

Likelihood-Maximizing Beamforming for Robust Hands-Free Speech Recognition

Seltzer, M.L.; Raj, B.; Stern, R.M.

TR2004-088 December 2004

Abstract

Speech recognition performance degrades significantly in distant-talking environments, where the speech signals can be severely distorted by additive noise and reverberation. In such environments, the use of microphone arrays has been proposed as a means of improving the quality of captured speech signals. Currently, microphone-array-based speech recognition is performed in two independent stages: array processing and then recognition. Array processing algorithms designed for signal enhancement are applied in order to reduce the distortion in the speech waveform prior to feature extraction and recognition. This approach assumes that improving the quality of the speech waveform will necessarily result in improved recognition performance and ignores the manner in which speech recognition systems operate. In this paper a new approach to microphone-array processing is proposed in which the goal of the array processing is not to generate an enhanced output waveform but rather to generate a sequence of features which maximizes the likelihood of generating the correct hypothesis. In this approach, called Likelihood-Maximizing Beamforming (LIMABEAM), information from the speech recognition system itself is used to optimize a filter-and-sum beamformer. Speech recognition experiments performed in a real distant-talking environment confirm the efficacy of the proposed approach.

IEEE Transactions on Speech and Audio Processing

This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Research Laboratories, Inc.; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Research Laboratories, Inc. All rights reserved.

Likelihood-Maximizing Beamforming for Robust Hands-Free Speech Recognition

Michael L. Seltzer, *Member, IEEE*, Bhiksha Raj, *Member, IEEE*,
and Richard M. Stern, *Member, IEEE*

Abstract

Speech recognition performance degrades significantly in distant-talking environments, where the speech signals can be severely distorted by additive noise and reverberation. In such environments, the use of microphone arrays has been proposed as a means of improving the quality of captured speech signals. Currently, microphone-array-based speech recognition is performed in two independent stages: array processing and then recognition. Array processing algorithms designed for signal enhancement are applied in order to reduce the distortion in the speech waveform prior to feature extraction and recognition. This approach assumes that improving the quality of the speech waveform will necessarily result in improved recognition performance and ignores the manner in which speech recognition systems operate. In this paper a new approach to microphone-array processing is proposed in which the goal of the array processing is not to generate an enhanced output waveform but rather to generate a sequence of features which maximizes the likelihood of generating the correct hypothesis. In this approach, called Likelihood-Maximizing Beamforming (LIMABEAM), information from the speech recognition system itself is used to optimize a filter-and-sum beamformer. Speech recognition experiments performed in a real distant-talking environment confirm the efficacy of the proposed approach.

Index Terms

microphone array processing, robust speech recognition, beamforming, adaptive filtering, distant-talking environments

EDICS Category: 1-RECO, 2-TRAN, 1-ENHA

Michael L. Seltzer*
Microsoft Research
Speech Technology Group

1 Microsoft Way
Redmond, WA 98052
tel: (425) 706-3763
fax: (425) 706-7329
email: mseltzer@microsoft.com

Bhiksha Raj
Mitsubishi Electric Research Labs
Cambridge Research Lab

201 Broadway
Cambridge, MA 02139
tel: (617) 621-7500
fax: (617) 621-7550
email: bhiksha@merl.com

Richard M. Stern
Carnegie Mellon University
Dept. of Electrical & Computer Engineering
and School of Computer Science

5000 Forbes Ave
Pittsburgh, PA 15213
tel: (412) 268-2535
fax: (412) 268-3890
email: rms@cs.cmu.edu

Likelihood-Maximizing Beamforming for Robust Hands-Free Speech Recognition

I. INTRODUCTION

State-of-the-art automatic speech recognition (ASR) systems are known to perform reasonably well when the speech signals are captured using a close-talking microphone worn near the mouth of the speaker. However, there are many environments where the use of such a microphone is undesirable for reasons of safety or convenience. In these settings, such as vehicles, meeting rooms, and information kiosks, a fixed microphone can be placed at some distance from the user. Unfortunately, as the distance between the user and the microphone grows, the speech signal becomes increasingly degraded by the effects of additive noise and reverberation, which in turn degrade speech recognition accuracy. The use of an array of microphones, rather than a single microphone, can compensate for this distortion in these distant-talking environments, by providing spatial filtering to the sound field, effectively focusing attention in a desired direction.

Many microphone array processing techniques which improve the quality of the output signal and increase the signal-to-noise ratio (SNR) have been proposed in the literature. The simplest and most common method is called delay-and-sum beamforming [1]. In this approach, the signals received by the microphones in the array are time-aligned with respect to each other in order to adjust for the path-length differences between the speech source and each of the microphones. The time-aligned signals are then weighted and added together. Any interfering signals that are not coincident with the speech source remain misaligned and are thus attenuated when the signals are combined. A natural extension of delay-and-sum beamforming is filter-and-sum beamforming, in which each microphone signal has an associated filter and the captured signals are filtered before they are combined.

In adaptive beamforming schemes such as the Generalized Sidelobe Canceller (GSC) [2], the array parameters are updated on a sample-by-sample or frame-by-frame basis according to a specified criterion. Typical criteria used in adaptive beamforming include a distortionless response in the look direction and/or the minimization of the energy from all directions not considered the look direction. In some cases, the array parameters can be calibrated to a particular environment or user prior to use, *e.g.* [3].

Adaptive filtering methods such as these generally assume that the target and jammer signals are uncorrelated. When this assumption is violated, as is the case for speech in a reverberant environment, the

methods suffer from signal cancellation because reflected copies of the target signal appear as unwanted jammer signals. While various methods have been proposed to mitigate this undesirable effect, *e.g.* [4], [5], signal cancellation nevertheless still arises in reverberant environments. As a result, conventional adaptive filtering approaches have not gained widespread acceptance for most speech recognition applications.

A great deal of recent research has focused specifically on compensating for the effects of reverberation. One obvious way to perform dereverberation is to invert the room impulse response. However, methods based on this approach have largely been unsuccessful because room impulse responses are generally non-minimum phase which causes instability in the inverse filters [6]. Rather than performing deconvolution, some researchers take a matched filter approach to dereverberation, *e.g.* [7], [8]. While there are theoretical benefits to such an approach in terms of improved SNR, matched filtering has been shown to provide only minimal improvement in speech recognition accuracy over conventional delay-and-sum processing, even if the room impulse responses are known *a priori* [9].

All of these microphone array processing methods were designed for signal enhancement, and as such, process incoming signals according to various signal-level criteria, *e.g.* minimizing the signal error, maximizing the SNR, or improving the perceptual quality as judged by human listeners. Conventional microphone-array-based speech recognition is performed by utilizing one of these algorithms to generate the best output waveform possible, which then gets treated as a single-channel input to a recognition system. This approach, shown in Figure 1(a), implicitly assumes that generating a higher quality output waveform will necessarily result in improved recognition performance. By making such an assumption, the manner in which speech recognition systems operate is ignored.

