

- [3] C. S. Burrus and T. W. Parks, "Time domain design of recursive digital filters," *IEEE Trans. Audio Electroacoust.*, vol. AU-18, pp. 137-141, June 1970.
- [4] F. N. Cornett, "First and second order techniques for the design of recursive filters," Ph.D. dissertation, Colorado State Univ., Ft. Collins, CO, 1976.
- [5] D. C. Farden and L. L. Scharf, "Statistical design of nonrecursive digital filters," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-22, pp. 188-196, June 1974.
- [6] E. Parzen, "Multiple time series Modeling," in *Multivariate Analysis II*, P. Krishnaiah, Ed. New York: Academic, pp. 389-409.
- [7] W. C. Kellogg, "Time domain design of nonrecursive least mean-square digital filters," *IEEE Trans. Audio Electroacoust.*, vol. AU-20, pp. 155-158, June 1972.
- [8] D. C. Farden and L. L. Scharf, "Authors reply to 'Comments on statistical design of nonrecursive digital filters,'" *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-23, pp. 495-496, Oct. 1975.
- [9] U. Grenander and G. Szegö, *Toeplitz Forms and Their Applications*. Berkeley, CA: Univ. California Press, 1958.
- [10] K. Steiglitz, "Computer-aided design of recursive digital filters," *IEEE Trans. Audio Electroacoust.*, vol. AU-18, pp. 123-129, June 1970.
- [11] P. Thajchayapong and P. J. W. Rayner, "Recursive digital filter design by linear programming," *IEEE Trans. Audio Electroacoust.*, vol. AU-21, pp. 107-112, Apr. 1973.
- [12] L. R. Rabiner, N. Y. Graham, and H. D. Helms, "Linear programming design of IIR digital filters with arbitrary magnitude functions," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-22, pp. 117-123, Apr. 1974.
- [13] D. E. Dudgeon, "Recursive filter design using differential correction," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-22, pp. 443-448, Dec. 1974.
- [14] H. Dubois and H. Leich, "On the approximation problem for recursive digital filters with arbitrary attenuation curve in the passband and the stopband," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-23, pp. 202-207, Apr. 1975.
- [15] C. Charalambous, "Minimax optimization of recursive digital filters using recent minimax results," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-23, pp. 333-346, Aug. 1975.
- [16] J. A. Cadzow, "Recursive digital filter synthesis via gradient based algorithms," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-24, pp. 349-356, Oct. 1976.
- [17] H. Clergeot and L. L. Scharf, "Connections between classical and statistical methods of FIR digital filter design," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-26, pp. 463-465, Oct. 1978.
- [18] K. Steiglitz and L. E. McBride, "A technique for identification of linear systems," *IEEE Trans. Automat. Contr.*, vol. AC-10, pp. 461-464, Oct. 1965.

Predictive Coding of Speech Signals and Subjective Error Criteria

BISHNU S. ATAL, SENIOR MEMBER, IEEE, AND MANFRED R. SCHROEDER, FELLOW, IEEE

Abstract—Predictive coding methods attempt to minimize the rms error in the coded signal. However, the human ear does not perceive signal distortion on the basis of rms error, regardless of its spectral shape relative to the signal spectrum. In designing a coder for speech signals, it is necessary to consider the spectrum of the quantization noise and its relation to the speech spectrum. The theory of auditory masking suggests that noise in the formant regions would be partially or totally masked by the speech signal. Thus, a large part of the perceived noise in a coder comes from frequency regions where the signal level is low. In this paper, methods for reducing the subjective distortion in predictive coders for speech signals are described and evaluated. Improved speech quality is obtained: 1) by efficient removal of formant and pitch-related redundant structure of speech before quantizing, and 2) by effective masking of the quantizer noise by the speech signal.

