

INNOVATION CODING WITH A CROSS-CORRELATED QUANTIZATION NOISE MODEL

Søren Vang Andersen, Morten Olesen, Søren Holdt Jensen, and Egon Hansen
 CPK, Aalborg University
 Fredrik Bajers Vej 7
 DK-9220 Aalborg Øst, Denmark
 e-mail: sva@cpk.auc.dk

ABSTRACT

We present the use of a cross-correlated quantization noise model in the recently proposed Kalman innovation speech coding scheme. Computer simulations and informal listening tests indicate that the incorporation of a cross-correlated noise model yields an improvement in both SNR and perceptual quality when compared to a uncorrelated noise model.

1 INTRODUCTION

The problem of efficiently coding a speech signal at low bit rates has lead to a variety of coding schemes. Many of these schemes achieve high coding efficiency through the use of statistical modeling of the speech signal. At low bit rates quantization noise becomes a non-negligible part of the decoded signal. Therefore statistical modeling of not only the speech signal but also the quantization noise may lead to novel, improved schemes.

The innovation speech coding scheme, proposed by Ramabadran and Sinha [4], is an example of a coding scheme with statistical modeling of quantization noise. To be specific, the decoder uses a Kalman estimator to estimate the speech signal present at the encoder and the estimation is based on innovations which are corrupted by quantization noise. The speech signal is modeled as an outcome of an autoregressive (AR) process, and the quantization noise is modeled as an outcome of a white process that is uncorrelated with other involved processes. This model for the quantization noise is here termed the uncorrelated model (UCM).

In this paper, we propose an extension of the innovation speech coding scheme. Our scheme uses a more comprehensive model for the quantization noise. This model incorporates the cross-correlation, which in fact exists, between the quantization noise process and the process for the input to the quantizer. Such a model is here termed a cross-correlated model (CCM).

2 INNOVATION CODING

The innovation coding scheme will be described briefly in the following. To generalize the description, we here

consider vector quantization; this follows as a direct generalization of the scalar scheme described in [4].

Let the AR model for the speech signal be given by the equation¹

$$X_{t+1} = \boldsymbol{\alpha}_t [X_t \quad X_{t-1} \quad \dots \quad X_{t-p+1}]^T + Y_t,$$

where t is the time index, p is the order of the model, $\boldsymbol{\alpha}_t$ is a row vector containing model coefficients, X_t is the process modeling the speech signal, and Y_t is a white process that is uncorrelated with other involved processes. A state space formulation of this AR model can be written as

$$\mathbf{S}_{t+j} = \boldsymbol{\Theta}_{t+j,t} \mathbf{S}_t + \mathbf{R}_{t+j,t}, \quad (1)$$

where j is the number of samples we advance in the signal for each iteration of the state space model. The state transition matrix $\boldsymbol{\Theta}_{t+j,t}$ is a square matrix defined by²

$$\boldsymbol{\Theta}_{t+j,t} \equiv \begin{cases} \begin{bmatrix} \boldsymbol{\alpha}_t \\ \mathbf{I} \quad \mathbf{0} \end{bmatrix}, & j = 1, \\ \boldsymbol{\Theta}_{t+j,t+j-1} \boldsymbol{\Theta}_{t+j-1,t}, & j > 1. \end{cases}$$

The state \mathbf{S}_t and state excitation $\mathbf{R}_{t+j,t}$ are column vectors defined by

$$\begin{aligned} \mathbf{S}_t &\equiv [X_t \quad X_{t-1} \quad \dots \quad X_{t-p+1}]^T \\ \mathbf{R}_{t+j,t} &\equiv \begin{cases} [Y_t \quad \mathbf{0}]^T, & j = 1, \\ \boldsymbol{\Theta}_{t+j,t+j-1} \mathbf{R}_{t+j-1,t} + \mathbf{R}_{t+j,t+j-1}, & j > 1. \end{cases} \end{aligned}$$

In the following, we will suppress time indexes to avoid an overcrowded notation. Instead we introduce a delay operator $D[\cdot]$ such that e.g. (1) becomes $\mathbf{S} = D[\boldsymbol{\Theta} \mathbf{S} + \mathbf{R}]$.

¹In this paper uppercase Arabic letters are used to denote processes, while corresponding lowercase letters denote their outcomes. Moreover $\mathbf{E}[\cdot]$, $\mathbf{Var}[\cdot]$, and $\mathbf{Corr}[\cdot]$ will be used to denote expectation, variance, and correlation respectively.

