



ELSEVIER

Speech Communication 34 (2001) 25–40

SPEECH
COMMUNICATION

www.elsevier.nl/locate/specom

Multi-stream adaptive evidence combination for noise robust ASR

Andrew Morris^{*,1}, Astrid Hagen¹, Hervé Glotin¹, Hervé Bourlard^{1,2}

Institut Dalle Molle d'Intelligence Artificielle Perceptive (IDIAP), Rue du Simplon 4, BP 592, CH-1920 Martigny, Switzerland

Abstract

In this paper, we develop different mathematical models in the framework of the multi-stream paradigm for noise robust automatic speech recognition (ASR), and discuss their close relationship with human speech perception. Largely inspired by Fletcher's "product-of-errors" rule (PoE rule) in psychoacoustics, multi-band ASR aims for robustness to data mismatch through the exploitation of spectral redundancy, while making minimum assumptions about noise type. Previous ASR tests have shown that independent sub-band processing can lead to decreased recognition performance with clean speech. We have overcome this problem by considering every combination of data sub-bands as an independent data stream. After introducing the background to multi-band ASR, we show how this "full combination" approach can be formalised, in the context of hidden Markov model/artificial neural network (HMM/ANN) based ASR, by introducing a latent variable to specify which data sub-bands in each data frame are free from data mismatch. This enables us to decompose the posterior probability for each phoneme into a reliability-weighted integral over all possible positions of clean data. This approach offers great potential for adaptation to rapidly changing and unpredictable noise. © 2001 Elsevier Science B.V. All rights reserved.

Zusammenfassung

In diesem Artikel werden im Rahmen des "multi-stream" Paradigmas für rausch-unempfindliche automatische Sprachverarbeitung (ASV) verschiedene mathematische Modelle entwickelt und ihr enger Zusammenhang mit der menschlichen Sprachwahrnehmung diskutiert. Zum grössten Teil inspiriert von Fletchers "Regel der Fehlermultiplikation" in der Psychoakustik, zielt die "multi-band" ASV auf eine erhöhte Rauschrobustheit durch die Nutzung spektraler Redundanz, wobei nur minimale Annahmen über das Rauschen aufgestellt werden müssen. Vorangegangene Experimente in ASV haben gezeigt, dass die unabhängige Verarbeitung einzelner (spektraler) Unterbänder ("sub-bands") zu Abnahme der Erkennungsleistung auf rausch-freiem Signal führen kann. Dieses Problem wird hier überwunden, indem jede Kombination von Unterbändern als ein unabhängiger Datenstrom miteinbezogen wird. Nach kurzer Einführung in die Hintergründe der "multi-stream" Sprachverarbeitung, wird gezeigt wie der sogenannte "full combination" Ansatz im Rahmen von HMM/ANN-basierter ASV formalisiert werden kann. Dies geschieht durch die Einführung einer unabhängigen Variablen, die beschreibt, welche Datenbänder im gegebenen Zeitfenster frei von Daten Uneinstimmigkeit sind. Dies ermöglicht es, die Posteriori-Wahrscheinlichkeit jedes Phonems in ein (zuverlässigkeits-)

^{*}Corresponding author.

E-mail addresses: morris@idiap.ch (A. Morris), hagen@idiap.ch (A. Hagen), glotin@idiap.ch (H. Glotin), bourlard@idiap.ch (H. Bourlard).

¹ <http://www.idiap.ch/>.

² Also at the Swiss Federal Institute of Technology, Lausanne, Switzerland.

gewichtetes Integral über alle möglichen Positionen rausch-freier Daten zu zerlegen. Dieser Ansatz bietet grosses Potenzial für die Anpassung an schnell veränderliches und unvorhersehbares Rauschvorkommen. © 2001 Elsevier Science B.V. All rights reserved.

Résumé

Dans cet article nous développons plusieurs modèles autour du paradigme “multi-stream” de la RAP (Reconnaissance Automatique de la Parole) robuste, et nous discutons de leurs relations avec la perception naturelle de parole. Fortement inspirée par la règle “produit des erreurs” de Fletcher, issue de la psychoacoustique, la reconnaissance “multi-bande” se veut être robuste à l’inadéquation des données par rapport aux conditions d’apprentissage, en exploitant la redondance spectrale, tout en faisant un minimum d’hypothèses sur la nature du bruit. Des études précédentes en RAP ont montré que le traitement indépendant des sous bandes peut diminuer le taux de reconnaissance en parole propre. Nous avons surmonté ce problème en considérant toutes les combinaisons de sous bandes comme des flux indépendants de données. Après un état de l’art sur la RAP multi-bandes, nous formalisons cette approche “full-combination” dans le contexte de la RAP HMM/ANN, en introduisant une variable latente qui spécifie à chaque trame de signal la combinaison de sous bandes la plus adéquate. Ceci nous permet de décomposer la probabilité a posteriori pour chaque phonème en une somme pondérée, sur toutes les positions possibles des données propres. Cette approche est prometteuse pour l’adaptation aux bruits imprévisibles et variant rapidement. © 2001 Elsevier Science B.V. All rights reserved.

Keywords: Noise robust ASR; Multi-band ASR; Expert combination; Noise adaptation; Latent variable decomposition

1. Introduction ³

Even low levels of mismatch between field data and data used in system training, due to noise or channel distortion, can still cause a sharp loss in automatic speech recognition (ASR) performance. As a result, present ASR technology is still inadequate for the majority of applications in which human-like performance is required in an ambient noisy environment ⁴. Conventional processing can deal with some kinds of data mismatch quite successfully. Slowly varying noise or channel mismatch effects can be removed by conventional robust preprocessing techniques, such as spectral subtraction for additive noise, or cepstral mean subtraction for convolutive noise. Stationary noise can also be reduced by noise modelling techniques, such as PMC (Varga and Moore,

1990; Gales and Young, 1993), and speaker variation can be effectively accommodated through a combination of training with large vocabulary multi-speaker databases, and speaker adaptation. However, the mismatch which remains after these techniques have been applied often causes an unacceptable loss in recognition performance.

One effect of the frequency analysis which is applied in the first stages of both natural and artificial speech processing, is that the level of redundancy in the resulting two-dimensional spectro-temporal representation is greatly increased. Another effect is that the data from different sources in the resulting “auditory scene” is largely separated. Conventional robust preprocessing can further increase this separation. Recognition experiments with band limited noise in both psychoacoustics (Fletcher, 1922; Allen, 1994) and ASR (Lippmann and Carlson, 1997; Morris et al., 1998) have shown that narrow sub-bands of clean data can often be sufficient for speech recognition. There is therefore great potential for improved noise robustness with any system, which is able to detect local data mismatch and focus recognition on the clean speech data, which remains.

³ This paper is a combined and extended version of two papers (Bourlard, 1999; Hagen et al., 1999) which were presented in the Tampere workshop “Robust Methods for Speech Recognition in Adverse Conditions”.

⁴ This is contrary to the surprisingly prevalent idea that “robust speech recognition is now a solved problem”, which the speech technology industry would of course like its clients to believe.

In Section 2, we introduce the advantages of multi-stream processing in general. In Section 3, we discuss evidence from auditory physiology and psychoacoustics for multi-channel processing in ASR. In Section 4, we briefly discuss the main historical approaches to multi-band ASR, and their associated limitations. One of these limitations is that independent sub-band processing reduces recognition performance in clean speech. In Section 5, we introduce the “full combination” approach to multi-band ASR, in which the assumption of sub-band independence is overcome by introducing a latent variable, which permits us to integrate over all possible positions of missing (i.e. strongly mismatching) data. In this way the full-band discriminant is decomposed into a mismatch weighted sum of clean-data discriminants over all sub-band combinations. The effectiveness of this approach is however dependent on the weights given to each sub-band combination discriminant and in Section 6 we present a number of different techniques for estimating these weights. In Section 7, we test some variations of this model on a free-format numbers recognition task, under different noise conditions. These results, and the problems, which they raise, are discussed in Section 8.