A speech recognition system does not interpret waveform-level information directly. It is a statistical pattern classifier that operates on a sequence of features derived from the waveform. We believe that this discrepancy between the waveform-based objective criteria used by conventional array processing algorithms and the feature-based objective criteria used by speech recognition systems is the key reason why sophisticated array processing methods fail to produce significant improvements in recognition accuracy over far simpler methods such as delay-and-sum beamforming. Speech recognition systems generate hypotheses by finding the word string that has the maximum likelihood of generating the observed sequence of feature vectors, as measured by statistical models of speech sound units. Therefore, an array processing scheme can only be expected to improve recognition performance if it generates a sequence of features which maximizes, or at least increases, the likelihood of the correct transcription, relative to other hypotheses.

In this paper, we present a new array processing algorithm called *Likelihood-MAXimizing BEAMforming*

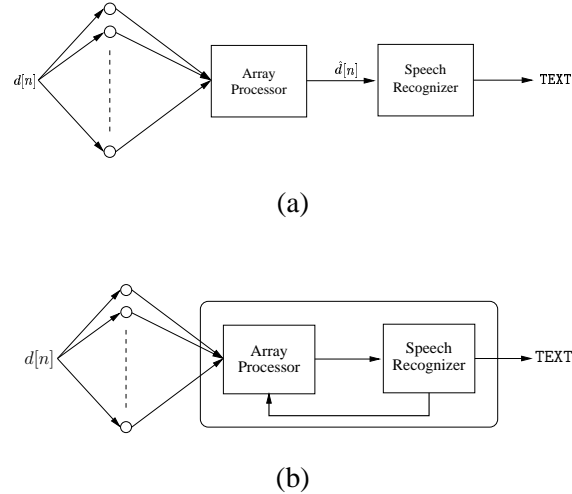


Fig. 1. (a) Conventional architecture used for speech recognition with a microphone array front-end. The objective of the array processor is to estimate the clean *waveform*. (b) An architecture for array processing optimized for speech recognition. The array processor and the speech recognizer are fully connected, allowing information from the recognizer to be used in the array processing. Note that the system no longer attempts to estimate the clean waveform.

(LIMABEAM), in which the microphone array processing problem is recast as one of finding the set of array parameters that maximizes the likelihood of the correct recognition hypothesis. The array processor and the speech recognizer are no longer considered two independent entities cascaded together, but rather two interconnected components of a single system, with the common goal of improved speech recognition accuracy, as shown in Figure 1(b). In LIMABEAM, the manner in which speech recognition systems process incoming speech is explicitly considered and pertinent information from the recognition engine itself is used to optimize the parameters of a filter-and-sum beamformer.

LIMABEAM has several advantages over current array processing methods. First, by incorporating the statistical models of the recognizer into the array processing stage, we ensure that the processing enhances those signal components important for recognition accuracy without undue emphasis on less important components. Second, in contrast to conventional adaptive filtering methods, no assumptions about the interfering signals are made. Third, the proposed approach requires no *a priori* knowledge of the room configuration, array geometry, or source-to-sensor room impulse responses. These properties enable us to overcome the drawbacks of previously-proposed array processing methods and achieve better recognition accuracy in distant-talking environments.

The remainder of this paper is organized as follows. In Section II, filter-and-sum beamforming is briefly reviewed. The LIMABEAM approach to microphone-array-based speech recognition is then described

in detail in Section III. In Section IV, two implementations of LIMABEAM are presented, one for use in situations in which the environmental conditions are stationary or slowly varying and one for use in time-varying environments. The performance of these two algorithms is evaluated in Section V through a series of experiments performed using a microphone-array-equipped Personal Digital Assistant (PDA). Some additional considerations for these algorithms are presented in Section VI. Finally, we present a summary of this work and some conclusions in Section VII.

II. FILTER-AND-SUM BEAMFORMING

In this work, we assume that filter-and-sum array processing can effectively compensate for the distortion induced by additive noise and reverberation. Assuming the filters have a finite impulse response (FIR), filter-and-sum processing is expressed mathematically as

$$y[n] = \sum_{m=0}^{M-1} \sum_{p=0}^{P-1} h_m[p] x_m[n - p - \tau_m] \quad (1)$$

where $h_m[p]$ is the p th tap of the filter associated with microphone m , $x_m[n]$ is the signal received by microphone m , τ_m is the steering delay induced in the signal received by microphone m to align it to the other array channels, and $y[n]$ is the output signal generated by the processing. M is the number of microphones in the array and P is the length of the FIR filters.

For notational convenience, we define ξ to be the vector of all filter coefficients for all microphones, as

$$\xi = [h_0[0], h_0[1], \dots, h_{M-1}[P-2], h_{M-1}[P-1]]^T \quad (2)$$

III. LIKELIHOOD-MAXIMIZING BEAMFORMING (LIMABEAM)

Conventionally, parameters of a filter-and-sum beamformer are chosen according to criteria designed according to the notion of a *desired signal*. In contrast, we consider the output waveform $y[n]$ to be incidental and seek the filter parameters that optimize recognition accuracy. Therefore, we forgo the notion of a desired signal, and instead focus on a *desired hypothesis*. In order to do so, we must consider both (1) the manner in which speech is input to the recognition system, *i.e.* the feature extraction process, and (2) the manner in which these features are processed by the recognizer in order to generate a hypothesis.

Speech recognition systems operate by finding the word string w most likely to generate the observed sequence of feature vectors $\mathcal{Z} = \{z_1, z_2, \dots, z_T\}$, as measured by the statistical models of the recognition system. When the speech is captured by a microphone array, the feature vectors are a function

of both the incoming speech and the array processing parameters. Recognition hypotheses are generated according to Bayes optimal classification as

$$\hat{\boldsymbol{w}} = \underset{\boldsymbol{w}}{\operatorname{argmax}} P(\mathcal{Z}(\boldsymbol{\xi})|\boldsymbol{w})P(\boldsymbol{w}) \quad (3)$$

where the dependence of the feature vectors \mathcal{Z} on the array processing parameters $\boldsymbol{\xi}$ is explicitly shown. The acoustic score $P(\mathcal{Z}(\boldsymbol{\xi})|\boldsymbol{w})$ is computed using the statistical models of the recognizer and the language score $P(\boldsymbol{w})$ is computed from a language model.

Our goal is to find the parameter vector $\boldsymbol{\xi}$ for optimal recognition performance. One logical approach to doing so is to choose the array parameters that maximize the likelihood of the correct transcription of the utterance that was spoken. This will increase the difference between the likelihood score of the correct transcription and the scores of competing incorrect hypotheses, and thus, increase the probability that the correct transcription will be hypothesized.

