



ELSEVIER

Speech Communication 37 (2002) 69–87

SPEECH
COMMUNICATION

www.elsevier.com/locate/specom

Automatic transcription of Broadcast News

S.S. Chen, E. Eide, M.J.F. Gales, R.A. Gopinath^{*}, D. Kanvesky, P. Olsen

IBM T.J. Watson Research Center, Yorktown Heights, NY 10598, USA

Abstract

This paper describes the IBM approach to Broadcast News (BN) transcription. Typical problems in the BN transcription task are segmentation, clustering, acoustic modeling, language modeling and acoustic model adaptation. This paper presents new algorithms for each of these focus problems. Some key ideas include Bayesian information criterion (BIC) (for segmentation, clustering and acoustic modeling) and speaker/cluster adapted training (SAT/CAT). © 2002 Elsevier Science B.V. All rights reserved.

Zusammenfassung

Dieser Beitrag beschreibt ein bei IBM entwickeltes Verfahren zur Transkription von Rundfunknachrichten. Zu den typischen Problemen der Transkription gehören die Segmentierung, Clusterung, akustische Modellierung, Sprachmodellierung und akustische Modelladaption. In diesem Beitrag werden neuartige Algorithmen für diese einzelnen Probleme präsentiert. Als einige der hier enthaltenen Schlüsselideen können das Kriterium der Bayes'schen Information (zur Segmentierung, Clusterung und akustischen Modellierung) sowie sprecher- und clusteradaptives Training genannt werden. © 2002 Elsevier Science B.V. All rights reserved.

Résumé

Cet article décrit l'approche d'IBM pour la transcription des nouvelles du journal télévisé. Les problèmes caractéristiques de cette tâche sont la segmentation, la classification automatique, la modélisation acoustique, l'estimation de modèles de langage et l'adaptation des modèles acoustiques. Cet article présente de nouveaux algorithmes pour chacun de ces problèmes. Parmi les idées clés, on peut citer le critère d'information Bayésien (pour la segmentation, la classification automatique et la modélisation acoustique) et l'apprentissage adapté pour chaque locuteur/classe de locuteurs. © 2002 Elsevier Science B.V. All rights reserved.

Keywords: Speech recognition; Broadcast news; Acoustic modeling; Adaptive training; Segmentation and clustering; Model selection

1. Introduction

Automatic transcription of Broadcast News (BN) is an extremely challenging and interesting large vocabulary continuous speech recognition

(LVCSR) task. The difficulty lies in an intermingling of speech and non-speech events such as music and background noise. Moreover, the speakers have different speaking styles and accents, degraded by microphone or channel conditions and with background noise and sometimes background speech. Any solution to this problem, depending on transcription accuracy and speed, enables a wide range of useful applications such as

^{*} Corresponding author.

E-mail address: rameshg@us.ibm.com (R.A. Gopinath).

close-captioning and audio-indexing of multimedia content.

The past few years have seen significant algorithmic advances in BN transcription in US English. This research effort has been sponsored by DARPA, through acoustic and textual data collections distributed by the linguistic data consortium (LDC) and through yearly evaluations conducted by NIST. About 160 h of acoustic training data and nearly 400 million words of textual training data in the BN domain is available from LDC.

Automatic transcription of BN is a multifaceted problem that can not reasonably be expected to be resolved with a single algorithmic technique. The effort at IBM has spanned a whole spectrum of algorithmic techniques and innovations. This paper reviews techniques developed by the authors over the last few years and deployed by IBM in the 1997–1999 Hub4 evaluations. The issues addressed span clustering and segmentation, acoustic modelling with and without channel adaptation, pronunciation modeling, language modeling as well as techniques to combine individual improvements (e.g., ROVER). While several of the techniques were developed in the context of BN transcription, they are general techniques that can be applied to other speech recognition tasks.

Segmentation of continuous audio into speaker/channel turns is useful in state-of-the-art speech recognition systems that employ speaker/channel adaptive schemes like cepstral mean normalization. Moreover, if the content is BN, one would like to *segment* the audio stream into homogeneous regions according to speaker identity, environmental condition and channel condition so that regions of different nature can be handled differently: for example, regions of pure music and noise can be rejected; also, one might design a separate recognition system for telephone speech. The various segmentation algorithms proposed in the literature fall into three categories: *decoder-guided* (Kubala et al., 1997; Woodland et al., 1997), *model-based* (Bakis et al., 1997a,b) and *metric-based* (Beigi and Maes, 1998; Gish and Schmidt, 1994; Siegler et al., 1997). The decoder-guided segmentation only places boundaries at silence locations, which in general has no direct connec-

tion with the acoustic changes in the data. Both the model-based and the metric-based segmentation schemes rely on thresholding of measurements which lack stability and robustness. More importantly, they do not generalize to unseen acoustic conditions. In contrast the so-called Bayesian information criterion (BIC) segmentation algorithm described in this paper is robust, threshold-free and generalizes well to unseen acoustic conditions and is based on detecting *any* change in the acoustic condition.

When speaker turns have been identified, it is useful to adapt the acoustic models to a given speaker for improved recognition performance. Since the same speaker may occur several times in a BN clip, it is necessary to cluster all segments from the same speaker to get the maximum gain from adaptation. Typically hierarchical clustering schemes are used for this purpose (Gish and Schmidt, 1994; Kubala et al., 1997; Siegler et al., 1997). In such schemes it is often difficult to determine the number of clusters. One can heuristically pre-determine the number of clusters or the minimum size of each cluster; accordingly, one can go down the tree to obtain desired clustering (Woodland et al., 1997). Another heuristic solution is to threshold the distance measures during the hierarchical process; the thresholding level is tuned on a training set (Siegler et al., 1997). Jin et al. (1997) shed some light on automatically choosing a clustering solution. In this paper, we view clustering as a model selection problem and show that the BIC criterion is an effective termination criterion for hierarchical clustering.

The diversity of BN speech data (several speakers, channel, and noise conditions) implies large variance in the trained acoustic models. By classifying acoustic training and test data into homogeneous conditions we may build separate acoustic models for each condition. This paper proposes two adaptive training methods: modified speaker adaptive training (SAT) and cluster adaptive training (CAT), both of which give significant performance improvements.

This paper also describes improvements to the modeling improvements that have a purely statistical basis. The classical modelling paradigm is expanded to include some modelling of covariance

structure and alternatives to the Gaussians density functions are suggested.

In speech recognition one models highly correlated features using diagonal Gaussians, even though full covariance Gaussians are more appropriate. Significant performance improvements are demonstrated using semi-tied covariances which model linearly transformed features using diagonal Gaussians.

A diagnostic study of histograms of any single dimension of the acoustic data assigned to a particular HMM state of a speech recognition system shows that these distributions are (a) skew-symmetric, (b) peakier around the mean and (c) have fatter tails than a typical Gaussian distribution. This paper proposes two efficient alternatives to model (b) and (c) that avoids dramatically increasing the number of Gaussian mixture components – Richter distributions and power-exponential distributions. Marginal performance improvements are obtained using these distributions.

Speech recognition systems typically use a hand-generated lexicon. Inconsistency between the pronunciations in the lexicon and the way a speaker pronounces the words is a major source of degradation in performance. A proposed resolution of this problem make use of network topologies for phonemes that incorporate common deviations in pronunciation from the lexicon.

Another problem area related to pronunciation deviation is fast speech. Fast speech is well known to significantly impair recognition accuracy (Sieger, 1995). This paper proposes two approaches – phone models with skip arcs, and faster-rate signal processing to handle fast speech.

In the BN domain there is a significant amount of training data. This paper shows that more accurate language models can be obtained by mixing language models built on subsets of the training corpus that span different periods in time rather than pooling all of the data and building a single model.

Schemes that combine hypothesis transcripts from several recognizers can greatly improve the accuracy. This is demonstrated using the ROVER algorithm (Fiscus, 1997) on a set of recognizers.

Our basic approach to BN transcription can be summarized as follows. First the continuous audio

stream is segmented into acoustically homogeneous clips of audio. Typically the audio in each segment is from a single speaker with a fixed background noise level and channel condition. A first-pass recognition of these segments is done using single canonical acoustic model built on all the available training data using adaptive training. Adaptive training takes into account the wide range of variability in the acoustic training data and produces a sharp canonical model. The first-pass recognition is followed by adaptation of the canonical model to each segment using the recognition transcripts as ground truth. This is followed by a second-pass recognition. The process of adapting the acoustic models for each segment and refining the transcription by recognition with the adapted models is iterated several times. The quality of the adaptation is linked to the amount of adaptation data in each segment. Since a typical broadcast show has several segments from the same speaker, the amount of adaptation data is increased by clustering the segments into similarity groups. Acoustic model adaptation is done on the clusters rather than on individual segments.

The recognizer used in all the experiments in this paper is the standard IBM rank-based single-pass stack decoder (see (Bahl et al., 1994, 1995) for details). The training data used for the experiments is the complete LDC Hub4 BN distributions (both for acoustic and language models). Depending on when the experiments were conducted different test data sets are used in the results reported in the paper. However, they are all drawn from either the 1996, 1997 or 1998 Hub4 evaluation data sets (or subsets thereof). Results are reported using the classification of the speech data along the so-called F-conditions (Pallet, 1997): prepared speech (F0), spontaneous speech (F1), low fidelity speech, including telephone channel speech (F2), speech in the presence of background music (F3), speech in the presence of background noise (F4), speech from non-native speakers (F5) and FX – all other speech.