²We use \mathbf{I} to denote an identity matrix and $\mathbf{0}$ to denote a zero vector or zero matrix. The dimensions of \mathbf{I} and $\mathbf{0}$ will follow from the context.

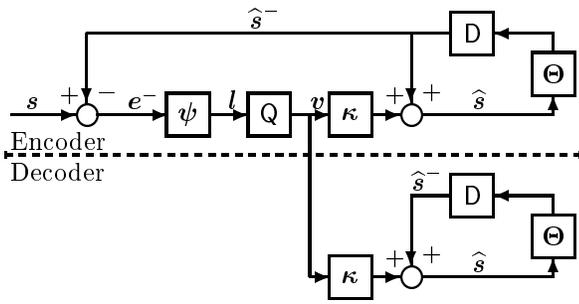


Figure 1: Block diagram of the innovation coding scheme.

With the state space formulation the innovation coding scheme can be illustrated as in Figure 1. The Kalman estimator provides an *a priori* estimate \hat{s}^- for s . From the *a priori* error e^- a vector of linear combinations $l = \psi e^-$ is obtained. This vector is quantized by the quantizer $Q[\cdot]$ to form the innovation v . The innovation is transmitted to the decoder and forms the basis for the estimation in both encoder and decoder. Pre-multiplying v with the Kalman gain matrix κ gives the term to update the *a priori* estimate \hat{s}^- into the *a posteriori* estimate \hat{s} . From this *a posteriori* estimate the next *a priori* estimate is obtained by pre-multiplication with Θ . The decoded speech signal is j consecutive elements of \hat{s} .

The influence of the quantizer is modeled as a zero-mean, additive process N , that is

$$V = Q[L] = L + N. \quad (2)$$

To determine the correlation matrix for N we restrict the treatment to schemes in which ψ is chosen such that $\text{Corr}[L]$ becomes positive definite and diagonal, and we assume the quantizer to be scaled using the diagonal elements of $\text{Corr}[L]$. It is then reasonable to assume the correlation matrix for N to be

$$\text{Corr}[N] = \Delta \text{Corr}[L]. \quad (3)$$

Here Δ is a diagonal matrix expressing the relative quantization noise variances. At this point we assume the UCM to hold. We can then use the standard Kalman estimator, see e.g. [3], and the encoding algorithm can be stated as a generalization of the algorithm given in [4].

UCM Encoding Algorithm:

$$\hat{s}^- = D[\Theta \hat{s}] \quad (4)$$

$$e^- = s - \hat{s}^- \quad (5)$$

$$\text{Corr}[E^-] = D[\Theta \text{Corr}[E] \Theta^T + \text{Corr}[R]] \quad (6)$$

$$l = \psi e^- \quad (7)$$

$$\text{Corr}[L] = \psi \text{Corr}[E^-] \psi^T \quad (8)$$

$$v = Q[l] \quad (9)$$

$$\text{Corr}[V] = (I + \Delta) \text{Corr}[L] \quad (10)$$

$$\kappa = \text{Corr}[E^-] \psi^T \text{Corr}[V]^{-1} \quad (11)$$

$$\hat{s} = \hat{s}^- + \kappa v \quad (12)$$

$$\text{Corr}[E] = (I - \kappa \psi) \text{Corr}[E^-]. \quad (13)$$

The decoding algorithm is identical except for the absence of equations (5), (7), and (9).

3 THE CROSS-CORRELATED QUANTIZATION NOISE MODEL

In this section we derive the modified encoding algorithm which incorporates the CCM. The basic assumption made is that $[\hat{S}^{-T} \ E^{-T}]^T$ is a zero-mean Gaussian process. Hence, for any two processes that are linear functions of $[\hat{S}^{-T} \ E^{-T}]^T$, uncorrelatedness will imply independence [5, p.51]. In various contexts we need this assumption in arguments of the following form:

Let A be a stochastic process that is uncorrelated with E^- , and for which $[A^T \ E^{-T}]^T$ is a zero-mean Gaussian process. Then A and E^- are independent. The quantization noise process N is a zero-mean process and a deterministic function of E^- . Thus independence of A and E^- implies independence of A and N . Because A and N are zero-mean and independent, they are uncorrelated. The argument above also holds for E^- replaced with L .

We refer to this argument as the *independence argument*.