For the time being, let us assume that the correct transcription of the utterance, which we notate as \boldsymbol{w}_C , is known. We can then maximize (3) for the array parameters $\boldsymbol{\xi}$. Because the transcription is assumed to be known *a priori*, the language score $P(\boldsymbol{w}_C)$ can be neglected. The maximum likelihood (ML) estimate of the array parameters can now be defined as the vector that maximizes the acoustic log-likelihood of the given sequence of words, expressed as

$$\hat{\boldsymbol{\xi}} = \underset{\boldsymbol{\xi}}{\operatorname{argmax}} \log (P(\mathcal{Z}(\boldsymbol{\xi})|\boldsymbol{w}_C)) \quad (4)$$

In an HMM-based speech recognition system, the acoustic likelihood $P(\mathcal{Z}(\boldsymbol{\xi})|\boldsymbol{w}_C)$ is computed as the total likelihood of all possible state sequences through the HMM for the sequence of words in the transcription \boldsymbol{w}_C . However, many of these sequences are highly unlikely. For computational efficiency, we assume that the likelihood of a given transcription is largely represented by the single most likely HMM state sequence. If \mathcal{S}_C represents the set of all possible state sequences through this HMM and \boldsymbol{s} represents one such state sequence, then the ML estimate of $\boldsymbol{\xi}$ can be written as

$$\hat{\boldsymbol{\xi}} = \underset{\boldsymbol{\xi}, \boldsymbol{s} \in \mathcal{S}_C}{\operatorname{argmax}} \left\{ \sum_i \log(P(\boldsymbol{z}_i(\boldsymbol{\xi})|s_i)) + \sum_i \log(P(s_i|s_{i-1}, \boldsymbol{w}_C)) \right\} \quad (5)$$

According to (5), in order to find $\hat{\boldsymbol{\xi}}$, the likelihood of the correct transcription must be *jointly optimized* with respect to both the array parameters and the state sequence. This joint optimization can be performed by alternately optimizing the state sequence and the array processing parameters in an iterative manner.

A. Optimizing the State Sequence

Given a set of array parameters ξ , the speech can be processed by the array and a sequence of feature vectors $\mathcal{Z}(\xi)$ produced. Using the features vectors and the transcription w_C , we want to find the state sequence $\hat{s} = \{\hat{s}_1, \hat{s}_2, \dots, \hat{s}_T\}$ such that

$$\hat{s} = \operatorname{argmax}_{s \in \mathcal{S}_C} \sum_i \log(P(s_i | s_{i-1}, w_C, \mathcal{Z}(\xi))) \quad (6)$$

This state sequence \hat{s} can be easily determined by forced alignment using the Viterbi algorithm [10].

B. Optimizing the Array Parameters

Given a state sequence, \hat{s} , we are interested in finding $\hat{\xi}$ such that

$$\hat{\xi} = \operatorname{argmax}_{\xi} \sum_i \log(P(z_i(\xi) | \hat{s}_i)) \quad (7)$$

This acoustic likelihood expression cannot be directly maximized with respect to the array parameters ξ for two reasons. First, the state distributions used in most HMM-based speech recognition systems are complicated density functions, *i.e.* mixtures of Gaussians. Second, the acoustic likelihood of an utterance and the parameter vector ξ are related through a series of linear and non-linear mathematical operations performed to convert a waveform into a sequence of feature vectors. Therefore, for a given HMM state sequence, no closed-form solution for the optimal value of ξ exists. As a result, non-linear optimization methods must be used.

We employ a gradient-based approach to finding the optimal value of ξ . For convenience, we define $\mathcal{L}(\xi)$ to be the total log likelihood of the observation vectors given an HMM state sequence. Thus,

$$\mathcal{L}(\xi) = \sum_i \log(P(z_i(\xi) | s_i)) \quad (8)$$

Using the definition of ξ given by (2), we define the gradient vector $\nabla_{\xi} \mathcal{L}(\xi)$ as

$$\nabla_{\xi} \mathcal{L}(\xi) = \left[\frac{\partial \mathcal{L}(\xi)}{\partial h_0[0]}, \frac{\partial \mathcal{L}(\xi)}{\partial h_0[1]}, \dots, \frac{\partial \mathcal{L}(\xi)}{\partial h_{M-1}[P-1]} \right]^T \quad (9)$$

Clearly, the computation of the gradient vector is dependent on the form of the HMM state distributions used by the recognition system and the features used for recognition. In the following sections, we derive the gradient expressions when the state distributions are modeled as Gaussian distributions or mixtures of Gaussians. In both cases, the features are assumed to be mel frequency cepstral coefficients (MFCC) or log mel spectra.

1) *Gaussian State Output Distributions*: We now derive the expression for $\nabla_{\xi}\mathcal{L}(\xi)$ for the case where the HMM state distributions are multivariate Gaussian distributions with diagonal covariance matrices. If we define μ_i and Σ_i to be the mean vector and covariance matrix, respectively, of the pdf of the most likely HMM state at frame i , the total log likelihood for an utterance can be expressed as

$$\mathcal{L}(\xi) = \sum_i \left\{ -\frac{1}{2} (z_i(\xi) - \mu_i)^T \Sigma_i^{-1} (z_i(\xi) - \mu_i) + \kappa_i \right\} \quad (10)$$

where κ_i is a normalizing constant. Using the chain rule, the gradient of $\mathcal{L}(\xi)$ with respect to ξ can be expressed as

$$\nabla_{\xi}\mathcal{L}(\xi) = - \sum_i \frac{\partial z_i(\xi)}{\partial \xi} \Sigma_i^{-1} (z_i(\xi) - \mu_i) \quad (11)$$

where $\partial z_i(\xi)/\partial \xi$ is the Jacobian matrix, composed of the partial derivatives of each element of the feature vector at frame i with respect to each of the array parameters. The Jacobian is of dimension $MP \times L$ where M is the number of microphones, P is the number of parameters per microphone, and L is the dimension of the feature vector.

It can be shown that for log mel spectral feature vectors, the elements of the Jacobian matrix can be expressed as

$$\frac{\partial z_i^l}{\partial h_m[p]} = 2 \frac{1}{M_i^l} \sum_{k=0}^{N/2} V^l[k] \Re (X_i^{mp}[k] Y_i^*[k]) \quad (12)$$

where $Y_i[k]$ is the DFT of frame i of the output signal $y[n]$, $X_i^{mp}[k]$ is the DFT of the signal captured by microphone m , beginning p samples prior to the start of frame i , $V^l[k]$ is the value of the l th mel filter applied to DFT bin k , and M_i^l is the l th mel spectral component in frame i . The size of the DFT is N and $*$ denotes complex conjugation. Note that in (12), we have assumed that Time-Delay Compensation (TDC) has already been performed and that the microphone signals $x_m[n]$ have already been time-aligned.

If optimization of the filter parameters is performed using MFCC features rather than log mel spectra, (12) must be modified slightly to account for the additional DCT operation. The full derivation of the Jacobian matrix for log mel spectral or cepstral features can be found in [11].

