

LIBRARY OF CONGRESS



0 020 814 014 0

DIGITAL SIGNAL PROCESSING 12 JAN 2002 1

TK 5102

.5

.D4463

Set 1

DIGITAL SIGNAL PROCESSING

12

JAN

2002

1

Digital Signal Processing

A Review Journal



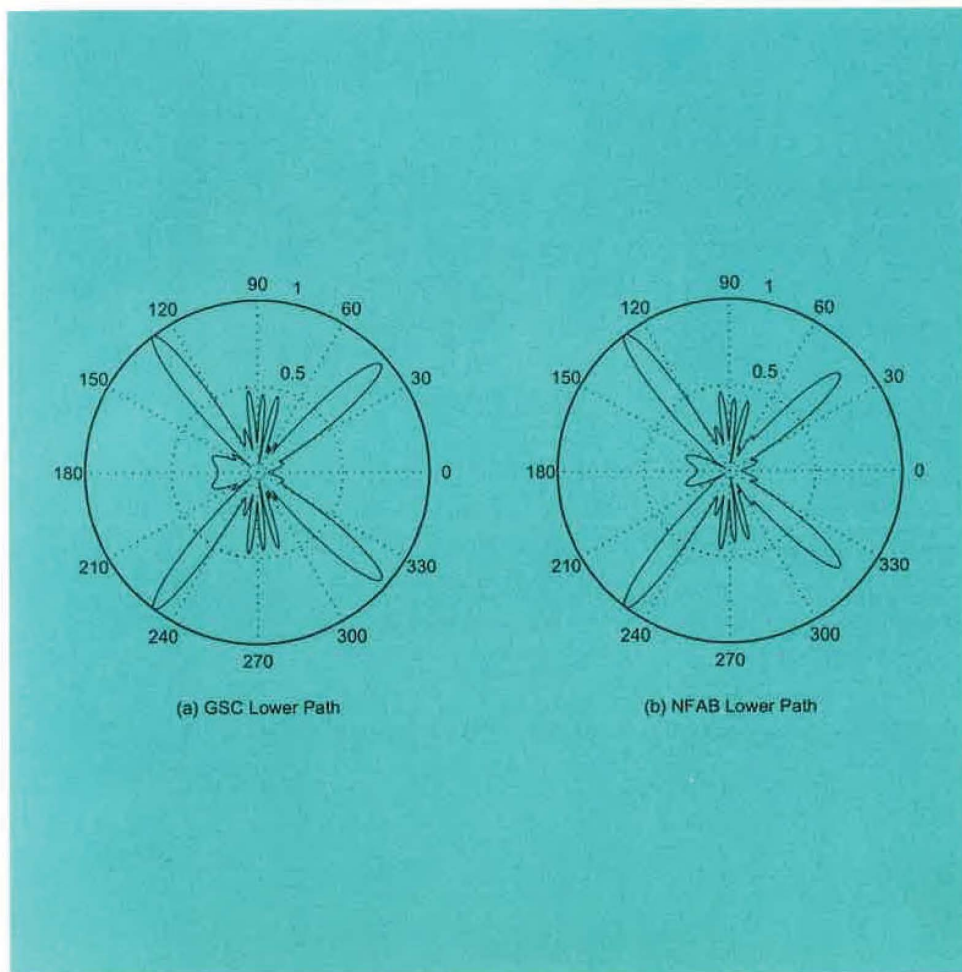
Volume 12, Number 1
January 2002

IDEAL[®] First
Articles published online first
<http://www.idealibrary.com>

TX 5-483-033



TX0005483033



Editors

Jim Schroeder
Joe Campbell

ISSN 1051-2004



**ACADEMIC
PRESS**

An Elsevier Science Imprint

Digital Signal Processing

A Review Journal

Editors

Jim Schroeder

*SPRI/CSSIP
Adelaide, SA, Australia
E-mail: schroeder@cssip.edu.au*

Joe Campbell

*M.I.T. Lincoln Laboratory
Lexington, Massachusetts
E-mail: j.campbell@ieee.org*

Editorial Board

Maurice Bellanger

*CNAM
Paris, France*

Robert E. Bogner

*University of Adelaide
Adelaide, SA, Australia*

Johann F. Böhme

*Ruhr-Universität Bochum
Bochum, Germany*

James A. Cadzow

*Vanderbilt University
Nashville, Tennessee*

G. Clifford Carter

*NUWC
Newport, Rhode Island*

A. G. Constantinides

*Imperial College
London, England*

Petar M. Djuric

*State University of New York
Stony Brook, New York*

Anthony D. Fagan

*University College Dublin
Dublin, Ireland*

Sadaoki Furui

*Tokyo Institute of Technology
Tokyo, Japan*

John E. Hershey

*General Electric Company
Schenectady, New York*

B. R. Hunt

*University of Arizona
Tucson, Arizona*

James F. Kaiser

*Duke University
Durham, North Carolina*

R. Lynn Kirlin

*University of Victoria
Victoria, British Columbia, Canada*

Ercan Kuruoğlu

*Istituto di Elaborazione della Informazione
Ghezzano, Italy*

Meemong Lee

*Jet Propulsion Laboratory
Pasadena, California*

Petre Stoica

*Uppsala University
Uppsala, Sweden*

Mati Wax

*Wavion, Ltd
Yoqneam, Isreal*

Rao Yarlagadda

*Oklahoma State University
Stillwater, Oklahoma*

Cover photo. Lower path directivity pattern at 5000 Hz. See the article by McCowan, Moore, and Sridharan in this issue.



Digital Signal Processing

Volume 12, Number 1, January 2002

© 2002 Elsevier Science (USA)

All Rights Reserved

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopy, recording, or any information storage and retrieval system, without permission in writing from the Publisher. *Exceptions:* Explicit permission from Academic Press is not required to reproduce a maximum of two figures or tables from an Academic Press article in another scientific or research publication provided that the material has not been credited to another source and that full credit to the Academic Press article is given. In addition, authors of work contained herein need not obtain permission in the following cases only: (1) to use their original figures or tables in their future works; (2) to make copies of their papers for use in their classroom teaching; and (3) to include their papers as part of their dissertations.

The appearance of the code at the bottom of the first page of an article in this journal indicates the Publisher's consent that copies of the article may be made for personal or internal use, or for the personal or internal use of specific clients. This consent is given on the condition, however, that the copier pay the stated per copy fee through the Copyright Clearance Center, Inc. (222 Rosewood Drive, Danvers, Massachusetts 01923), for copying beyond that permitted by Sections 107 or 108 of the U.S. Copyright Law. This consent does not extend to other kinds of copying, such as copying for general distribution, for advertising or promotional purposes, for creating new collective works, or for resale. Copy fees for pre-2002 articles are as shown on the article title pages; if no fee code appears on the title page, the copy fee is the same as those for current articles.

1051-2004/02 \$35.00

MADE IN THE UNITED STATES OF AMERICA

This journal is printed on acid-free paper.



DIGITAL SIGNAL PROCESSING (ISSN 1051-2004)

Published quarterly by Elsevier Science.
Editorial and Production Offices: 525 B Street, Suite 1900, San Diego, CA 92101-4495
Accounting and Circulation Offices: 6277 Sea Harbor Drive, Orlando, FL 32887-4900
2002: Volume 12, Price \$343.00 U.S.A. and Canada; \$374.00 all other countries
All prices include postage and handling

Information concerning personal subscription rates may be obtained by writing to the Publishers. All correspondence, permission requests, and subscription orders should be addressed to the office of the Publishers at 6277 Sea Harbor Drive, Orlando, FL 32887-4900 (telephone: 407-345-2000). Send notices of change of address to the office of the Publishers at least 6 to 8 weeks in advance. Please include both old and new addresses. POSTMASTER: Send changes of address to *Digital Signal Processing*, 6277 Sea Harbor Drive, Orlando, FL 32887-4900.