I. INTRODUCTION

FOR autocorrelated signals, such as speech, predictive coding [1]-[4] is an efficient method of encoding the signal into digital form. The coding efficiency is achieved by quan-

tizing and transmitting only the signal which cannot be predicted from the already coded signal. In predictive coders, the power of the quantizer noise is proportional to the power of the prediction error. Thus, efficient prediction is important for minimizing the quantizer error. Small quantization error, however, does not ensure that the distortion in the speech signal is *perceptually* small; it is necessary to consider the *spectrum* of the quantization noise and its relation to the speech spectrum. The theory of auditory masking suggests that noise in the formant regions would be partially or totally masked by the speech signal. Thus, a large part of the perceived noise in a coder comes from the frequency regions where the signal level is low. Moreover, we can tolerate more distortion in the transitional segments in speech (where rapidly changing formants produce wider formant regions) in comparison to the steady segments.

In this paper, we discuss methods for modifying the spectrum of the quantization noise in a predictive coding system for speech to reduce the perceptible distortion introduced by such coders. The proper spectral shaping is realized by controlling the frequency response of the feedback network in the predictive coder independently of the predictor. The methods permit adaptive adjustment of the noise spectrum dependent

Manuscript received July 18, 1978; revised November 28, 1978.

B. S. Atal is with Bell Laboratories, Murray Hill, NJ 07974.

M. R. Schroeder is with the Drittes Physikalisches Institut, University of Göttingen, Göttingen, Germany, and Bell Laboratories, Murray Hill, NJ 07974.

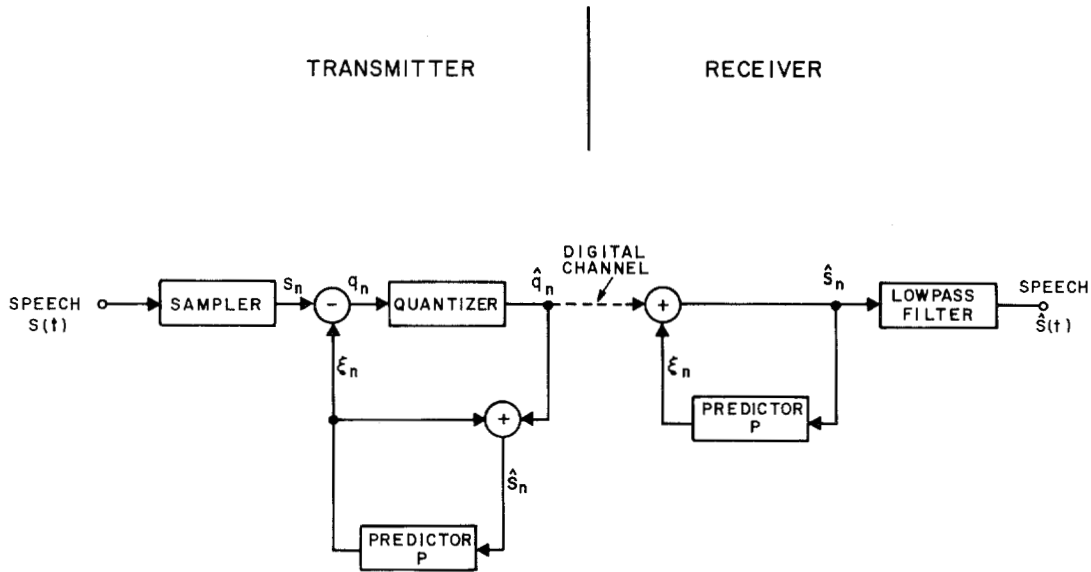


Fig. 1. Block diagram of a predictive coder.

on the time-varying speech spectrum. Improved speech quality is obtained both by exploiting auditory masking and by efficient prediction of formant and pitch-related redundancies in speech before quantizing.