The starting point of the derivation is the following important statement:

If the quantizer satisfies the centroid condition for a squared error distance measure, then the following relation holds:

$$E[LN^T] = -\text{Corr}[N]. \quad (14)$$

This follows using the same method of proof as in [2, p.357]. Because $\text{Corr}[L]$ is diagonal, the *independence argument* can be used to generalize the above statement to quantizers constructed as the product of lower dimensional quantizers that all satisfy the centroid condition. We are now ready to modify the Kalman estimator to incorporate relation (14) into the noise model.

The Kalman gain for the purpose of updating the *a priori* estimate to the *a posteriori* estimate is given in [3, p.311] as

$$\kappa = E[SV^T] \text{Corr}[V]^{-1}. \quad (15)$$

If we use the general property that \hat{S}^- and E^- are uncorrelated, and we use the *independence argument*, then (15) may be written as

$$\begin{aligned} \kappa &= \text{Corr}[E^-] \psi^T \text{Corr}[V]^{-1} \\ &\quad + E[E^- N^T] \text{Corr}[V]^{-1}. \end{aligned} \quad (16)$$

In the Kalman estimator for the UCM the second term in (16) is a zero matrix. For the CCM this, however, cannot be true. To see this, use (7) to obtain the following relation:

$$\mathbf{E}[\mathbf{L}\mathbf{N}^T] = \boldsymbol{\psi}\mathbf{E}[\mathbf{E}^-\mathbf{N}^T]. \quad (17)$$

From (14) we see that the left-hand side of (17) is not a zero matrix which means that $\mathbf{E}[\mathbf{E}^-\mathbf{N}^T]$ cannot be a zero matrix. Because $\boldsymbol{\psi}$ is chosen such that $\text{Corr}[\mathbf{L}]$ becomes positive definite, $\boldsymbol{\psi}$ cannot have more rows than columns and the rank of $\boldsymbol{\psi}$ will equate its number of rows. Hence, if $\boldsymbol{\psi}$ is square, then it has a unique inverse, and $\mathbf{E}[\mathbf{E}^-\mathbf{N}^T]$ can be determined from (17). On the other hand, if $\boldsymbol{\psi}$ is not square, some more care must be taken. One solution for this case is outlined in the following paragraph.

Consider the expansion of $\boldsymbol{\psi}$ into a square full rank matrix $\boldsymbol{\psi}_E \equiv [\boldsymbol{\psi}^T \quad \boldsymbol{\psi}_A^T]^T$ such that the expansion of \mathbf{L} ,

$$\mathbf{L}_E \equiv \boldsymbol{\psi}_E \mathbf{E}^- = [\mathbf{L}^T \quad \mathbf{L}_A^T]^T, \quad (18)$$

will have positive definite and diagonal correlation matrix $\text{Corr}[\mathbf{L}_E]$. Note that $\boldsymbol{\psi}_E$ will exist provided that $\text{Corr}[\mathbf{E}^-]$ is positive definite. Now, suppose no information about the extra linear combinations \mathbf{l}_A is passed through the quantizer. Then the extra linear combinations have not changed the available information and consequently they have not changed the estimate. Formally this can be expressed by the expanded innovation process $\mathbf{V}_E \equiv [\mathbf{V}^T \quad \mathbf{0}]^T$, where the zero information part is assigned the *a priori* expected value for \mathbf{l}_A , i.e. $\mathbf{0}$. The expanded quantization noise process thus becomes $\mathbf{N}_E \equiv \mathbf{V}_E - \mathbf{L}_E = [\mathbf{N}^T \quad -\mathbf{L}_A^T]^T$. We observe that \mathbf{L} and \mathbf{L}_A are uncorrelated, zero-mean, Gaussian processes and use the *independence argument* to obtain

$$\mathbf{E}[\mathbf{L}_E \mathbf{N}_E^T] = -\text{Corr}[\mathbf{L}_E] \Delta_E, \quad (19)$$

where Δ_E is defined as

$$\Delta_E \equiv \begin{bmatrix} \Delta & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

Now (18) and (19) leads to

$$\begin{aligned} \boldsymbol{\psi}_E \mathbf{E}[\mathbf{E}^-\mathbf{N}_E^T] &= \mathbf{E}[\mathbf{L}_E \mathbf{N}_E^T] \\ &= -\text{Corr}[\mathbf{L}_E] \Delta_E \\ &= -\boldsymbol{\psi}_E \text{Corr}[\mathbf{E}^-] \boldsymbol{\psi}_E^T \Delta_E. \end{aligned}$$

Since $\boldsymbol{\psi}_E$ has a unique inverse, we have

$$\mathbf{E}[\mathbf{E}^-\mathbf{N}_E^T] = -\text{Corr}[\mathbf{E}^-] \boldsymbol{\psi}_E^T \Delta_E.$$

Equating first columns in this matrix expression yields the unknown cross-correlation

$$\mathbf{E}[\mathbf{E}^-\mathbf{N}^T] = -\text{Corr}[\mathbf{E}^-] \boldsymbol{\psi}^T \Delta. \quad (20)$$

Note that this is the result which we can obtain directly from (17) in the case that $\boldsymbol{\psi}$ is square.