This material may be protected by Copyright law (Title 17 U.S. Code)

Near-field Adaptive Beamformer for Robust Speech Recognition

Iain A. McCowan, Darren C. Moore, and S. Sridharan

Speech Research Laboratory, RCSAVT, School of EESE, Queensland University of Technology, GPO Box 2434, Brisbane QLD 4001, Australia

E-mail: iain@ieee.org; moore@idiap.ch; s.sridharan@qut.edu.au

McCowan, I. A., Moore, D. C., and Sridharan, S., Near-field Adaptive Beamformer for Robust Speech Recognition, *Digital Signal Processing* 12 (2002) 87–106.

This paper investigates a new microphone array processing technique specifically for the purpose of speech enhancement and recognition. The main objective of the proposed technique is to improve the low frequency directivity of a conventional adaptive beamformer, as low frequency performance is critical in speech processing applications. The proposed technique, termed *near-field adaptive beamforming* (NFAB), is implemented using the standard generalized sidelobe canceler (GSC) system structure, where a near-field superdirective (NFSD) beamformer is used as the fixed upper-path beamformer to improve the low frequency performance. In addition, to minimize signal leakage into the adaptive noise canceling path for near-field sources, a compensation unit is introduced prior to the blocking matrix. The advantage of the technique is verified by comparing the directivity patterns with those of conventional filter-sum, NFSD, and GSC systems. In speech enhancement and recognition experiments, the proposed technique outperforms the standard techniques for a near-field source in adverse noise conditions. © 2002 Elsevier Science (USA)

Key Words: microphone array; beamforming; near-field; adaptive; superdirectivity; speech recognition.

1. INTRODUCTION

Currently, much research is being undertaken to improve the robustness of speech recognition systems in real environments. This paper focuses on the use of a microphone array to enhance the noisy input speech signal prior to recognition. While the use of microphone arrays for speech recognition has been studied for some time by a number of researchers, a persistent problem has been the poor low frequency directivity of conventional beamforming techniques with



practical array dimensions. Low frequency performance is critical for speech processing applications, as significant speech energy is located below 1 kHz.

By explicitly maximizing the array gain, superdirective beamforming techniques are able to achieve greater directivity than conventional techniques with closely spaced sensor arrays [1]. This directivity generally comes at the expense of a controlled reduction in the white noise gain of the array. Recent work has demonstrated the suitability of superdirective beamforming for speech enhancement and recognition tasks [2, 3]. By employing a spherical propagation model in its formulation, rather than assuming a far-field model, *near-field superdirectivity* (NFSD) succeeds in achieving high directivity at low frequencies for near-field speech sources in diffuse noise conditions [4]. In previous work, near-field superdirectivity has been shown to lead to good speech recognition performance in high noise conditions for a near-field speaker [5].

Superdirective techniques are typically formulated assuming a diffuse noise field. While this is a good approximation to many practical noise conditions, further noise reduction would result from a more accurate model of the actual noise conditions during operation. Adaptive array processing techniques continually update their parameters based on the statistics of the measured input noise. The *generalized sidelobe canceler* (GSC) [6] presents a structure that can be used to implement a variety of adaptive beamformers. A block diagram of the basic GSC system is shown in Fig. 1. The GSC separates the adaptive beamformer into two main processing paths—a standard fixed beamformer, \mathbf{w} , with L constraints on the desired signal response, and an adaptive path, consisting of a blocking matrix, \mathbf{B} , and a set of adaptive filters, \mathbf{a} . As the desired signal has been constrained in the upper path, the lower path filters can be updated using an unconstrained adaptive algorithm, such as the least-mean-square (LMS) algorithm.

While the theory of adaptive techniques promises greater signal enhancement, this is not always the case in real situations. A common problem with the GSC system is leakage of the desired signal through the blocking matrix, resulting in signal degradation at the beamformer output. This is particularly problematic for broadband signals, such as speech, and especially for speech recognition applications where signal distortion is critical.

In this paper we propose a system that is suited to speech enhancement in a practical near-field situation, having both the good low frequency performance of near-field superdirectivity and the adaptability of a GSC system, while taking

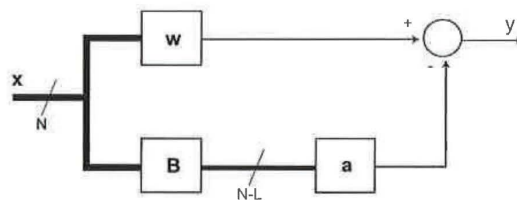


FIG. 1. Generalized sidelobe canceler structure.

care to minimize the problem of signal degradation for near-field sources. We begin by formulating a concise model for near-field sound propagation in Section 2. This model is then used in Section 3 to develop the proposed *near-field adaptive beamforming* (NFAB) technique. To demonstrate the benefit of the technique over existing methods, an experimental evaluation assessing directivity patterns, speech enhancement performance, and speech recognition performance is detailed in Sections 4 and 5.

2. NEAR-FIELD SOUND PROPAGATION MODEL

In sensor array applications, a succinct means of characterizing both the array geometry and the location of a signal source is via the *propagation vector*. The propagation vector concisely describes the theoretical propagation of the signal from its source to each sensor in the array. In this section, we develop an expression for the propagation vector of a sound source located in the near-field of a microphone array using a spherical propagation model. This expression is then used in the formulation of the proposed near-field adaptive beamformer in the following sections.

Many microphone array processing techniques assume a planar signal wavefront. This is reasonable for a far-field source, but when the desired source is close to the array a more accurate spherical wavefront model must be employed. For a microphone array of length L , a source is considered to be in the near-field if $r < 2L^2/\lambda$, where r is the distance to the source and λ is the wavelength.

We define the reference microphone as the origin of a 3-dimensional vector space, as shown in Fig. 2. The position vector for a source in direction (θ_s, ϕ_s) , at distance r_s from the reference microphone, is denoted \mathbf{p}_s and is given by:

$$\mathbf{p}_s = r_s [\hat{\mathbf{x}}, \hat{\mathbf{y}}, \hat{\mathbf{z}}] \begin{bmatrix} \cos \theta_s \sin \phi_s \\ \sin \theta_s \sin \phi_s \\ \cos \phi_s \end{bmatrix}. \quad (1)$$

The microphone position vectors, denoted as \mathbf{p}_i ($i = 1, \dots, N$), are similarly defined. The distance from the source to microphone i is thus

$$d_i = \|\mathbf{p}_s - \mathbf{p}_i\|, \quad (2)$$

where $\|\cdot\|$ is the Euclidean vector norm.