II. SPECTRUM OF QUANTIZING NOISE IN PREDICTIVE CODERS

Fig. 1 shows a schematic diagram of a predictive coder. Its operation can be summarized as follows. The input speech signal is sampled to produce a sequence of sample values $s_0, s_1, \dots, s_n, \dots$. The predictor P forms a linear estimate of each sample value based on the previously decoded output sample values (reconstructed speech samples). This estimate ξ_n is subtracted from the actual sample value s_n and the resulting difference q_n is quantized and transmitted to the receiver. The quantized difference \hat{q}_n is added back to the predicted value ξ_n both at the transmitter and at the receiver to form the next output sample \hat{s}_n . It is readily seen in Fig. 1 that, in the absence of channel errors, the coder output \hat{s}_n is given by

$$\begin{aligned} \hat{s}_n &= \xi_n + \hat{q}_n \\ &= \xi_n + s_n - \xi_n + \delta_n \\ &= s_n + \delta_n \end{aligned} \quad (1)$$

where δ_n is the error introduced by the quantizer (difference between the output and the input of the quantizer) at the n th sample. Thus, the difference between the output and the input speech sample values is identical to the error introduced by the quantizer. Assuming that the spectrum of the quantizer error is white (a reasonable assumption, particularly if the prediction error is white), the noise at the output of the coder is also white.

III. GENERALIZED PREDICTIVE CODER

The quantizer input q_n in Fig. 1 can be written as

$$\begin{aligned} &= s_n - \sum_{k=1}^m (s_{n-k} + \delta_{n-k}) a_k \\ &= s_n - \sum_{k=1}^m s_{n-k} a_k - \sum_{k=1}^m \delta_{n-k} a_k \end{aligned} \quad (2)$$

where the predictor P is represented by a transversal filter with m delays and m gains a_1, a_2, \dots, a_m . The quantizer input thus consists of two parts: 1) the prediction error based on the prediction of the input s_n from its own past, and 2) the filtered error signal obtained by filtering of the quantizer error through the filter P . An easy way of modifying the spectrum of the quantizing noise is to use a filter F different from P for filtering the quantizer error [5], [6]. Fig. 2 shows such a generalized predictive coder. The quantizer input q_n is now given by

$$q_n = s_n - \sum_{k=1}^m s_{n-k} a_k - \sum_{k=1}^m \delta_{n-k} b_k \quad (3)$$

where b_1, b_2, \dots, b_m are m' gain coefficients of the transversal filter F . We will assume that both $1 - F$ and $1 - P$ have their roots inside the unit circle. The coder output is now given as

$$\hat{s}_n = s_n + \sum_{k=1}^m (\hat{s}_{n-k} - s_{n-k}) a_k + \delta_n - \sum_{k=1}^m \delta_{n-k} b_k. \quad (4)$$

Representing Fourier transforms by upper case letters, (4) can be written in frequency-domain notations as

$$\hat{S} - S = \Delta \frac{1 - F}{1 - P}. \quad (5)$$

For $F = P$, the output noise is the same as the quantizer noise, and the two coders shown in Figs. 1 and 2 are identical. However, with $F \neq P$, the coder of Fig. 2 allows greater flexibility

