

Article

Mutual Information, the Linear Prediction Model, and CELP Voice Codecs

Jerry Gibson

¹ Department of Electrical and Computer Engineering, University of California, Santa Barbara;
gibson@ece.ucsb.edu

Abstract: We write the mutual information between an input speech utterance and its reconstruction by a Code-Excited Linear Prediction (CELP) codec in terms of the mutual information between the input speech and the contributions due to the short term predictor, the adaptive codebook, and the fixed codebook. We then show that a recently introduced quantity, the log ratio of entropy powers, can be used to estimate these mutual informations in terms of bits/sample. A key result is that for many common distributions and for Gaussian autoregressive processes, the entropy powers in the ratio can be replaced by the corresponding minimum mean squared errors. We provide examples of estimating CELP codec performance using the new results and compare to the performance of the AMR codec and other CELP codecs. Similar to rate distortion theory, this method only needs the input source model and the appropriate distortion measure.

Keywords: Autoregressive models; entropy power; linear prediction model; CELP voice codecs; mutual information

1. Introduction

Voice coding is a critical technology for digital cellular communications, voice over Internet protocol (VoIP), voice response applications, and videoconferencing systems [1,2]. Code-Excited Linear Prediction (CELP) is the most widely deployed method for speech coding today, serving as the primary speech coding method in the Adaptive Multirate (AMR) codec [3] and in the more recent Enhanced Voice Services (EVS) codec [4], both of which are used in cell phones and VoIP. At a high level, a CELP Codec consists of a linear prediction model excited by an adaptive codebook and a fixed codebook. The linear prediction model captures the vocal tract dynamics (short term memory) and the adaptive codebook folds in the long term periodicity due to speaker pitch. The fixed codebook tries to represent the excitation for unvoiced speech and any remaining excitation energy not modelled by the adaptive codebook [2,5].

The encoder uses an analysis-by-synthesis paradigm to select the best fixed codebook excitation based on minimizing a frequency-weighted perceptual error signal. However, the linear prediction coefficients and the long term predictor parameters are calculated to minimize the unweighted mean squared error and then substituted into the codec structure prior to the analysis-by-synthesis optimization.

No characterization of the relative contributions of the several CELP codec components to the overall CELP codec perceptual performance is known. That is, given a segment of an input speech signal, what is the relative importance of the short term predictor, the long term predictor, and the fixed codebook to the perceptual quality achieved by the CELP codec? The mean squared prediction error can be calculated for the first two, and a translation of mean squared error to perceptual quality has been devised for calculating rate distortion bounds for speech coders in Gibson and Hu [6]. An

35 indication of the relative reduction in bits/sample provided by each component for a given quality
 36 would be very useful in optimizing codec performance and for investigating new adaptive codec
 37 structures.

38 In this paper, we develop and extend an approach suggested by Gibson [5] and decompose the
 39 mutual information between the input speech and the speech reconstructed by a CELP codec into
 40 the sum of unconditional and conditional mutual informations between the input speech and the
 41 linear prediction component, the adaptive codebook, and the fixed codebook, and show that this
 42 decomposition can be used to predict CELP codec performance based upon analyzing the input speech
 43 utterance, without actually implementing the CELP codec and processing the speech. We present
 44 examples comparing the estimated CELP codec performance and the actual performance achieved
 45 by CELP codecs. The approach is highly novel and the agreement between the actual and estimated
 46 performance is surprising.

47 The paper is outlined as follows. Section 2 provides an overview of the principles of Code-Excited
 48 Linear Prediction (CELP) needed for the remainder of the paper, while Sec. 3 develops the
 49 decomposition of the mutual information between the input speech and the speech reconstructed by
 50 the CELP codec. The concept of entropy power as defined by Shannon [7] is presented in Sec. 4, and
 51 the ordering of mutual information as a signal is passed through a cascaded signal processing system
 52 is stated in Sec. 5. The recently proposed quantity, the log ratio of entropy powers, is given in Sec. 6,
 53 where it is shown that the mean squared estimation errors can be substituted for the entropy powers
 54 in the ratio for an important set of probability densities [5,8,9]. The mutual information between the
 55 input speech and the short term prediction sequence is discussed in Sec. 7 and the mutual information
 56 provided by the adaptive and fixed codebooks about the input speech is developed in Sec. 8. The
 57 promised analysis of CELP codecs using these prior mutual information results based on only the input
 58 speech model and the distortion measure is presented in Sec. 9. The final section contains conclusions
 59 drawn from the results in the paper.

60 2. Code-Excited Linear Prediction (CELP)

61 Block diagrams of a Code-Excited Linear Prediction (CELP) encoder and decoder are shown in
 62 Figs. 1 and 2, respectively [1,2].