In such a model, the differences in distance to each sensor can be significant for a near-field source, resulting in phase misalignment across sensors. The difference in propagation time to each microphone with respect to the reference microphone ($i = 1$) is given by

$$\tau_i = \frac{d_i - d_1}{c}, \quad (3)$$

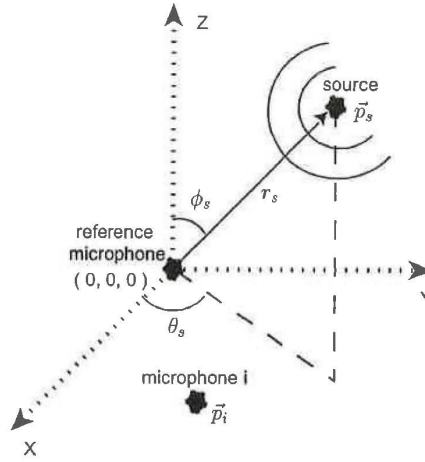


FIG. 2. Near-field propagation model.

where $c = 340 \text{ ms}^{-1}$ for sound. In addition, the wavefront amplitude decays at a rate proportional to the distance traveled. The resulting amplitude differences across sensors are negligible for far-field sources, but can be significant in the near-field case. The microphone attenuation factors, with respect to the amplitude on the reference microphone, are given by

$$\alpha_i = \frac{d_1}{d_i}. \quad (4)$$

Thus, if $x_1(f)$ is the desired source at the reference microphone, the signal on the i th microphone is given by

$$x_i(f) = \alpha_i x_1(f) e^{-j2\pi f \tau_i}. \quad (5)$$

Consequently, we define the near-field propagation vector for a source at distance r and direction (θ, ϕ) as

$$\mathbf{d}(f, r, \theta, \phi) = [\alpha_1 e^{-j2\pi f \tau_1} \dots \alpha_i e^{-j2\pi f \tau_i} \dots \alpha_N e^{-j2\pi f \tau_N}]^T. \quad (6)$$

3. NEAR-FIELD ADAPTIVE BEAMFORMING

The proposed system structure is shown in Fig. 3. The objective of the proposed technique is to add the benefit of good low frequency directivity to a standard adaptive beamformer, as low frequency performance is critical in speech processing applications. The upper path consists of a fixed near-field superdirective beamformer, while the lower path contains a near-field compensation unit, a blocking matrix and an adaptive noise canceling filter. The principal components of the system are discussed in the following sections.

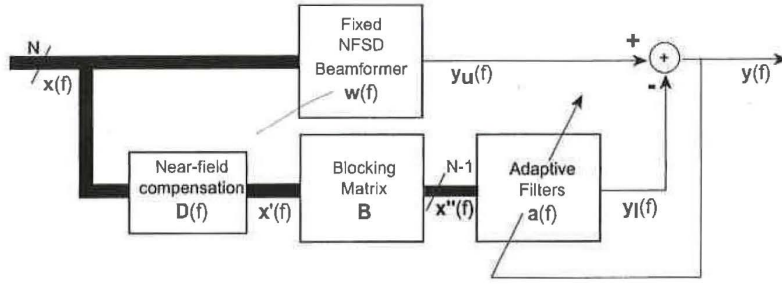


FIG. 3. Near-field adaptive beamformer.

Section 3.1 gives an explanation of the near-field superdirective beamformer. Section 3.2 proposes the inclusion of a near-field compensation unit in the adaptive sidelobe canceling path and examines its effect on reducing signal distortion at the output. Once this near-field compensation has been performed, a standard generalized sidelobe canceling blocking matrix and adaptive filters can be applied to reduce the output noise power, as discussed in Section 3.3.

3.1. Near-field Superdirective Beamformer

Superdirective beamforming techniques are based upon the maximization of the array gain, or directivity index. The array gain is defined as the ratio of output signal-to-noise ratio to input signal-to-noise ratio and for the general case can be expressed in matrix notation as [1]

$$G(f) = \frac{\mathbf{w}(f)^H \mathbf{P}(f) \mathbf{w}(f)}{\mathbf{w}(f)^H \mathbf{Q}(f) \mathbf{w}(f)}, \quad (7)$$

where $\mathbf{w}(f)$ is a column vector of channel gains,

$$\mathbf{w}(f) = [w_1(f) \dots w_i(f) \dots w_N(f)]^T, \quad (8)$$

$(\)^H$ is the complex conjugate transpose operator, and $\mathbf{P}(f)$ and $\mathbf{Q}(f)$ are the cross-spectral density matrices of the signal and noise respectively. In practical speech processing applications the form of the signal and noise cross-spectral density matrices is generally unknown and must be estimated, either from mathematical models (fixed beamformers) or from the statistics of the multichannel inputs (adaptive beamformers). Superdirective beamformers are calculated based on assumed mathematical models for the $\mathbf{P}(f)$ and $\mathbf{Q}(f)$ matrices.

When the desired signal is known to emanate from a single source at location (r_s, θ_s, ϕ_s) , the signal cross-spectral matrix \mathbf{P} simplifies to the propagation vector of the source, and the array gain can be expressed as

$$G(f) = \frac{|\mathbf{w}(f)^H \mathbf{d}(f, r_s, \theta_s, \phi_s)|^2}{\mathbf{w}(f)^H \mathbf{Q}(f) \mathbf{w}(f)}, \quad (9)$$

where $\mathbf{d}(f, r, \theta, \phi)$ is the propagation vector for the desired source, as defined in Eq. (6).

A diffuse (spherically isotropic) noise field is often a good approximation for many practical situations, particularly in reverberant closed spaces, such as in a car or an office [7, 8]. For diffuse noise, the noise cross-spectral density matrix \mathbf{Q} can be formulated as

$$\mathbf{Q}(f) = \frac{1}{4\pi} \int_{\phi} \int_{\theta} \mathbf{d}(f, \theta, \phi) \mathbf{d}(f, \theta, \phi)^H \sin \theta d\theta d\phi, \quad (10)$$

where $\mathbf{d}(f, \theta, \phi)$ is the propagation vector of a far-field noise source ($r \gg 2L^2/\lambda$) in direction (θ, ϕ) .

The superdirectivity problem is thus formulated as:

$$\max_{\mathbf{w}(f)} \frac{|\mathbf{w}(f)^H \mathbf{d}(f, r_s, \theta_s, \phi_s)|^2}{\mathbf{w}(f)^H \mathbf{Q}(f) \mathbf{w}(f)}. \quad (11)$$

By using a spherical propagation model to formulate the propagation vector, \mathbf{d} , the standard superdirective formulation can be optimized for a near-field source [9, 4]. As such, the only difference in the calculation of the standard and near-field superdirective channel filters is the form of the propagation vector, \mathbf{d} . For a near-field source, the assumption of plane wave (far-field) propagation leads to errors in the array response to the desired signal due to curvature of the direct wavefront. A thorough discussion of the use of a near-field model for superdirective microphone arrays is given by Ryan and Goubran [9].

Cox [10] gives the general superdirective filter solution subject to

1. L linear constraints, $\mathbf{C}(f)^H \mathbf{w}(f) = \mathbf{g}(f)$ (explained below); and
2. a constraint on the maximum white noise gain, $\mathbf{w}(f)^H \mathbf{w}(f) = \delta^{-2}$, where δ^2 is the desired white noise gain.

as

$$\mathbf{w}(f) = \{\mathbf{Q}(f) + \epsilon \mathbf{I}\}^{-1} \mathbf{C}(f) \{\mathbf{C}(f)^H [\mathbf{Q}(f) + \epsilon \mathbf{I}]^{-1} \mathbf{C}(f)\}^{-1} \mathbf{g}(f), \quad (12)$$

where ϵ is a Lagrange multiplier that is iteratively adjusted to satisfy the white noise gain constraint. The white noise gain is the array gain for spatially white (incoherent) noise; that is, $\mathbf{Q}(f) = \mathbf{I}$. A constraint on the white noise gain is necessary as an unconstrained superdirective solution will in fact result in significant gain to any incoherent noise, particularly at low frequencies. Cox [10] states that the technique of adding a small amount to each diagonal matrix element prior to inversion is in fact the optimum means of solving this problem. A study of the relationship between the multiplier ϵ and the desired white noise gain δ^2 , shows that the white noise gain increases monotonically with increasing ϵ . One possible means of obtaining the desired value of ϵ is thus an iterative technique employing a binary search algorithm between a specified minimum and maximum value for ϵ . The computational expense of the iterative procedure is not critical, as the beamformer filters depend only on the source

location and array geometry, and thus must only be calculated once for a given configuration.