2. Automatic segmentation and clustering

The key idea is to view segmentation and clustering as model selection problems. This allows us

to use the popular BIC model selection criterion for segmentation and clustering. In segmentation the goal is to detect changes in speaker identity, environmental condition and channel condition without knowing the class of changes a priori. The paper proposes a maximum likelihood approach to detect turns of a Gaussian process; the decision of a turn is based on the BIC, a model selection criterion. The paper also proposes applying BIC as a termination criterion in the hierarchical methods for speaker clustering. Experimental results indicate that the segmentation algorithm successfully detects acoustic changes and the clustering algorithm can produce clusters with high purity and enhance un-supervised adaptation as much as the ideal clustering by the true speaker identities.

2.1. Model selection

The problem of model selection is to choose one among a set of candidate parametric models to describe a given data set. Models with too few parameters will not be accurate (under-fitting), while those with too many parameters will not generalize (over-fitting). There are two popular methods for model selection: cross-validation (non-parametric) and BIC (parametric) (Schwarz, 1978).

Consider modeling the data $\mathcal{X} = \{x_i : i = 1, \dots, N\}$ using one model from a set of parametric models $\mathcal{M} = \{M_i : i = 1, \dots, K\}$. For any given model M assume that the $\#(M)$ parameters are chosen to maximize the likelihood and let $L(\mathcal{X}, M)$ denote this maximum value. BIC is a likelihood criterion penalized by the number of parameters in the model and is defined as

$$\text{BIC}(M) = \log L(\mathcal{X}, M) - \lambda \frac{1}{2} \#(M) \log(N), \quad (1)$$

where the penalty weight $\lambda = 1$. One chooses the model M with maximal BIC value. This procedure can be shown to be a large-sample version of the Bayesian procedure for the case of independent, identically distributed observations and linear models (Schwarz, 1978). By varying the penalty weight $\lambda > 0$, one can tradeoff the relative importance of likelihood and model complexity, al-

though only $\lambda = 1$ corresponds to the strict definition of BIC.

2.2. Change detection via BIC

In this section, we describe a maximum likelihood (ML) approach for acoustic change detection based on BIC. The feature sequence $\mathbf{x} = \{x_i \in \mathcal{R}^d, i = 1, \dots, N\}$ extracted from the continuous audio stream (e.g., mel-cepstral features) is assumed to be drawn from independent multivariate Gaussian process: $x_i \sim N(\mu_i, \Sigma_i)$, where μ_i and Σ_i are the mean and covariance, respectively.

2.2.1. Detecting one change point

Consider a simplified problem with at most one change point. A change at time i is resolved with the following hypothesis test: $H_0 : x_1 \cdots x_N \sim N(\mu, \Sigma)$ versus $H_1 : x_1 \cdots x_i \sim N(\mu_1, \Sigma_1); x_{i+1} \cdots x_N \sim N(\mu_2, \Sigma_2)$. If Σ, Σ_1 and Σ_2 are the sample covariances from all the data, from $\{x_1, \dots, x_i\}$ and from $\{x_{i+1}, \dots, x_N\}$, respectively, then, the ML ratio statistics are

$$R(i) = N \log |\Sigma| - N_1 \log |\Sigma_1| - N_2 \log |\Sigma_2| \quad (2)$$

and the ML estimate of the changing point is $\hat{i} = \arg \max_i R(i)$.

The above test can be viewed as a model selection problem. The data is modeled as one Gaussian or two Gaussians. The difference in BIC values of these two models is

$$\text{BIC}(i) = R(i) - \lambda P, \quad (3)$$

where the penalty $P = \frac{1}{2}(d + \frac{1}{2}d(d + 1)) \log N$, and the penalty weight $\lambda = 1$; d is the dimension of the space. If (3) is positive, the model with two Gaussians is favored. A change has occurred if $\{\max_i \text{BIC}(i)\} > 0$. Clearly the ML changing point can also be expressed as $\hat{i} = \arg \max_i \text{BIC}(i)$.

Compared with the metric-based segmentation methods (Siegler et al., 1997; Beigi and Maes, 1998) the proposed scheme is robust. The main reason appears to be that these schemes use the distance between *fixed-size* windows to the left and right of proposed change point which use a limited number of samples. In contrast, the BIC criterion uses all the samples on either side of the proposed change point, and is hence robust. Experimental

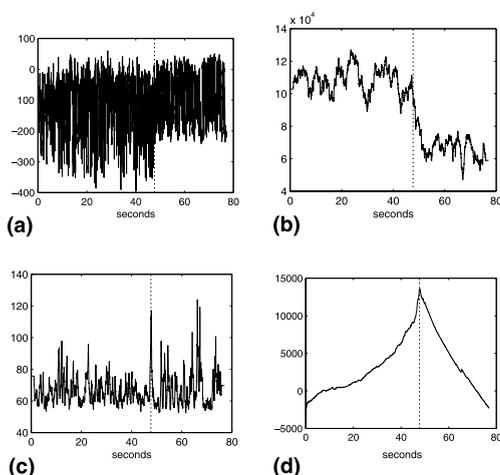


Fig. 1. Detecting one changing point. (a) First cepstral dimension, (b) log likelihood distance, (c) KL2 distance, (d) BIC criterion.

results on 77 s of speech from two speakers is shown in Fig. 1. Panel (a) plots the first dimension of the cepstral vectors; the dotted line indicates the location of the change. Panels (b) and (c) show the log likelihood ratio and symmetric Kullback–Liebler distances, respectively, computed with 1 s windows (100 frames) on either side of a hypothesized change point. In both cases a local maximum is attained at the location of the change. However, there are several other maxima which do not correspond to any change points. Panel (d) displays the BIC criterion that clearly predicts the change point.

2.2.2. Detecting multiple changing points

The proposed single change point detection scheme can be generalized to detecting multiple change points using the following algorithm:

- (1) Initialize the interval $[a, b]$: $a = 1; b = 2$;
- (2) Detect if there is one changing point in $[a, b]$ via BIC;
- (3) if (no change in $[a, b]$)
 - let $b = b + 1$;
 - else
 - let \hat{t} be the changing point detected;
 - set $a = \hat{t} + 1; b = a + 1$;
 - end

(4) go to (2).

By expanding the window $[a, b]$, the final decision is made based on as many samples as possible. Finally note that the BIC scheme while inherently threshold-free can also be viewed as a dynamic thresholding scheme on the log likelihood distance $(\lambda \frac{1}{2}(d + \frac{1}{2}d(d + 1)) \log N)$. The accuracy of this procedure depends on the detectability of the true change points. If $\mathcal{T} = \{t_i\}$ are the true change points, the detectability is defined as $D(t_i) = \min(t_i - t_{i-1} + 1, t_{i+1} - t_i + 1)$. Low detectability leads to missed change points which can contaminate the next Gaussian model and thus affect the detection of the next change point. While in principle this algorithm has a quadratic complexity, one can dramatically reduce the cost by performing a crude search with further refinement. An efficient implementation is described in (Tritschler and Gopinath, 1999).

2.2.3. Change point detection on 1997 Hub4 evaluation data

Table 1 compares the results of applying this segmentation procedure on the 1997 Hub4 evaluation data (3 h long) to the hand-segmentation provided by NIST. The scheme detected 462 changes of which 19 (4.1%) were spurious (i.e., occurred within a speaker turn). About 20 (4.3%) detected points had biases as shown in panel (a) in Fig. 2. However, these biases are small (< 1 s).

The NIST segmentation has 620 changes of which 207 (33.4%) are missed. Of these 154 (25.0%) were short turns (< 2 s). Panel (b) shows the histogram of the detectability of the true changes; there were 223 true changes with detectability less than 2 s. Panel (c) shows the histogram of the detectability of the true changes which were missed in the detection; it is clear that most of the errors came from low detectabilities. Panel (d) describes the Type-II error rates according to different degrees of

Table 1
Change detection error rates

Type-I Error	4.1%		
Type-II Error	33.4%	≤ 2 s	25.0%
		> 2 s	8.4%

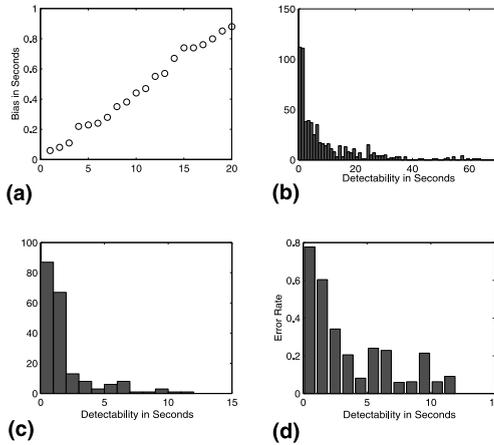


Fig. 2. Error analysis of change detection. (a) Biases, (b) histogram: all true changes, (c) histogram: missed true changes, (d) type-II errors analysis.