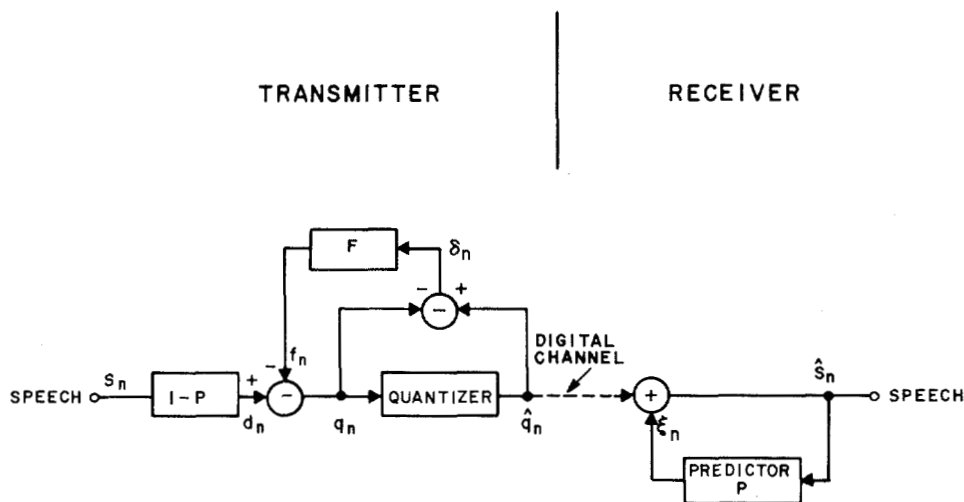


Fig. 2. Block diagram of a generalized predictive coder with adjustable noise spectrum.

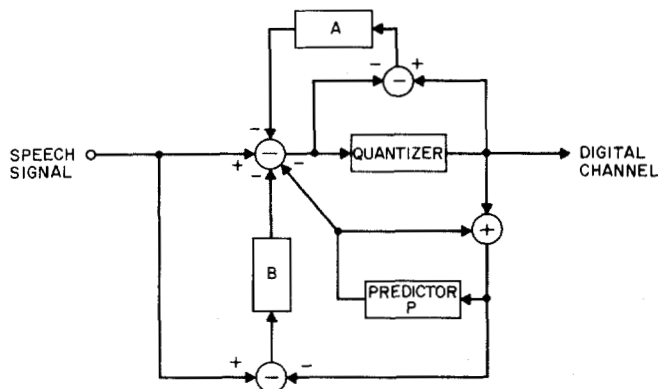


Fig. 3. Another configuration for the generalized predictive coder with adjustable noise spectrum.

choice of the feedback filter F .¹ Under the assumption that the quantizer noise is white, the spectrum of the coder output noise is determined only by the factor $(1 - F)/(1 - P)$ as shown in (5). Let the squared magnitude of this factor at a frequency f be $\Gamma(f)$. Then

$$\Gamma(f) = \left| \frac{1 - F(e^{2\pi jfT})}{1 - P(e^{2\pi jfT})} \right|^2 \quad (6)$$

where T is the sampling interval. Equation (6) implies an important constraint on the average value of $\log \Gamma(f)$, that is,

$$\frac{1}{f_s} \int_0^{f_s} \log \Gamma(f) df = 0 \quad (7)$$

where f_s is the sampling frequency. Expressed on a decibel scale, the average value of $\log \Gamma(f)$ is 0 dB. The proof of (7) is relatively straightforward [7] and is outlined below.

Consider the function $1 - F$ which is expressed in the z -

transform notation as

$$1 - F(z) = 1 - \sum_{k=1}^{m'} b_k z^{-k} = \prod_{k=1}^{m'} (1 - z_k z^{-1}) \quad (8)$$

where z_k is the k th root of $1 - F(z)$. The function $\log [1 - F(z)]$ is given by

$$\log [1 - F(z)] = \sum_{k=1}^{m'} \log (1 - z_k z^{-1}). \quad (9)$$

Since $|z_k| < 1$ (all the zeros of $1 - F(z)$ are inside the unit circle), the right side of (9) can be expressed as a polynomial function of z^{-1} . Therefore,

$$\log [1 - F(z)] = \sum_{n=1}^{\infty} c_n z^{-n} = \sum_{n=1}^{\infty} c_n e^{-2\pi jfTn} \quad (10)$$

where $c_n = \sum_{k=1}^{m'} z_k^n$. The integral of $\log [1 - F(z)]$ over the frequency range from 0 to f_s is then given by

$$\int_0^{f_s} \log [1 - F(e^{2\pi jfT})] df = \sum_{n=1}^{\infty} c_n \int_0^{f_s} e^{-2\pi jfTn} df = 0. \quad (11)$$

¹A somewhat different configuration of a predictive coder for controlling the spectrum of the output noise is shown in Fig. 3. It is easily verified that the Fourier transform of the output noise for this coder is given by $\Delta(1 - A)/(1 - B)$. In principle, the coder of Fig. 2 can be made equivalent to the coder of Fig. 3 by appropriate choice of the filter F . However, as a practical matter, one or the other coder may be simpler to implement depending on the choice for the spectrum of the