We provide a brief description of the various blocks in Figs. 1 and 2 to begin. The CELP encoder is
 an implementation of the Analysis-by-Synthesis (AbS) paradigm [1]. CELP, like most speech codecs in
 the last 45 years, is based on the linear prediction model for speech, wherein the speech is modeled as

$$s(k) = \sum_{i=1}^N a_i s(k-i) + Gw(k) \quad (1)$$

63 where we see that the current speech sample at time instant k is represented as a weighted linear
 64 combination of N prior speech samples plus an excitation term at the current time instant. The weights,
 65 $a_i, i = 1, 2, \dots, N$, are called the linear prediction coefficients. The Synthesis Filter in Fig. 1 has the form
 66 of this linear combination of past outputs and the Fixed and Adaptive Codebooks model the excitation,
 67 $w(k)$. The LP Analysis block calculates the linear prediction coefficients, and we see that the block also
 68 quantizes the coefficients so that encoder and decoder use exactly the same coefficients.

69 The adaptive codebook is used to capture the long term memory due to the speaker pitch and
 70 the fixed codebook is selected to be an algebraic codebook, which has mostly zero values and only a
 71 relatively few nonzero pulses. The Pitch Analysis block calculates the Adaptive Codebook long term
 72 memory. The process is AbS in that for a block of (say) M input speech samples, the linear prediction
 73 coefficients and long term memory are calculated and a perceptual weighting filter is constructed
 74 using the linear prediction coefficients. Then, for every length M sequence (codevector) in the Fixed
 75 Codebook, (say) there are L code vectors in the Fixed Codebook, a synthesized sequence of speech
 76 samples are produced. This is the Fixed Codebook Search block. The best codevector out of the L in

characteristic of CELP codecs that is well known is that those speech attributes not captured by the short term predictor must be accounted for, as best as possible, by the excitation codebooks, but an objectively meaningful measure of the individual component contributions is yet to be advanced.

In the next section, we propose a decomposition in terms of the mutual information and conditional mutual information with respect to the input speech provided by each component in the CELP structure that appears particularly useful and interesting for capturing the performance and the tradeoffs involved.

3. A Mutual Information Decomposition

Gibson [5] proposed the following decomposition of the several contributions to the synthesized speech by the CELP codec components. In particular, letting X represent a frame of input speech samples, we define X_R , X_N , and X_C as the reconstructed speech, the prediction component, and the combined fixed and adaptive codebook components, respectively. Then we can write the mutual information between the input speech and the reconstructed speech as

$$I(X; X_R) = I(X; X_N, X_C) = I(X; X_N) + I(X; X_C | X_N) \quad (2)$$

This expression states that the mutual information between the original speech X and the reconstructed speech X_R equals the mutual information between X and X_N , the N th order linear prediction of X , plus the mutual information between X and the combined codebook excitations X_C conditioned on X_N . Thus, to achieve or maintain a specified mutual information between the original speech and the reconstructed speech, any change in X_N must be offset by an adjustment of X_C . This fits what is known experimentally and that was alluded to earlier. If we define X_A to represent the Adaptive codebook contribution and X_F to represent the Fixed codebook contribution, we can further decompose $I(X; X_C | X_N)$ as

$$\begin{aligned} I(X; X_C | X_N) &= I(X; X_A, X_F | X_N) \\ &= I(X; X_A | X_N) + I(X; X_F | X_N, X_A) \\ &= I(X; X_F | X_N) + I(X; X_A | X_N, X_F) \end{aligned} \quad (3)$$

where we have used the chain rule for mutual information [16].

While these expressions are interesting, the challenge that remains is to characterize each of these mutual informations without actually having to calculate them directly from data, which is a difficult problem in and of itself [17].

An interesting quantity introduced and analyzed by Gibson in a series of papers is the log ratio of entropy powers [5,8,9]. Specifically, the log ratio of entropy powers is related to the difference in mutual information, and further, in many cases, the entropy powers can be replaced with the minimum mean squared prediction error (MMSPE) in the ratio. Using the MMSPE, the difference in mutual informations can be easily calculated. The following sections develop these concepts before we apply them to an analysis of the CELP structure.

4. Entropy Power/Entropy Rate Power

Given a random variable X with probability density function $p(x)$, we can write the differential entropy $h(X) = -\int_{-\infty}^{\infty} p(x) \log p(x) dx$ where the variance $\text{var}(X) = \sigma^2$. Since the Gaussian distribution has the maximum differential entropy of any distribution with mean zero and variance σ^2 [16],

$$h(X) \leq \frac{1}{2} \log 2\pi e \sigma^2 \quad (4)$$

from which we obtain

$$Q = \frac{1}{(2\pi e)} \exp 2h(X) \leq \sigma^2 \quad (5)$$

111 where Q was defined by Shannon to be the *entropy power* associated with the differential entropy of
 112 the original random variable [7]. In addition to defining entropy power, this equation shows that
 113 the entropy power is the *minimum variance* that can be associated with the not-necessarily-Gaussian
 114 differential entropy $h(X)$.