The constraint matrix, $\mathbf{C}^H(f)$, is of order $L \times N$, where there are L linear constraints being applied, and the vector $\mathbf{g}(f)$ is a length- L column vector of constraining values. The constraints generally include one specifying unity response for the desired signal, $\mathbf{d}^H(f)\mathbf{w}(f) = 1$, and where this is the sole constraint the above solution can be simplified by substituting $\mathbf{C}(f) = \mathbf{d}(f)$ and $\mathbf{g}(f) = 1$, giving

$$\mathbf{w}(f) = \frac{[\mathbf{Q}(f) + \epsilon\mathbf{I}]^{-1}\mathbf{d}(f)}{\mathbf{d}(f)^H[\mathbf{Q}(f) + \epsilon\mathbf{I}]^{-1}\mathbf{d}(f)}. \quad (13)$$

Once the optimal filters $\mathbf{w}(f)$ have been calculated, the near-field superdirective beamformer output is calculated as

$$y_u(f) = \mathbf{w}(f)^H \mathbf{x}(f), \quad (14)$$

where $\mathbf{x}(f)$ is the N -channel input column vector

$$\mathbf{x}(f) = [x_1(f) \dots x_i(f) \dots x_N(f)]^T. \quad (15)$$

3.2. Near-field Compensation Unit

The first element in the adaptive path of standard GSC is the blocking matrix [6]. Its purpose is to block the desired signal from the adaptive noise estimate. To ensure complete blocking, the desired signal must both be time aligned and have equal amplitudes across all channels. If this is the case, cancellation occurs if each row of the blocking matrix sums to zero, and all rows are linearly independent.

For a near-field desired source, to align the desired signal on all channels, a near-field compensation must first be applied to the input channels prior to blocking. To ensure full cancellation we need to compensate for both phase misalignment and amplitude scaling of the desired signal across sensors. We define the diagonal matrix

$$\mathbf{D}(f) = [\text{diag}(\mathbf{d}(f))]^{-1}, \quad (16)$$

where $\mathbf{d}(f)$ is the near-field propagation vector from Eq. (6). In this paper we define the diagonal operator, $\text{diag}(\cdot)$, to produce a diagonal matrix from a vector parameter. Conversely, if invoked with a matrix parameter, it produces a row vector corresponding to the matrix diagonal. The near-field compensation can be applied as

$$\mathbf{x}'(f) = \mathbf{D}(f)\mathbf{x}(f). \quad (17)$$

Once this near-field compensation has been performed, a standard GSC blocking matrix can be employed to block the desired signal from the adaptive path.

The inclusion of this compensation unit is critical for a near-field desired signal. Without compensation for both phase and amplitude differences between sensors, blocking of the desired signal will not be ensured, leading to signal

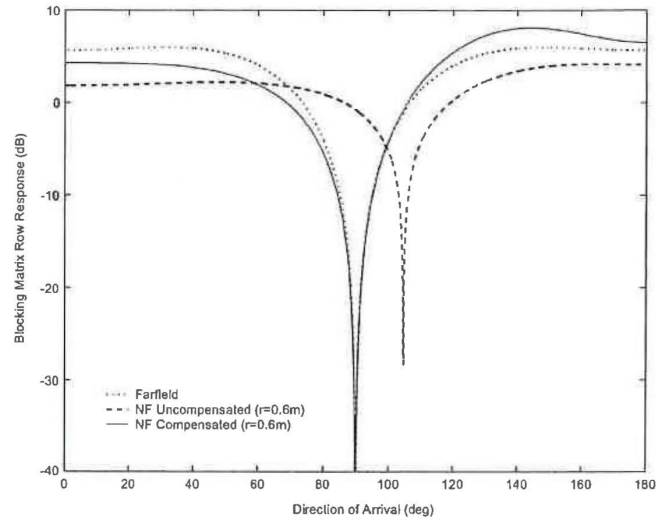


FIG. 4. Comparison of blocking matrix row beam-patterns.

cancellation at the output. The near-field compensation effectively ensures that a true null exists in the beam-pattern of each blocking matrix row in the direction *and* distance corresponding to the desired source. To illustrate, Fig. 4 shows the directivity pattern at 2 kHz for the first row in the blocking matrix using the array shown in Fig. 5, with the desired source directly in front of the center microphone at a distance of 0.6 m. The figure shows the compensated response in the far- and near-fields, as well as the uncompensated near-field response. It is clear that the uncompensated system will allow a high degree of signal leakage into the adaptive path as it blocks noise sources rather than the desired signal.

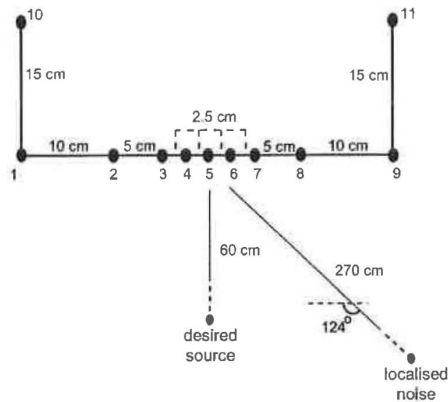


FIG. 5. Experimental configuration.

3.3. Blocking Matrix and Adaptive Noise Canceling Filter

The blocking matrix and adaptive noise canceling filters are taken from the standard GSC technique [6]. The order of the blocking matrix is $N \times (N - L)$, where there are L constraints applied in the fixed upper path beamformer. Generally only a unity constraint on the desired signal is specified, and the standard $N \times (N - 1)$ Griffiths–Jim blocking matrix is used:

$$\mathbf{B} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ -1 & 1 & \vdots & \vdots & \vdots \\ 0 & -1 & \ddots & 0 & 0 \\ 0 & 0 & \ddots & 1 & 0 \\ \vdots & \vdots & \vdots & -1 & 1 \\ 0 & 0 & \dots & 0 & -1 \end{bmatrix}. \quad (18)$$

The output of the blocking matrix is calculated as

$$\mathbf{x}''(f) = \mathbf{B}^H \mathbf{x}'(f), \quad (19)$$

where $\mathbf{x}''(f)$ is an $(N - 1)$ -length column vector. Defining the $(N - 1)$ -length adaptive filter column vector as

$$\mathbf{a}(f) = [a_1(f) \dots a_i(f) \dots a_{N-1}(f)]^T, \quad (20)$$

the output of the lower path is given as

$$y_l(f) = \mathbf{a}(f)^H \mathbf{x}''(f). \quad (21)$$

The NFAB output is then calculated from the upper and lower path outputs as

$$y(f) = y_u(f) - y_l(f) \quad (22)$$

and the adaptive filters are updated using the standard unconstrained LMS algorithm

$$\mathbf{a}_{k+1}(f) = \mathbf{a}_k(f) + \mu \mathbf{x}''_k(f) y_k(f), \quad (23)$$

where μ is the adaptation step size and k denotes the current frame.