2. Multi-stream processing

Multi-expert systems arise in many different fields of data classification and function approximation in general. These systems have a number of proven theoretical and practical advantages, of which the following are of particular relevance to ASR:

- *Hierarchical systems of experts reduce problem perplexity.* Unsupervised-training can be used to train a hierarchical system of experts together with a gating network for expert selection. The gating network may be trained to use large scale features for expert selection, so that each expert is trained on a sub-region of the input data space, with correspondingly reduced perplexity (Jordan and Jacobs, 1994).
- *Linear combination of multiple experts can improve generalisation.* Linear combinations of

estimators have been studied and used by the statistics community for a long time. When expert outputs are linearly combined (even as a simple average), the expected committee error will decrease, both in theory (Bishop, 1995), and in practice (Raviv and Intrator, 1996). This error will also decrease further if the spread of the experts’ predictions can be increased without increasing the expected errors of the individual members.⁵ Different experts can be obtained by varying the parametric function used to model each expert, and/or by varying the data which is used to train each expert, either by using different features from the same data, or different noise added to the same data.

3. Multi-stream processing in human speech recognition

In any recognition process, it is advantageous to constructively combine as many sources of information⁶ as are available (Morgan et al., 1998). An example from human perception is that the auditory system is “wired” to combine visual with acoustic information even at the sub-conscious level, so that perceived phoneme category is directly influenced by lip movements (McGurk and McDonald, 1976; Moore, 1997). Experiments in ASR have demonstrated that combining mouth shape with acoustic data can strongly improve recognition performance with noisy speech (Tomlinson et al., 1996; Dupont and Luetin, 1998; Girin et al., 1998; Westphal and Waibel, 1999).

Further evidence for the use of multiple experts in the mammalian auditory system is seen at the first stage of central auditory processing. In the cochlear nucleus, each fibre in the auditory nerve splits and carries the same data through about seven different types of specialised nerve cell, each having a very different characteristic response. The outputs from these cells are recombined at higher levels of processing (Pickles, 1988).

⁵ The optimal linear expert combination can be obtained as a function of the error covariance matrix.

⁶ Sometimes also known as “multiple cues” (Moore, 1997).

While investigating the effects of band limited noise on human hearing, Fletcher (1922) found that the error rate for human phoneme perception using the full frequency range was approximately equal to the product of the error rate using high-pass filtered speech, with the error rate using low-pass filtered speech at the same cut-off frequency. Furthermore, this error rate was independent of the cut-off frequency used. A generalisation of this rule to more than two sub-bands was more recently popularised by Allen (1994), and is now commonly known as the Fletcher–Allen principle, or “product-of-errors” rule (PoE rule):

In human perception, the error rate for full-band perception is equal to the product of the sub-band error rates obtained through perception of each sub-band on its own.

$$P(\text{error}) = \prod_i P(\text{error}_i). \quad (1)$$

Under the assumption of sub-band error independence, it follows from this rule that

Full-band classification is incorrect if and only if classification is incorrect in every sub-band

or equivalently:

Full-band classification is correct if and only if classification is correct in any sub-band

Fig. 1 shows how the probability of correct classification, under the PoE rule, is distributed as a function of the probability of correct classification in each of two sub-bands. The PoE rule serves as proof of existence for a system, which combines multiple guesses at the speech information with an infallible mechanism for selecting the correct guess when it is present. Although the accuracy of the PoE rule has more recently been questioned (Steeneken and Houtgast, 1999), this powerful recognition paradigm has strongly motivated the development of multi-band ASR. Of course, as the number of sub-bands is increased each sub-band gets narrower and the error rate for recognition within each isolated sub-band will increase. However, in Table 1 it is shown that, if a correct answer can always be spotted when present, then the increased number of guesses can easily offset the penalty of increased sub-band error rate, and the optimal number of sub-bands may be considerably greater than the four which we use throughout this

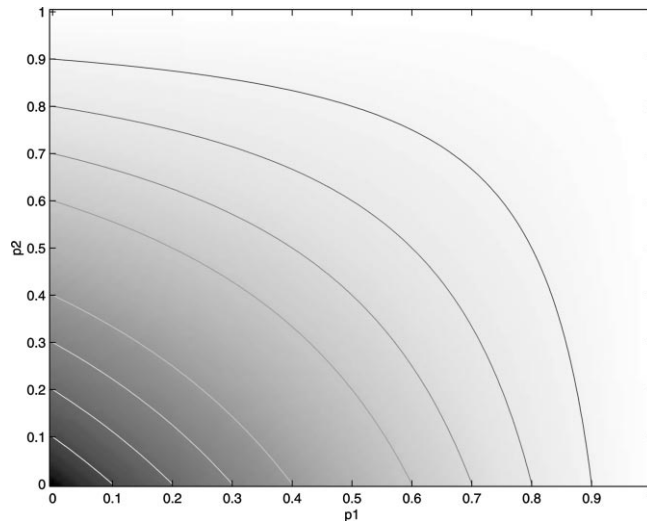


Fig. 1. Probability of correct classification when classification error follows the PoE rule, for two (independent) observation channels (white corresponding to maximum probability). Classifiers yield probabilities $p_1 = \hat{P}(q_k|x_1)$ and $p_2 = \hat{P}(q_k|x_2)$. $P(\text{correct}) = 1 - (1 - p_1)(1 - p_2) = p_1 + p_2 - p_1p_2$. Contours show lines of equal correct recognition probability (contour height is value on contour-axis intercept). $E[P(\text{correct})]$ is maximum when the expected values of p_1 and p_2 are maximum, and their covariance is zero.

Table 1

Sub-band word error rates for the Numbers95 database of connected digits are shown, from left to right, for when the full frequency range is divided into 1, 2 and 4 sub-bands, respectively^a

WER	Band 1	PoE WER	Band 1	Band 2	PoE WER	Band 1	Band 2	Band 3	Band 4	PoE WER
Clean	7.1	7.1	16.1	32.3	5.20	33.1	29.7	38.9	55.1	2.11
SNR 12	15.6	15.6	22.6	46.7	10.55	39.6	41.1	50.7	70.8	9.58
SNR 0	50.0	50.0	54.4	81.6	44.39	70.7	75.4	80.2	88.9	38.01

^a The PoE WER values would theoretically result if sub-band error rates combined according to the product-of-errors rule. In this case, although WER increases in each sub-band as the number of sub-bands is increased, the larger number of sub-bands would show considerable advantage. (See Section 7 for more details on sub-band frequency ranges and database.)

paper. It is interesting to note here that this trend has recently been followed in (Hermansky and Sharma, 1999).

4. Multi-band processing in artificial speech recognition

While there are many possibilities in ASR for combining evidence from different data streams, such as vision with acoustics, or acoustic features from different time scales (Greenberg, 1997; Kingsbury et al., 1998; Wu et al., 1998), in the present work we are primarily concerned with the combination of experts operating on different sub-divisions of the acoustic frequency spectrum.