2) *Mixture of Gaussians State Output Distributions*: Most state-of-the-art recognizers do not model the state output distributions as single Gaussians but rather as mixtures of Gaussians. It can be shown [11] that when the HMM state distributions are modeled as mixtures of Gaussians, the gradient expression can be expressed as

$$\nabla_{\xi}\mathcal{L}(\xi) = - \sum_i \sum_{k=1}^K \gamma_{ik}(\xi) \frac{\partial z_i(\xi)}{\partial \xi} \Sigma_{ik}^{-1} (z_i(\xi) - \mu_{ik}) \quad (13)$$

where $\gamma_{ik}(\xi)$ represents the *a posteriori* probability of the k th mixture component of state s_i , given $z_i(\xi)$.

Comparing (11) and (13), it is clear that the gradient expression in the Gaussian mixture case is simply a weighted sum of the gradients of each of the Gaussian components in the mixture, where the weight on each mixture component represents its *a posteriori* probability of generating the observed feature vector.

C. Optimizing Log Mel Spectra vs. Cepstra

Array parameter optimization is performed in the log mel spectral domain, rather than the cepstral domain. Because the log mel spectra are derived from the energy using a series of triangular weighting functions of unit area, all components of the vectors have approximately the same magnitude. In contrast, the magnitude of cepstral coefficients decreases significantly with increasing cepstral order. When there is a large disparity in the magnitudes of the components of a vector, the larger components dominate the objective function and tend to be optimized at the expense of smaller components in gradient-descent-based optimization methods. Using log mel spectra avoids this potential problem.

In order to perform the array parameter optimization in the log mel spectral domain but still perform decoding using mel frequency cepstral coefficients (MFCC), we employ a parallel set of HMMs trained on log mel spectra, rather than cepstra. To obtain parallel models, we employed the Statistical Re-estimation (STAR) algorithm [12], which ensures that the two sets of models have identical frame-to-state alignments.

These parallel log mel spectral models were trained without feature mean normalization, since mean normalization is not incorporated into the optimization framework (we will revisit this issue in Section VI).

D. Gradient-based Array Parameter Optimization

Using the gradient vector defined in either (11) or (13), the array parameters can be optimized using conventional gradient descent [13]. However, improved convergence performance can be achieved by other methods, *e.g.* those which utilize estimates of the Hessian. In this work, we perform optimization using the method of conjugate gradients [14], using the software package found in [15]. In this method, the step size varies with each iteration and is determined by the optimization algorithm itself.

IV. LIMABEAM IN PRACTICE

In the previous section, a new approach to microphone array processing was presented in which the array processing parameters are optimized specifically for speech recognition performance using information from the speech recognition system itself. Specifically, we showed how the parameters of a filter-and-sum beamformer can be optimized to maximize the likelihood of a known transcription. Clearly, we are faced with a paradox: prior knowledge of the correct transcription, is required in order

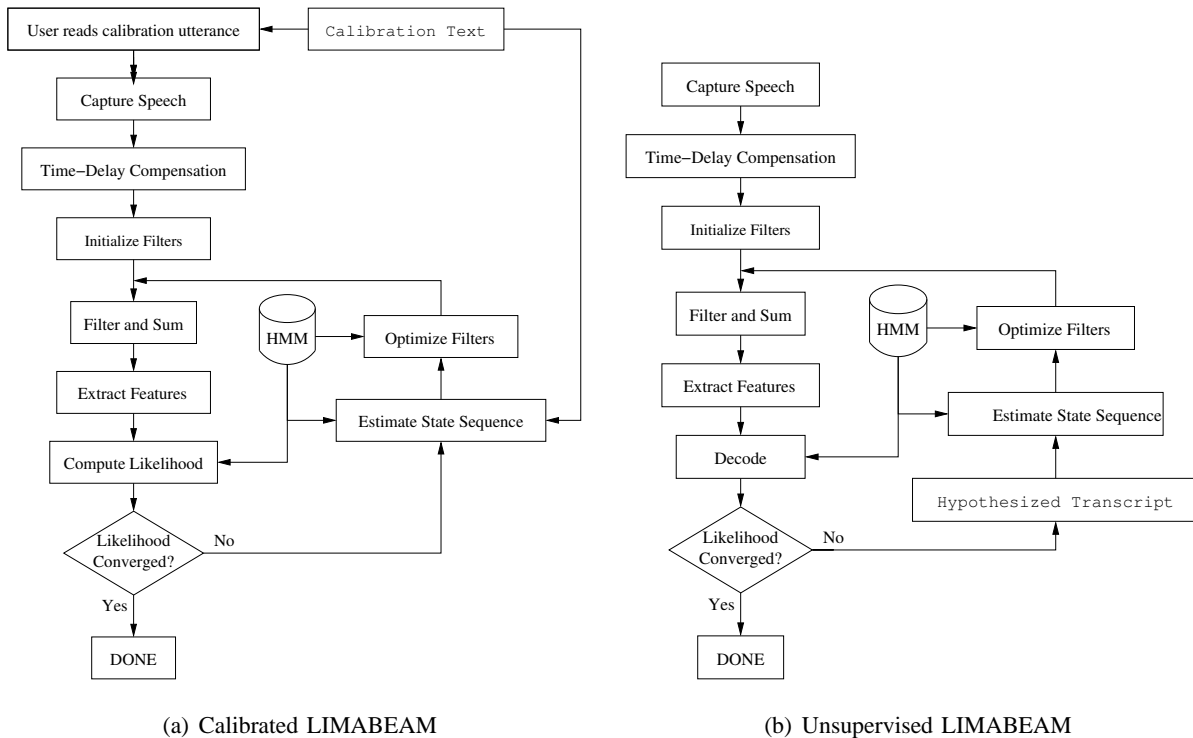


Fig. 2. Flowcharts of Calibrated LIMABEAM and Unsupervised LIMABEAM. In Calibrated LIMABEAM, the parameters of the filter-and-sum beamformer are optimized using a calibration utterance with a known transcription and then fixed for future processing. In Unsupervised LIMABEAM, the parameters are optimized for each utterance independently using hypothesized transcriptions.

to maximize its likelihood, but if we had such knowledge, there would be no need for recognition in the first place. In this section we present two different implementations of LIMABEAM as solutions to this paradox. The first method is appropriate for situations in which the environment and the user’s position do not vary significantly over time, such as in a vehicle or in front of a desktop computer terminal, while the second method is more appropriate for time-varying environments.

A. Calibrated LIMABEAM

In this approach, the LIMABEAM algorithm is cast as a method of microphone array calibration. In the calibration scenario, the user is asked to speak an enrollment utterance with a *known* transcription. An estimate of the most likely state sequence corresponding to the enrollment transcription is made via forced alignment using the features derived from the array signals. These features can be generated using an initial set of filters, *e.g.* from a previous calibration session or a simple delay-and-sum configuration.

Using this estimated state sequence, the filter parameters can be optimized. Using the optimized filter

parameters, a second iteration of calibration can be performed. An improved set of features for the calibration utterance is generated and used to re-estimate the state sequence. The filter optimization process can then be repeated using the updated state sequence. The calibration process continues in an iterative manner until the overall likelihood converges. Once convergence occurs, the calibration process is complete. The resulting filters are now fixed and used to process future incoming speech to the array. Because the array parameters are calibrated to maximize the likelihood of the enrollment utterance, we refer to this method as *Calibrated LIMABEAM*. A flowchart of the calibration algorithm is shown in Figure 2(a).

