive noise in HMM-based text-dependent speaker verification. Speech and noise data were obtained from the YOHO [6] and NOISEX-92 databases [13] respectively. Speaker recognition Equal Error Rates (EER) are presented for noise-contaminated speech at different signal-to-noise ratios (SNRs) and different noise sources. For example, average EER for speech in operations room noise at 6dB SNR dropped from approximately 20% un-compensated to less than 5% using PMC. Finally, it is shown that speaker recognition performance is relatively insensitive to the exact value of the parameter that determines the relative amplitudes of the speech and noise components of the PMC model.

## 1. INTRODUCTION

It is well known that noise contamination of speech signals results in increased speech and speaker recognition errors, due to the consequent mismatch between training and test conditions, and loss of information. Hence considerable effort has been applied to the development of robust noise compensation techniques. These techniques generally fall into two categories, speech pre-processing and adaptation of the recognition stage. The first class of methods attempt to pre-process the corrupted speech such that the resulting parameters are representative of clean speech. Techniques in this category include spectral subtraction [2], and spectral mapping [11]. Examples of techniques that modify the recognition stage include noise masking [8], HMM decomposition [12] and Parallel Model Combination (PMC) [5]. Methods based on pre-processing are often computationally simpler, but have the disadvantage that any relevant information that is discarded in the pre-processing stage cannot be recovered for use in recognition. The work in this paper focuses on the second approach, and in particular on the use of Parallel Model Combination (PMC) in which speech and noise HMMs are compiled into a single composite HMM prior to recognition.

PMC has been applied successfully to HMM-based automatic speech recognition, where it has been shown to improve recognition performance on speech contaminated with additive noise [5]. This paper presents the results of experiments to investigate the utility of PMC for noise-robust HMM based text-dependent speaker verification. Clean speaker verification data from the YOHO database [6], designated as test data, was mixed with two different types of noises from the NOISEX-92 database [13] at a range of different signal-to-noise ratios. This data was then processed using clean speech models combined with the appropriate noise models at the appropriate mixing levels using PMC. The results show that PMC gives significant reductions in equal error rates. For example, un-compensated equal error rates of 20%, 33% and 45% (at +6dB, 0dB and –6dB respectively) are reduced to 5%, 14% and 32% respectively using PMC. Further experiments were conducted to measure the sensitivity of the PMC error rate to changes in the value of the mixing level. The results show that performance is relatively insensitive to mixing level, and that restriction to seven different mixing levels only results in a small increase in error relative to the experiments with correctly matched mixing levels.

## 2. NOISE ROBUST SPEAKER VERIFICATION

The goal of speaker verification is to confirm the claimed identity of a subject by exploiting individual differences in their speech. It is useful to distinguish between text-dependent speaker verification, where the decision is made using speech corresponding to known text, and text-independent speaker verification, where the speech is unconstrained [4]. The present study is concerned with the former. Text-dependent verification is clearly the simpler problem, and is amenable to word-level acoustic modelling techniques from automatic speech recognition, and in particular the use of HMMs. Once the decision to use HMM techniques has been made, it is natural to ask whether HMM noise compensation techniques from speech recognition can also be applied successfully. The noise compensation technique that is the subject of the present study is PMC [5].

## 3. PARALLEL MODEL COMBINATION (PMC)

PMC is based on the premise that noise compensation should occur during the pattern processing stage of speech or speaker recognition and not during parameterisation. In particular, the decision that a component of the data is 'noise' should emerge

from the recognition process, rather than precede it. In this way, all of the information contained in the speech signal is retained and exploited for correct verification. Although similar to HMM decomposition [12], PMC has the advantage that it is able to operate in the cepstral domain and therefore inherits the advantages of parameter decorrelation.

Let $\sigma_S$ and $\sigma_N$ be single Gaussian states of a 'clean' speech HMM and a noise HMM respectively, in the cepstral domain. Suppose that the means and variances of $\sigma_S$ and $\sigma_N$ are denoted by $\{\mu_S{}^c, \Sigma_S{}^c\}$ and $\{\mu_N{}^c, \Sigma_N{}^c\}$. PMC creates a combined cepstral state $\sigma_S \otimes \sigma_N$ by inverse-transforming $\{\mu_S{}^c, \Sigma_S{}^c\}$ and $\{\mu_N{}^c, \Sigma_N{}^c\}$ into the linear, spectral domain (via the log spectral domain), combining them appropriately into a single distribution, and then mapping this distribution back into the cepstral domain. During this process, PMC makes a number of assumptions. It is assumed that speech and noise are independent, that they are additive, and, during the combination process, that the sum of two log-normal distributions is log normal [5].