3.4. Summary of Technique

In summary, the proposed NFAB technique is characterized by the series of equations

$$y_u(f) = \mathbf{w}(f)^H \mathbf{x}(f) \quad (24a)$$

$$\mathbf{x}''_k = \mathbf{B}^H \mathbf{D}(f) \mathbf{x}(f) \quad (24b)$$

$$y_l(f) = \mathbf{a}(f)^H \mathbf{x}_k''(f) \quad (24c)$$

$$y(f) = y_u(f) - y_l(f) \quad (24d)$$

$$\mathbf{a}_{k+1}(f) = \mathbf{a}_k(f) + \mu \mathbf{x}_k''(f) y_k(f), \quad (24e)$$

where all terms have been defined in the preceding discussion.

4. EXPERIMENTAL CONFIGURATION

For the experimental evaluation in this paper, we used the 11 element array shown in Fig. 5. The array consists of a nine element broadside array, with an additional two microphones situated directly behind the end microphones. The total array is 40 cm wide and 15 cm deep in the horizontal plane. The broadside microphones are arranged according to a standard broadband subarray design, where different subarrays are used for different frequency ranges for the fixed upper path beamformer. The two endfire microphones are included for use by the near-field superdirective beamformer in the low frequency range. The four subarrays are thus

- ($f < 1$ kHz): microphones 1–11;
- ($1 \text{ kHz} < f < 2$ kHz): microphones 1, 2, 5, 8, and 9;
- ($2 \text{ kHz} < f < 4$ kHz): microphones 2, 3, 5, 7, and 8; and
- ($4 \text{ kHz} < f < 8$ kHz): microphones 3–7.

The array was situated in a computer room, with different sound source locations, as shown in Fig. 5. The two sound sources were

1. the desired speaker situated 60 cm from the center microphone, directly in front of the array; and
2. a localized noise source at an angle of 124° and a distance of 270 cm from the array.

Impulse responses of the acoustic path between each source and microphone were measured from multichannel recordings made in the room with the array using the maximum length sequence technique detailed in Rife and Vanderkooy [11]. As the impulse responses were calculated from real recordings made simultaneously across all input channels, they take into account the real acoustic properties of the room and the array. The multichannel desired speech and localized noise microphone inputs were then generated by convolving the original single-channel speech and noise signals with these impulse responses. In addition, a real multichannel background noise recording of normal operating conditions was made in the room with other workers present. This recording is referred to in the experiments as the ambient noise signal and is approximately diffuse in nature. It consists mainly of computer noise, a variable level of background speech, and noise from an air-conditioning unit. The ambient noise effectively represents a diffuse noise field, while the localized noise represents a coherent noise source. In this paper, we specify the levels of the two different noise sources independently, as the signal to ambient-noise ratio (SANR) and

signal to localized-noise ratio (SLNR). These values are calculated as the average segmental SNR from the speech and noise input, as measured at the center microphone of the array.

In this way, realistic multichannel input signals can be simulated for specified levels of ambient and localized noise. As well as facilitating the generation of different noise conditions, simulating the multichannel inputs using the impulse response method is more practical than making real recordings for speech recognition experiments, as existing single channel speech corpora may be used.

5. EXPERIMENTAL RESULTS

This section presents the results of the experimental evaluation. The proposed NFAB technique is compared to a conventional fixed filter-sum beamformer, a fixed near-field superdirective beamformer, and a conventional GSC adaptive beamformer. These beamformers are specified in Table 1.

The techniques are first assessed in terms of the directivity pattern in order to demonstrate the advantage of the proposed NFAB over conventional beamforming techniques, particularly at low frequencies. Following this, the techniques are evaluated for speech enhancement in terms of the improvement in signal to noise ratio and the log area ratio. Finally, the techniques are compared in a hands-free speech recognition task in noisy conditions using the TIDIGITS database [12].

5.1. Directivity Analysis

As has been stated, the main objective of the proposed technique is to produce an *adaptive beamformer* that exhibits good *low frequency* performance for *near-field* speech sources. To assess the effectiveness of the proposed technique in achieving this objective, in this section we analyze the horizontal directivity pattern. The directivity of a filter-sum beamformer is expressed in matrix notation as

$$h(f, r, \theta, \phi) = \mathbf{w}_o(f)^H \mathbf{d}(f, r, \theta, \phi), \quad (25)$$

where \mathbf{w}_o is the length N channel filter vector

$$\mathbf{w}_o(f) = [w_{o,1}(f) \dots w_{o,i}(f) \dots w_{o,N}(f)]^T, \quad (26)$$

TABLE 1
Beamforming Techniques in Evaluation

Technique	Description	Filters
FS	Conventional FS beamformer	$\mathbf{w}_o(f) = [\text{diag}(\mathbf{D}(f))]^H$
NFSD	Near-field superdirective beamformer	$\mathbf{w}_o(f) = \mathbf{w}(f)$
GSC	GSC system with FS fixed upper path beamformer	$\mathbf{w}_o(f) = [\text{diag}(\mathbf{D}(f))]^H - \mathbf{D}(f)\mathbf{B}\mathbf{a}(f)$
NFAB	Near-field adaptive beamformer	$\mathbf{w}_o(f) = \mathbf{w}(f) - \mathbf{D}(f)\mathbf{B}\mathbf{a}(f)$

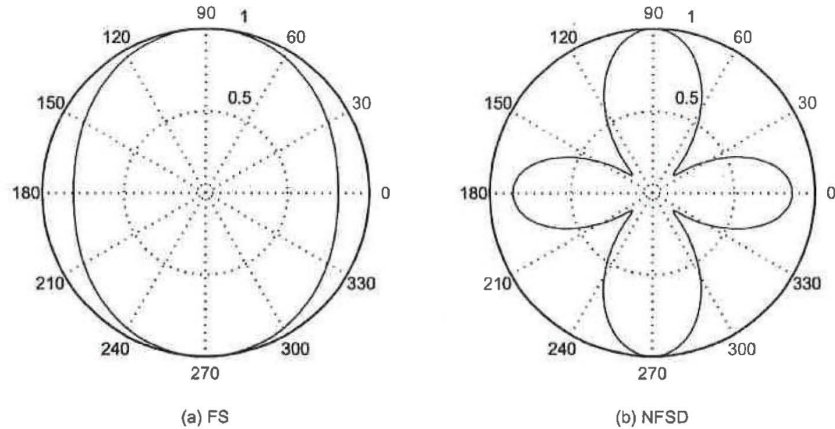


FIG. 6. Upper path directivity pattern at 300 Hz.

5.1.1. Upper path directivity. First, we seek to demonstrate the directivity improvement that NFSD achieves at low frequencies compared to a conventional filter-sum (FS) beamformer. For the FS beamformer, a common solution is to choose $\mathbf{w}_o(f) = [\text{diag}(\mathbf{D}(\mathbf{f}))]^H$. This effectively ensures that the desired signal is aligned for phase and amplitude across sensors using a spherical propagation model. For NFSD, we use the filter vector $\mathbf{w}(f)$ described in Section 3.1. Figure 6 shows the near-field directivity pattern at 300 Hz for the FS and NFSD. From these figures, it is clear that the NFSD technique results in greater directional discrimination at low frequencies compared to a conventional beamformer. At higher frequencies ($f > 1$ kHz), conventional beamformers offer reasonable directivity, and so the FS and NFSD techniques give comparable performance.