B. *Unsupervised LIMABEAM*

In order for the proposed calibration algorithm to be effective, the array parameters learned during calibration must be valid for future incoming speech. This implies that there will not be any significant changes over time to the environment or the user's position. While this is a reasonable assumption for several situations, there are several applications in which either the environment or the position of the user do vary over time. In these cases, filters obtained from calibration may no longer be valid. Furthermore, there may be situations in which requiring the user to speak a calibration utterance is undesirable. For example, a typical interaction at an information kiosk is relatively brief and requiring the user to calibrate the system will significantly increase the time it takes for the user to complete a task.

In these situations, it is more appropriate to optimize the array parameters more frequently, *i.e.* on an utterance-by-utterance basis. However, we are again faced with the paradox discussed earlier. In order to maximize the likelihood of the correct transcription of the test utterances, we require *a priori* knowledge of the very transcriptions that we desire to recognize. In this case, where the use of a calibration utterance is no longer appropriate, we solve this dilemma by estimating the transcriptions and using them in an *unsupervised* manner to perform the array parameter optimization.

In *Unsupervised LIMABEAM*, the filter parameters are optimized on the basis of a hypothesized transcription, generated from an initial estimate of the filter parameters. Thus, this algorithm is a multi-pass algorithm. For each utterance or series of utterances, the current set of filter parameters are used to generate a set of features for recognition which in turn, are used to generate a hypothesized transcription. Using the hypothesized transcription and the associated feature vectors, the most likely state sequence is estimated using Viterbi alignment as before. The filters are then optimized using the estimated state sequence, and a second pass of recognition is performed. This process can be iterated until the likelihood converges. A flowchart of the algorithm is shown in Figure 2(b).

V. EXPERIMENTAL EVALUATION

In order to evaluate the proposed Calibrated LIMABEAM and Unsupervised LIMABEAM algorithms, we employed the CMU WSJ PDA corpus, recorded at CMU. This corpus was recorded using a PDA mockup, created with a Compaq iPaq outfitted with 4 microphones using a custom-made frame attached to the PDA. The microphones were placed in a $5.5 \text{ cm} \times 14.6 \text{ cm}$ rectangular configuration, as shown in Figure 3. The 4 microphones plus a close-talking microphone worn by the user were connected to a digital audio multi-track recorder. The speech data was recorded at a sampling rate of 16 kHz.

Recordings were made in a room approximately $6.0 \text{ m} \times 3.7 \text{ m} \times 2.8 \text{ m}$. The room contained several desks, computers and a printer. The reverberation time of the room was measured to be approximately 270 ms. Users read utterances from the Wall Street Journal (WSJ) test set [16] which were displayed on the PDA screen. All users sat in a chair in the same location in the room and held the PDA in whichever hand was most comfortable. No instructions were given to the user about how to hold the PDA. Depending on the preference or habits of the user, the position of the PDA could vary from utterance-to-utterance or during a single utterance.

Two separate recordings of the WSJ test set were made with 8 different speakers in each set. In the first set, referred to as *PDA-A*, the average SNR of the array channels is approximately 21 dB. In the second recording session, a humidifier was placed near the user to create a noisier environment. The SNR of the second set, referred to as *PDA-B*, is approximately 13 dB.

Speech recognition was performed using Sphinx-3, a large-vocabulary HMM-based speech recognition system [17]. Context-dependent 3-state left-to-right HMMs with no skips (8 Gaussians/state) were trained using the speaker-independent WSJ training set, consisting of 7000 utterances. The system was trained with 39-dimensional feature vectors consisting of 13-dimensional MFCC parameters, along with their delta and delta-delta parameters. A 25-ms window length and a 10-ms frame shift were used. Cepstral mean normalization (CMN) was performed in both training and testing.

A. Experiments Using Calibrated LIMABEAM

The first series of experiments were performed to evaluate the performance of Calibrated LIMABEAM algorithm. In these experiments, a single calibration utterance for each speaker was chosen at random from utterances at least 10 s in duration. For each speaker, delay-and-sum beamforming was performed on the calibration utterance and recognition features were generated from the delay-and-sum output signal. These features and the known transcription of the calibration utterance were used to estimate the most likely state sequence via forced alignment. Using this state sequence, a filter-and-sum beamformer with



Fig. 3. A 4-microphone PDA mockup used to record the CMU WSJ PDA corpus.

20 taps per filter was optimized. In all cases, the steering delays were estimated using the PHAT method [18] and the filters were initialized to a delay-and-sum configuration for optimization. The filters obtained were then used to process all remaining utterances for that speaker.

Experiments were performed using both 1 Gaussian per state and 8 Gaussians per state in the log-likelihood expression used for filter optimization. The results of these experiments are shown in Figure 4(a) and Figure 4(b) for the PDA-A and PDA-B test sets, respectively. For comparison, the results obtained using only a single microphone from the array and using conventional delay-and-sum beamforming are also shown. The Generalized Sidelobe Canceller (GSC) algorithm ([2]) with parameter adaptation during the non-speech regions only (as per [4]) was also performed on the PDA data. The recognition performance was significantly worse than delay-and-sum beamforming and therefore the results are not reported here.

As the figures show, the calibration approach is, in general, successful at improving the recognition accuracy over delay-and-sum beamforming. On the less noisy PDA-A test data, using mixtures of Gaussians in the likelihood expression to be optimized resulted in a significant improvement over conventional delay-and-sum processing, whereas the improvement using single Gaussians is negligible. On the other hand, the improvements obtained in the noisier PDA-B set are substantial in both cases and the performance is basically the same. While it is difficult to compare results across the two test sets directly because the speakers are different in each set, the results obtained using single Gaussians versus mixtures of Gaussians generally agree with intuition. When the test data are well matched to the training data, *i.e.* same domain and distortion, using more descriptive models is beneficial. As the mismatch between the