Both the static and dynamic (velocity and acceleration) parameters, which are normally included in a cepstrum-based representation for speech recognition can be compensated using PMC. However, in this study, only static parameters were compensated, as it was believed that compensating for the delta and acceleration parameters would only bring marginal improvements in performance as most of the speaker dependent characteristics will be removed by the differencing process. The speech mean and covariance matrix in the log spectral domain are given by:

$$\mu_S^l = C^{-1} \mu_S^c$$

$$\Sigma_S^l = C^{-1} \Sigma_S^c (C^{-1})^T$$

where C is the cosine transform defined by the matrix

$$C_{ij} = \cos(i(j - 0.5)\pi / B)$$

Similar expressions exist for the parameters of the noise state. The mean and covariance matrices for the state $\sigma_S \otimes \sigma_N$ are then given by:

$$\mu_{S \otimes N}^l = \log(\exp(\mu_S^l) + g \exp(\mu_N^l))$$

$$\Sigma_{S \otimes N}^l = \log(\exp(\Sigma_S^l) + g \exp(\Sigma_N^l))$$

where $g$ is a gain matching term which determines the signal-to-noise ratio.

## 4.    EXPERIMENTAL DATA

The experiments used speech data from three corpora: TIMIT, YOHO [6] and NOISEX-92 [13]. TIMIT was used to initialise a set of phoneme-level HMMs, because no phoneme transcriptions were provided with YOHO.

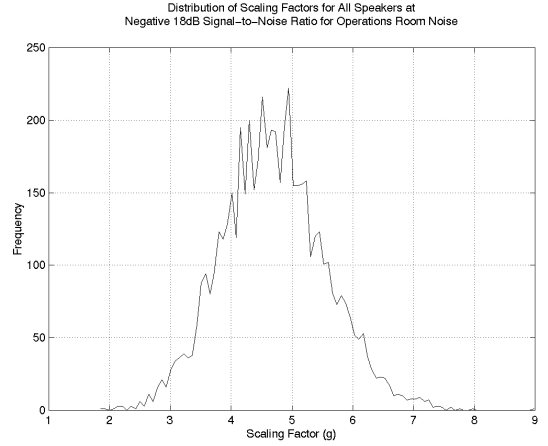Speech in the YOHO database is sampled at 8k samples per second. The corpus consists of recordings of 138



*Figure 1  Distribution of the scaling factor*

subjects speaking connected digit phrases in an office environment. It was chosen because of its established use in evaluation of speaker recognition systems, and enables the results obtained from the present experiment to be compared with those from other laboratories [3].

The NOISEX-92 database contains a range of different types of noise, including speech noise, car noise, and operations room noise. It has been used previously to investigate noise compensation techniques for automatic speech recognition [5].

## 5.    EXPERIMENTAL METHOD

### 5.1    Mixing Speech and Noise

The speech level was measured using a software implementation of the procedure described in [7], which accounted for the silence intervals present in speech and only calculated the level for voiced segments. Before applying this process, the speech data was amplitude scaled to ¼ of the maximum 16-bit integer value to prevent saturation when noise was added. The speech signal was then added to the scaled version of the noise signal to give noise contaminated speech signals at–18dB to +18dB at 6dB SNR intervals.

The actual scaling factor used to obtain a particular SNR varied between speakers and speech files. Hence for each SNR the average scaling factor g was computed across all speakers. This scaling factor was used in the construction of the PMC models for that SNR. An example of the distribution of the scaling factor for a single SNR and different speakers and files is shown in figure 1.

Since experiments were performed on scaled data, the enrolment data was also scaled to prevent any mismatch in the baseline results.

### 5.2    HMM System

HMM training and recognition used the Hidden Markov Model Toolkit (HTK) system, which was easily modified to accommodate PMC. Speech was parameterised into 39 dimensional representation based on Mel-Frequency Cepstral Coefficients (MFCC), using a 25ms Hamming window. This

included the $0^{th}$ order energy term, plus velocity ($\Delta$) and acceleration ($\Delta^2$) of each coefficient.

Initially, a total of 57 3-state, 4-component Gaussian mixture monophone HMMs were created using the TIMIT database. These were later expanded to 78 tied-state triphones to cover all possible pronunciations of connected digit pairs, e.g. 29_30_31. The triphones were subsequently used to create speaker dependent models using speech from the YOHO database. For each speaker, all 24 utterances in each of 4 enrolment sessions were used for training.

## 5.3    Recognition and Scoring

Of the 138 speakers in YOHO, 118 were used as authorised speakers and 20 were used to train a general speaker model (GSM) [10]. Furthermore, all 4 utterances in the 10 verification sessions were used to calculate the False Reject Rate (FRR). To calculate the False Accept Rate (FAR) for each speaker, 40 utterances were randomly chosen from the authorised speaker set except the speaker on test to form an impostor set. The authorised speaker model (AS) was then used to recognise utterances from the impostor set.