5.1.2. Lower path directivity. Second, we wish to demonstrate the effect of the noise canceling path. The directivity of the noise canceling filters can be obtained by using the channel filters $\mathbf{w}_o(f) = \mathbf{D}(f)\mathbf{B}\mathbf{a}(f)$. The blocking matrix and adaptive filters essentially implement a conventional (nonsuperdirective) beamformer that adaptively focuses on the major sources of noise. To examine the directivity of the lower path filters, the beamformer was run on an input speech signal with a white localized noise source (at the location shown in Fig. 5) added at an SLNR of 0 dB and a low level of ambient noise (SANR = 20 dB). The steady-state adaptive filter vector, $\mathbf{a}(f)$, was written to file for both the proposed NFAB technique and the conventional GSC beamformer. The near-field directivity patterns of the lower path filters are plotted in Figs. 7 and 8 for 300 and 5000 Hz, respectively. We see that the lower path adaptive filters for both beamformers converge to similar solutions in terms of directivity, producing a main lobe in the direction of the coherent noise source ($\approx 124^\circ$ from Fig. 5), as well as a null in the location of the desired speaker. As expected, the directivity of the adaptive path is poor at low frequencies, as seen in Fig. 7.

5.1.3. Overall beamformer directivity. Finally, we examine the directivity pattern of the overall beamformer for the NFAB and conventional adaptive

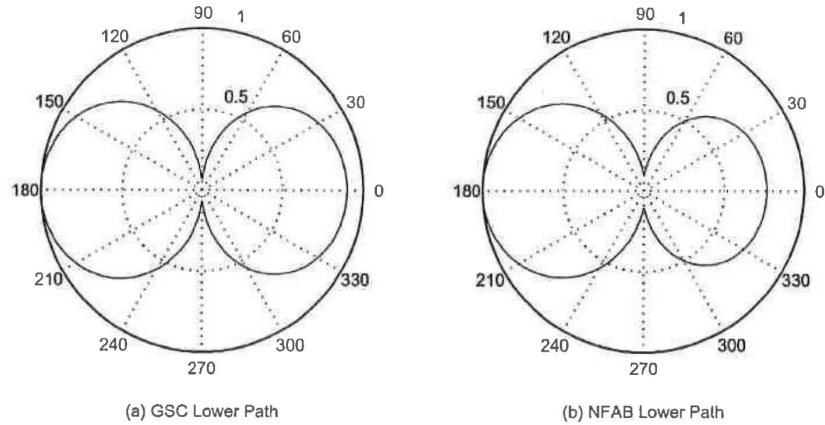


FIG. 7. Lower path directivity pattern at 300 Hz.

systems. The near-field directivity patterns at 300 Hz are shown in Fig. 9. We see that the directivity pattern of the NFAB system exhibits a true null in the direction and at the distance of the noise source, while the directivity of the conventional beamformer is too poor to significantly attenuate the noise at this frequency. At frequencies above 1 kHz the directivity performance of both techniques is comparable.

5.1.4. *Summary of beamformer directivity.* In summary we see that, in terms of directivity, the proposed NFAB system:

- outperforms the conventional FS system in terms of low frequency performance and the ability to attenuate coherent noise sources,
- outperforms the NFSD system due to the ability to attenuate coherent noise sources, and

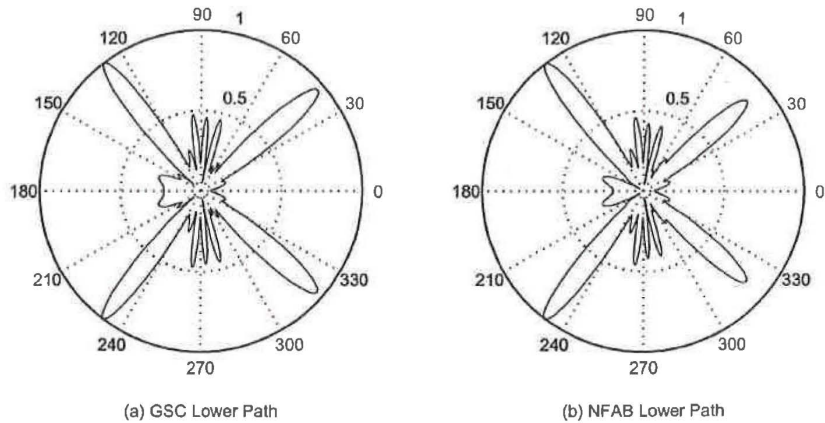


FIG. 8. Lower path directivity pattern at 5000 Hz.

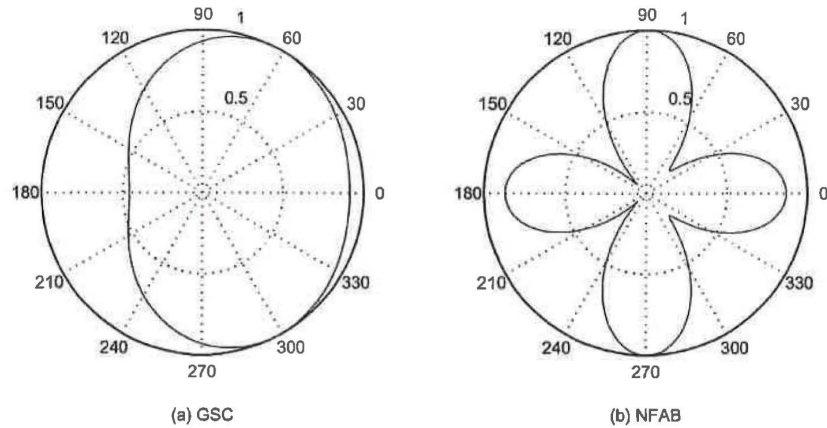


FIG. 9. Overall beamformer directivity pattern at 300 Hz.

- outperforms the conventional GSC system in terms of low frequency performance.

In this way, we see that the proposed system succeeds in meeting the stated objectives and should therefore demonstrate improved performance in speech processing applications.

5.2. Speech Enhancement Analysis

The signal plots in Fig. 10 give an indication of the level of enhancement achieved by the NFAB technique. For the desired speech signal, we used a segment of speech from the TIDIGITS database corresponding to the digit sequence *one-nine-eight-six*. Ambient noise was added at an SANR level of 10 dB, and a localized white noise signal was added at an SLNR level of 0 dB. The plots indicate that NFAB succeeds in reducing the noise level with negligible distortion to the desired signal.

To better measure the level of enhancement, objective speech measures were used to compare the different techniques. Two measures were used, these being the SNR improvement and the log area ratio distortion measure. The *SNR improvement* is defined as the difference in SNR at the array output and input. As the true SNR cannot be measured, it is estimated as the average segmental signal-plus-noise to noise ratio. While the signal to noise ratio is a useful measure for assessing noise reduction, it does not necessarily give a good indication of how much distortion has been introduced to the desired speech signal. The *log area ratio* (LAR) measure of speech quality is more highly correlated with perceptual intelligibility in humans [13]. The log area ratio measure for a frame of speech is calculated as

$$\text{LAR}(n) = \left| \frac{1}{P} \sum_{i=1}^P \left[\log \frac{1+r_o(i)}{1-r_o(i)} - \log \frac{1+r_p(i)}{1-r_p(i)} \right] \right|^{1/2}, \quad (27)$$

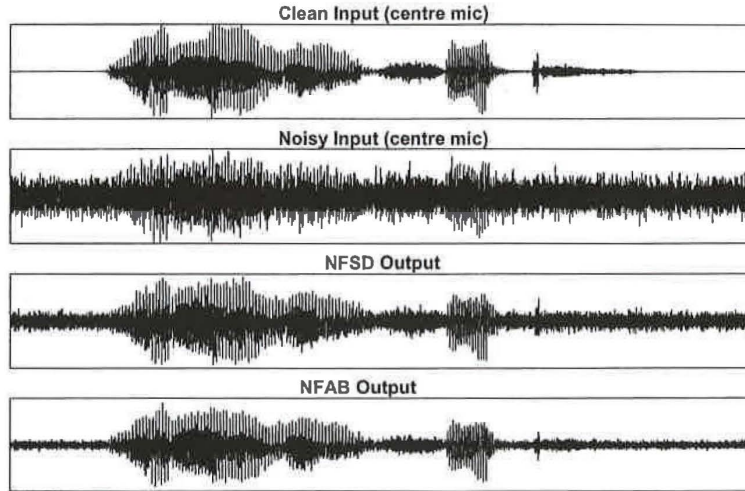


FIG. 10. Sample enhanced signal.

where n is the frame number, and r_o and r_p are the original and processed P th-order linear predictive coefficients of the n th frame, respectively. The overall log area ratio distortion measure for the signal is calculated as the average distortion over all input frames.