Although the main motivation behind the multi-band approach is to exploit spectral data redundancy in a way which reflects the PoE rule for human speech perception, further potential advantages of the multi-band approach include the following:

- *Channel specific processing.* Different recognition strategies might ultimately be applied in each sub-band. For example, higher frequencies could use greater time resolution, and lower frequencies greater frequency resolution. It would also be possible to use sub-band specific speech sub-units.
- *Sub-unit specific expert combination.* It is sometimes possible to weight each expert according to the speech sub-unit, which is being distinguished. Consonants, for example, might give more weight to high frequency sub-bands.
- *Channel asynchrony.* Models discussed here use the same phoneme set for each expert, and force

synchrony between experts, but it would be possible to permit some level of sub-band asynchrony⁷ (Bourlard and Dupont, 1996; Hermansky et al., 1996; Tomlinson et al., 1997; Mirghafori, 1999), (Fig. 2), and to use speech sub-units specific to each frequency sub-range (Mirghafori, 1999).

The first multi-band ASR systems were based on the hidden Markov model/artificial neural network (HMM/ANN) model (Bourlard and Morgan, 1994; Bourlard and Morgan, 1997). In standard full-band ASR, an multi-layer perceptron (MLP) is first used to transform the acoustic data into posterior phoneme probabilities,^{8,9} $P(q_k|x^n; \Theta)$, for each word sub-unit, q_k , and data frame, x^n (Fig. 3). The posterior probabilities from the MLP are then passed as scaled likelihoods (Hennebert et al., 1997) into an HMM for decoding. In early multi-band processing (Bourlard and Dupont,

⁷ When streams are not frame synchronous the complexity of the decoding algorithm required may be considerably greater than for a standard recogniser. Results to date have indicated that allowing asynchrony between streams does not give any significant performance improvement (Mirghafori, 1999).

⁸ Any sufficiently flexible parametric function which is trained to minimise the sum of square errors against a target classification function (with one output per class) will generate output class posteriors which are optimal against the distribution of the training data (Richard and Lippmann, 1991). As the number of hidden units increases from zero, MLP performance increases rapidly to the Bayesian optimum, when the posterior probability for data belonging to each class is equivalent to the probability that classification is correct if this class is selected (Duda and Hart, 1993). Fig. 3 shows that this equivalence is very close for a typical MLP used for phoneme classification in continuous speech.

⁹ See Nomenclature section for full definition of all mathematical symbols used.

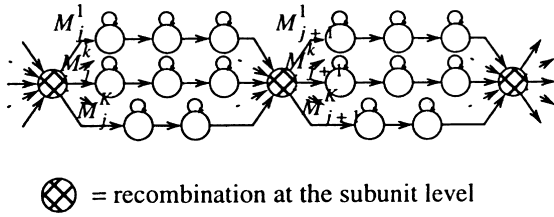


Fig. 2. General form of K -streams recogniser with anchor points between speech units (to force synchrony between different streams). Note that the model topology is not necessarily the same for the different sub-systems.

1996; Okawa et al., 1998; Mirghafori, 1999), the aim was to divide the frequency range into a number of sub-bands, x_i , and to process each of these independently, so one MLP was trained for each frequency sub-band, x_i . The posterior probabilities, $P(q_k|x_i; \Theta_i)$, were then combined, usually at the frame level, before decoding as for the full-band HMM/ANN. A number of different combination strategies have been proposed:

1. the linear weighted average:

$$P(q_k|x) \cong \sum_{i=1}^d w_i P(q_k|x_i; \Theta_i), \quad (2)$$

2. the geometric weighted average:

$$P(q_k|x) \cong \prod_{i=1}^d P(q_k|x_i; \Theta_i)^{w_i}, \quad (3)$$

3. MLP combination.

In this context, $w_i \in [0, 1]$ is a fixed weight for each sub-band expert which should reflect the extent to which it contains features which are useful for phoneme discrimination. Best results were obtained using the more flexible of these functions, MLP combination (Mirghafori, 1999).

This approach has a number of limitations. If it was required to adapt weights to changing noise conditions, then supervised weight training could not be used. This problem could be overcome using methods similar to those described in Section 6. A more important limitation of this sub-band approach is that *independent sub-band processing loses all joint spectral information*, such as the shape of the spectral envelope. The PoE rule tells us that intelligent combination of a number of narrow band recognisers should be able to equal or outperform the full-band expert, but experimental results have repeatedly shown that independent sub-band processing leads to a significant decrease in performance with clean speech.

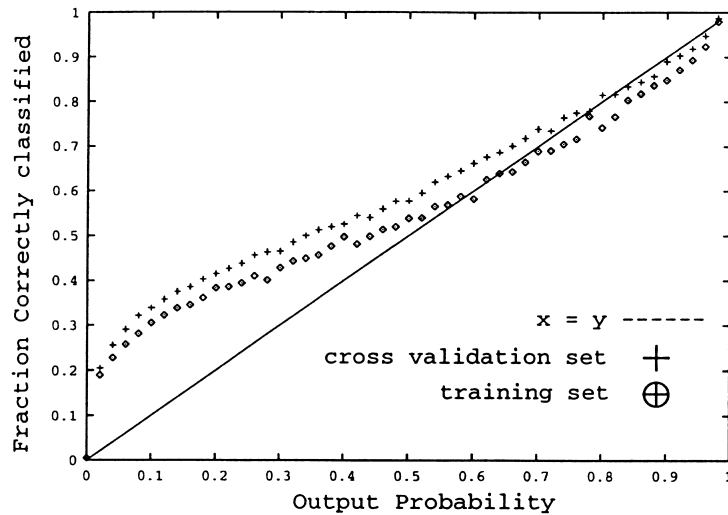


Fig. 3. It is possible to generate “good” posterior probabilities from a neural network, and these are indeed good measures of the probability being correct. This plot was generated from real speech data by collecting statistics over the acoustic parameters from 1750 resource management speaker-independent training sentences and 500 cross-validation sentences (not used for training, but for which correct classification was known) (Boulard and Morgan, 1994).

5. Multi-band ASR with latent variables

The first (rule of the method which leads to the truth) is to accept nothing as true which you do not clearly recognise to be so. (Descartes, 1619).

It is now known that the PoE rule is only approximately true (Steeneken and Houtgast, 1999). Even in human perception, where the error rate in each sub-band is far lower than with sub-band ASR (especially in noise), some use is made of joint sub-band information. One of the first methods, which were attempted to overcome this problem, was to combine the four sub-band experts with a full-band expert. This led to performance improvement towards, but not significantly surpassing, that of the full-band expert on its own (Mirghafori, 1999). In (Hermansky et al., 1996) tests were made with a seven sub-band system in which

an expert was used for each possible sub-band combination.

Several combination strategies were tested. Simple linear combination showed satisfactory performance in clean speech, while hand selection of the correct expert (when present) showed the potential of this system for very strong robustness to band limited noise. This confirmed the following result, which had previously been demonstrated in the context of missing feature theory (MFT) (Lippmann and Carlson, 1997; Morris et al., 1998): an effective strategy for reducing the effects of data mismatch is to detect and simply ignore strongly mismatching data.

At any point in time there are a number of sub-bands whose mismatch level is such that recognition will improve if data in these sub-bands is ignored.

An obvious problem with the strategy of ignoring unreliable data is that estimating the local data mismatch level in each sub-band is not easy. At best we can assign a probability to each sub-band being “clean” (i.e. having mismatch below some threshold level) at time step n . In the case of d sub-bands, there are 2^d possible sub-band combinations (if we include both the full and empty sets).

