# Mel-Frequency Cepstral Coefficient-Based Bandwidth Extension of Narrowband Speech

*Amr H. Nour-Eldin, Peter Kabal*

Department of Electrical & Computer Engineering
McGill University, Montréal, Québec, Canada
`amr.nour-eldin@mail.mcgill.ca, peter.kabal@mcgill.ca`

## Abstract

We present a novel MFCC-based scheme for the Bandwidth Extension (BWE) of narrowband speech. BWE is based on the assumption that narrowband speech (0.3–3.4 kHz) correlates closely with the highband signal (3.4–7 kHz), enabling estimation of the highband frequency content given the narrow band. While BWE schemes have traditionally used LP-based parametrizations, our recent work has shown that MFCC parametrization results in higher correlation between both bands reaching twice that using LSFs. By employing high-resolution IDCT of highband MFCCs obtained from narrowband MFCCs by statistical estimation, we achieve high-quality highband power spectra from which the time-domain speech signal can be reconstructed. Implementing this scheme for BWE translates the higher correlation advantage of MFCCs into BWE performance superior to that obtained using LSFs, as shown by improvements in log-spectral distortion as well as Itakura-based measures (the latter improving by up to 13%).

**Index Terms**: Bandwidth extension, high-resolution IDCT, highband certainty, mutual information, source-filter model

## 1. Background and introduction

In traditional telephone networks, speech bandwidth is limited to the 0.3–3.4 kHz range. As a result, narrowband speech has sound quality inferior to its wideband counterpart and it shows reduced intelligibility especially for consonant sounds. Wideband speech reconstruction through Bandwidth Extension (BWE) attempts to regenerate the low (20–300 Hz) and highband (3.4–7 kHz) signals lost during the filtering processes employed in traditional networks, thereby providing backward compatibility with existing networks.

Traditionally, BWE research efforts have primarily used linear predictive (LP) techniques. By using such techniques, the reconstruction problem is divided into two separate tasks; forming a highband residual error (excitation) signal, and, recreating a set of higband linear predictive coefficients (LPCs). Once these two components have been generated, the highband residual excites the highband LP synthesis filter to regenerate the missing highband signal that can then be added to the available narrowband signal to generate wideband speech. The problem of reconstructing highband features from the corresponding narrowband ones has been addressed using two approaches; codebook mapping and statistical estimation. The underlying assumption is that narrowband speech correlates closely with the highband signal, and hence, the higher frequency speech content can be estimated from the narrowband signal.

In contrast to the ample research published on BWE techniques, the correlation assumption between the narrow- and high-band spectral envelopes has received little attention. Certainty about the high band given the narrow band was quantified in [1] as the ratio of Mutual Information (MI) between the two bands to the discrete entropy of the high band. The authors show that this ratio (representing correlation between the two bands) is quite low. Accordingly, it was concluded that existing BWE schemes perform reasonably, not because they accurately predict the true high band, but rather by extending the narrow band such that the overall wideband signal sounds pleasant.

More recently, we investigated in [2] the effect of the type of speech parametrization on the resulting correlation between narrow and high frequency bands (quantified by highband certainty). In particular, we considered Mel-Frequency Cepstral Coefficients (MFCCs) as well as Line Spectral Frequencies (LSFs). We showed that, for similar dimensionalities, MFCCs result in highband certainties that can reach almost twice as those resulting from LSFs. We argued that this higher correlation can be attributed to the Discrete Cosine Transform (DCT) employed in MFCC generation. DCT results in a decorrelation of cepstral coefficients leading to higher separability between different speech classes, with the advantage that MFCCs result in feature space modelling more discriminative of these classes. These results are confirmed by the findings of [3] which show MFCCs to have the highest speech class separability and second highest MI content among several speech parametrizations. LSFs, on the other hand, are widely used in speech coding, and are particularly attractive for BWE for their quantization error resilience and perceptual significance properties (where properties of formants and valleys can be related to LSF pairs). More importantly, LSFs have the important advantage of being easily convertible into LP coefficients, and hence, coupled with an excitation estimate, can be directly used for BWE. In contrast, reconstruction of the time-domain speech signal from MFCCs is more difficult at best. Based on these results for the correlation between speech frequency bands, we concluded in [2] that—notwithstanding the LSF advantage of straightforward speech reconstruction—MFCC-based BWE is inherently better.