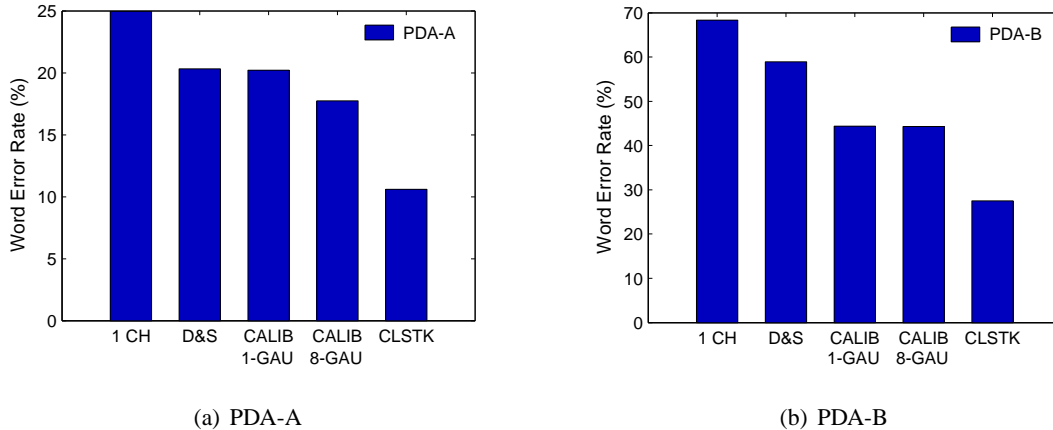


Fig. 4. Word Error Rate obtained using Calibrated LIMABEAM on the CMU WSJ PDA-A and PDA-B corpora. The figures show the performance obtained using a single microphone, delay-and-sum beamforming, and the proposed Calibrated LIMABEAM method with 1 Gaussian per state or 8 Gaussians per state in the optimization. The performance obtained using a close-talking microphone is also shown.

training and test data increases, *e.g.* the SNR decreases, more general models give better performance.

Comparing Figure 4(a) and Figure 4(b), there is a significant disparity in the relative improvement obtained on the PDA-A test set compared with the PDA-B test set. A 12.7% relative improvement over delay-and-sum beamforming was obtained in the PDA-A test set, while the improvement on PDA-B was 24.6%. As described above, the users were not told to keep the PDA in the same position from utterance to utterance. The users in this corpus each read approximately 40 utterances while holding the PDA in their hand. Therefore, we can expect some movement will naturally occur. As a result, the filter parameters obtained from calibration using an utterance chosen at random may not be valid for many of the utterances from that user. We re-examine this hypothesis in the next section, where the parameters are adjusted for each utterance individually using the unsupervised approach.

Finally, the experimental procedure described constitutes a single iteration of the Calibrated LIMABEAM algorithm. Performing additional iterations did not result in any further improvement.

B. Experiments Using Unsupervised LIMABEAM

A second series of experiments was performed to evaluate the performance of the Unsupervised LIMABEAM algorithm. In this case, the filter parameters were optimized for each utterance individually in the following manner. Delay-and-sum beamforming was used to process the array signals in order to generate an initial hypothesized transcription. Using this hypothesized transcription and the features derived from the delay-

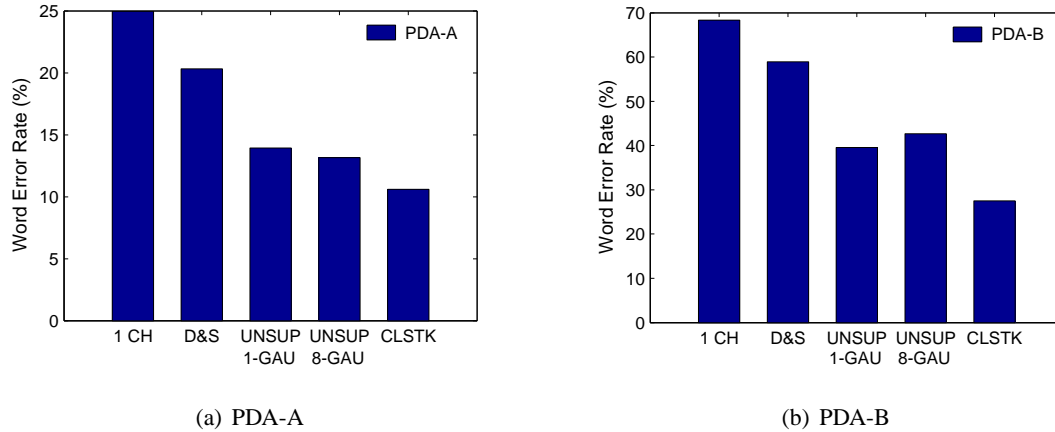


Fig. 5. Word Error Rate obtained using Unsupervised LIMABEAM on the CMU WSJ PDA-A and PDA-B corpora. The figures show the performance obtained using a single microphone, delay-and-sum beamforming, and the proposed Unsupervised LIMABEAM method with 1 Gaussian per state or 8 Gaussians per state in the optimization. The performance obtained using a close-talking microphone is also shown.

and-sum output, the state sequence was estimated via forced alignment. Using this state sequence, the filter parameters were optimized. As in the calibrated case, 20 taps were estimated per filter and the filters were initialized to a delay-and-sum configuration. We again compared the recognition accuracy obtained when optimization is performed using HMM state output distributions modeled as Gaussians or mixtures of Gaussians. The results are shown in Figure 5(a) and Figure 5(b) for PDA-A and PDA-B, respectively.

There is sizable improvement in recognition accuracy over conventional delay-and-sum beamforming in both test sets. Using Unsupervised LIMABEAM, an average relative improvement of 31.4% was obtained over delay-and-sum beamforming over both test sets. It is interesting to note that by comparing Figure 4(a) and Figure 5(a), we can see a dramatic improvement in performance using the unsupervised method, compared to that obtained using the calibration algorithm. This confirms our earlier conjecture that the utterance used for calibration was not representative of the data in the rest of the test set, possibly because the position of the PDA with respect to the user varied over the course of the test set.

Additionally, we can also see that the effect of optimizing using Gaussian mixtures versus single Gaussians in the unsupervised case is similar to that seen in the calibration experiments. As the mismatch between training and testing conditions increases, better performance is obtained from the more general single Gaussian models. As in the calibration case, these results were obtained from only a single iteration of Unsupervised LIMABEAM, and additional iterations did not improve the performance further.

TABLE I

WER OBTAINED ON PDA-A USING UNSUPERVISED LIMABEAM AND AN OPTIMIZED SINGLE-CHANNEL POST-FILTER APPLIED TO THE OUTPUT OF A DELAY-AND-SUM BEAMFORMER.

Processing Method	Filter Length	WER (%)
Delay-and-sum	–	20.3
Unsupervised ML post-filter	20	20.1
Unsupervised ML post-filter	80	18.4
Unsupervised LIMABEAM	20	13.9

C. LIMABEAM vs. Sum-and-Filter Processing

There is another class of methods for microphone array processing which can be referred to as *sum-and-filter* methods. In such methods, the array signals are processed using conventional delay-and-sum beamforming or another array processing algorithm and the single-channel output signal is then passed through a *post-filter* for additional spectral shaping and noise removal [19], [20].

We performed a series of experiments to compare the performance of the proposed maximum likelihood filter-and-sum beamformer to that of a single-channel post-filter optimized according to the same maximum likelihood criterion and applied to the output of a delay-and-sum beamformer. For these experiments, the parameters of both the filter-and-sum beamformer and the single-channel post-filter were optimized using Unsupervised LIMABEAM with single-Gaussian HMM state output distributions. For the filter-and-sum beamformer, 20-tap filters were estimated as before. For the post-filtering, filters with 20 taps and 80 taps were estimated, the latter being the same number of total parameters used in the filter-and-sum case.