Let s_i denote the i th combination of sub-bands from x . In estimating the full-band posterior $P(q_k|x)$, we would like to select the expert i , which gives the “best” estimate for this phoneme. In the presence of band limited noise, we can reasonably assume that the best combination will be the largest set of clean sub-bands. Define an indicator latent variable c_i to represent the event that s_i is the largest set of clean sub-bands. Providing events s_i include the empty set, the set of events c_i are mutually exclusive and exhaustive. Therefore,

$$P(q_k|x) = P\left(q_k \cap \bigcup_{i=1}^{2^d} c_i \middle| x\right), \quad (4)$$

because the set of events c_i are exhaustive,

$$= \sum_i P(q_k \cap c_i|x), \quad (5)$$

because c_i are mutually exclusive,

$$= \sum_i P(c_i|x)P(q_k|c_i \cap x). \quad (6)$$

The condition “ c_i true” tells us that x can be partitioned into certain and uncertain parts (s_i, s'_i). If nothing is known about the uncertain data, then this data is simply “missing”, or “not given”, so that

$$P(q_k|c_i \cap x) = P(q_k|s_i). \quad (7)$$

Provided the number of combinations is not too large (d not greater than about 7), we can train an MLP expert on clean data from each combination¹⁰ to output phoneme posterior estimates $P(q_k|s_i; \Theta_i)$.¹¹ It is given that s_i in Eq. (7) is clean, so $P(q_k|s_i) \cong P(q_k|s_i; \Theta_i)$. Therefore,

¹⁰ It is important to note that data should be orthogonalised within each sub-band (Bourlard and Dupont, 1996; Okawa et al., 1998). Orthogonalisation across the full data vector would usually spread noise between sub-bands, after which the noise would not be band limited and the advantage of sub-band processing would be lost. Alternative means of spectral data orthogonalisation have been proposed (Nadeu et al., 1995) which mix features across time rather than frequency, but this would also have a tendency to mix clean with noisy data.

¹¹ We distinguish here between the true probabilities $P(q_k|s_i)$ and their estimates $P(q_k|s_i; \Theta_i)$, which are parametric functions (MLPs) trained on clean speech, and therefore can only be assumed to be accurate when it is known that the speech data in sub-band subset s_i is clean.

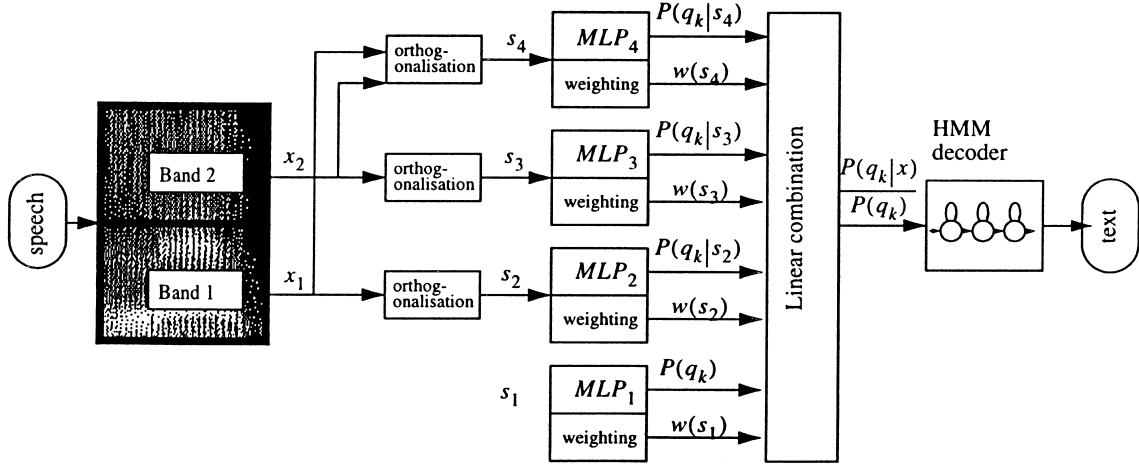


Fig. 4. Full Combination multi-band ASR HMM/ANN system with two sub-bands. Each possible combination of frequency sub-bands is first orthogonalised and then input to an MLP which has been trained for this sub-band combination, and to a weighting function. Posterior phoneme probabilities from each MLP expert are then combined in a linear weighted sum to give a single posterior for each phoneme. These posteriors are then converted to scaled likelihoods before input to a fixed parameter HMM for Viterbi decoding.

$$P(q_k|x) \cong \sum_i P(c_i|x)P(q_k|s_i; \Theta_i). \quad (8)$$

Eq. (8) gives us a factorisation of the full-band posterior into a composition of combination weights and clean combination posteriors.¹² We call this the “full combination” (FC) multi-band approach. The utility of this method depends directly on the accuracy with which combination weights $P(c_i|x)$ can be estimated and adapted to changing noise conditions. Various methods for combination weight estimation are described in Section 6. The full-combination ASR system used in the experiments reported here was implemented in an HMM/ANN based system with four sub-bands. A two sub-band full-combination multi-band model (FCM) system is illustrated in Fig. 4.

Full-combination approximation. For d much greater than about seven it becomes impractical

to train a separate MLP for all 2^d sub-band combinations. One way to alleviate this problem may be to prune out combinations, which have a priori negligible weight for a particular application. Alternatively, although we must avoid the assumption $p(x_i|x_j) \cong p(x_i)$ of full sub-band independence, it is easy to show (Morris et al., 1999) that under the weaker assumption of *conditional* sub-band independence, $p(x_i|x_j \cap q_k) \cong p(x_i|q_k)$, all 2^d combination posteriors can be expressed in terms of the d single sub-band posteriors. If $|s_i|$ is the number of sub-bands in combination s_i , then the clean combination posteriors $P(q_k|s_i)$ in Eq. (8) can be approximated as follows:

$$P(q_k|s_i) \cong \frac{\bar{P}_{ki}}{\sum_m \bar{P}_{mi}}, \quad (9)$$

where $\bar{P}_{ki} = P^{1-|s_i|}(q_k) \prod_{x_j \in s_i} P(q_k|x_j)$.

It is shown in Section 7 that the full-combination approximation based on Eq. (9) consistently outperforms the assumption of full independence that is implicit in the simple sub-band combination of Eq. (2).

¹² Likelihood based models (such as HMMs) can be decomposed in a similar way, but some complications can arise. More details on both posteriors and likelihood based decomposition are discussed in (Morris, 1999).

6. Full-combination expert weighting

We describe here a number of approaches to both fixed and adaptive weighting.

If the latent variable c in Eq. (8) is not an indicator for each sub-band combination being the largest clean combination, then we are no longer justified in assuming that $P(q_k|s_i) \cong P(q_k|s_i; \Theta_i)$. On the other hand, if we can assume that the data is clean, and we are interested in exploiting the possibility that data in some sub-bands should carry more weight than data in others, then we can replace c by the latent variable b , which indicates whether each sub-band combination is “the best” or most useful. As b_i are also mutually exclusive and exhaustive, the same working used to derive Eq. (8) also gives us ¹³

$$P(q_k|x) \cong \sum_i P(b_i|x)P(q_k|s_i; \Theta_i). \quad (10)$$

Although the weights $P(b_i|x)$ in Eq. (10) depend on x , if we are interested in using fixed weights, then we must ignore x and assume that $P(b_i|x) \cong P(b_i)$.

6.1. Fixed weight estimation using linear and non-linear LMSE

Fixed weights $w_i = P(b_i)$, $i = 1, \dots, 2^d$, can be estimated using the (supervised) least mean square error (LMSE) criterion,

$$w_{\text{LMSE}} = \arg \min_w \sum_{n=1}^N \sum_{k=1}^{2^d} (y_k^n(w) - t_k^n)^2, \quad (11)$$

where $y_k^n(w)$ is the combined posterior for q_k at frame n , given fixed weights w , and $t_k^n = P(q_k|x^n)$ is the target posterior, $= 1$ if target class for frame n is class q_k , else $= 0$. For *linear* LMSE

$$y_k^n(w) = P(q_k|x^n; w, \Theta) = \sum_{i=1}^{2^d} w_i P(q_k|s_i^n; \Theta_i). \quad (12)$$

In this case, the resulting LMSE “normal equations” are linear and can be solved directly, but in

this case the weights cannot be constrained to be positive or sum to 1 across all experts. For non-linear LMSE, when an MLP is trained using “back error propagation” (a particular case of gradient descent), weights can be constrained to sum to 1, if required, by using the softmax activation function in the output layer.