Despite MFCCs' advantages over LSFs in terms of speech class separability, the difficulty of synthesizing speech from MFCCs has restricted their use to fields that do not require inverting MFCC vectors back into the original time-domain speech signals, e.g., automatic speech recognition. This difficulty arises from the non-invertibility of several steps employed in MFCC generation; using the magnitude of the complex spectrum, mel-scale filterbank *binning* and higher-order cepstral coefficient truncation. Consequently, all BWE techniques encountered in the literature are based on LP representations of the wideband (or highband) signals from which the wideband (or highband) frequency content is reconstructed (and added to the

September 22–26, Brisbane Australia

narrowband signal). The availability of the narrowband signal, however, has allowed researchers to investigate the effect of several types of narrowband parametrizations on increasing the correlation between narrowband feature vectors and LP-based wideband (or highband) feature vectors. Examples include [4] whose narrowband feature vectors consist of a mixture of auto-correlation coefficients, zero-crossing rate, normalized frame-energy, gradient index, local kurtosis, and the spectral centroid. A rare use of MFCCs in BWE is that of [5] which employs a Vector Quantization (VQ) codebook to map MFCC-parametrized narrowband signals to LSF wideband signals. Informal listening tests in [5] show clear preference for wideband speech reconstructed using the narrowband MFCC representation compared to that of the conventional LP-based representation, despite the reported increase in Log-Spectral Distortion (LSD). Despite the BWE performance improvements resulting from such alternative narrowband parametrizations, these improvements are limited by the wideband (or highband) LP-based representation. This limitation arises from the lower correlation between the alternative narrowband features and the LP-based highband ones (e.g., narrowband MFCCs correlate less with highband LSFs than with highband MFCCs).

In the work presented herein, we exploit the superiority of MFCCs over LSFs in terms of frequency bands' correlation by using MFCCs to represent both narrow- and high-band spectral envelopes for BWE (rather than limiting their use to the narrow band only as in [5]). To reconstruct highband speech from MFCCs (obtained by Gaussian Mixture Model (GMM) statistical estimation from input narrowband MFCCs), we employ high-resolution inverse DCT (IDCT) similar to that of [6] resulting in fine mel-scale cepstra, from which the linear power spectra can be recreated. The high-resolution IDCT effectively uses cosine functions to interpolate between mel-filterbank log-energies to reconstruct the cepstrum with finer detail (otherwise lost due to mel-filterbank binning). As in [7], we use a source-filter model to reconstruct speech from the estimated power spectra through inverse Fourier transform to obtain auto-correlation coefficients, to which the Levinson-Durbin recursion is applied. From the LPCs thus obtained, speech is synthesized by exciting the corresponding LP synthesis filters by an enhanced excitation signal [8] obtained from the narrow band. This MFCC inversion scheme thus eliminates the requirements of pitch estimation and voicing decisions of the more complex sinusoidal model techniques (employed in the field of distributed speech recognition), such as that of [9].

In contrast to [5], our MFCC-based BWE technique shows an LSD objective quality improvement of about 0.14 dB (about 3%) compared to LSF-based BWE with the same GMM complexity. More importantly, by using two variants of the more subjectively correlated Itakura-Saito distortion, we find a 7.5% improvement in highband spectral shape reconstruction due to the use of MFCCs rather than LSFs, reaching 13.2% when normalization for the effect of the reconstructed highband gain is applied. These results demonstrate the superiority of MFCC-based BWE over conventional LP-based schemes.