Table 2
Segmentation and clustering in Hub4 1997 task

	Error rate (%)
NIST hand-segmentation	19.8
IBM segmentation	20.4

detectability: when detectability is below 1 s, as the Type-II error rate is 78%, most such changing points were missed; as the detectability increases, the Type-II error drops. Table 2 compares recognition performance of our segmentation with manual segmentation provided by NIST.

2.3. Clustering via BIC

Clustering can also be viewed as model selection. We have a set of segments with associated data samples. Let $\mathcal{C}_k = \{c_i : i = 1, \dots, k\}$ be a clustering with k clusters, each c_i having n_i samples. Modeling each c_i as $N(\mu_i, \Sigma_i)$, the BIC value of \mathcal{C}_k is

$$\text{BIC}(\mathcal{C}_k) = \sum_{i=1}^k \left\{ -\frac{1}{2} n_i \log |\Sigma_i| \right\} - \lambda P, \quad (4)$$

where $P = \frac{1}{2}(d + \frac{1}{2}d(d+1)) \log N$ and $\lambda = 1$. One chooses the clustering with maximal BIC value.

2.3.1. Hierarchical clustering via greedy BIC

Finding the global best clustering using BIC as above is expensive. However, in hierarchical clustering BIC can be optimized greedily. Consider binary bottom-up clustering with some distance metric. Let $\mathcal{C} = \{c_1, \dots, c_k\}$ and $\mathcal{C}' = \{c, c_3, \dots, c_k\}$ be two clusterings, where, in the latter clusters c_1 and c_2 are merged into, say, cluster c . Modeling each c_i as $N(\mu_i, \Sigma_i)$ and c as $N(\mu, \Sigma)$ the increase in BIC value by going from \mathcal{C} to \mathcal{C}' is

$$\text{BIC} = n \log |\Sigma| - n_1 \log |\Sigma_1| - n_2 \log |\Sigma_2| - \lambda P, \quad (5)$$

where $n = n_1 + n_2$, $P = \frac{1}{2}(d + \frac{1}{2}d(d+1)) \log N$ and $\lambda = 1$. Two nodes are not to be merged if (5) is negative. Since the BIC value is increased at each merge, we are searching for an “optimal” clustering tree by optimizing the BIC criterion in a greedy fashion. Note that BIC is used as a termination criterion. The distance measure in the bottom-up process could be BIC or any other distance measure. Clearly this approach also works for top-down clustering.

2.3.2. Speaker clustering on the Hub4 1996 evaluation data

We experimented with this algorithm on the 1996 Hub4 evaluation data (Bakis et al., 1997a,b) that had 824 segments from 28 speakers. Bottom-up clustering with log likelihood ratio distance and maximum linkage is used with the BIC termination criterion (5). Our algorithm produces 31 clusters with high purity. Recognition results with un-supervised speaker adaptation using MLLR on these clusters is shown in Table 3. There is negligible degradation in performance relative to ideal speaker clustering.

Table 3
MLLR adaptation enhanced by BIC clustering

	Prepared	Spontaneous
Baseline	18.8%	27.0%
MLLR w/o clustering	18.7%	26.9%
MLLR w/ideal clustering	17.5%	24.8%
MLLR w/BIC clustering	17.5%	24.6%

3. Acoustic modeling

3.1. Bayesian information criterion for Gaussian mixtures

HMM states in a recognizer are typically modeled using Gaussian mixtures. The number of Gaussians per state is either chosen based on the total number of Gaussians the system can support and/or number of samples in training associated with a particular state. Since too few Gaussians leads to under-fitting and too many to over-fitting, BIC is a natural choice for determining the number of Gaussians for each state. For each state and integer n we have a model with n Gaussians and associated BIC value. One chooses the model with maximal BIC value.

Fig. 3 illustrates how this procedure works for a particular HMM state. The horizontal axis represents the number of Gaussians. The vertical axis represents the log-likelihood in Panel (a) and the

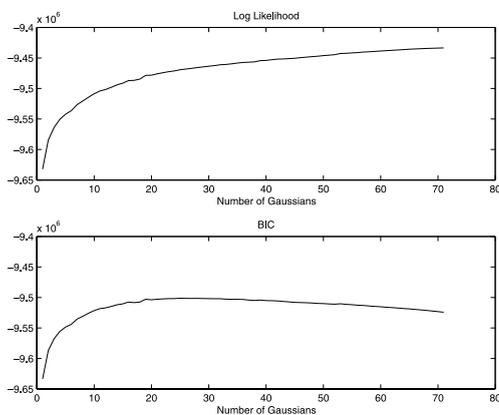


Fig. 3. Choosing the number of Gaussians by maximizing the BIC criterion.

BIC value in Panel (b). As the number of Gaussians increases, the likelihood always improves, whereas the BIC value first increases and then decreases. The optimal BIC value attained at $n = 27$. The total number of Gaussians in the system (for all the states) is automatically determined with the above procedure. A modified BIC scheme with penalty weight $\lambda \neq 1$ can be used to give models different number of Gaussians for all the states; smaller values of λ lead to larger models.

To compare above procedure to heuristic thresholding a collection of acoustic models of varying (total) sizes ranging from 90 to 289 K Gaussians were built using different values of λ . A 90 K Gaussian model was also built using heuristic thresholding. Recognition results on a subset of the 1997 Hub4 evaluation test set with manual segmentation from NIST is shown in Table 4. For the same total number of Gaussians, the BIC procedure gives a better model. Moreover, BIC gives a sequence of better and larger models for decreasing values of λ . In our experiments, compared to thresholding, the BIC procedure tends to favor more Gaussians for complex sounds (vowels) and favor less Gaussians for simple sounds (fricatives).

3.2. Semi-tied covariance matrices

Modeling highly correlated features accurately with Gaussian distributions requires full covariances. This requires a large number of parameters and thus restricts the number of Gaussian distributions that can be robustly estimated from a given data-set. This section describes an alternative to full covariance matrices where the covariance

Table 4
Comparison of the BIC approach with the thresholding approach on the 1997 evaluation subset

	# Gaussians (K)	All	F0	F1	F2	F3	F4	F5	FX
Standard	90	26.0	11.9	23.5	31.7	28.4	28.5	22.3	42.3
$\lambda = 1.00$	90	25.2	11.6	23.1	30.5	27.7	26.2	20.5	41.8
$\lambda = 0.80$	135	24.7	11.2	21.2	29.5	29.0	26.8	21.6	41.2
$\lambda = 0.65$	178	24.2	10.7	21.5	29.3	26.5	25.9	21.4	40.3
$\lambda = 0.54$	237	23.8	10.7	21.6	29.3	26.5	24.2	19.7	39.6
$\lambda = 0.45$	289	23.5	10.5	21.5	28.9	24.4	24.6	20.7	39.0

matrix is now split into two elements, one full and one diagonal, which may be tied at separate levels. The covariance matrix is of the form

$$\Sigma^{(m)} = \mathbf{H}^{(r)\text{T}} \Sigma_{\text{diag}}^{(m)} \mathbf{H}^{(r)}, \quad (6)$$

where r is the *semi-tied class* of Gaussian component m . Though the use of multiple semi-tied classes has been shown to give some gains on large vocabulary tasks, it was decided to restrict the number of classes to one. This greatly simplifies adaptation of these models (though semi-tied transforms may be adapted as described in (Gales, 1999))

Semi-tied covariances can also be viewed as class-dependent linear feature space transformations and this viewpoint leads to an efficient implementation of semi-tied covariances (Gales, 1999; Gopinath, 1998). In fact, if there is a single semi-tied class, then it reduces to standard diagonal Gaussian models in the transformed space. Thus if the original feature-space (in our case determined by LDA (Bahl et al., 1995)) was \mathbf{o} , then the effective space in which models are built is

$$\hat{\mathbf{o}} = \mathbf{H}^{-1} \mathbf{o} = \mathbf{A} \mathbf{o}. \quad (7)$$

The task is now to estimate \mathbf{A} in an ML fashion.

3.2.1. Training semi-tied covariance matrices

A couple of schemes have been used to obtain the semi-tied transform and related model parameters. The one used in this work is to maximize the following expression with respect to \mathbf{A} .

$$\begin{aligned} \mathcal{Q}(\mathcal{M}, \hat{\mathcal{M}}) \\ = \sum_{\tau, m} \gamma_m(\tau) \log \left(\frac{|\mathbf{A}|^2}{|\text{diag}(\mathbf{A} \mathbf{W}^{(m)} \mathbf{A}^{\text{T}})|} \right), \end{aligned} \quad (8)$$

where

$$\mathbf{W}^{(m)} = \frac{\sum_{\tau} \gamma_m(\tau) (\mathbf{o}(\tau) - \mu^{(m)}) (\mathbf{o}(\tau) - \mu^{(m)})^{\text{T}}}{\sum_{\tau} \gamma_s(\tau)}, \quad (9)$$

$$\mu^{(m)} = \frac{\sum_{\tau} \gamma_m(\tau) \mathbf{o}(\tau)}{\sum_{\tau} \gamma_m(\tau)}. \quad (10)$$

For very large numbers of Gaussian components it is not possible to store separate within class co-

Table 5

Semi-tied covariances: % WER on the 1996 Hub4 evaluation data: (a) baseline, (b) single semi-tied class, (c) four semi-tied classes

Expt	F0 (planned)	F1 (spontaneous)
Baseline	21.1	29.1
1 Semi-tied class	19.3	28.4
4 Semi-tied classes	19.4	29.0

variance matrices for all Gaussian components.¹ For this reason, covariance matrices were stored at the state level, rather than the Gaussian component level. This was found to have very little effect on recognition performance. Having optimized this with respect to \mathbf{A} , the mean and the variance of the Gaussian components associated with the HMMs are then found using

$$\hat{\mu}^{(m)} = \frac{\sum_{\tau} \gamma_s(\tau) \mathbf{A} \mathbf{o}(\tau)}{\sum_{\tau} \gamma_m(\tau)} \quad (11)$$

and

$$\hat{\Sigma}^{(m)} = \text{diag}(\mathbf{A} \mathbf{W}^{(m)} \mathbf{A}^{\text{T}}). \quad (12)$$

After obtaining the new feature-space transformation the HMM parameters are updated using Baum–Welch re-estimation. Table 5 shows experimental results on the planned speech (F0) and spontaneous speech (F1) portions of the 1996 DARPA Hub4 evaluation test. The first experiment used a single semi-tied class, while the second used four classes; one each for the HMM states of the following sounds: (a) stop-consonants and flaps, (b) fricatives, (c) vowels and diphthongs, and (d) nasals, glides and silence.