6.2. Fixed weight estimation using maximum likelihood

The fixed weights, $w_i = P(b_i)$, can be estimated using relative-frequencies ¹⁴ as follows:

$$w_i = P(b_i) \cong n_i/n, \quad (13)$$

where n_i is the number of frames of training data for which expert i has the largest posterior, across all experts, for the target phoneme (and therefore has the smallest Kullback–Leibler distance from the target probability distribution), and n is the number of frames of training data.

6.3. Adaptive weighting using estimated level of local data mismatch

It is usually reasonable to assume that sub-band combination reliability can be estimated from the reliability of each of its component sub-bands, and that sub-band reliabilities are independent. The local data mismatch level in each sub-band can be estimated from measures $\gamma(x_j)$ of *likeness to speech data*, such as local signal to noise ratio (SNR) ¹⁵ (Hirsch and Ehrlicher, 1995; Morris et al., 1999), harmonicity index (Berthommier and Glotin, 1999) or, in the case of stereo data, interaural time delay (Glotin et al., 1999). By definition of c_i the adaptive weights ($P(c_i|x)$) can then be obtained from these sub-band reliabilities as follows:

¹⁴ Relative-frequency is the maximum likelihood (ML) estimate for a Bernoulli probability.

¹⁵ SNR is not directly related to mismatch level and should be used with care. Low noise energy implies low data mismatch, but low noise energy combined with zero speech energy gives infinitely low SNR.

¹³ It is also possible to combine the latent variables b_i and c_j by using a double summation (Morris, 1999).

$$\begin{aligned}
c_i &\iff (\text{reliable}(x_j) \forall x_j \in s_i) \cap (\neg \text{reliable}(x_j) \forall x_j \notin s_i) \\
&\Rightarrow P(c_i|x) = \prod_{x_j \in s_i} P(\text{reliable}(x_j)) \prod_{x_j \notin s_i} P(\neg \text{reliable}(x_j)).
\end{aligned} \tag{14}$$

If the measure $\gamma(x_j)$ is increasing with level of mismatch, then we can model $P(\text{reliable}(x_j))$ as $P(\gamma(x_j) < \varepsilon_j)$, where the threshold ε_j for each sub-band is estimated as the average level of $\gamma(x_j)$, over the training set, above which recognition improves when band x_j is ignored. If we assume that the estimated mismatch value $\hat{\gamma}(x_j)$ is its true value plus a zero mean Gaussian error, with variance σ_j^2 , then we obtain

$$\begin{aligned}
P(\text{reliable}(x_j)) &\cong P(\gamma(x_j) < \varepsilon_j) \\
&= \Phi\left(\frac{\varepsilon_j - \hat{\gamma}(x_j)}{\sigma_j}\right),
\end{aligned} \tag{15}$$

where σ_j^2 are parameters which reflect our confidence in the estimation precision for each sub-band, which could be tuned, for example, to maximise performance in clean speech.

7. Experimentation

Databases. Speech was taken from the Numbers95 database of multi-speaker US English free-format numbers¹⁶ telephone speech, with 30 words and 33 phonemes (Cole et al., 1995). Factory noise from the Noisex92 database (Varga et al., 1992), and car noise from an in-house database, were added at varying SNR levels relative to the average signal energy in each utterance (excluding non-speech periods). Multi-band ASR is best suited to band limited noise, but the noises used in these initial tests was not band limited (Fig. 5).

Acoustic features. Tests were made using PLP coefficients (Hermansky, 1990; Rao and Pearlman, 1996; Morgan et al., 1998), with an analysis window size of 25 ms, and 12.5 ms shift. Tests were also made using J-Rasta-PLP features (Hermansky and Morgan, 1994; Morgan et al., 1998), which result from PLP features after an adaptive

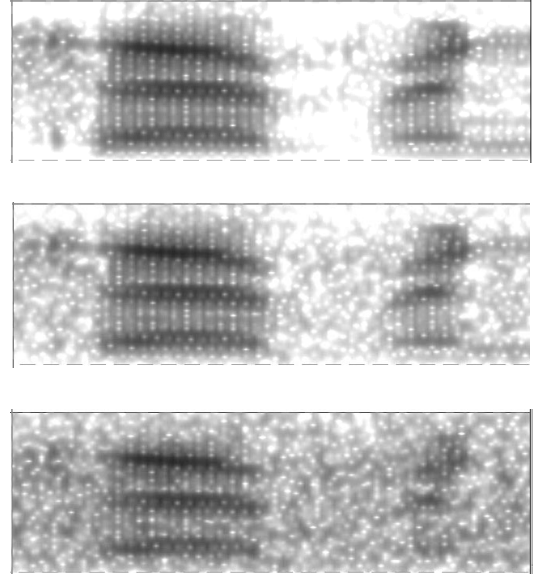


Fig. 5. Signal “seven” from Numbers95 (top) mixed with Daimler–Chrysler car noise at SNR = 12 (middle) and 0 dB (bottom).

cancellation of slowly changing additive and convolutive noises. In order to compete with the state-of-the-art baseline system used here, data features had to be near orthogonal. Orthogonalisation was applied within each data stream, so as not to spread noise between streams.

Recognition systems. Baseline tests were made with a full-band HMM/ANN hybrid (Bourlard and Morgan, 1997) in which the ANN was an MLP with one hidden layer of 1000 units, trained to input 9 consecutive data frames and output the posterior probabilities $P(q_k|x)$ for each of 33 phonemes q_k and each data frame, x . Training used the full Numbers95 training set (Cole et al., 1995). In recognition, posterior probabilities from the MLP, divided by their priors¹⁷, were passed as

¹⁷ Although we should *theoretically* divide posteriors by priors to get the (scaled) likelihoods which are required by the HMM for decoding, in practice the division by estimated priors is only helpful if these estimates accurately represent the priors in the test data. Inaccurate priors can result if some of the phonemes in the training set were undersampled, or if the true priors in the training set do not match the priors in the test set. Better results were obtained in the experiments reported here when the posteriors were *not* divided by the priors. This is equivalent to assuming that all priors are equal.

¹⁶ E.g. “two hundred eleven”.

scaled likelihood to a fixed parameter HMM for decoding. The HMM for each phoneme used a one to three repeated-state model. No language model was used.

For the full-combination multi-band system a separate MLP (of identical design to the full-band MLP) was trained on clean data for each sub-band combination. Multiple MLP outputs were then merged at the frame level (which here was also the state level), using Eq. (8), to give a single posterior probability for each class, before these were passed as scaled likelihood to the same HMM as used by the full-band system.

Sub-band definition. The number of sub-bands used in sub-band ASR, and their frequency ranges, is somewhat arbitrary. For these experiments we have chosen to work with four sub-bands. The original reason for this choice was based on having one sub-band for each formant. Frequency ranges were chosen accordingly and are displayed in Table 2. Although it would be instructive to test with different sub-band configurations, we kept the sub-band specification constant so that test results could be compared over a large number of different experimental configurations.

Recognition tests. Recognition tests were made with the two different noise types as described, at SNR levels clean (45 dB SNR), 12 dB and 0 dB SNR. Tests used the full Numbers95 development test set. Results are summarised in Fig. 6 for PLP features, and Fig. 7 for J-Rasta-PLP features.