## 2. Review of highband certainty results

As stated above, highband certainty (representing the correlation between narrow and high frequency bands) is defined as the ratio of mutual information to discrete highband entropy. As described in [2], we estimate the mutual information, $I(X, Y)$, between narrow- and high-band feature vectors ($X$ and $Y$, respectively) using GMMs to model the marginal and joint dis-

Table 1: *Information measures (in bits) and highband certainty.*

| | $\mathrm{Dim}(X,Y)$ | $I(X;Y)$ | $H(Y)$ | $\frac{I(X;Y)}{H(Y)}$ |
|---|---|---|---|---|
| MFCCs | (10,6) | 1.59 | 7.82 | 20.3% |
| | (5,3) | 1.48 | 7.87 | 18.8% |
| LSFs | (10,6) | 0.84 | 7.88 | 10.7% |
| | (5,3) | 1.18 | 7.90 | 14.9% |

tributions of both sets of vectors, while the discrete highband entropy, $H(Y)$, is estimated using VQ of the highband vectors[1].

Table 1 shows the information measure results obtained for both MFCCs and LSFs for two different narrow- and high-band dimensionalities. We observe that despite the differing highband dimensionality or type of parametrization used, our discrete highband entropy estimates, $H(Y)$, are almost equal[2]. This confirms the convergence of our VQ estimates to the true highband entropy. With highband entropy estimates being equal, both MI and highband certainty figures show that MFCCs outperform LSFs in terms of capturing information mutual to both bands. This observation is further confirmed by the fact that both MI and highband certainty increase with increasing MFCC dimensionality, in contrast to their decrease using LSFs.

## 3. MFCC parametrization

Our application of the well-known MFCC parametrization of speech for the narrowband (0–4 kHz) and highband (4–8 kHz) signals (obtained by filtering the wideband speech to be used in BWE GMM training), is summarized as follows:

1. *Pre-emphasis:* A single-pole (at $z = -0.97$) high-pass filter is used to emphasize the highband formants of amplitudes lower than those of narrowband formants.
2. *Windowing:* A Hamming window is used to mitigate the edge effect of discontinuities due to framing. We use 20 ms frames with 50% overlap.
3. *Magnitude spectrum:* FFT (Fast Fourier transform) is applied followed by a magnitude operation.
4. *Mel-scale filterbank binning:* Mel-scale triangular filters are applied to the magnitude spectrum with FFT coefficients within each filter squared and summed resulting in mel-scale filterbank energies. We use 15 filters for the 0–4 kHz narrow band and 6 for the 4–8 kHz high band.
5. *Log operation:* Filterbank log-energies are obtained.
6. *DCT:* Type III DCT of the log-energies is applied per

$$c_n = \sqrt{\frac{2}{N}} \sum_{k=0}^{N-1} (\log Y_k) \cos\left(\frac{(2k+1)n\pi}{2N}\right),$$

where $c_n$ is the $n$th MFCC ($0 \leq n \leq N-1$), $N$ is the number of mel-scale filters, and $Y_k$ (or $X_k$) is the $k$th highband (or narrowband) mel-scale filter energy. Table 1 indicates that correlation between the two MFCC-parametrized frequency bands is almost twice that when using LSFs at $\mathrm{Dim}(X,Y) = (10,6)$. Accordingly, we use these dimensionalities for the parameters of our BWE scheme to emphasize the performance improvement using MFCCs versus LSFs, with $c_0$ included since the ratio of band energies represents an important measure of dependence between both bands.

---

[1]Please refer to [2] for complete details on: (a) the estimation of $I$ and $H$, and (b), the training and testing data sets used for Table 1.

[2]The discrete highband entropy estimates of Table 1 in [2] were improved upon by better implementation of the LBG training algorithm, resulting in the estimates of Table 1 above.