3.3. Pronunciation networks

Typically in a large vocabulary speech recognition system a word is represented as a sequence of its constituent phonemes as defined by a hand-generated lexicon, or as several alternative se-

¹ An alternative optimization does allow training of these very large systems.

quences if multiple pronunciations of the word are allowed.

A major factor in the performance of a speech recognizer is the consistency between the way the speaker pronounces the words and the way that the pronunciation is specified in the dictionary; when there is a mismatch recognition errors are likely to occur.

One approach to bridging the gap between lexicon and pronounced speech is to alter the dictionary to reflect common pronunciation variations. Alternatively, one may choose to leave the dictionary unchanged and reflect the pronunciation deviations in the acoustic model topology. In this paper we report on an implementation of the latter, an automatic method for discovering an appropriate model for each context-dependent phoneme which allows for such phenomena as reduced pronunciations and substituted phonemes where warranted by observation on training data.

3.3.1. *Generating observations of strings*

For the purposes of discussion, the processing required to build context-dependent pronunciation networks may be divided into two parts: a sequence of “pre-processing” steps resulting in sets of observations from which to train the pronunciation networks and the actual training of those networks. The material described in this section constitutes the preprocessing steps; the discussion in the next section describes the building of the pronunciation networks given the observations.

The observations from which the networks are built are sequences of label strings; each observed label sequence corresponds to the output of a constrained phoneme recognizer for the portion of the waveform aligned to a given phoneme by a Viterbi alignment.

The first step towards the generation of the sets of observation label strings from which each pronunciation network is built, then, is to perform a Viterbi alignment of the text to the speech waveform for each of the training utterances. The canonical pronunciation of each word as defined by the lexicon, with the usual three-state, left-to-right model is used in this computation. The alignment step provides a labeled segmentation of each utterance into phonetic regions, with the labels cor-

responding to the states in the model for each phoneme.

The second phase of preprocessing is to perform some type of phoneme recognition on the training data. We have chosen to use the same alphabet as in the Viterbi alignment, i.e., thirds-of-phones. Our phoneme recognizer consists of the same acoustic models as the baseline speech recognition system along with a set of transition probabilities among context-dependent phones. Finding the most-likely sequence of phones yields quite clean label sequences. The resulting label sequences are segmented according to the subphonetic boundaries calculated in the Viterbi alignment and pooled according to the state label associated with the alignment.

The third and final phase of preprocessing is to partition the set of observation sequences for each phoneme into context-dependent units. For this we use a decision tree which asks questions about phonetic context and measures the goodness of a split in terms of the reduction in entropy in the distribution over labels. Each leaf in this tree represents a context-dependent unit for which we will build a pronunciation network using the technique outlined in the previous section. Dropping the training data down the tree yields a set of label sequences for each context-dependent unit from which to build the network which will characterize it.

Note that this tree is distinct from that which defines the context-dependent units for which acoustic models are built as will be discussed in the following section.

3.3.2. *Pronunciation models given observed strings*

The discussion in this section assumes that we have a set of observations in the form of sequences of phonemic labels and that the observations have been partitioned based on phonetic context. These preprocessing steps were described in the previous section. Having obtained the observations, we would like to discover from them the observations a network of states which represents well the collection of observed label sequences for each context; the procedure by which we build a network to represent the observations for each context is the subject of this section.

This procedure is similar in spirit to that described in (Stolcke and Omohundro, 1994) but differs from their implementation in order to reduce the required computations.

For each context, we first discard all sequences containing a symbol which has not been observed at least N times in the pool of training sequences for this context, where N is a hand-set threshold. We then divide the remaining training observations into two sets, a “development” set for building the initial networks and a second “held-out” set for reducing the number of states in the network.

Next, we build a large initial network which can explain all the sequences in the development training set. We have investigated initializing this network as null and as the default three state left-to-right network and have found the latter to provide slightly better recognition results. Against the starting point of null or the default network, we check the first observation to see if it may be explained by the existing network. If so, we update transition probabilities only. If the observation does not align to the existing network, we add a parallel path consisting of one state each time the phoneme in the sequence differs from its left neighbor, with transition probabilities derived from the number of repetitions of the phonemes within the observation. For each remaining observation in the development training pool we repeat the procedure, checking whether it may be explained by the existing network; if so we update transition probabilities and if not we add a parallel path to the network. After all observations have been incorporated into the network, we prune all branches with transition probability into the branch less than ϵ , where ϵ is a hand-determined threshold.

Having built the initial network, we begin collapsing states where such merges are favorable on our held-out training data. Each state in the network is labeled from the alphabet of subphonemic units, e.g. AA_1 . Only states which carry the same label are considered for merge. The probability distribution over labels for each state is concentrated entirely on the associated label; the likelihood of the held-out data is computed before and after a given merge. If the likelihood increases, or

if the number of observations which can be modeled by the resulting network is larger than was the case prior to the merge, the merge is retained; otherwise the two states are kept distinct. Because states are collapsed whenever an advantageous merge is found, the order in which states are evaluated for merging impacts the final network. We have somewhat arbitrarily started with the first two states having the same label and hold the first state fixed, sequentially evaluating the second state until a good merge is found. Once a merge is retained, we reset the starting state to be the next state in the network from the current starting state and iterate until no more good merges exist or until there are fewer than S states in the network, where S is a threshold set by hand.

Merging of two states consists of creating a new state whose parents are the union of the parents of the two states and whose children are itself plus the children of the two states excluding those states themselves (self-loops) and deleting the two unmerged states and their incoming and outgoing arcs.

Finally, two iterations of the EM algorithm are run to estimate transition probabilities for the network.

3.3.3. Using pronunciation networks in recognition

The error rates resulting from using the pronunciation networks to define the context-dependent phoneme model topologies are compared with the errors resulting from a baseline three-state, left-to-right model topology which has no skips. We report the error rates on both the Wall Street Journal and BN databases. On WSJ, the baseline of 8.0% error fell to 7.4% when the pronunciation network topologies were used in lieu of the default networks. On the BN data the error rate was again reduced by using the pronunciation networks, as detailed in Table 6. Appealingly, the largest reduction in error occurs in the case of spontaneous speech, a condition where we would expect that pronunciations variations from the lexical representation might be larger than what would be observed in the case of planned speech, and therefore a condition for which this type of modeling would be potentially valuable.

Table 6
Error rates on Broadcast News data

	F0	F1	F2	F3	F4	F5	FX
A	9.1	20.8	28.0	25.1	24.4	19.6	37.1
B	8.9	20.1	27.8	25.0	24.4	19.5	37.4

Condition A is the baseline. Condition B uses pronunciation networks. F0 = Clean/Planned Speech. F1 = Clean Spontaneous Speech. F2 = Speech on Telephone Channels. F3 = Speech with Background Music. F4 = Speech in Degraded Acoustics. F5 = Non-native speakers. FX = Combinations of F1–F5.

3.4. Tail distribution modeling

An examination of the histograms of any single dimension of the acoustic data assigned to a particular HMM state shows that the distributions are (a) skew-symmetric, (b) peakier than typical Gaussians and (c) have tails that taper off at a slower rate than for a Gaussian tail. Despite these limitations, mixtures of Gaussians perform well in speech recognition experiments. Addressing these issues may lead to improved performance. An option (that increases the number of parameters dramatically is) to increase the number of Gaussian components in the mixtures. This section describes two efficient alternatives that address (b) and (c) – Richter distributions and power-exponential distributions. A Richter distribution (Brown, 1987; Richter, 1986) is a mixture of Gaussians with the same mean and covariances that are multiples of each other. A power exponential distribution is obtained by raising the exponent in a Gaussian distribution to a power that is possibly different from that of the Gaussian. For large powers they are similar to a uniform distribution whereas for small powers they have sharp peaks and heavy tails (Basu et al., 1999).