The results of these experiments are shown in Table I. As the table shows, jointly optimizing the parameters of a filter-and-sum beamformer provides significantly better speech recognition performance compared to optimizing the parameters of a single-channel post-filter.

D. Incorporating TDC Into LIMABEAM

In the filter-and-sum equation shown in (1), the steering delays were shown explicitly as τ_m and in the experiments performed thus far, we estimated those delays and performed TDC prior to optimizing the filter parameters of the beamformer. Thus, at the start of LIMABEAM, the microphone array signals are all in phase. However, because TDC is simply a time-shift of the input signals, it can theoretically be incorporated into the filter optimization process. Therefore it is possible that we can do away with

TABLE II

WER OBTAINED ON PDA-B USING BOTH LIMABEAM METHODS WITH AND WITHOUT TIME-DELAY COMPENSATION (TDC) PRIOR TO BEAMFORMER OPTIMIZATION.

Array Processing Method	TDC	WER (%)
Delay-and-sum	Yes	58.9
Calibrated LIMABEAM	Yes	44.4
Calibrated LIMABEAM	No	45.7
Unsupervised LIMABEAM	Yes	39.5
Unsupervised LIMABEAM	No	39.8

the TDC step and simply let the LIMABEAM algorithm implicitly learn the steering delays as part of the filter optimization process.

To test this, we repeated the Calibrated LIMABEAM and Unsupervised LIMABEAM experiments on the PDA-B test set. We compared the performance of both algorithms with and without TDC performed prior to filter parameter optimization. In the case where TDC was not performed, the initial set of features required by LIMABEAM for state-sequence estimation was obtained by simply averaging the array signals together without any time alignment. The results of these experiments are shown in Table II.

As the results in the table show, there is very little degradation in performance when the TDC is incorporated into the filter optimization process. It should be noted that in these experiments, because the users held the PDA, they were never significantly off-axis to the array. Therefore there was not a significant difference between the initial features obtained from delay-and-sum and those obtained from averaging the unaligned signals together. In situations where the user is significantly off-axis, initial features obtained from simple averaging without TDC may be noisier than those obtained after TDC. This may degrade the quality of the state-sequence estimation, which may, in turn, degrade the performance of the algorithm. In these situations, performing TDC prior to filter parameter optimization is preferable.

VI. OTHER CONSIDERATIONS

A. Combining the LIMABEAM Implementations

In situations where Calibrated LIMABEAM is expected to generate improved recognition accuracy, the overall performance can be improved further by performing Calibrated LIMABEAM and Unsupervised LIMABEAM sequentially. As is the case with all unsupervised processing algorithms, the performance of Unsupervised LIMABEAM is dependent on the accuracy of the data used for adaptation.

By performing Calibrated LIMABEAM prior to Unsupervised LIMABEAM, we can use the calibration method as a means of obtaining more reliable transcriptions to use in the unsupervised optimization.

To demonstrate the efficacy of this approach, an experiment was performed using the PDA-B test set in which Unsupervised LIMABEAM was performed using transcriptions generated by Calibrated LIMABEAM rather than by delay-and-sum beamforming as before. Recalling Figure 4(b), Calibrated LIMABEAM generated a 24.6% relative improvement over delay-and-sum processing on the PDA-B test set. By performing Unsupervised LIMABEAM using the transcriptions generated by the calibrated beamformer rather than by delay-and-sum beamforming, the Word Error Rate (WER) was reduced from 42.8% to 37.9%. For comparison, the WER obtained from delay-and-sum beamforming was 58.9%.

B. Data Sufficiency for LIMABEAM

One important factor to consider when using either of the two LIMABEAM implementations described is the amount of speech data used in the filter optimization process. If too little data are used for optimization or the data are unreliable, then the filters produced by the optimization process will be sub-optimal and could potentially degrade recognition accuracy.

In Calibrated LIMABEAM, we are attempting to obtain filters that generalize to future utterances using a very small amount of data (only a single utterance). As a result, if the beamformer contains too many parameters, the likelihood of overfitting is quite high. For example, experiments performed on an 8-channel microphone array in [11] showed that a 20-tap filter-and-sum beamformer can be reliably calibrated with only 3-4 s of speech. However, when the filter length is increased to 50 taps, overfitting occurs and recognition performance degrades. When 8-10 s of speech data are used, the 50-tap filters can be calibrated successfully and generate better performance than the 20-tap filters calibrated on the same amount of data.

For Unsupervised LIMABEAM to be successful, there has to be a sufficient number of correctly labeled frames in the utterance. Performing unsupervised optimization on an utterance with too few correctly hypothesized labels will only degrade performance, propagating the recognition errors further.

C. Incorporating Feature Mean Normalization

Speech recognition systems usually perform better when mean normalization is performed on the features prior to being processed by the recognizer, both in training and decoding. Mean normalization can easily be incorporated into the filter parameter optimization scheme by performing mean normalization on both the features and the Jacobian matrix in the likelihood expression and its gradient.

TABLE III

WER OBTAINED BY APPLYING UNSUPERVISED MLLR AFTER ARRAY PROCESSING ON THE PDA-B TEST SET.

Processing Method	WER (%)
Delay-and-sum	58.9
Delay-and-sum + MLLR	44.7
Calibrated LIMABEAM	44.4
Calibrated LIMABEAM + MLLR	40.8

However, we found no additional benefit to incorporating feature mean normalization into the array parameter optimization process. We believe this is because the array processing algorithm is already attempting to perform some degree of channel compensation for both the room response *and* the microphone channel, as it is impossible to separate the two.

D. Applying Additional Robustness Techniques

There is a vast literature of techniques designed to improve speech recognition accuracy under adverse conditions, such as additive noise and/or channel distortion. These algorithms typically operate in the feature space, *e.g.* Codeword-Dependent Cepstral Normalization (CDCN) [21], or the model space, *e.g.* Maximum Likelihood Linear Regression (MLLR) [22].

We have found that applying such techniques after LIMABEAM results in further improvements in performance. For example, Table III shows the WER obtained when batch-mode unsupervised MLLR with a single regression class is applied after delay-and-sum beamforming and after Calibrated LIMABEAM for the PDA-B test set.

As the table shows, performing unsupervised MLLR after delay-and-sum beamforming results in recognition accuracy that is almost as good as Calibrated LIMABEAM alone. However, when MLLR is applied after Calibrated LIMABEAM, an additional 10% reduction in WER is obtained. It should also be noted that in this experiment, the MLLR parameters were estimated using the entire test set, while the parameters estimated by Calibrated LIMABEAM were estimated from only a single utterance. Furthermore, by comparing these results to those shown in Table II, we can see that the performance obtained by applying MLLR to the output of delay-and-sum beamforming is still significantly worse than that obtained by Unsupervised LIMABEAM alone.