## 4. Highband speech synthesis

Two of the six steps of MFCC generation involve non-invertible loss of information; discarding phase information in Step 3 and the many-to-one mapping of the mel-scale filterbank binning of Step 4. The DCT of Step 6 also involves potential loss of information depending on whether the MFCC vectors are truncated.

Starting with narrowband speech input (sampled at 8 kHz), we recover from the lost information (after upsampling to 16 kHz and lowpass filtering with $f_c = 4$ kHz) as follows:

### 4.1. High-resolution IDCT

Since the narrow band is available as BWE input, no inversion is needed for narrowband MFCCs. These are calculated only to be used as inputs to the maximum-likelihood estimation of highband parameters from the trained GMMs described in Section 5.1. Since no truncation was applied to highband MFCCs in Step 6 above, the highband log-energies can be perfectly reconstructed by simple IDCT from the highband MFCCs estimated from the GMMs. These 6 log-energies can be viewed as scaled samples of the cepstrum at the center frequencies of the mel-scale filters, insufficient to recreate a spectrum. Finer cepstral detail can, however, be obtained by interpolating from these log-energies by increasing the resolution of the IDCT per

$$\log \hat{Y}_{k'} = \sqrt{\frac{2}{N}} \sum_{n=0}^{N-1} c_n \cos \left( \frac{(2k'+1)n\pi}{2iN} \right),$$

where $0 \leq k' \leq iN-1$, $N = 6$, and $i$ is an interpolation factor. Thus, The total number of log-energies (or cepstral samples) to be estimated in the 4–8 kHz range is $iN$. The interpolation factor, $i$, is determined by the desired mel-scale resolution. Using the frequency linear to mel-scale conversion; $f_{mel} = 2595 \log_{10} \left( 1 + f_{Hz}/700 \right)$, we obtain for 1 mel resolution in the $f_{Hz_1} = 4$ to $f_{Hz_2} = 8$ kHz band

$$i = \left\lceil \frac{f_{mel_2} - f_{mel_1}}{N+1} \right\rceil = 100,$$

resulting in a fine 600-sample cepstral resolution in the highband range. Thus, we effectively interpolate between the mel-frequency band centers using the DCT basis functions themselves as the interpolating functions [6].

### 4.2. MFCC preservation

Given a fixed highband MFCC dimensionality ($n_{max}$=5 corresponding to 6 cepstral coefficients as noted in Step 6), the choice for the number, $N$ (where $6 \leq N \leq 600$), of highband mel-scale filters is influenced by two opposing distortions. As $N$ increases, there will be more cepstral samples (log-energies) to interpolate from, resulting in fewer intermediate samples to interpolate (and hence, lower distortion due to interpolation errors). However, an increasing $N$ also involves truncation of an increasing number of higher-order cepstral coefficients, which in turn translates into IDCT distortion since the truncated cepstral coefficients are assumed to equal zero. By measuring the Euclidean distances between original cepstral samples (obtained just prior to the DCT of Step 6) and those samples resulting from high-resolution IDCT involving truncation for various values of $N$, we were able to conclude empirically that IDCT distortions resulting from MFCC truncation exceed those due to errors of interpolation from fewer log-energies. In fact, the interpolation performed implicitly by the high-resolution IDCT is quite accurate, leading to our choice in Step 6 to preserve cepstral coefficients by setting $N$ equal to MFCC dimensionality.

### 4.3. Highband LP synthesis

By exponentiation of the interpolated cepstra followed by mel-to-linear conversion, we obtain highband power spectra. Computing the inverse Fourier transform of the two-sided power spectra results in the auto-correlation coefficients, which can then be used to solve the Yule-Walker equations by means of the Levinson-Durbin recursion. Thus, we obtain highband LPCs minimizing the forward predictor mean-square-error. These LPCs represent the coefficients of the all-pole vocal tract filter.