3.4.1. Richter distributions

Richter distributions are described by

$$f(\mathbf{o}; \mu, \Sigma, p(v)) = \int \mathcal{N}(\mathbf{o}; \mu, v^2 \Sigma) p(v) dv, \quad (13)$$

where $p(v)$ is a probability density function and includes Gaussian and Cauchy distributions as special cases. For example, $p(v) = \delta_1(v)$ (where $\delta_{v_r}(v)$ is the Kronecker delta function) corresponds to a standard Gaussian distribution. By appropriately choosing $p(v)$ the tails and the peakness of the distribution can be controlled. An EM scheme

is described for ML estimates of μ and Σ that does not require explicitly obtaining the distribution $p(v)$ is described in (Liporace, 1982). A discrete version of (13) is described in (Richter, 1986), where $p(v) = \sum_r w_r \delta_{v_r}(v)$ with $w_r > 0$, $\sum_r w_r = 1$. The paper also gives formulae for ML estimates of μ , Σ and w_r . This form of distribution was used in (Brown, 1987) for discrete speech modeling, though in the experiments described the discrete distribution of v was determined a priori rather than trained from the data.

This paper proposes using Richter mixtures for modeling HMM states in LVCSR systems.

$$\mathcal{L}(\mathbf{o}) = \sum_m \sum_r w_r^{(m)} \mathcal{N}(\mathbf{o}; \mu^{(m)}, v_r^{(m)2} \Sigma^{(m)}). \quad (14)$$

It is possible to tie the Richter distribution parameters $w^{(m)}$ and $v^{(m)}$ over many Richter distributions. In our experiments globally tied Richter distribution parameters were obtained using the Hub4 training data. The tails of the Richter distribution are longer than those of the Gaussian distribution, indicating that, in a likelihood sense, the Gaussian components are sub-optimal.

The following re-estimation formulae (modified versions of those in Brown, 1987) are used:

$$\hat{\mu}^{(m)} = \frac{\sum_{r,\tau} \frac{\gamma_r^{(m)}(\tau)}{\hat{v}_r^{(pm)2}} \mathbf{o}(\tau)}{\sum_{r,\tau} \frac{\gamma_r^{(m)}(\tau)}{\hat{v}_r^{(pm)2}}}, \quad (15)$$

$$\hat{\Sigma}^{(m)} = \frac{\sum_{r,\tau} \frac{\gamma_r^{(m)}(\tau)}{\hat{v}_r^{(pm)2}} \hat{\mathbf{W}}^{(m)}(\tau)}{\sum_{r,\tau} \gamma_r^{(m)}(\tau)},$$

$$\hat{v}_r^{(p)2} = \frac{\sum_{M^{(p)},\tau} \gamma_r^{(m)}(\tau) \hat{q}^{(m)}(\tau)}{d \sum_{M^{(p)},\tau} \gamma_r^{(m)}(\tau)}, \quad (16)$$

$$\hat{w}_r^{(m)} = \hat{c}^{(m)} \frac{\sum_{\tau} \gamma_r^{(m)}(\tau)}{\sum_{M^{(p)},r,\tau} \gamma_r^{(m)}(\tau)},$$

where

$$\hat{c}^{(m)} = \frac{\sum_{r,\tau} \gamma_r^{(m)}(\tau)}{\sum_{S^{(m)},r,\tau} \gamma_r^{(m)}(\tau)},$$

$$\hat{q}^{(m)}(\tau) = \left(\mathbf{o}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right)^T \hat{\boldsymbol{\Sigma}}^{(m)-1} \left(\mathbf{o}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right)$$

and

$$\hat{\mathbf{W}}^{(m)}(\tau) = \left(\mathbf{o}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right) \left(\mathbf{o}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right)^T. \quad (17)$$

$M^{(p)}$ is the set of components sharing the same Richter parameters, p_m is the Richter class of component m , d is the dimensionality of the observation vector $\mathbf{o}(\tau)$ and $\gamma_r^{(m)}(\tau)$ is the posterior probability of being in Richter component r of component m at time τ and $S^{(m)}$ is the set of components in the same state as m . Formulae (16) and (15) yield an iterative estimation scheme since the mean and the variance are a function of $\hat{\mathbf{v}}^{(r)}$, which itself is a function of the estimates of the mean and variance. The sufficient statistics for this operation are the occupancy, sum and sum squared of the feature vector for each Richter distribution of each component. Thus if there are M components and R Richter distributions per component, the equivalent of $M \times R$ components must be stored. An alternative to this and the one used in this section is to either update the Richter distribution parameters or the means and variances. In this case it is only necessary to store parameters at the Richter tying level or the component level.

One reason for using Richter distributions rather than additional Gaussian components is the efficiency of the likelihood calculation. Indeed,

$$\mathcal{L}(\mathbf{o}(\tau); m) = \sum_{m,r} b_r^{(m)} \exp \left(- \frac{q^{(m)}(\tau)}{2v_r^{(m)2}} \right),$$

where $q^{(m)}(\tau)$ is a function of the component, m , and observation

$$q^{(m)}(\tau) = \left(\mathbf{o}(\tau) - \boldsymbol{\mu}^{(m)} \right)^T \boldsymbol{\Sigma}^{(m)-1} \left(\mathbf{o}(\tau) - \boldsymbol{\mu}^{(m)} \right) \quad (18)$$

and $b_r^{(m)}$ is a function of the Richter component r , but independent of the observation

$$b_r^{(m)} = \frac{1}{\sqrt{2^d \pi^d |\det \boldsymbol{\Sigma}^{(m)}|}} \frac{w_r^{(m)}}{\sqrt{v_r^{(m)2d}}}.$$

The additional cost is in the log-add over the Richter components. This may be ignored if a max of the components is taken, rather than the sum.

It is also common to use linear transformations to adapt model parameters to be more representative of a particular speaker, or acoustic environment. A variety of linear transformations and re-estimation formulae are described in (Gales, 1998a). Modifying these formulae to handle Richter distributions is trivial. The main modification is to deal with $\gamma_r^{(m)}(\tau)/v_r^{(m)2}$ rather than the standard posterior component probability. As an example the estimation formulae for the transform $\hat{\mathbf{A}}$ in ML linear regression, where Richter components are used, is considered. Here

$$\hat{\boldsymbol{\mu}}^{(m)} = \mathbf{A} \boldsymbol{\mu}^{(m)} + \mathbf{b} = \mathbf{W} \boldsymbol{\zeta}^{(m)}, \quad (19)$$

where $\boldsymbol{\zeta}^{(m)}$ is the extended mean vector $[\boldsymbol{\mu}^{(m)T} \ 1]^T$. The estimation formula for row i of the transformation matrix is

$$\hat{w}_i = \mathbf{k}^{(i)} \mathbf{G}^{(i)-1}, \quad (20)$$

where

$$\mathbf{G}^{(r)} = \sum_{M,\tau,r} \frac{\gamma_r^{(m)}(\tau)}{v_r^{(m)2} \sigma_i^{(m)2}} \boldsymbol{\zeta}^{(m)} \boldsymbol{\zeta}^{(m)T}$$

and

$$\mathbf{k}^{(i)} = \sum_{M,\tau,r} \frac{\gamma_r^{(m)}(\tau)}{v_r^{(m)2} \sigma_i^{(m)2}} \boldsymbol{\zeta}^{(m)T} \mathbf{o}_i(\tau). \quad (21)$$

Similarly modifications to the variance adaptation formulae are possible.

3.4.2. Power exponentials

Consider the class of densities

$$f(\mathbf{o}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \alpha) = \rho_\alpha |\det \boldsymbol{\Sigma}|^{-1/2} \exp \left(- (\gamma_\alpha q)^\alpha \right), \quad (22)$$

where

$$q = (\mathbf{o} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{o} - \boldsymbol{\mu}), \quad (23)$$

$$\gamma_\alpha = \frac{\Gamma(3/\alpha)}{\Gamma(1/\alpha)} \quad \text{and} \quad \rho_\alpha = \frac{\alpha\Gamma^{1/2}(3/\alpha)}{2\Gamma^{3/2}(1/\alpha)}. \quad (24)$$

This class was recently suggested and studied in (Basu et al., 1999). The one-dimensional case appears to have first been suggested by Subbotin (1923). The class (22) will be referred to as the power exponential distribution. It is also known as the error function, p-Gaussians or as α -Gaussians.

Following (14) a model is considered where each state in the system is modelled by a mixture of power exponential distribution, i.e.,

$$\mathcal{L}(\mathbf{o}) = \sum_m w^{(m)} f(\mathbf{o}; \mu^{(m)}, \Sigma^{(m)}, \alpha^{(m)}). \quad (25)$$

It is worth noticing that the class of functions described in (22) is not a subset of the class described in (13). Power exponential distributions can not in general be modelled with Richter distributions. This fact can be verified by noticing that functions in the class (13) are all log concave, whereas the power exponentials are not log concave for $0 < \alpha < 1$. This makes the framework of (Liporace, 1982) unsuitable for parameter update for $0 < \alpha < 1$.