A set of experiments was conducted in which the localized white noise was replaced with a localized speech-like noise source taken from the NOISEX database [14]. This is essentially a white noise signal that has been shaped with a speech-like spectral envelope and thus represents a more realistic noise scenario than white noise. The signal to localized noise ratio (SLNR) was varied from 20 to 0 dB, with the ambient noise present at a constant SANR level of 10 dB. The output signal to noise ratio improvement and log area ratios are given in Tables 2 and 3 for the different enhancement techniques.¹ The measures have been averaged over 10 randomly chosen speech segments taken

TABLE 2
Signal to Noise Ratio Improvement (SANR = 10 dB)

Technique	SLNR (dB)				
	0	5	10	15	20
FS	0.5	0.4	0.3	0.1	0.1
NFSD	1.4	1.6	1.1	0.5	0.2
GSC	1.6	1.8	1.9	2.5	3.3
NFAB	5.5	5.9	6.4	7.5	7.9

¹ Sample sound files are also available at <http://www.speech.qut.edu.au/pages/people/mccowani>.

TABLE 3
Log Area Ratio: (SANR = 10 dB)

Technique	SLNR (dB)				
	0	5	10	15	20
Noisy input	3.6	3.3	2.9	2.6	2.5
FS	3.1	2.7	2.5	2.4	2.3
NFSD	2.8	2.3	2.1	2.0	1.9
GSC	2.6	2.1	1.8	1.7	1.6
NFAB	2.5	1.9	1.6	1.6	1.4

from different speakers in the TIDIGITS database. These results are plotted in Fig. 11.

The SNR results show that the proposed NFAB gives considerably greater noise reduction compared to the FS, NFSD, and GSC techniques, providing approximately 6–8 dB of SNR improvement compared to the noisy input signal. Even with a relatively low level of localized noise (high SLNR), the NFAB technique offers significantly greater noise reduction than these other methods. In addition, the proposed technique gives less distortion than the other techniques, as measured by the LAR. As would be expected, the fixed NFSD technique gives slightly less distortion than the adaptive GSC technique.

From these results we see that, in a high level of diffuse and coherent noise with a near-field desired speech source, the proposed NFAB technique succeeds in significantly reducing the noise level, while minimizing the distortion to the speech signal.

5.3. Speech Recognition Analysis

The same noise scenario was used for experiments in robust speech recognition. The training and test data for the experiments were taken from the male adult portion of the TIDIGITS database. Tied-state triphone hidden Markov models and standard MFCC parameterization with energy, delta, and acceler-

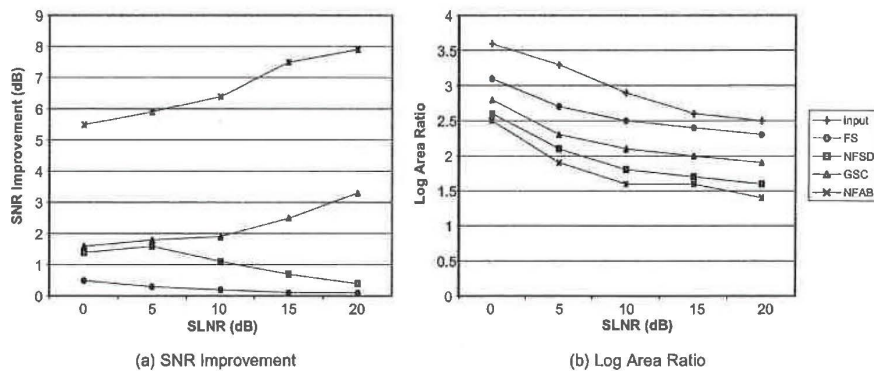


FIG. 11. Speech enhancement measures: (a) SNR improvement and (b) LAR.

TABLE 4
Speech Recognition Results: Word Recognition Rates

Technique	SLNR (dB)			
	10	5	0	-5
Noisy input	86.8	65.9	23.2	13.1
FS	89.2	81.7	62.9	36.4
NFSD	97.7	93.2	77.2	45.4
GSC	88.8	83.8	73.8	56.8
NFAB	98.2	96.7	91.1	76.7

ation coefficients were used. The models were trained with the clean input to the center microphone and then refined using MAP adaptation to better match the noisy environment. The noise segments used in the adaptation process were taken from a separate recording made in the room. The recognition results are given as percentage word recognition rates in Table 4 and shown graphically in Fig. 12.

The results clearly show that NFAB gives excellent robustness to adverse noise conditions in a near-field speech recognition application. The results at low noise levels show that the baseline recognition system is already quite robust to noise, due to the use of MAP adaptation. At more realistic noise levels, however, unenhanced performance is clearly unsatisfactory. For example, at an SLNR of 0 dB and SANR of 10 dB, the word *error* rate for the unprocessed input is 76.8%. While standard GSC and NFSD are able to reduce this to 26.2 and 22.8%, respectively, the proposed NFAB technique succeeds in reducing the error rate to 8.9%. As would be expected, the figure shows that NFAB offers

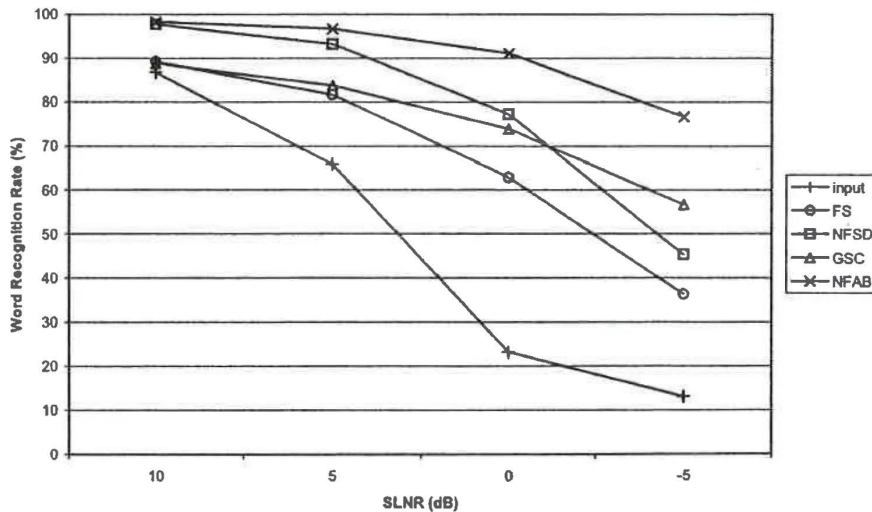


FIG. 12. Speech recognition results: Word recognition rates.

similar performance to NFSD when the noise is approximately diffuse (high SLNR) and demonstrates improved ability to attenuate any coherent noise sources (low SLNR) due to the GSC-style adaptive noise canceling path. The recognition performance of NFAB is seen to be similar to NFSD at low levels of coherent noise (high SLNR) and degrades at a rate comparable to GSC with increasing levels of coherent noise.