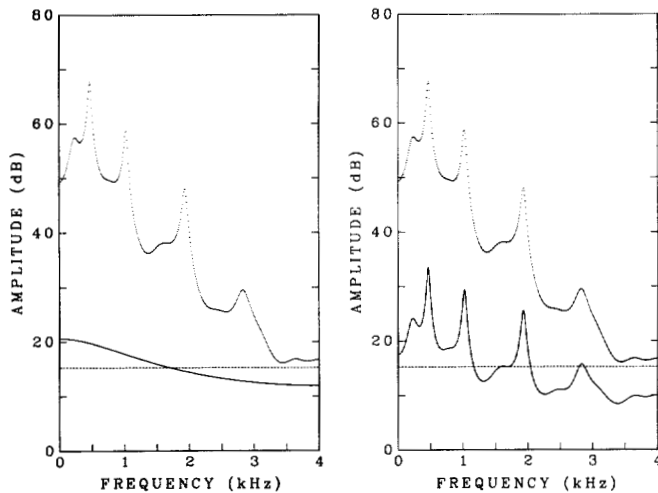


Fig. 4. Two possible shapes for the spectrum of output noise (solid curve) in the coder shown in Fig. 2. The average level of the logarithmic spectrum (shown as a dashed line) is the same in both cases. The speech spectrum is shown by the dotted curve.

$$\int_0^{f_s} \log [1 - P(e^{2\pi ifT})] df = 0. \quad (12)$$

Equation (7) follows directly from (11) and (12).

Assuming that the power of the quantizer noise δ_n is not changed significantly by the feedback loop—a desirable condition for satisfactory operation of the coder—the average value of \log power spectrum of output noise is then determined solely by the quantizer and is not altered by the choice of the filter F or the predictor P . The filter F , however, redistributes the noise power from one frequency to another. Thus, reduction in quantizer noise at one frequency can be obtained only at the expense of increasing the quantizer noise at another frequency. Since a large part of perceived noise in a coder comes from the frequency regions where the signal level is low, the filter F can be used to reduce the noise in such regions while increasing the noise in the formant regions where the noise could be effectively masked by the speech signal. Some examples of the possible shapes for the spectrum of the quantizing noise together with the speech spectrum are illustrated in Fig. 4. In each case, the logarithmic spectrum of the quantizing noise has equal area above and below the average level shown by the dashed line.

IV. APPLICATION TO SPEECH SIGNALS

A. Selection of Predictor

Linear prediction is a well-known method of removing the redundancy in a signal. For speech, the prediction is done most conveniently in two separate stages [4], [8]: a first prediction based on the short-time spectral envelope of speech, and a second prediction based on the periodic nature of the spectral fine structure. The short-time spectral envelope of speech is determined by the frequency response of the vocal tract and for voiced speech also by the spectrum of the glottal pulse. The spectral fine structure arising from the quasi-periodic nature of voiced speech is determined mainly by the pitch period. The fine structure for unvoiced speech is an

Prediction Based on Spectral Envelope

Prediction based on the spectral envelope involves relatively short delays. The predictor can be characterized in the z -transform notation as

$$P_s(z) = \sum_{k=1}^p a_k z^{-k} \quad (13)$$

where z^{-1} represents a delay of one sample interval and a_1, a_2, \dots, a_p are p predictor coefficients. The value of p typically is 10 for speech sampled at 8 kHz. A higher value may often be desirable.