The estimation formula for $w^{(m)}$ is identical to the standard HMM re-estimation formulae. Update formulae for $\mu^{(m)}$ and $\Sigma^{(m)}$ are suggested in (Basu et al., 1999):

$$\hat{\mu}^{(m)} = \frac{\sum_\tau \gamma^{(m)}(\tau) (q^{(m)}(\tau))^{\alpha^{(m)}/2-1} \mathbf{o}(\tau)}{\sum_\tau \gamma^{(m)}(\tau) (q^{(m)}(\tau))^{\alpha^{(m)}/2-1}} \quad (26)$$

and

$$\hat{\Sigma}^{(m)} = \frac{\sum_\tau \gamma^{(m)}(\tau) (q^{(m)}(\tau))^{\alpha^{(m)}/2-1} \hat{\mathbf{W}}^{(m)}(\tau)}{\sum_\tau \gamma^{(m)}(\tau)}, \quad (27)$$

where $q^{(m)}(\tau)$ is defined in Eq. (18), $\hat{\mathbf{W}}^{(m)}(\tau)$ in Eq. (17) and $\gamma^{(m)}(\tau)$ is defined to be the posterior probability of being in the power exponential component m at time τ . It is not known that the overall likelihood is guaranteed to increase with the update given by (26) and (27), but numerical evidence suggests that this is true.

Special consideration for $0 < \alpha < 1$ is suggested in (Basu et al., 1999). The powers $\alpha^{(m)}$ can either be

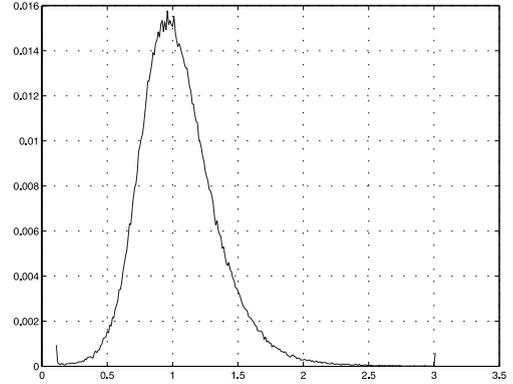


Fig. 4. The distribution of powers, α , after training using (28) on the Hub4 1997 data.

fixed on a global level or they can be updated according to the formula given in (Basu et al., 1999):

$$\hat{\alpha}^{(m)} = \operatorname{argmax}_\alpha \sum_\tau \gamma^{(m)}(\tau) \times \log \left(f(\mathbf{o}(\tau); \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}, \alpha) \right). \quad (28)$$

With this update of $\alpha^{(m)}$ the likelihood is guaranteed to increase. Fig. 4 shows the distribution of α estimated on a per component case. The mean of the values of α is approximately one. It is interesting to note that the Gaussian component equivalent of power exponential components, $\alpha = 2$, occurs infrequently. Again, this indicates that Gaussian components are sub-optimal in a likelihood sense.

Currently adaptation of power exponentials have not been investigated.

3.4.3. Tail-distribution modeling – results

The two forms of modified tail distribution modeling were investigated on the 1998 Hub4 partitioned evaluation test set. The baseline system for the Richter components had a total of about 135,000 components. A 4 distribution Richter component system ($R = 4$) was initialised using the means and variances of the baseline system. The Richter parameters were tied at the state level. Table 7 shows the comparison of a Richter system and the equivalent baseline system. The adaptation scheme used in both was a global mean

Table 7
Results on the Hub4 1997 partitioned evaluation test set

	Error rate (%)		
	F0	F1	Avg
<i>Richter system</i>			
Base	11.6	18.5	18.7
Base + Adapt	10.1	17.0	16.4
Richter	11.3	18.1	18.4
Richter + Adapt	10.1	16.9	16.3
<i>Power exponential system</i>			
Base ($\alpha = 2$)	11.8	22.9	26.1
$\alpha = 1$	11.5	23.0	25.5
EM update for α	11.9	22.6	25.4

and full variance transform described in (Gales, 1998a). This was applied in an unsupervised batch adaptation mode. Using Richter components showed a small gain in performance over the standard Gaussian components. After adaptation the performance of the two systems are indistinguishable.

The experiments using power exponential components used a modified baseline system consisting of approximately 120,000 Gaussians. The test was performed on a subset of the 1997 partitioned evaluation that was used for development (Chen et al., 1999). Finally a smaller language model than for that of the tests with the Richter distribution were used, thus degrading the performance for the spontaneous speech category, F1, and for some of the more difficult conditions, F2–FX. Two power exponential systems were built. The first used a fixed value of $\alpha^{(m)} = 1$ for all components, motivated by Fig. 4. The second system used a per-component value of $\alpha^{(m)}$ obtained using Eq. (28).

Table 7 shows the performance of the various power exponential systems. Again only small reductions in word error rate were observed using the improved tail modeling.

4. Adaptive training

Acoustic models trained on non-homogeneous data have large variance since they model both

inter-speaker and *intra-speaker* variability. Assuming some form of adaptation is used during testing, it is preferable to adapt a canonical model that just represents the intra-speaker variability. This paper addresses two forms of adaptive training. The first is a modified form of the so-called SAT scheme (Anastasakos et al., 1996), where a constrained-model space transform is used instead of MLLR for adaptation. This yields simple estimation formulae for the canonical model. The second form of adaptive training uses a simple form of transform for adaptation with a more complex canonical model.

As with all adaptive training schemes the training procedure is as follows:

1. Partition the training data into “appropriate” groups.
2. Estimate the adaptation transform for each partitioned subset of the training data.
3. Estimate the canonical model given the partition adaptation transforms.
4. Goto (2) unless convergence criteria are satisfied.

Each of the estimation stages is described for both forms of adaptive training considered.

4.1. Modified speaker adaptive training

Modified speaker adaptive training (MSAT) uses a constrained model-space transform. Here

$$\hat{\boldsymbol{\mu}}^{(sm)} = \mathbf{A}^{(r)} \boldsymbol{\mu}_a^{(m)} + \mathbf{b}^{(r)}, \quad (29)$$

$$\hat{\boldsymbol{\Sigma}}^{(m)} = \mathbf{A}^{(r)} \boldsymbol{\Sigma}^{(m)} \mathbf{A}^{(r)T}. \quad (30)$$

Since the majority of speech recognition systems, and the ones considered in this paper, use diagonal covariance matrices, this form of transform may be efficiently implemented as multiple transformations of the feature-space. Thus,

$$\begin{aligned} \mathcal{L}(\mathbf{o}; \boldsymbol{\mu}_a^{(m)}, \boldsymbol{\Sigma}^{(m)}, \mathbf{A}^{(r)}) \\ = \frac{1}{|\mathbf{A}^{(r)}|^2} \mathcal{N}(\mathbf{A}^{(r)} \mathbf{o} + \mathbf{b}^{(r)}; \boldsymbol{\mu}_a^{(m)}, \boldsymbol{\Sigma}^{(m)}), \end{aligned} \quad (31)$$

where $\mathbf{A}^{(r-1)} = \mathbf{A}^{(r)}$ and $\mathbf{b}^{(r)} = -\mathbf{b}^{(r)}$.

4.1.1. Adaptation transform

The estimation of the adaptation transform for a particular partition of the training data² is obtained in a simple iterative fashion.

$$\mathbf{w}_i = (\alpha \mathbf{p}_i + \mathbf{k}^{(i)}) \mathbf{G}^{(i-1)}, \quad (32)$$

where \mathbf{p}_i is the extended cofactor row vector $[0 \ c_{i1} \ \cdots \ c_{in}]$ ($c_{ij} = \text{cof}(\mathbf{A}_{ij})$),

$$\mathbf{G}^{(i)} = \sum_m \frac{1}{\sigma_i^{(m)2}} \sum_{\tau} \gamma_m(\tau) \zeta(\tau) \zeta(\tau)^{\text{T}}, \quad (33)$$

$$\mathbf{k}^{(i)} = \sum_m \frac{1}{\sigma_i^{(m)2}} \mu_{\text{ai}}^{(m)} \sum_{\tau} \gamma_m(\tau) \zeta(\tau)^{\text{T}} \quad (34)$$

and α satisfies a simple quadratic expression given in (Gales, 1998a).

4.1.2. Canonical model

In (Gales, 1998a) the estimation of the canonical model is given by

$$\hat{\boldsymbol{\mu}}_a^{(m)} = \frac{\sum_{s,\tau} \gamma_m^{(s)}(\tau) \hat{\boldsymbol{\theta}}^{(rs)}(\tau)}{\sum_{s,\tau} \gamma_m^{(s)}(\tau)} \quad (35)$$

and

$$\begin{aligned} \hat{\boldsymbol{\Sigma}}^{(m)} &= \frac{\sum_{s,\tau} \gamma_m^{(s)}(\tau) \left(\hat{\boldsymbol{\theta}}^{(rs)}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right) \left(\hat{\boldsymbol{\theta}}^{(rs)}(\tau) - \hat{\boldsymbol{\mu}}^{(m)} \right)^{\text{T}}}{\sum_{s,\tau} \gamma_m^{(s)}(\tau)}, \end{aligned} \quad (36)$$

where

$$\hat{\boldsymbol{\theta}}^{(rs)}(\tau) = \mathbf{A}^{(rs)} \mathbf{o}(\tau) + \mathbf{b}^{(rs)} \quad (37)$$

and $m \in M^{(r)}$. Thus the sufficient statistics required to estimate the canonical model have the same complexity as the standard Gaussian component estimation schemes.

² During testing exactly the same procedure is used for the particular partition of the test data. Of course during testing the canonical model is not updated.