115 5. Cascaded Signal Processing

Figure 3 shows a cascade of N signal processing operations with the Estimator blocks at the output of each stage as studied by Messerschmitt [18]. He used the conditional mean at each stage and the corresponding conditional mean squared errors to obtain a representation of the distortion contributed by each stage. We analyze the cascade connection in terms of information theoretic quantities, such as mutual information, differential entropy, and entropy rate power. Similar to Messerschmitt, we consider systems that have no hidden connections between stages other than those explicitly shown. Therefore, we conclude directly from the Data Processing Inequality [16] that

$$I(X; Y_1) \geq \dots \geq I(X; Y_{N-1}) \geq I(X; Y_N) \geq I(X; \hat{X}) \quad (6)$$

Since $I(X; Y_n) = h(X) - h(X|Y_n)$, it follows from Eq. (6) that for non-negative $h(\cdot)$,

$$h(X|Y_1) \leq \dots \leq h(X|Y_{N-1}) \leq h(X|Y_N) \leq h(X) \quad (7)$$

116 For the optimal estimators at each stage, the basic Data Processing Inequality also yields $I(X; Y_n) \geq I(X; \hat{X}_n)$ and thus $h(X|Y_n) \leq h(X|\hat{X}_n)$. These are the fundamental results that additional processing

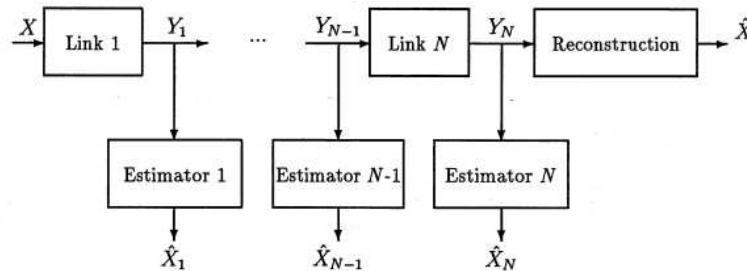


Figure 3. N-Link System Block Diagram (Adapted from [18])

117 cannot increase the mutual information.
 118

119 Now we notice that the series of inequalities in Eq. (7) along with the entropy power expression
 120 in Eq. (5) gives us the series of inequalities in terms of entropy power at each stage in the cascaded
 121 signal processing operations

$$Q_{X|Y_1} \leq \dots \leq Q_{X|Y_{N-1}} \leq Q_{X|Y_N} \leq Q_X \quad (8)$$

We can also write that

$$Q_{X|Y_n} \leq Q_{X|\hat{X}_n} \quad (9)$$

122 In the context of Eq. (8), the notation $Q_{X|Y_n}$ denotes the minimum variance when reconstructing
 123 an approximation to X given the sequence at the output of stage n in the chain [8].

124 6. Log Ratio of Entropy Powers

We can use the definition of the entropy power in Eq. (5) to express the logarithm of the ratio of two entropy powers in terms of their respective differential entropies as [8]

$$\log \frac{Q_X}{Q_Y} = 2[h(X) - h(Y)] \quad (10)$$

We can write a conditional version of Eq. (5) as

$$Q_{X|Y_n} = \frac{1}{(2\pi e)} \exp 2h(X|Y_n) \leq \text{Var}(X|Y_n) \quad (11)$$

and from which we can express Eq. (10) in terms of the entropy powers at successive stages in the signal processing chain, Fig. 3, as

$$\frac{1}{2} \log \frac{Q_{X|Y_n}}{Q_{X|Y_{n-1}}} = h(X|Y_n) - h(X|Y_{n-1}) \quad (12)$$

If we add and subtract $h(X)$ to the right hand side of Eq. (12), we then obtain an expression in terms of the difference in mutual information between the two stages as

$$\frac{1}{2} \log \frac{Q_{X|Y_n}}{Q_{X|Y_{n-1}}} = I(X; Y_{n-1}) - I(X; Y_n) \quad (13)$$

125 From the series of inequalities on the entropy power in Eq. (8), we know that both expressions in Eqs.
126 (12) and (16) are greater than or equal to zero.

127 These results are from [8] and extend the Data Processing Inequality by providing a new
128 characterization of the information loss between stages in terms of the entropy powers of the two
129 stages. Since differential entropies are difficult to calculate, it would be particularly useful if we could
130 obtain expressions for the entropy power at two stages and then use Eqs. (12) and (16) to find the
131 difference in differential entropy and mutual information between these stages.