Finally, the Kalman gain $\boldsymbol{\kappa}$ can be obtained by inserting (20) into (16) which yields

$$\boldsymbol{\kappa} = \text{Corr}[\mathbf{E}^-] \boldsymbol{\psi}^T (\mathbf{I} - \Delta) \text{Corr}[\mathbf{V}]^{-1}.$$

The rest of the derivation follows steps parallel to those for a Kalman estimator for uncorrelated noise, see e.g. [3].

CCM Encoding Algorithm:

The CCM algorithm is obtained from the UCM algorithm by deleting (10), and replacing (11) and (13) with the following:

$$\boldsymbol{\kappa} = \text{Corr}[\mathbf{E}^-] \boldsymbol{\psi}^T \text{Corr}[\mathbf{L}]^{-1} \quad (21)$$

$$\text{Corr}[\mathbf{E}] = (\mathbf{I} - \boldsymbol{\kappa} (\mathbf{I} - \Delta) \boldsymbol{\psi}) \text{Corr}[\mathbf{E}^-]. \quad (22)$$

4 SIMULATION RESULTS

To investigate the impact of the CCM on the coder performance, simulations on speech signals were carried out. These simulations are described in the following.

The original speech signal was sampled at 8 kHz and linearly quantized using 16 bits per sample. The AR model for the speech signal was a two step filter model consisting of a long delay filter followed by a short delay filter. The short delay filter was a 10'th order all-pole filter. Coefficients for the short delay filter were determined for 10 ms frames using the autocorrelation method with a 30 ms Hamming window covering the last frame, the current frame, and the next frame. The long delay filter was of the form described in [1] using 3 nonzero coefficients. Coefficients for the long delay filter as well as $\text{Var}[Y]$ were estimated for each 10 ms frame. The coefficient vector $\boldsymbol{\alpha}$ was obtained through convolution of the long and short delay inverse filters as described in [4]. The order of this filter was fixed at $p = 120$. Innovations were quantized using the product of 1-bit scalar quantizers optimized for a Gaussian process. The simulations were focused on the issue of quantizing innovations and therefore no quantization of the AR model was included.

Four simulations were made. Results from these simulations are listed in Table 1. In simulations A and B, $\boldsymbol{\psi}$

Table 1: Simulation results

| ID. | Model | j | SEG-SNR [dB] | | | | Comb. |
|-----|-------|---|--------------|------|------|------|-------|
| | | | F1 | F2 | M1 | M2 | |
| A | UCM | 1 | 17.4 | 15.8 | 16.8 | 10.1 | 14.9 |
| B | CCM | 1 | 19.3 | 20.0 | 19.5 | 17.9 | 19.2 |
| C | UCM | 2 | 18.4 | 18.2 | 17.7 | 14.3 | 17.2 |
| D | CCM | 2 | 19.1 | 19.8 | 19.3 | 17.8 | 19.0 |

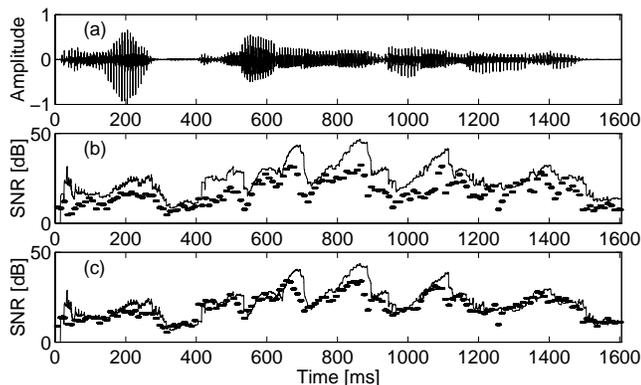


Figure 2: (a) The original speech signal (F2). (b) Time averaged SNR estimate (points) and SNR estimate obtained from the signal model and the *a posteriori* error correlation matrix (solid line), from simulation A with the UCM. (c) Like (b) but from simulation B with the CCM.