The input to the quantizer in Fig. 2 consists of two parts: 1) the prediction error d_n based on the prediction of the input s_n from its own past, and 2) the filtered error signal f_n obtained by filtering of the quantizer noise δ_n through the filter F . Under the assumption that the quantizer noise is uncorrelated with the prediction error, the total power ϵ_q at the input to the quantizer is the sum of the powers in the prediction error d_n and the filtered noise f_n . That is,

$$\epsilon_q = \epsilon_p + \epsilon_f, \quad (14)$$

where ϵ_p is the power in the prediction error and ϵ_f is the power in the filtered noise. It is our experience that, for satisfactory operation of the coder, ϵ_f should be less than ϵ_p . The power in the filtered noise is determined both by the power in the quantizer error δ_n and the power gain G of the filter F . The power gain G equals the sum of the squares of the filter coefficients. For $F = P_s$, the power gain is usually large and can often exceed 200. Such a high power gain causes excessive feedback of the noise power to the quantizer input, particularly for coarse quantizers, resulting in poor performance of the coder. Excessive feedback can be prevented by requiring that the power gain of the filter $P_s(z)$ is not large. The reason for the high power gain is as follows.

The spectrum of the filter $1 - P_s(z)$ is approximately the reciprocal of the speech spectrum (within a scaling constant). The low-pass filter used in the analog-to-digital conversion of the speech signal forces the reciprocal spectrum (and thus $|1 - P_s(z)|$) to assume a high value in the vicinity of the cutoff frequency of the filter. The power gain, which is equal to the integral of the power spectrum $|1 - P_s(e^{2\pi ifT})|^2$ with respect to the frequency variable f , thus also becomes large.

These artificially high power gains will not arise if the low-pass filter used in the sampling process was an ideal low-pass filter with a cutoff frequency exactly equal to the half the sampling frequency. The amplitude-versus-frequency response of a practical low-pass filter falls off gradually. The computed covariance matrix used in LPC analysis therefore has missing components corresponding to the speech signal rejected by the low-pass filter. The missing high-frequency components produce artificially low eigenvalues of the covariance matrix corresponding to eigenvectors related to such components. The high power gain of P_s is precisely caused by the small eigenvalues. The covariance matrix of the low-pass filtered speech is nearly singular thereby resulting in a nonunique solution of the predictor coefficients. Thus a variety of dif

spectrum equally well in the passband of the low-pass filter. We wish to avoid solutions which lead to high power gains of the predictor P_s .

The ill-conditioning of the covariance matrix can be avoided by adding to the covariance matrix another matrix proportional to the covariance matrix of high-pass filtered white noise. We define a new covariance matrix $\hat{\Phi}$ (with its (ij) term represented by $\hat{\phi}_{ij}$) and a new correlation vector \hat{c} (with its i th term represented by \hat{c}_i) by the equations

$$\hat{\phi}_{ij} = \phi_{ij} + \lambda \epsilon_{\min} \mu_{i-j} \tag{15}$$

and

$$\hat{c}_i = c_i + \lambda \epsilon_{\min} \mu_i \tag{16}$$

where

$$\phi_{ij} = \langle s_{n-i} s_{n-j} \rangle,$$

$$c_i = \langle s_n s_{n-i} \rangle,$$

λ is a small constant (suitable values are in the range 0.01–0.10), ϵ_{\min} is the minimum value of the mean-squared prediction error, μ_i is the autocorrelation of the high-pass filtered white noise at a delay of i samples, and $\langle \rangle$ indicates averaging over the speech samples contained in the analysis segment. Ideally, the high-pass filter should be the filter complimentary to the low-pass filter used in the sampling process. We have obtained reasonably satisfactory results with the high-pass filter $[\frac{1}{2}(1 - z^{-1})]^2$. For this filter, the autocorrelations are $\mu_0 = \frac{3}{8}$, $\mu_1 = -\frac{1}{4}$, $\mu_2 = \frac{1}{16}$, and $\mu_k = 0$ for $k > 2$. By making the scale factor on the noise covariance matrix in (15) and (16) proportional to the mean-squared prediction error, we find that it is possible to use a fixed value of λ . The results are not very sensitive to small variations in the value of λ . The minimum value of the mean-squared prediction error is determined by the Cholesky decomposition [8] of the original covariance matrix $[\phi_{ij}]$. A modified form [9] of the covariance method is used to determine the predictor coefficients from the new covariance matrix $\hat{\Phi}$. The first two steps in this modified procedure are identical to the usual covariance method [8]. That is, the matrix $\hat{\Phi}$ is expressed as a product of a lower triangular matrix L and its transpose L^t by Cholesky decomposition and a set of linear equations $Lq = \hat{c}$ is solved. The partial correlation at a delay m is obtained from