### 4.4. Highband excitation signal

Rather than using a voicing-based model (based on the pitch extracted from the narrowband signal for voiced segments) as in [7] and [9], we use the narrowband signal equalized in the 3.4–4 kHz band to provide the excitation signal. As shown in [8], Gaussian noise modulation by the 3–4 kHz signal envelope (containing pitch harmonics) results in a superior excitation signal that is robust to differences in phonemes and speaker gender, leading to excellent highband signal reconstruction. In contrast, [7] uses a simple series of pitch pulses or white noise as the excitation for voiced and unvoiced segments, respectively. The loss of phase information (in Step 3 of MFCC calculation) is thus partially mitigated by using the equalized narrowband signal for excitation generation (thereby using phase information in the 3–4 kHz band to reconstruct highband phase). Moreover, the unimportance of phase for speech intelligibility [10] makes the accurate estimation of phase unwarranted.

## 5. MFCC-based BWE

### 5.1. System description

We implement MFCC-based BWE by modifying the *dual-mode* system detailed in [2] to incorporate MFCC parametrization and inversion as described in Sections 3 and 4. Shown in Figure 1, our system is based on that of [8] which exploits equalization to extend the bandwidth of narrowband speech up to 4 kHz. Besides being more accurate than any estimation algorithm in this frequency range, equalization up to 4 kHz also allows extraction of the enhanced excitation signal described in Section 4.4. GMM statistical estimation is used to generate the complementary spectrum, represented by LSFs/MFCCs, in the 4–8 kHz band. The estimated highband parameters, converted to LPCs, are then used together with the estimated excitation signal to reconstruct highband speech through LP synthesis. In addition, an excitation gain, $g$, is used to scale the synthesized highband components such that their energy is equal to that of the corresponding frequency band in the original wideband speech used for GMM training. Being a perceptual property, this gain improves the subjective quality of the extended speech. As $g$ is assumed to be correlated with the narrowband spectrum, it can also be statistically estimated from narrowband parameters.

### 5.2. Results and analysis

We evaluate BWE performance in the missing 4–7 kHz band (thus, only considering the effect of GMM modelling without the equalization effects in the 3.4–4 kHz band) by LSD (dB);

$$d_{LSD}^2 = \frac{1}{\pi} \int_{\omega_l}^{\omega_h} \left( 20 \log_{10} \frac{g}{|Y(e^{j\omega})|} - 20 \log_{10} \frac{\hat{g}}{|\hat{Y}(e^{j\omega})|} \right)^2 d\omega,$$

where $\omega_l$ and $\omega_h$ are the cutoff frequencies of the missing high band, $g$ and $Y(e^{j\omega})$ are the highband gain and frequency spectrum of the original wideband signal, respectively, while $\hat{g}$ and
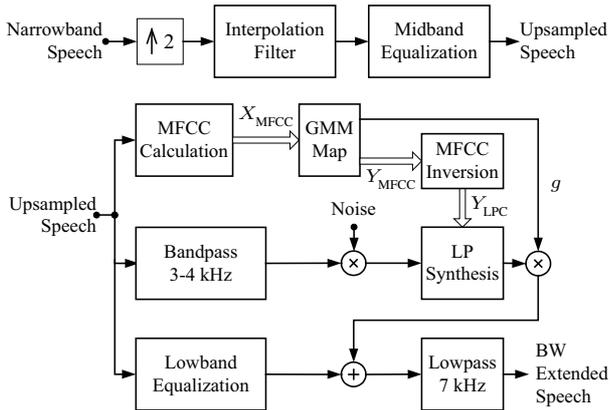
Figure 1: *Dual-mode MFCC-based BWE system.*

$\hat{Y}(e^{j\omega})$ are those of the GMM-estimated reconstructed signal.