For both FRR and FFA, the GSM was used to normalise scores from the speaker dependent models [1]. The decision rule used, for a particular threshold $t$, is as follows:

$$\text{If} \quad \frac{P(X \mid S)}{P(X \mid GSM)} \geq t \quad \text{then 'Accept', else 'Reject'}.$$

The Equal Error Rate (EER) corresponds to the value of $t$ for with FRR = FAR.

## 6.    EVALUATION

Four sets of experiments were carried out. First, the performance of the baseline system was tested using clean speech. The second experiment investigated the degradation in performance when noise contaminated speech was used in verification without PMC. Next, the same experiment was performed using PMC. The second and third experiments considered both operations room noise (as an example of 'typical' noise) and speech noise (as 'worst case'). Finally, the dependency of performance on the gain matching term, $g$, was investigated.

### 6.1    Baseline System Performance

The EER achieved using clean verification speech was 0.57%. The target, therefore, was to get the EER as close to this value for all SNRs.

### 6.2    Un-Compensated System Performance

The results without noise compensation for SNRs varying between −18dB and +18dB are shown in figure 2. The EER deteriorated at an average rate of 10% for every 6dB reduction in SNR. This trend continued until performance levelled off at 50% EER ('random' error for a two class problem). As expected, the performance degradations were much worse for 'speech' noise than for 'operations room' noise.
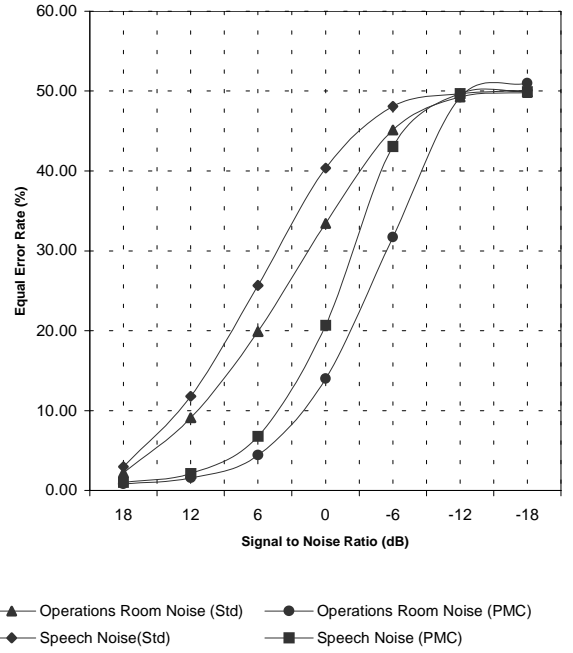


*Figure 2  Graph of compensated results against un-compensated results*

### 6.3    Compensated System Performance

PMC compensation was performed using the average scaling factor obtained during speech and noise mixing to set the gain matching term g (section 5.1). With PMC, performance at 18dB was close to the baseline result. Although performance degradation was still observed as SNR was reduced, this was not as severe as the un-compensated models. An average of 50% performance increase was observed between standard and PMC recognition. However, this was only true until -12dB where the performance curve of PMC rejoins that of standard recognition.

### 6.4    Dependency on g

To investigate the dependency of the EER on the gain matching term $g$, experiments were conducted using operations room noise in which the value of $g$ was held constant while the SNR was varied between +18dB and −18dB. Values of $g$ which were considered correspond to +18, +12, +6, 0, -6, -12 and −18dB (figure 3). As one would expect, the results show that sensitivity to the absolute value of g increases as SNR increases. However, the results show that performance of the system is relatively insensitive to the exact value of $g$, and that for an SNR of $R$dB, good performance is obtained with values of $g$ corresponding to SNRs of $R\pm6$dB. This suggests that good speaker verification performance can be obtained over SNRs ranging from +18dB to −18dB using PMC models with $g$ corresponding to no more that 7 different SNRs. For example, at −6dB the best performance is actually achieved using the value of $g$ which matches 0dB SNR, and the error rates for values of $g$ corresponding to −12dB, 0dB and +6dB are all within 10% of the error rate for the 'correct' value of $g$ for −6dB.