was a single row vector chosen to be an eigenvector corresponding to the largest eigenvalue of $\text{Corr}[\mathbf{E}^-]$. This choice was proposed in [4]. In simulations C and D, $j = 2$ and ψ was chosen with two rows equating eigenvectors corresponding to the two largest eigenvalues of $\text{Corr}[\mathbf{E}^-]$. Simulation A and C used the UCM, whereas simulations B and D used the CCM. Simulations were carried out on speech uttered by two female (F) and two male (M) speakers. The combined length of the speech from all four speakers was 10 seconds. As an objective performance indicator segmental SNR (SEG-SNR) was used. The SEG-SNR was calculated using 10 ms frames. Frames with signal power more than 40 dB below the average signal power were left out of the calculation. From Table 1, we see that the CCM for the combined data improved the SEG-SNR with 4.3 dB for the case $j = 1$ and 1.8 dB for the case $j = 2$. This result was supported by an improvement in perceptual quality when judged in informal listening test. In some instants this improvement was substantial.

To describe further the difference between the UCM and the CCM, we investigated the adjustment of the Kalman estimator and quantizer. The signal model provides an estimate of the signal variance and $\text{Corr}[\mathbf{E}^-]$ gives an estimate of the error variance. Using these two estimates, an estimate for the SNR of the decoded speech signal can be made. In Figure 2 this SNR estimate is plotted together with an SNR estimate obtained by time averaging over 10 ms frames. It is seen that these two estimates are more alike when the CCM is used. This observation was the same for all four speakers which we take as an indication that the CCM models the quantization noise more accurately than the UCM. The scaled input to the quantizer is expected to be a unit variance, zero-mean Gaussian process. Based on all 10 seconds of speech in simulations A and B normal probability plots for this process were made. The results are shown

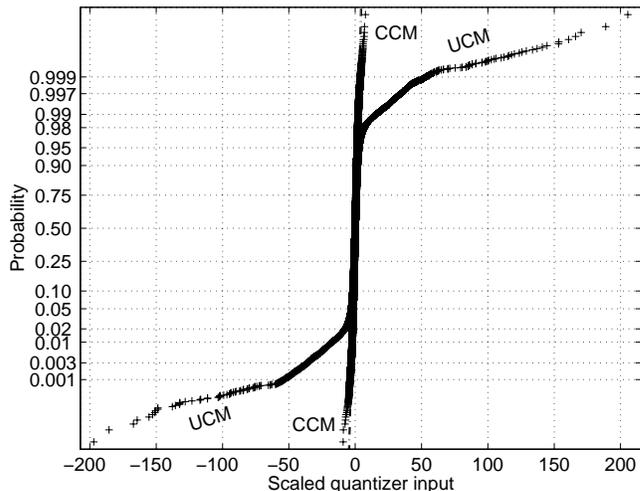


Figure 3: Normal probability plot for the scaled quantizer input. Outcomes of simulation A with the UCM and simulation B with the CCM are plotted. The dash dot line indicates an exact Gaussian process.

in Figure 3. It is clearly seen that when the UCM is used, the scaled quantizer input has a larger amount of outliers than when the CCM is used. Indeed, when the CCM is used the process is close to Gaussian. This result agrees with the assumption made in Section 3.

5 CONCLUSION

The quantization noise modeling in the innovation coding scheme was improved. This improvement was obtained by incorporating the cross-correlation between the quantization noise and the input to the quantizer. In simulations on speech signals the cross-correlated model resulted in a significant increase in the segmental SNR. This increase was supported by an improvement in perceptual quality when judged in informal listening test. In addition to this the adjustment of the estimation statistics and quantizer scaling improved in accuracy.

REFERENCES

- [1] B. S. Atal. Predictive Coding of Speech at Low Bit Rates. *IEEE Trans. Comm.*, 30(4):600–614, 1982.
- [2] A. Gersho and R. M. Gray. *Vector Quantization and Signal Compression*. Kluwer Academic Publishers, 1992.
- [3] S. Haykin. *Adaptive Filter Theory*. Printice Hall, 3rd edition, 1996.
- [4] T. V. Ramabadran and D. Sinha. Speech Data Compression Through Sparse Coding of Innovations. *IEEE Trans. Speech, Audio.*, 2(2):274–284, 1994.
- [5] K. S. Shanmugan and A. M. Breipohl. *Random Signals: Detection, Estimation and Data Analysis*. John Wiley and Sons, Inc., 1988.