It is apparent from these results that NFAB is an enhancement technique that is well suited to speech recognition. The experimental results for both speech enhancement and recognition demonstrate that for an adaptive beamformer to be applicable in speech processing applications, it should exhibit good directivity at low frequencies and take care to minimize any signal degradation.

6. CONCLUSIONS

A new microphone array processing technique designed specifically for near-field speech processing applications has been proposed, termed near-field adaptive beamforming (NFAB). The technique incorporates a fixed near-field superdirective beamformer into a GSC-style adaptive beamforming structure and as such exhibits the benefits of good low frequency performance and the ability to adaptively attenuate coherent noise signals. Distortion due to the adaptive noise canceling path is minimized by the introduction of a near-field compensation unit.

Two major problems with common conventional microphone array techniques are their poor low frequency performance and the introduction of signal distortion in adaptive techniques. By taking care to address both these issues, the proposed NFAB technique succeeds in significantly outperforming conventional beamforming techniques in terms of objective speech quality measures and speech recognition results in both diffuse and coherent noise.

Speech enhancement results indicate that NFAB succeeds in significantly reducing the output noise power, while also minimizing the distortion to the desired signal. These characteristics make it ideal as an enhancement technique for robust speech recognition. In a high noise configuration, with a signal to localized noise ratio of 0 dB, and a signal to ambient noise ratio of 10 dB, the proposed technique succeeds in increasing the recognition rate from 23.2 to 91.1%. For the same configuration, near-field superdirectivity and conventional GSC only achieve 77.2 and 73.8%, respectively.

In summary, near-field adaptive beamforming has been shown to be a speech enhancement technique that produces a high quality, highly intelligible signal for applications requiring hands-free speech acquisition where the desired speaker is in the array's near-field.

ACKNOWLEDGMENTS

The near-field superdirective technique was developed at France Telecom R&D [4]. The initial application of these techniques in a speech recognition application was researched by the author during a six month research period at France Telecom R&D in 1999 [5], under the supervision of Yannick Mahieux. The authors specifically wish to acknowledge Claude Marro for his helpful and insightful review of this paper.

REFERENCES

1. Cox, H., Zeskind, R., and Kooij, T., Practical supergain. *IEEE Trans. Acoustics, Speech Signal Process.* **34** (1986), 393–398.
2. Bitzer, J., Simmer, K. U., and Kammeyer, K., Multi-microphone noise reduction techniques for hands-free speech recognition—a comparative study. In *Robust Methods for Speech Recognition in Adverse Conditions (ROBUST-99)*, Tampere, Finland, May 1999, pp. 171–174.
3. Doerbecker, M., Speech enhancement using small microphone arrays with optimized directivity. In *Proc. Int. Workshop on Acoustic Echo and Noise Control*, September 1997, pp. 100–103.
4. Täger, W., Near field superdirectivity (NFS). In *Proceedings of ICASSP'98*, 1998, pp. 2045–2048.
5. McCowan, I., Marro, C., and Mauuary, L., Robust speech recognition using near-field superdirective beamforming with post-filtering. In *Proceedings of ICASSP'2000*, 2000, Vol. 3, pp. 1723–1726.
6. Griffiths, L. and Jim, C., An alternative approach to linearly constrained adaptive beamforming. *IEEE Trans. Antennas Propagation* **30** (1982), 27–34.
7. Meyer, J. and Simmer, K. U., Multi-channel speech enhancement in a car environment using wiener filtering and spectral subtraction. In *Proceedings of ICASSP'97*, 1997, Vol. 2, pp. 1167–1170.
8. Bitzer, J., Simmer, K. U., and Kammeyer, K., Theoretical noise reduction limits of the generalized sidelobe canceller (gac) for speech enhancement. In *Proceedings of ICASSP'99*, 1999, Vol. 5, pp. 2965–2968.
9. Ryan, J. G. and Goubran, R. A., Near-field beamforming for microphone arrays. In *Proceedings of ICASSP'97*, 1997, pp. 363–366.
10. Cox, H., Zeskind, R., and Owen, M., Robust adaptive beamforming. *IEEE Trans. Acoustics, Speech Signal Process.* **35** (1987), 1365–1376.
11. Rife, D. and Vanderkooy, J., Transfer-function measurement with maximum-length sequences. *J. Audio Eng. Soc.* **37** (1989), 419–444.
12. Texas Instruments and NIST. Studio quality speaker-independent connected-digit corpus (TIDIGITS). CD-ROM, February 1991. NIST Speech Discs 4-1, 4-2, and 4-3.
13. Quackenbush, S. R., Barnwell, T. P., and Clements, M. A., *Objective Measures of Speech Quality*. Prentice-Hall, Englewood Cliffs, NJ, 1988.
14. Defence Research Agency Speech Research Unit. NOISEX-92. CD-ROM, June 1992.

IAIN A. MCCOWAN received the B.Eng (Hons) and B.InfoTech from the Queensland University of Technology, Brisbane, in 1996. In February 1998 he joined the Research Concentration in Speech, Audio, and Video Technology at the Queensland University of Technology where he is currently completing his Ph.D. His main research interests are in the fields of robust speech recognition and speech enhancement using microphone arrays. He is a student member of the Institute of Electrical and Electronic Engineers.

DARREN C. MOORE received the B.Eng (Hons) and B.InfoTech from the Queensland University of Technology, Brisbane, in 1997. In February 1998 he joined the Research Concentration in Speech, Audio, and Video Technology at the Queensland University of Technology, where he is currently completing a M.Eng. His professional interests lie in the field of speech enhancement using microphone arrays and in the implementation of real-time DSP solutions.

S. SRIDHARAN obtained his B.Sc (electrical engineering) and M.Sc (communication engineering) from the University of Manchester Institute of Science and Technology, United Kingdom and Ph.D (signal processing) from the University of New South Wales, Australia. Dr. Sridharan is senior member of the IEEE, USA and a corporate member of IEE, United Kingdom and IEAust of Australia. He is currently a professor in the School of Electrical and Electronic Systems Engineering of the Queensland University of Technology (QUT) and is also the Head of the Research Concentration in Speech, Audio, and Video Technology at QUT.



Digital Signal Processing

A Review Journal

Volume 12, Number 1, January 2002

Contents

- 1 Rapid Power-Line Frequency Monitoring
 Ronald M. Adelson

- 12 A Constrained Block Iterative Algorithm for Multiple-Access
 Interference Suppression
 John F. Doherty and Yu-Tsun Hsieh

- 21 An Adaptive Combined Classifier System for Invariant
 Face Recognition
 G. A. Khuwaja

- 47 Representation and Recognition of 3D Free-Form Objects
 George Mamic and Mohammed Bennamoun

- 77 A New Algorithm for 16QAM Carrier Phase Estimation Using
 QPSK Partitioning
 Feng Rice, Mark Rice, and Bill Cowley

- 87 Near-field Adaptive Beamformer for Robust Speech
 Recognition
 Iain A. McCowan, Darren C. Moore, and S. Sridharan

- 107 Multiplierless Adaptive Filtering
 **Tamal Bose, Anand Venkatachalam,
 and Ratchaneekorn Thamvichai**

- 119 Nonlinear Least l_p -Norm Filters for Nonlinear Autoregressive
 α -Stable Processes
 Ercan E. Kuruoğlu

The homepage for Digital Signal Processing can be found at <http://www.academicpress.com/dsp>. All articles are available online at <http://www.idealibrary.com>.



1051-2004(200201)12:1;1-E