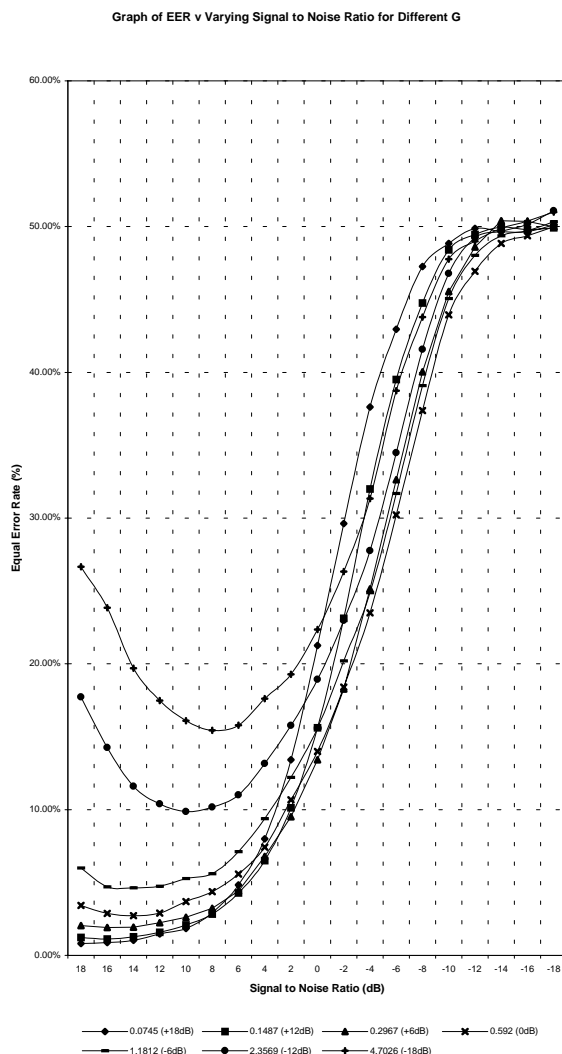Graph of EER v Varying Signal to Noise Ratio for Different G



*Figure3   Dependency of the EER on g*

## 7.   CONCLUSION

This paper reports the results of experiments, which investigate the utility of Parallel Model Combination (PMC) for noise-robust HMM-based text-dependent speaker verification. Speech data was taken from the YOHO corpus, and noise data from NOISEX-92. The system used context-sensitive phoneme-level HMMs with 4-component Gaussian mixture states. The results show that the baseline EER of 0.57% for 'clean' speech degrades rapidly as a consequence of contamination with 'speech' or 'operations room' noise. For example the clean speech EER of 0.57% drops to 33% and 41% respectively at 0dB SNR. Using PMC with the correctly matched gain factor *g*, the corresponding figures are 14% and 21%, and both EERs are below 10% at +6dB SNR. Finally, it has been shown that performance is relatively insensitive to the exact value of the gain factor *g*, provided that it corresponds to a SNR that is within ±6dB of the true SNR.

## 8.   REFERENCES

[1]    R. Auckenthaler, M. Carey and H. Lloyd-Thomas, "Score Normalization for Text-Independent Speaker Verification Systems," Digital Signal Processing, Vol. 10, pp. 42-54, 2000.

[2]    S.F. Boll, "Suppression of Acoustic Noise in Speech Using Spectral Subtraction," IEEE Transactions on Acoustics Speech and Signal Processing, Vol. ASSP-27, 1979.

[3]    J. P. Campbell Jr., "Speaker Recognition: A Tutorial," Proceedings of the IEEE, Vol. 85, No. 9, September 1997.

[4]    S. Furui, "An Overview of Speaker Recognition Technology," ESCA workshop on Automatic Speaker Recognition, Identification and Verification, pp. 1-9, 1994.

[5]    M.J. Gales and S.J. Young, "Robust Continuous Speech Recognition Using Parallel Model Combination," IEE Transactions on Speech and Audio Processing, Vol. 4, No. 5, pp. 352-359, September 1996.

[6]    A. Higgins, "YOHO speaker verification," Presented at the Speech Research Symposium, Baltimore, MD, 1990.

[7]    ITU-T Recommendation, "Objective Measurement of Active Speech Level," pp. 56, September 1993.

[8]    D.H. Klatt, "A digital Filter Bank for Spectral Matching," Proceedings of the ICASSP, pp. 573-576 , 1979.

[9]    T. Matsui, T. Kanno, S. Furui, "Speaker Recognition Using HMM Composition in Noisy Environments," Computer Speech and Language, Vol. 10, pp. 107-116, 1996.

[10]   D.A. Reynolds, T.F. Quatieri and R.B. Dunn, "Speaker Verification Using Adapted Gaussian Mixture Models," Digital Signal Processing, Vol. 10: (1-3), pp. 19-41, 2000.

[11]   H.B.D Sorensen, "A Cepstral Noise Reduction Molti-Layer Neural Network," Proceedings of the ICASSP, S14.14, pp. 933-936, 1991.

[12]   A.P. Varga and R.K. Moore, "Hidden Markov Model Decomposition of Speech and Noise," Proceedings of the ICASSP 90, 1990.

[13]   A.P. Varga, H.J.M. Steeneken, M. Tomlinson and D. Jones, "The NOISEX-92 Study on the Effect of Additive Noise on Automatic Speech Recognition," Technical Report, DRA Speech Research Unit, Malvern, England, 1992.