While LSD is widely used for evaluating spectral envelope degradation due to its tractability and historic value, it does not take into account the perceptual importance of some aspects of the LP speech spectrum representation (e.g., LSD weights bandwidth differences for formants and valleys equally). In contrast, the Itakura-Saito distortion [11] has some perceptual relevance in that it weights differences in the LP spectra more heavily for peaks (which generally occur ar formant locations) than for valleys. However, due to its sensitivity to LP gain, a gain-optimized variant; the Itakura distortion [11], was derived by finding the LP model gains that minimize the Itakura-Saito distortion, thus rendering it gain-independent. This variant was shown in [12] to have a correlation of 0.73 with the subjective Diagnostic Acceptability Measure (versus 0.63 for LSD).

In our context, the Itakura-Saito distortion is given by

$$d_{\mathrm{IS}}\!\left(\tfrac{g^2}{|Y|^2}, \tfrac{\hat{g}^2}{|\hat{Y}|^2}\right) = \frac{1}{2\pi}\int_{-\pi}^{\pi}\left[\frac{g^2/|Y|^2}{\hat{g}^2/|\hat{Y}|^2} - \log\frac{g^2/|Y|^2}{\hat{g}^2/|\hat{Y}|^2} - 1\right]\mathrm{d}\omega.$$

The Itakura-Saito distortion does not fulfill the symmetry condition for distance metrics. However, a symmetrized version; the COSH measure, can be constructed by the arithmetic mean;

$$d_{\mathrm{COSH}} = \frac{1}{2}\left[d_{\mathrm{IS}}\!\left(\tfrac{g^2}{|Y|^2}, \tfrac{\hat{g}^2}{|\hat{Y}|^2}\right) + d_{\mathrm{IS}}\!\left(\tfrac{\hat{g}^2}{|\hat{Y}|^2}, \tfrac{g^2}{|Y|^2}\right)\right].$$

In a similar manner, we symmetrize the nonsymmetric gain-optimized Itakura distortion, given by

$$d_{\mathrm{It}}\!\left(\tfrac{g^2}{|Y|^2}, \tfrac{\hat{g}^2}{|\hat{Y}|^2}\right) \triangleq \min_{\hat{g}>0} d_{\mathrm{IS}}\!\left(\tfrac{g^2}{|Y|^2}, \tfrac{\hat{g}^2}{|\hat{Y}|^2}\right) = \log\left(\frac{\hat{\boldsymbol{y}}^T \boldsymbol{R}_Y \hat{\boldsymbol{y}}}{g^2}\right),$$

by the arithmetic mean;

$$d_{\mathrm{I}} = \frac{1}{2}\left[d_{\mathrm{It}}\!\left(\tfrac{g^2}{|Y|^2}, \tfrac{\hat{g}^2}{|\hat{Y}|^2}\right) + d_{\mathrm{It}}\!\left(\tfrac{\hat{g}^2}{|\hat{Y}|^2}, \tfrac{g^2}{|Y|^2}\right)\right],$$

where $\hat{\boldsymbol{y}}^T$ is the reconstructed LPC vector, and $\boldsymbol{R}_Y$ is the Toeplitz autocorrelation matrix of the original signal LP model.

Table 2 shows the distortion results obtained for MFCC- and LSF-based BWE. In contrast to [5], our MFCC-based scheme does improve $d_{\mathrm{LSD}}$. Although minor in comparison to the highband certainty gains of Table 1, it is important to note that the 0.14 dB $d_{\mathrm{LSD}}$ improvement was achieved with no increase in GMM complexity or data requirements of the LSF-based system. In comparison, the earlier version [13] of the dual-mode system in [8] achieves a highband $d_{\mathrm{LSD}}$ reduction of 0.96 dB by employing equalization and GMM statistical estimation compared to VQ codebook mapping. The significance of

Table 2: *LSD and Itakura-based distortion results.*

|  | $d_{\mathrm{LSD}}$ (dB) | $d_{\mathrm{COSH}}$ | $d_{\mathrm{I}}$ |
|---|---|---|---|
| LSFs | 4.74 | 7.95 | 0.79 |
| MFCCs | 4.60 | 7.35 | 0.69 |
| Improvement | 0.14 (2.9%) | 0.60 (7.5%) | 0.10 (13.2%) |

our $d_{\mathrm{LSD}}$ result further becomes evident by a comparison to that of [4], whose considerably more complex speaker-independent HMM-based system requires increasing the number of HMM states from 16 to 64 to achieve similar $d_{\mathrm{LSD}}$ reduction.