4.2. Cluster adaptive training

In CAT the set of Gaussian component means associated with a particular partition of the training data is given by a linear interpolation of a set of P cluster means (Gales, 1998b). Thus CAT is defined as follows. For a particular Gaussian component, $m \in M_w^{(r_m)}$, the mean for speaker s is given by

$$\boldsymbol{\mu}^{(sm)} = \mathbf{M}^{(m)} \boldsymbol{\lambda}^{(sr_m)}, \quad (38)$$

where $\mathbf{M}^{(m)}$ is the matrix of P cluster mean vectors for component m ,

$$\mathbf{M}^{(m)} = [\boldsymbol{\mu}_c^{(m1)} \ \cdots \ \boldsymbol{\mu}_c^{(mP)}], \quad (39)$$

where $\boldsymbol{\mu}_c^{(mp)}$ is the mean of Gaussian component m associated with cluster p and the cluster weight vector, $\boldsymbol{\lambda}^{(sr_m)}$, for speaker s is given by (assuming that a bias cluster is used)

$$\boldsymbol{\lambda}^{(sr_m)} = [\lambda_1^{(sr_m)} \ \cdots \ \lambda_{P-1}^{(sr_m)} \ 1]^{\text{T}} \quad (40)$$

and r_m is the cluster weight class of Gaussian component m .

4.2.1. Adaptation transform

$$\boldsymbol{\lambda}^{(sr)} = \mathbf{G}_w^{(sr)-1} \mathbf{k}_w^{(sr)}, \quad (41)$$

where

$$\mathbf{G}_w^{(sr)} = \sum_{m \in M_w^{(r)}, \tau} \gamma_m(\tau) \mathbf{M}^{(m)\text{T}} \boldsymbol{\Sigma}^{(m)-1} \mathbf{M}^{(m)}, \quad (42)$$

$$\mathbf{k}_w^{(sr)} = \sum_{m \in M_w^{(r)}} \mathbf{M}^{(m)\text{T}} \boldsymbol{\Sigma}^{(m)-1} \sum_{\tau} \gamma_m(\tau) \mathbf{o}(\tau). \quad (43)$$

4.2.2. Canonical model

The canonical model associated with CAT has a different form to the canonical models used in SAT or MSAT. Here

$$\mathcal{M} = \{ \{ \mathbf{M}^{(1)}, \dots, \mathbf{M}^{(M)} \}, \{ \boldsymbol{\Sigma}^{(1)}, \dots, \boldsymbol{\Sigma}^{(M)} \} \}. \quad (44)$$

These parameters are estimated using

$$\mathbf{M}^{(m)\text{T}} = \mathbf{G}^{(m)-1} \mathbf{K}^{(m)}, \quad (45)$$

$$\Sigma^{(m)} = \text{diag} \left(\frac{\mathbf{L}^{(m)} - \mathbf{M}^{(m)} \mathbf{K}^{(m)}}{\sum_{s,\tau} \gamma_m(\tau)} \right), \quad (46)$$

where

$$\mathbf{G}^{(m)} = \sum_{s,\tau} \gamma_m(\tau) \lambda^{(sr_m)} \lambda^{(sr_m)\text{T}}, \quad (47)$$

$$\mathbf{K}^{(m)} = \sum_{s,\tau} \gamma_m(\tau) \lambda^{(sr_m)} \mathbf{o}(\tau) \mathbf{o}(\tau)^\text{T}, \quad (48)$$

$$\mathbf{L}^{(m)} = \sum_{s,\tau} \gamma_m(\tau) \mathbf{o}(\tau) \mathbf{o}(\tau)^\text{T}. \quad (49)$$

4.3. Adaptive training – results

The standard baseline system for the adaptive training experiments were used. In order to get a valid set of comparisons a set of initial experiments were carried out using this initial model set and the standard linear adaptation schemes. Table 8 shows the results on the Hub4 1997 partitioned evaluation test set. The baseline recognition rate was 18.7% and this was used to give the adaptation hypothesis and initial alignments to obtain the transforms. Using a single MLLR transform, which has 3660 parameters, the error rate was reduced by about 8%. Similar performance was obtained using the constrained model-space transform. Using MLLR in conjunction with a full variance transform an additional 5% reduction in word error rate was achieved. It is worth noting that this transform is a generalization of both MLLR and the constrained model-space transform.

For the adaptive training experiments the partitioning of the training data was as follows. For the training data labeled with both speaker and

focus conditions the partitions were on in terms of these speaker/focus pairings. Where the focus condition was not given, data was simply grouped according to the speaker.

For the MSAT experiments a single full transformation matrix with bias was estimated for all partitions with greater than 1000 frames. For those partitions having less than a 1000 frames a diagonal transformation matrix with bias was estimated. The motivation for this is that the single global semi-tied transform will roughly de-correlate the data for all components. During recognition a single full transform with bias was estimated for each test partition. Table 9 shows the recognition performance of adaptation of MSAT system. Using a single full constrained model-space transform a 12% reduction in word error rate was achieved compared to baseline and a 5% reduction compared to using the same adaptation scheme on the baseline SI model. Rather than using a constrained model-space transform the MLLR plus full variance transform could also be used. This gave an overall 15% reduction over the baseline system. However the gain from adapting an adaptively trained system over adapting the baseline SI system was only 3%. This is expected, since as the complexity, and power, of the transform increases the gain from using adaptive training is reduced.

For the CAT experiments four initial clusters were generated. These were based on the data with speaker and focus condition information and were grouped as *Male/Female* with *Clean/Noisy*. MLLR was used to transform the baseline SI model set to each of the groupings. A single interpolation weight was then computed for each training set partition using the alignments from the baseline SI model set. The canonical model parameters were then updated. During recognition a single set of

Table 8
Results for adapting the standard SI system on the Hub4 1997 partitioned evaluation test set

System	Error rate (%)		
	F0	F1	Avg
Base	11.6	18.5	18.7
Base + MLLR	10.7	18.0	17.2
Base + Constr	10.6	17.0	17.2
Base + MLLR + Cov	10.1	17.0	16.4

Table 9
Results for the MSAT system on the Hub4 1997 partitioned evaluation test set

System	Error rate (%)		
	F0	F1	Avg
MSAT + Constr	10.1	16.3	16.4
MSAT + MLLR + Cov	9.9	16.2	15.9

Table 10
Results for CAT system on the Hub4 1997 partitioned evaluation test set

System	Error rate (%)		
	F0	F1	Avg
CAT	11.0	17.3	17.6
CAT + MLLR + Cov	10.0	16.1	15.7

interpolation weight was calculated for each test partition. Table 10 shows the performance of the CAT system on the Hub4 1997 partitioned evaluation test set. Using only CAT, requiring estimation of only four parameters, a 6% reduction in word error rate was achieved. Though quite small it is only about 2% worse than using MLLR whilst requiring a factor of about 1000 fewer parameters to be estimated. Rather than using CAT alone it is possible to combine it with other more powerful adaptation schemes. Using MLLR plus a full covariance transform in addition to CAT yielded a 16% reduction in word error rate over the baseline and a 4% reduction over using the same transform on the baseline SI model.

5. Language modeling

The LM training corpus for BN provided by LDC has billions of words spanning several years. Cleaning up this data for building LMs requires filters that depend on the source and period the data spans. Moreover, processing time of LMs may depend non-linearly on the length of the corpus (e.g., LMs with classes based on MMI criteria (Brown et al., 1992)). Our approach therefore was to partition the training corpus into groups, build several LMs for each group, and integrate these LMs using a mixture model. The corpus was partitioned into four groups: text from 1996, 1997 and 1998 and the text of the acoustic training data. Three types of LMs – 3- and 4-g language models with deleted interpolation and a maximum entropy class based language model (Jelinek, 1997) – were built on each partition. The mixture weights were obtained based on either word-error-rate or perplexity minimization on a development test set. Clearly the word-error-rate

Table 11
Mixture LM and single LM

Lang. model	Set1	Set2
6-LM tun	18.8	15.7
8-LM tun	18.0	14.9
8-LM perp	18.0	14.9
14-LM perp	17.8	14.7
4g-LM 96	20.2	17.3

minimization technique is more expensive and was done in a greedy fashion by starting with two components and gradually increasing the number of LM components; however, both techniques lead to LMs with similar performance on test data. The results with several mixtures sizes based on these two types of optimization (WER and perplexity) are shown in Table 11, where, for comparison, results with the single best LM component – a 4-g LM built on the 1996 training text is also shown.

6. Some additional techniques

6.1. Rover

Recently a scheme that combines hypothesis transcripts from several recognizers was shown to greatly improve the accuracy (Fiscus, 1997). The scheme – *recognition output voting error reduction* or ROVER – aligns a set of hypothesis transcripts using dynamic programming and picks the most frequent word (weighted with confidence scores if available) among words aligned to each other; ties are resolved arbitrarily. If the hypothesis transcripts are ordered (e.g., the “best” transcript first), the induced order can be used to resolve the tie among words.