Top figures show results using different weighting measures, car noise (left) and factory noise (right), for the following:

1. full-band baseline HMM/ANN hybrid ASR system,
2. full combination multi-band hybrid, with four sub-bands, using equal weights, “FCM, equal wts”,
3. FCM using relative-frequency, “FCM, RF wts”,

4. FCM with cheating, “FCM, cheating”, i.e. correct phoneme selected whenever selected by any expert.

Bottom figures show results using different multi-band models, all using equal weights, for the following:

1. full-band baseline HMM/ANN hybrid ASR system,
2. full combination multi-band hybrid, using equal weights, “FCM, equal wts”,
3. full combination approximation, using Eq. (9), using equal weights, “approximate full combination multi-band model (AFC), equal wts”,
4. early multi-band approach (one expert per sub-band), with equal weights, “early mband, eq wts”.

Results. As usual, J-Rasta-PLP features show a strong improvement over PLP features (which give similar performance to MFCCs). Note that the benefits of robust preprocessing are largely complementary to those offered by multi-band processing. Of the two noise types tested, the more impulsive factory noise has a worse effect on performance than car noise, which is near stationary. The “cheating” performance, which would obtain if the PoE rule was in effect (i.e. if there was correct phoneme recognition whenever any sub-band combination expert had correct recognition), is very robust down to about 12 dB SNR. None of the methods tested has performance anywhere near as good as cheating, and only one (FCM with fixed RF weights and Rasta features) has performance, which improves over the baseline in noise. Equal weights perform almost as well as FC weights, especially with Rasta features. The failure of any of the multi-band methods tested to significantly improve over the full-band baseline may be because both noise types also affected every sub-band, while the advantage of multi-band processing is strongest when one or more sub-bands are often clean. However, the

Table 2
Sub-bands definition

Sub-bands def.	Band 1	Band 2	Band 3	Band 4
Frequency range/Hz	216–778	707–1632	1506–2709	2122–3769
#coeffs, 1pc order	5, 3	5, 3	3, 2	3, 2

Multi-stream adaptive evidence combination for noise robust ASR

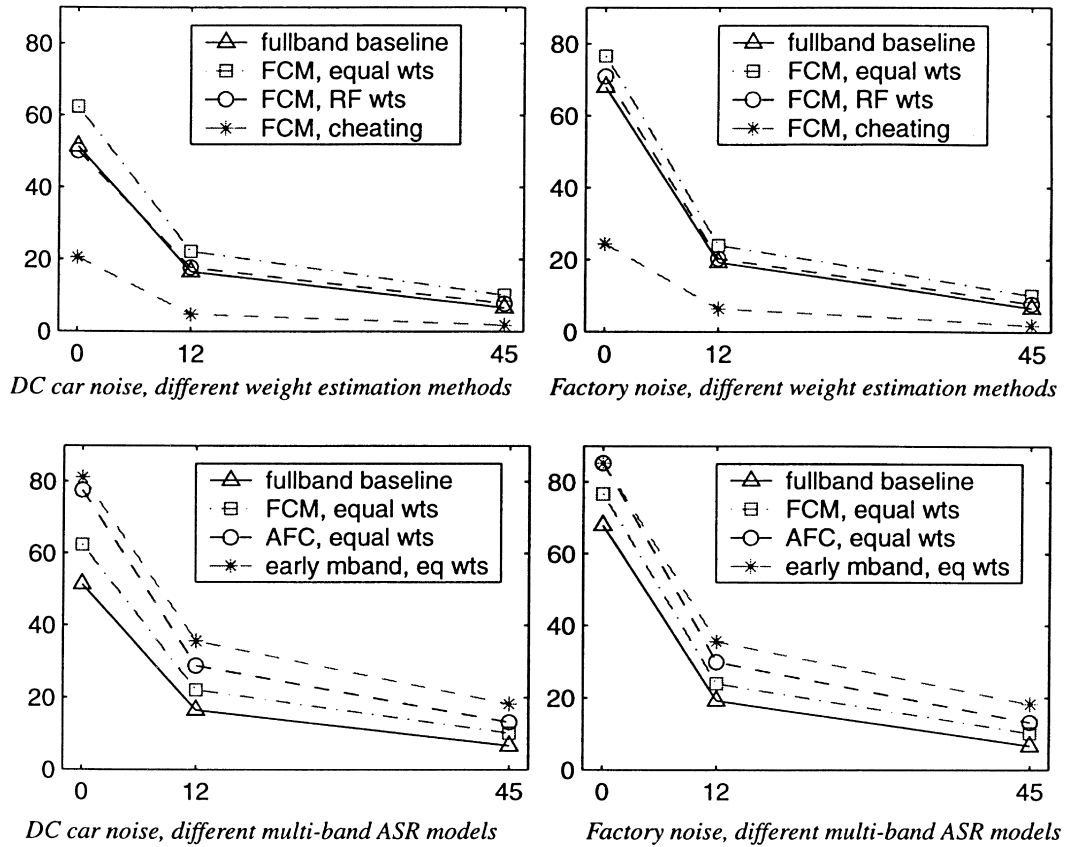


Fig. 6. WER (vertical axis) against SNR with PLP features.

full-combination method consistently outperforms the full-combination approximation method (using just one expert per sub-band), which in turn significantly outperforms the “early sub-band” approach (which also used one expert per sub-band, but combines them in a less principled way).

8. Summary and conclusion

After an introduction to the general advantages of multi-stream processing, and to the evidence for its use in human perception, we recounted how Fletcher’s PoE rule for independent sub-band processing in speech perception, as well as a

number of other factors, motivated the development of multi-band processing for the purpose of robustness to band-limited noise in ASR. We then argued that the original formulation of multi-band ASR was disadvantaged by independent sub-band processing, and described how this problem could be theoretically overcome by performing separate recognition on every possible combination of sub-bands, which enables us to integrate over all possible positions of noisy or mismatching data. A derivation of this “full combination” multi-band approach was presented in the context of HMM/ANN based ASR in which an MLP outputs phone posterior probabilities for each time frame and for each sub-band combination. An indicator