VII. SUMMARY AND CONCLUSIONS

In this paper, we introduced Likelihood-Maximizing Beamforming (LIMABEAM), a novel approach to microphone array processing designed specifically for improved speech recognition performance. This method differs from previous array processing algorithms in that no waveform-level criteria are used to optimize the array parameters. Instead, the array parameters are chosen to maximize the likelihood of the correct transcription of the utterance, as measured by the statistical models used by the recognizer itself. We showed that finding a solution to this problem involves jointly optimizing the array parameters and the most likely state sequence for the given transcription and described a method for doing so.

We then developed two implementations of LIMABEAM which optimized the parameters of a filter-and-sum beamformer. In the first method, called Calibrated LIMABEAM, an enrollment utterance with a known transcription is spoken by the user and used to optimize the filter parameters. These filter parameters are then fixed and used to process future utterances. This algorithm is appropriate for situations in which the environment and the user's position do not vary significantly over time. For time-varying environments, we developed an algorithm for optimizing the filter parameters in an unsupervised manner. In Unsupervised LIMABEAM, the optimization is performed on each utterance independently using a hypothesized transcription obtained from an initial pass of recognition.

The performance of these two LIMABEAM methods was demonstrated using a microphone-array-equipped PDA. In the Calibrated LIMABEAM method, we were able to obtain an average relative improvement of 18.6% over conventional beamforming, while the average relative improvement obtained using Unsupervised LIMABEAM was 31.4%. We were able to improve performance further still by performing Calibrated LIMABEAM and Unsupervised LIMABEAM in succession, and also by applying HMM adaptation after LIMABEAM.

The experiments performed in this paper showed that we can obtain significant improvements in recognition accuracy over conventional microphone array processing approaches in environments with moderate reverberation over a range of SNRs. However, in highly reverberant environments, an increased number of parameters is needed in the filter-and-sum beamformer to effectively compensate for the reverberation. As the number of parameters to optimize increases, the data insufficiency issues discussed in Section VI begin to emerge more significantly, and the performance of LIMABEAM suffers. To address these issues and improve speech recognition accuracy in highly reverberant environments, we have begun developing a subband filtering implementation of LIMABEAM.

ACKNOWLEDGMENTS

This research was sponsored by the Space and Naval Warfare Systems Center, San Diego, under Grant No. N66001-99-1-8905. The content of the information in this publication does not necessarily reflect the position or the policy of the U. S. Government, and no official endorsement should be inferred. Michael L. Seltzer was also supported by a Microsoft Research Graduate Fellowship.

The authors would like to thank Dr. Yasunari Obuchi of the Hitachi Central Research Lab, for recording and preparing the PDA speech data used in this work.

REFERENCES

- [1] D. H. Johnson and D. E. Dudgeon, *Array Signal Processing*. New Jersey: Prentice Hall, 1993.
- [2] L. J. Griffiths and C. W. Jim, "An alternative approach to linearly constrained adaptive beamforming," *IEEE Transactions on Antennas and Propagation*, vol. AP-30, no. 1, pp. 27–34, January 1982.
- [3] S. Nordholm, I. Claesson, and M. Dahl, "Adaptive microphone array employing calibration signals," *IEEE Transactions on Speech and Audio Processing*, vol. 7, no. 3, pp. 241–252, May 1999.
- [4] D. Van Compernelle, "Switching adaptive filters for enhancing noisy and reverberant speech from microphone array recordings," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 2, Albuquerque, NM, April 1990, pp. 833–836.
- [5] O. Hoshuyama, A. Sugiyama, and A. Hirano, "A robust adaptive beamformer for microphone arrays with a blocking matrix using constrained adaptive filters," *IEEE Transactions on Signal Processing*, vol. 47, no. 10, pp. 2677–2684, October 1999.
- [6] S. Neely and J. Allen, "Invertibility of a room impulse response," *Journal of the Acoustical Society of America*, vol. 66, no. 1, pp. 165–169, July 1979.
- [7] J. L. Flanagan, A. C. Surendran, and E. E. Jan, "Spatially selective sound capture for speech and audio processing," *Speech Communication*, vol. 13, no. 1-2, pp. 207–222, October 1993.
- [8] S. Affes and Y. Grenier, "A signal subspace tracking algorithm for microphone array processing of speech," *IEEE Transactions on Speech and Audio Processing*, vol. 5, no. 5, pp. 425–437, September 1997.
- [9] B. Gillespie and L. E. Atlas, "Acoustic diversity for improved speech recognition in reverberant environments," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, Orlando, FL, May 2002, pp. 557–560.
- [10] A. J. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE Transactions on Information Theory*, vol. 13, no. 2, pp. 260–269, April 1967.
- [11] M. L. Seltzer, "Microphone array processing for robust speech recognition," Ph.D. dissertation, Carnegie Mellon University, Pittsburgh, PA, July 2003. [Online]. Available: <http://www.cs.cmu.edu/~robust/Thesis/mseltzer-phdthesis.pdf>
- [12] P. Moreno, B. Raj, and R. M. Stern, "A unified approach for robust speech recognition," in *Proceedings of Eurospeech*, vol. 1, Madrid, Spain, September 1995, pp. 481–485.
- [13] S. Haykin, *Adaptive Filter Theory*. New Jersey: Prentice Hall, 2002.
- [14] J. Nocedal and S. Wright, *Numerical Optimization*. New York: Springer, 1999.
- [15] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes in C: The Art of Scientific Computing*. New York: Cambridge University Press, 1988.

- [16] D. B. Paul and J. M. Baker, "The design of the Wall Street Journal-based CSR corpus," in *Proceedings of the ARPA Speech and Natural Language Workshop*, Harriman, NY, February 1992, pp. 357–362.
- [17] P. Placeway, S. Chen, M. Eskenazi, U. Jain, V. Parikh, B. Raj, M. Ravishankar, R. Rosenfeld, K. Seymore, M. Siegler, R. Stern, and E. Thayer, "The 1996 Hub-4 Sphinx-3 system," in *Proceedings of the DARPA Speech Recognition workshop*. DARPA, February 1997.
- [18] C. H. Knapp and C. Carter, "The generalized correlation method for estimation of time delay," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. ASSP-24, no. 4, pp. 320–327, August 1976.
- [19] R. Zelinski, "A microphone array with adaptive post-filtering for noise reduction in reverberant rooms," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 5, New York, NY, May 1988, pp. 2578–2581.
- [20] C. Marro, Y. Mahieux, and K. U. Simmer, "Analysis of noise reduction and dereverberation techniques based on microphone arrays with postfiltering," *IEEE Transactions on Speech and Audio Processing*, vol. 6, no. 3, pp. 240–259, May 1998.
- [21] A. Acero, *Acoustical and Environmental Robustness in Automatic Speech Recognition*. Boston, MA: Kluwer Academic Publishers, 1993.
- [22] C. J. Leggetter and P. C. Woodland, "Speaker adaptation of HMMs using linear regression," Cambridge University, Cambridge, UK, Tech. Rep. CUED/F-INFENG/ TR. 181, June 1994.



Michael L. Seltzer Biography text here.



Bhiksha Raj Biography text here.



Richard M. Stern Biography text here.