An interesting fact – that deletions are more probable than substitution by a particular word – can be used to improve ROVER. To see this, consider an enumeration of possible alignments from a best-first ordered set of three hypothesis transcripts: (w_1, w_1, w_1) , (w_1, w_1, w_2) , (w_1, w_2, w_1) , (w_2, w_1, w_1) , (w_1, w_1, \emptyset) , (w_1, \emptyset, w_1) , (\emptyset, w_1, w_1) , $(w_1, \emptyset, \emptyset)$, $(\emptyset, w_1, \emptyset)$, $(\emptyset, \emptyset, w_1)$, (w_1, w_2, \emptyset) , (w_1, \emptyset, w_2) , (\emptyset, w_2, w_1) and (w_1, w_2, w_3) . In all but three of the cases viz., (w_1, w_2, \emptyset) , (w_1, \emptyset, w_2) and (\emptyset, w_2, w_1) , ROVER picks the best alternative. However, in

Table 12

Enumeration of individual systems used in the 1998 Hub4 evaluation and some corresponding rover combinations

System #	Short description	WER I	WER II
1	289 K baseline	15.7%	13.3%
2	Smaller set of phonemes	16.3%	14.7%
3	Telephone models	18.4%	16.6%
4	Left context models	17.5%	14.9%
5	Power exponential densities	18.9%	16.0%
6	Speaker adapted training	15.5%	12.8%
7	Cluster adapted training	15.4%	13.1%
	ROVER (7, 6, 1, 2, 3, 4, 5, \emptyset , \emptyset)	14.5%	12.4%
	ROVER (7, 2, 3, \emptyset)	14.3%	12.5%
	ROVER (7, 2, 3)	14.8%	12.7%
	ROVER (7, 6, 2, 3, \emptyset , \emptyset)	14.4%	12.4%

WER I and WER II refers to the error rates of the two individual test sets in the 1998 Hub4 evaluation.

these three cases, a deletion is to be preferred by the fact above. One way to accomplish this in ROVER is to add a fourth \emptyset script in the alignment process. Depending on the number of hypothesis transcripts one or more \emptyset hypothesis can be added to improve performance. Table 12 shows the performance of several systems used by IBM in the 1998 Hub4 evaluation along with ROVER numbers on subsets of these systems using this idea.

6.2. Fast speaking rate

Fast speech is well known to significantly impair recognition accuracy (Sieger, 1995). This paper proposes two approaches – phone models with skip arcs, and faster-rate signal processing to handle fast speech. In the former, standard 3-state left-to-right phone models are modified allowing a skip from the first to the third state. In the latter, the frame rate is increased to 110 from 100 per second. On test data segments are labeled as fast or normal, with the fast segments transcribed using modified phones at the fast rate.

Speaking rate is determined using the average duration of phones. On the training data using Viterbi alignment this information is obtained for both fast and normal speakers. A subset of phones where there is marked difference in duration was identified (say \mathcal{P}) and this subset included almost all the consonants and some vowels. During testing each segment of speech is aligned to an initial pass transcription. Each phone in \mathcal{P} is labeled fast

or normal by comparing its duration to that in training. Let N^n and N^f denote the number of phones in \mathcal{P} labeled normal and fast, respectively. Then, for an integer s we define the metric M_s such that an entire segment is labeled fast or normal based on $N^n - N^f < s$. The metric M_s , given the difficulty of estimating speaking rate (e.g., Monkowski et al., n.d.; Burshtein, 1995), is motivated solely by the fact that it led to good performance improvements on speech recognition using modified phone models and faster signal processing. Table 13 shows experimental results using this approach with and without un-supervised adaptation (MLLR) for two metrics viz., M_0 and M_2 on Set1 in the 1998 Hub4 evaluation test. Rows in the Table with M_0 correspond to word error rates on that subset of segments of Set1 that were labeled as fast; similarly for M_2 . The last row shows that on the entire Set1 test set about 0.3% absolute improvement is obtained using this technique.

Table 13
Decoding results for M_0 and M_2 speech metrics

Decoding	All	F0	F1
M_0 + Baseline	21.2	13.0	34.8
M_0 + Modified	19.9	12.0	31.1
M_0 + MLLR	20.5	12.5	33.9
M_0 + MLLR + Modified	19.1	11.4	31.1
M_2 + Baseline	23.3	13.0	35.0
M_2 + MD	21.5	12.0	31.2
M_2 + MLLR + Modified	20.9	10.7	30.4
Set1 + Baseline	18.0	8.9	19.9
Set1 + M_2 + Modified	17.7	8.7	19.1

7. Summary

Transcription of BN is an extremely challenging problem. This paper has described several techniques that are useful in building a state-of-the-art speech recognition system for BN. Specifically the paper proposes BIC for segmentation, clustering and choosing the number of Gaussians in a mixture model, pronunciation modeling for conversational speech, semi-tied covariance modeling, tail-distribution modeling, adaptive training (modified SAT and CAT) to handle the heterogeneity in BN data.

Acknowledgements

The authors would like to thank M. Monkowski and M. Franz for many valuable discussions on fast speech problems and language modeling.

References

- Anastasakos, T. et al., 1996. A compact model for speaker-adaptive training. In: Proc. ICSLP-96.
- Bahl, L.R. et al., 1994. Robust methods for using context-dependent features and models in a continuous speech recognizer. In: Proc. ICASSP.
- Bahl, L.R. et al., 1995. Performance of the IBM large vocabulary continuous speech recognition system on the ARPA Wall Street Journal task. In: Proc. ICASSP, pp. 41–44.
- Bakis, R. et al., 1997a. Transcription of BN Shows with the IBM LVCSR system. In: Proc. DARPA Sp. Reco. Workshop.
- Bakis, R. et al., 1997b. Transcription of broadcast news shows with the IBM large vocabulary speech recognition system. In: Proc. Speech Recognition Workshop, pp. 67–72.
- Basu, S. et al., 1999. Power exponential densities for the training and classification of acoustic feature vectors in speech recognition. Research report, T.J. Watson Research Center.
- Beigi, H., Maes, S., 1998. Speaker, channel and environment change detection. In: Proc. World Congress on Automation.
- Brown, P., 1987. The acoustic-modeling problem in automatic speech recognition. Ph.D. Thesis, IBM T.J. Watson Research Center.
- Brown, P.F. et al., 1992. Class-based n -gram models of natural language. *Comput. Linguist.* 18 (4), 467–479.
- Burshtein, D., 1995. Robust parametric modelling of durations in hidden Markov models. ICASSP, 548–551.
- Chen, S.S. et al., 1999. Recent improvements to IBM's speech recognition system for automatic transcription of broadcast news. In: Proc. Broadcast News Transcription and Understanding Workshop.
- Fiscus, J.G., 1997. A post-processing system to yield reduced word error rates: recognizer output voting error reduction (rover). In: Proc. IEEE ASRU Workshop, Santa Barbara, pp. 347–352.
- Gales, M.J.F., 1998a. Maximum likelihood linear transformations for HMM-based speech recognition. *Comput. Speech Language* 12, 75–98.
- Gales, M.J.F., 1998b. Cluster adaptive training for speech recognition. In: Proc. ICSLP, pp. 1783–1786.
- Gales, M.J.F., 1999. Semi-tied covariance matrices for hidden Markov models. *IEEE Trans. Speech Audio Process.* 7, 272–281.
- Gish, H., Schmidt, N., 1994. Text-independent speaker identification. *IEEE Signal Process. Mag.*, 18–21.
- Gopinath, R.A., 1998. Maximum likelihood modeling with Gaussian distributions for classification. In: Proc. ICASSP.
- Jelinek, F., 1997. *Statistical Methods for Speech Recognition*. MIT Press, Cambridge.
- Jin, H., Kubala, F., Schwartz, R., 1997. Automatic speaker clustering. In: Proc. Speech Recognition Workshop, pp. 108–111.
- Kubala, F. et al., 1997. The 1996 BBN Byblos Hub-4 transcription system. In: Proc. Speech Recognition Workshop, pp. 90–93.
- Liporace, L.A., 1982. Maximum likelihood estimation for multivariate observations of Markov sources. *IEEE Trans. Inform. Theory* 28, 729–734.
- Monkowski, M.D. et al., n.d. Context dependent phonetic duration models for decoding conversational speech.
- Pallet, D., 1997. Overview of the 1997 DARPA speech recognition workshop. In: Proc. DARPA Speech Recognition Workshop, 2–5 February, Chantilly, VA.
- Richter, A.G., 1986. Modelling of continuous speech observations. In: *Advances in Speech Processing Conference*. IBM Europe Institute.
- Schwarz, G., 1978. Estimating the dimension of a model. *Ann. Stat.* 6, 461–464.
- Sieger, M.A., 1995. Measuring and compensating for the effects of speech rate in large vocabulary continuous speech recognition. Thesis, Carnegie Mellon University.
- Siegler, M., Jain, U., Ray, B., Stern, R., 1997. Automatic segmentation, classification and clustering of broadcast news audio. In: Proc. Speech Recognition Workshop, pp. 97–99.
- Stolcke, A., Omohundro, S., 1994. Best-first model merging for hidden Markov model induction. International Computer Science Institute Report TR-94-003.
- Subbotin, M., 1923. On the law of frequency of errors. *Mathematicheskii Sbornik* 31, 296–301.
- Tritschler, A., Gopinath, R.A., 1999. An improved segmentation and clustering scheme using the Bayesian Information Criterion. In: Proc. Eurospeech 1999, Budapest, Hungary.
- Woodland, P., Gales, M., Pye, D., Young, S., 1997. The Development of the 1996 HTK broadcast news transcription system. In: Proc. Speech Recognition Workshop, pp. 73–78.