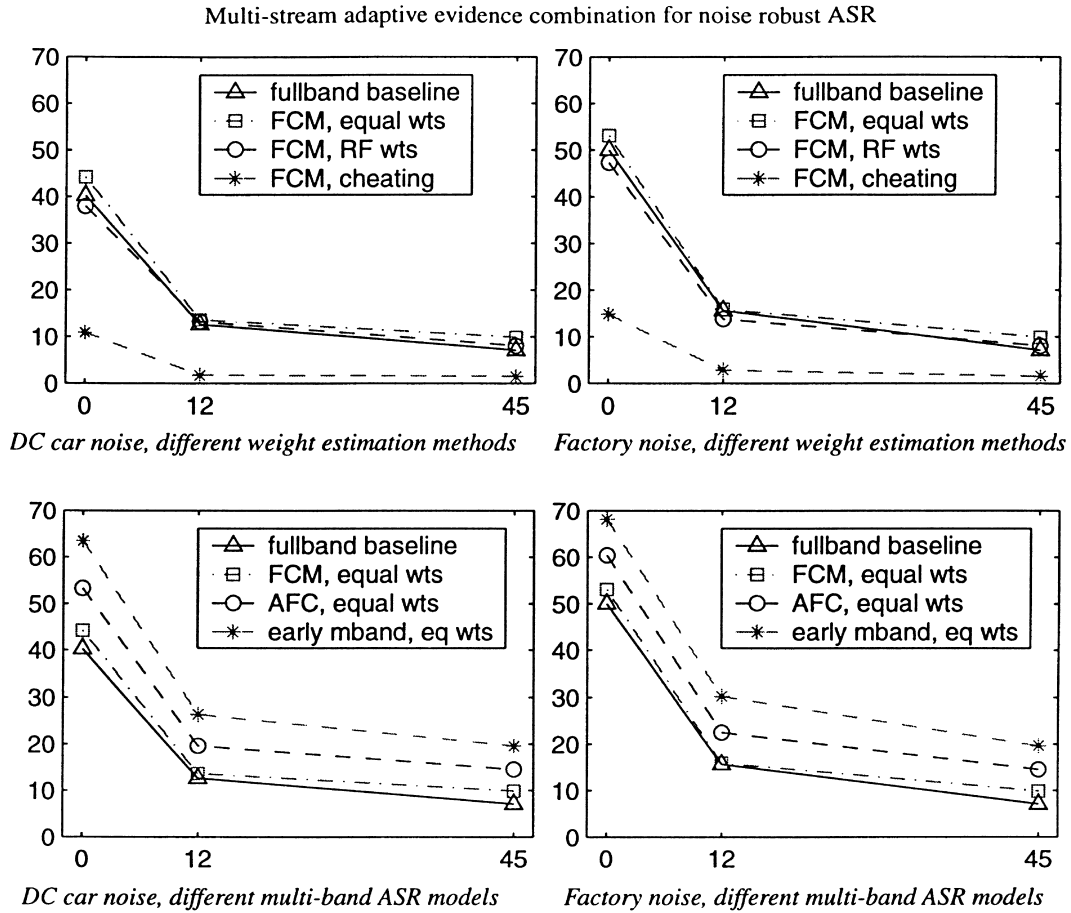


Fig. 7. WER (vertical axis) against SNR with J-Rasta-PLP features.

latent variable was used to specify whether each sub-band combination is the largest set of clean sub-bands. In this way, the full-band phoneme posterior for each phoneme was decomposed into reliability weighted sum of corresponding clean sub-band-combination posteriors. While theoretically attractive, especially in the context of time varying and band limited noise, the effectiveness of this approach depends strongly on the method used for reliability weighting. A number of weighting methods was presented. Some of these methods are only suitable for estimating static weights, while others can be used in the more general context where mismatch level varies over time and frequency.

The full-combination technique was tested on a speaker independent continuous speech free-for-

mat numbers recognition task, for clean speech and for speech plus car or factory noise. Tests compared the full-combination technique with the early sub-band approach, and also with the full-combination approximation. Results showed that, in the case of four sub-bands,

- both the full-combination method and its approximation outperform the early sub-band approach, even when equal weights are used for each combination expert,
- relative frequency based FCM weighting marginally outperforms the full-band baseline in noise when using J-Rasta-PLP features.

These results, for static weighting, have not shown a clear advantage for multi-band ASR over the full-band baseline. However, the main advantage of adaptive weighting in multi-band ASR is with

mismatch which is band limited and/or changing, while the initial experiments which we report here use only noises which are not band limited and are near constant. Tests were also made using adaptive weights based on estimated local SNR (Hirsch and Ehrlicher, 1995), but these weights did not perform better than equal weights so results were not shown.

The full advantage of the full-combination approach will not come into play unless mismatch is unequally distributed over time and/or frequency, with some proportion of data remaining clean. It is therefore imperative that effective noise robust preprocessing is used to remove as many constants or slowly changing noise as possible prior to multi-band processing. The FCM results reported here do not much improve over the baseline system, but they do show some increased advantage when noise robust J-Rasta-PLP features are used in place of simple PLP features, even when the noise is not band limited. It is also clear that expert weighting (mismatch estimation) should be as accurate and rapidly adapting as possible.

Future directions. The static weighting methods reported in Section 6 are clearly not optimal in the presence of noise. The techniques for adaptive mismatch estimation outlined in Section 6.3 assume that training data is “speech like”, and mismatching data is non speech like. However, mismatch concerns only how field data differs from the data population used in training, and not how it differs from “clean speech”. As the training population is defined by its pdf or “likelihood function”, it can be argued that mismatch estimation should be based, directly or indirectly, on use of the training data pdf. One such approach would be to base mismatch detection on the assumption that clean and mismatching data have different (distributions of) likelihood values. Another would be to use weights, which maximise the local clean-data likelihood (although this method could be applied only in the context of likelihood based FCM decomposition). We have developed both of these approaches in more detail in (Morris, 1999), although they have not yet been fully tested. Other approaches making more or less direct use of the training data pdf to downweight mismatching sub-bands are described in (de Veth et al., 1999; Ming and Smith, 1999).

Finally, it should be noted that although the “early” multi-band approach tested here shows significantly reduced performance in relation to the full-combination approach, a lot of variants of independent sub-band processing have reported positive performance results (Bourlard and Dupont, 1996; Kingsbury et al., 1998; Mirghafori, 1999). With the present full-combination approach we are theoretically limited to linear expert combination. However, some of the best results for independent sub-band processing were obtained using MLP combination, which is highly non-linear. There are also a number of other combination strategies which worked well for independent sub-band processing which we have not yet tested with FCM, such as weighting by some function of the entropy in the posteriors output from each expert (Hermansky et al., 1996; Berthommier and Glotin, 1999; Mirghafori, 1999).

Nomenclature

ASR	automatic speech recognition
ANN	artificial neural network
MLP	multi-layer perceptron
HMM	hidden Markov model
ML	maximum likelihood
SNR	signal to noise ratio
WER	word error rate
FCM	full-combination multi-band model
AFC	approximate full combination multi-band model
$P(x)$	probability of “event x ” occurring
$p(x)$	probability density at x of a continuous value
$P(x; \Theta)$	function with parameters Θ used to estimate $P(x)$
$\hat{P}(x)$	function (parameters unspecified) used to estimate $P(x)$
$P(a \cap b)$	probability of events a and b occurring (same as $P(a, b)$)
$P(a \cup b)$	probability of events a or b occurring
$\neg a$	not a , the negation of condition a
q_k	speech unit whose presence is being estimated, or event that data x is from this class
x, x^n	data window vector at time step n
d	number of spectral sub-bands
x_i	i th sub-band of x , for $i = 1, \dots, d$

s_i, s'_i	i th sub-band combination of x , and its complement, for $i = 1, \dots, 2^d$
x clean	x is from data population used in model training (i.e. x has no data mismatch)
b_i	s_i is best subset of x for estimating speech unit posteriors
c_i	s_i is largest clean subset of x
b_{ik}	s_i is best subset of x for estimating speech unit posteriors when true unit is q_k
w	vector of sub-band expert weights
$\Phi(x)$	cumulative density function for the standard Gaussian pdf

Acknowledgements

This work was supported by the Swiss Federal Office for Education and Science (OFES) in the framework of both the EC/OFES SPHEAR (SPeech, HEARing and Recognition) project and the EC/OFES RESPITE project (REcognition of Speech by Partial Information TEchniques).