Finally, by employing $d_{\mathrm{COSH}}$ and $d_{\mathrm{I}}$ to compare MFCC-based BWE performance to that of LSFs, not only do we gain a more subjectively correlated measure of performance improvement, but we also gain a separation of the highband spectral shape-related improvement from errors in highband gain reconstruction (by exploiting the gain-insensitivity of $d_{\mathrm{I}}$). Accordingly, Table 2 shows a higher perceptually-relevant improvement of 7.5% due to improved highband spectral shape reconstruction, further reaching 13.2% when the effect of inaccuracies in highband gain reconstruction is eliminated.

# 6. References

[1] M. Nilsson, H. Gustafsson, S. V. Andersen and W. B. Kleijn, "Gaussian mixture model based mutual information estimation between frequency bands in speech", in *Proc. ICASSP*, vol. 1, Orlando, FL, USA, pp. 525–528, 2002.

[2] A. H. Nour-Eldin and P. Kabal, "Objective analysis of the effect of memory inclusion on bandwidth extension of narrowband speech", in *Proc. InterSpeech*, Antwerp, Belgium, pp. 2489–2492, 2007.

[3] P. Jax and P. Vary, "Feature selection for improved bandwidth extension of speech signals", in *Proc. ICASSP*, vol. 1, Montreal, QC, Canada, pp. 697–700, 2004.

[4] P. Jax and P. Vary, "On artificial bandwidth extension of telephone speech", *Signal Processing*, vol. 83, no. 8, pp. 1707–1719, 2003.

[5] N. Enbom and W. B. Kleijn, "Bandwidth expansion of speech based on vector quantization of the mel-frequency cepstral coefficients", in *Proc. IEEE Workshop on Speech Coding*, Porvoo, Finland, pp. 171–173, 1999.

[6] T. Ramabadran, J. Meunier, M. Jasiuk and B. Kushner, "Enhancing distributed speech recognition with back-end speech reconstruction", in *Proc. EuroSpeech*, Aalborg, Denmark, pp. 1859-1862, 2001.

[7] B. Milner and X. Shao, "Speech reconstruction from mel-frequency cepstral coefficients using a source-filter model", in *Proc. ICSLP*, Denver, CO, USA, pp. 2421–2424, 2002.

[8] Y. Qian and P. Kabal, "Combining equalization and estimation for bandwidth extension of narrowband speech", in *Proc. ICASSP*, vol. 1, Montreal, QC, Canada, pp. 713–716, 2004.

[9] D. Chazan, R. Hoory, G. Cohen and M. Zibulski, "Speech reconstruction from Mel frequency cepstral coefficients and pitch frequency", in *Proc. ICASSP*, vol. 3, Istanbul, Turkey, pp. 1299–1302, 2000.

[10] D. L. Wang and J. S. Lim, "The unimportance of phase in speech enhancement", *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 30, no. 4, pp. 679–681, 1984.

[11] R. M. Gray, A. Buzo, A. H. Gray, Jr. and Y. Matsuyama, "Distortion measures for speech processing", *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 28, no. 4, pp. 367–376, 1980.

[12] S. R. Quackenbush, T. P. Barnwell III and M. A. Clements, *Objective Measures of Speech Quality*, Englewood Cliffs, NJ: Prentice-Hall, 1988.

[13] Y. Qian and P. Kabal, "Dual-Mode wideband speech recovery from narrowband speech", in *Proc. EuroSpeech*, Geneva, Switzerland, pp. 14331437, 2003.