$$r_m = \frac{q_m}{\left[\langle s_n^2 \rangle - \sum_{j=1}^{m-1} q_j^2 \right]^{1/2}} \tag{17}$$

where q_m is the m th component of q . The partial correlations are transformed to predictor coefficients using the well-known relation between the partial correlations and the predictor coefficients for all-pole filters [7, p. 110]. The modified procedure ensures that all of the zeros of the polynomial $1 - P_s(z)$ are inside the unit circle. Using the above procedure, reasonable power gains are realized without introducing significant bias in the spectrum at lower frequencies.² Examples of spectral envelopes of speech computed for $\lambda = 0$ (uncorrected) and for $\lambda = 0.05$ with the high-pass filter $[\frac{1}{2}(1 - z^{-1})]^2$ are

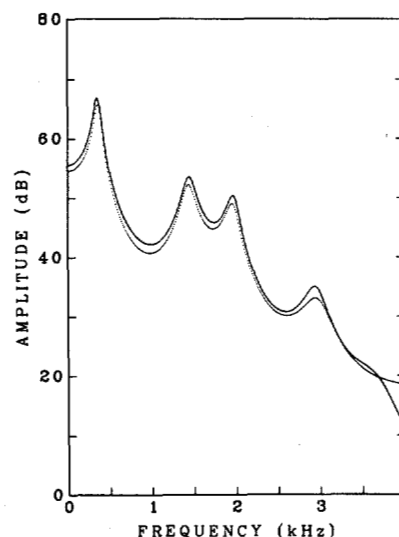


Fig. 5. Spectral envelopes of speech based on LPC analysis with high-frequency correction for $\lambda = 0$ (solid curve) and for $\lambda = 0.05$ (dotted curve). The power gains in the two cases are 204.6 and 12.6, respectively.

204.6 for $\lambda = 0$ to 12.6 for $\lambda = 0.05$. The speech signal was sampled at a rate of 10 kHz. The anti-aliasing filter had attenuation of 3 dB at 4.2 kHz and more than 40 dB at 5 kHz.

Prediction Based on Spectral Fine Structure

Adjacent pitch periods in voiced speech show considerable similarity. The quasi-periodic nature of the signal is present—although to a lesser extent—in the difference signal obtained after prediction on the basis of spectral envelope. The periodicity of the difference signal can be removed by further prediction. The predictor for the difference signal can be characterized in the z -transform notation by

$$P_d(z) = \beta_1 z^{-M+1} + \beta_2 z^{-M} + \beta_3 z^{-M-1} \tag{18}$$

where M represents a relatively long delay in the range 2 to 20 ms. In most cases, this delay would correspond to a pitch period (or possibly, an integral number of pitch periods). The degree of periodicity in the difference signal varies with frequency. The three amplitude coefficients β_1 , β_2 , and β_3 provide a frequency-dependent gain factor in the pitch-prediction process. We found it necessary to use at least a third-order predictor for pitch prediction. The difference signal after prediction based on spectral envelope has a nearly flat spectrum up to half the sampling frequency. Due to a fixed sampling frequency unrelated to pitch period, the individual samples of the difference signal do not show a high period-to-period correlation. The third-order pitch predictor provides an interpolated value with a much higher correlation than the individual samples. Higher order pitch predictors provide even better improvement in the prediction gain.

Let the n th sample of the difference signal after the first (“formant”) prediction be given by

²Another possible solution, namely, undersampling of the speech signal, for avoiding excessive power gains was suggested by one of the

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.