References

- Allen, J.B., 1994. How do humans process and recognise speech?. *IEEE Trans. Speech Signal Process.* 2 (4), 567–576.
- Berthommier, F., Glotin, H., 1999. A new SNR-feature mapping for robust multi-stream speech recognition. In: *Proc. ICPHS'99*, pp. 711–715.
- Bishop, C. (Ed.), 1995. *Neural Networks for Pattern Recognition*. Clarendon Press, Oxford, pp. 365–368.
- Bouclard, 1999. Non-stationary multi-channel (multi-stream) processing towards robust and adaptive ASR. In: *Proceedings of Tampere Workshop on Robust Methods for Speech Recognition in Adverse Conditions*, pp. 1–10.
- Bouclard, H., Dupont, S., 1996. A new ASR approach based on independent processing and recombination of partial frequency bands. In: *Proc. ICSLP'96*, Philadelphia, pp. 422–425.
- Bouclard, H., Morgan, N., 1994. *Connectionist speech recognition – a hybrid approach*. Kluwer Academic Publishers, Dordrecht.
- Bouclard, H., Morgan, N., 1997. Hybrid HMM/ANN systems for speech recognition: overview and new research directions. In: *Proceedings of International School on Neural Nets: Adaptive Processing of Temporal Information*.
- Cole, R.A., Noel, T., Lander, L., Durham, T., 1995. New telephone speech corpora at CSLU. In: *Proceedings of European Conference on Speech Communication and Technology*, Vol. 1, pp. 821–824.
- de Veth, J., de Wet, F., Cranen, B., Boves, L., 1999. Missing feature theory in ASR: make sure you missing the right type of features. In: *Proceedings of Workshop on Robust Methods for Speech Recognition in Adverse Conditions*, pp. 231–234.
- Duda, R.O., Hart, P.E., 1993. *Pattern Classification and Scene Analysis*. Wiley, New York.
- Dupont, S., Luetttin, J., 1998. Using the multi-stream approach for continuous audio-visual speech recognition: experiments on the M2VTS database. In: *Proc. ICSLP'98*, pp. 1283–1286.
- Fletcher, H., 1922. The nature of speech and its interpretation. *J. Franklin Inst.* 193 (6), 729–747.
- Gales, M.J.F., Young, S.J., 1993. HMM recognition in noise using parallel model combination. In: *Proc. Eurospeech'93*, pp. 837–840.
- Girin, L., Feng, G., Schwartz, J.-L., 1998. Fusion of auditory and visual information for noisy speech enhancement: a preliminary study of vowel transitions. In: *Proc. ICASSP'98*, pp. 1005–1008.
- Glottin, H., Berthommier, F., Tessier, E., 1999. A CASA-labelling model using the localisation cue for robust cocktail-party speech recognition. In: *Proc. Eurospeech'99*, pp. 2351–2354.
- Greenberg, S., 1997. On the origins of speech intelligibility in the real world. In: *Proceedings ESCA Workshop on Robust Speech Recognition for Unknown Communication Channels*, pp. 23–32.
- Hagen, A., Morris, A.C., Bouclard, H., 1999. Different weighting schemes in the full combination sub-bands approach for noise robust ASR. In: *Proceedings Tampere Workshop on Robust Methods for Speech Recognition in Adverse Conditions*, pp. 199–202.
- Hennebert, J., Ris, C., Bouclard, H., Renals, S., Morgan, N., 1997. Estimation of global posteriors and forward-backward training of hybrid systems. In: *Proc. Eurospeech'97*, pp. 1951–1954.
- Hermansky, H., 1990. Perceptual linear predictive (PLP) analysis of speech. *J. Acoust. Soc. Am.* 87 (4), 1738–1752.
- Hermansky, H., Morgan, N., 1994. RASTA processing of speech. *IEEE Trans. Speech Audio Process.* 2 (4), 578–589.
- Hermansky, H., Sharma, S., 1999. Temporal patterns (TRAPS) in ASR noisy speech. In: *Proc. ICASSP'99*, pp. 298–292.
- Hermansky, H., Tibrewela, S., Pavel, M., 1996. Towards ASR on partially corrupted speech. In: *Proc. ICSLP'96*, pp. 462–465.
- Hirsch, H.G., Ehrlicher, C., 1995. Noise estimation techniques for robust speech recognition. In: *ICASSP'95*, pp. 153–156.
- Jordan, M.I., Jacobs, R.A., 1994. Hierarchical mixtures of experts and the EM algorithm. *Neural Comput.* 6, 181–214.
- Kingsbury, B., Morgan, N., Greenberg, S., 1998. Robust speech recognition using the modulation spectrogram. *Speech Communication* 25 (1–3), 117–132.

- Lippmann, R.P., Carlson, B.A., 1997. Using missing feature theory to actively select features for robust speech recognition with interruptions, filtering and noise. In: *Proc. Eurospeech'97*, pp. 37–40.
- McGurk, H., McDonald, J., 1976. Hearing lips and seeing voices. *Nature* 264, 746–748.
- Ming, J., Smith, F.J., 1999. Union: a new approach for combining sub-band observations for noisy speech recognition. In: *Proceedings of Workshop on Robust Methods for Speech Recognition in Adverse Conditions*, pp. 175–178.
- Mirghafori, N., 1999. A multi-band approach to automatic speech recognition. PhD Dissertation, University of California at Berkeley, December 1998. Reprinted as ICSI Technical Report, ICSI TR-99-04.
- Moore, B.C.J., 1997. *An Introduction to the Psychology of Hearing*, 4th edition. Academic Press, New York.
- Morgan, N., Boulard, H., Hermansky, H., 1998. Automatic speech recognition: an auditory perspective. Research Report IDIAP-RR 98-17.
- Morris, A.C., 1999. Latent variable decomposition for posteriors or likelihood based sub-band ASR. Research Report IDIAP-Com 99-04.
- Morris, A.C., Cooke, M., Green, P., 1998. Some solutions to the missing feature problem in data classification, with application to noise robust ASR. In: *Proc. ICASSP'98*, pp. 737–740.
- Morris, A.C., Hagen, A., Boulard, H., 1999. The full-combination sub-bands approach to noise robust HMM/ANN based ASR. In: *Proc. Eurospeech'99*, pp. 599–602.
- Nadeu, C., Hernando, J., Gorricho, M., 1995. On the decorrelation of filterbank energies in speech recognition. In: *Proc. Eurospeech'95*, pp. 1381–1384.
- Okawa, S., Boccheri, E., Potamianos, A., 1998. Multi-band speech recognition in noisy environment. In: *Proc. ICASSP'98*, pp. 641–644.
- Pickles, J.O., 1988. *An Introduction to the Physiology of Hearing*. Academic Press, New York.
- Rao, S., Pearlman, W.A., 1996. Analysis of linear prediction, coding and spectral estimation from sub-bands. *IEEE Trans. Inf. Theory* 42, 1160–1178.
- Raviv, Y., Intrator, N., 1996. Bootstrapping with noise: an effective regularisation technique. *Connection Sci.*, Special Issue on Combining Estimators, 8, 356–372.
- Richard, M.D., Lippmann, R.P., 1991. Neural network classifiers estimate Bayesian a-posteriori probabilities. *J. Neural Comput.* 3 (4), 461–483.
- Steeneken, H.J.M., Houtgast, T., 1999. Mutual dependence of the octave-band weights in predicting speech intelligibility. *Speech Communication* 28 (2), 109–123.
- Tomlinson, J., Russel, M.J., Brooke, N.M., 1996. Integrating audio and visual information to provide highly robust speech recognition. In: *Proc. ICASSP'96*, pp. 821–824.
- Tomlinson, J., Russel, M.J., Moore, R.K., Bucklan, A.P., Fawley, M.A., 1997. Modelling asynchrony in speech using elementary single-signal decomposition. In: *Proc. ICASSP'97*, pp. 1247–1250.
- Varga, A., Moore, R., 1990. Hidden Markov model decomposition of speech and noise. In: *Proc. ICASSP'90*, pp. 845–848.
- Varga, A., Steeneken, H.J.M., Tomlinson, M., Jones, D., 1992. The Noisex-92 study on the effect of additive noise on automatic speech recognition. Technical Report DRA Speech Research Unit.
- Westphal, M., Waibel, A., 1999. Towards spontaneous speech recognition for on-board car navigation and information systems. In: *Proc. Eurospeech'99*, pp. 1955–1958.
- Wu, S.-L., Kingsbury, B.E., Morgan, N., Greenberg, S., 1998. Performance improvements through combining phone and syllable scale information in automatic speech recognition. In: *Proc. ICASSP'98*, pp. 459–462.