132 We are interested in studying the change in the differential entropy and mutual information
133 brought on by different signal processing operations by investigating the log ratio of entropy powers.

In the following we highlight several cases where Eq. (10) holds with equality when the entropy powers are replaced by the corresponding variances. The Gaussian and Laplacian distributions often appear in studies of speech processing and other signal processing applications[10,15,19], so we show that substituting the variances for entropy powers in the log ratio of entropy powers for these distributions satisfies Eq. (10) exactly. For two i.i.d. Gaussian distributions with zero mean and variances σ_X^2 and σ_Y^2 , we have directly that $Q_X = \sigma_X^2$ and $Q_Y = \sigma_Y^2$, so

$$\frac{1}{2} \log \frac{Q_X}{Q_Y} = \frac{1}{2} \log \frac{\sigma_X^2}{\sigma_Y^2} = [h(X) - h(Y)] \quad (14)$$

134 which satisfies Eq. (10) exactly. Of course, since the Gaussian distribution is the basis for the definition
135 of entropy power, this result is not surprising.

For two i.i.d. Laplacian distributions with variances λ_X^2 and λ_Y^2 [20], their corresponding entropy powers $Q_X = 2e\lambda_X^2/\pi$ and $Q_Y = 2e\lambda_Y^2/\pi$, respectively, so we form

$$\begin{aligned} \frac{1}{2} \log \frac{Q_X}{Q_Y} &= \frac{1}{2} \log \frac{\lambda_X^2}{\lambda_Y^2} \\ &= [h(X) - h(Y)]. \end{aligned} \quad (15)$$

136 since $h(X) = \ln(2e\lambda_X)$ so the Laplacian distribution also satisfies Eq. (10) exactly [5]. We thus conclude
 137 that we can substitute the variance, or for zero mean Laplacian distributions, the mean squared value
 138 for the entropy power in Eq. (10) and the result is the difference in differential entropies.

139 Interestingly, using mean squared errors or variances in Eq. (10) is accurate for many other
 140 distributions as well. It is straightforward to show that Eq. (10) holds with equality when the entropy
 141 powers are replaced by mean squared error for the logistic, uniform, and triangular distributions
 142 as well. Further, the entropy powers can be replaced by the ratio of the squared parameters for the
 143 Cauchy distribution.

144 Therefore, the satisfaction of Eq. (10) with equality occurs in more than one or two special cases.
 145 The key points are first that the entropy power is the smallest variance that can be associated with
 146 a given differential entropy, so the entropy power is some fraction of the mean squared error for
 147 a given differential entropy. Second Eq. (10) utilizes the ratio of two entropy powers, and thus, if
 148 the distributions corresponding to the entropy powers in the ratio are the same, the scaling constant
 149 (fraction) multiplying the two variances will cancel out. Therefore, we are not saying that the mean
 150 squared errors equal the entropy powers in any case but for Gaussian distributions. It is the new
 151 quantity, the log ratio of entropy powers that enables the use of the mean squared error to calculate the
 152 loss in mutual information at each stage.

153 7. Mutual Information in the Short Term Prediction Component

154 Gibson [5] used Eq. (16) to investigate the change in mutual information as the predictor order,
 155 denoted in the following by N , is increased for different speech frames. Based on several analyses
 156 of the MMSPE and the fact that the log ratio of entropy powers can be replaced with the log ratio of
 157 MMSPE's for several different distributions, as outlined in Sec. 6 and in [9], we can use the expression

$$\frac{1}{2} \log \frac{MMSPE(X, X_{N-1})}{MMSPE(X, X_N)} = I(X; X_N) - I(X; X_{N-1}) \quad (16)$$

158 as in [5] to associate a change in mutual information with a change in the predictor order. Figure
 159 4 (bottom) shows 160 time domain samples from a narrowband (200 to 3400 Hz) speech sentence
 160 sampled at 8,000 samples/sec, and the top plot is the magnitude of the spectral envelope calculated
 161 from the linear prediction model using Eq. (1). We show the MMSPE and the corresponding change in
 162 mutual information for predictor orders $N = 1, 2, \dots, 10$, in Table 1. We see that the mutual information
 163 between the input speech frame and a 10th order predictor is 1.52 bits/sample. We can examine the
 164 mutual information between the input speech and a 10th order linear predictor for other frames in the
 165 same speech utterance.

To categorize for easy reference the differences in the speech frames, we use a common indicator
 of predictor performance, the Signal-to-Prediction Error (SPER) in dB [21], also called the Prediction
 Gain [15], defined as

$$SPER(dB) = 10 \log_{10} \frac{MSE(X)}{MMSPE(X, X_{10})} \quad (17)$$

166 where $MSE(X)$ is the average energy in the utterance and $MMSPE(X, X_{10})$ is the minimum mean
 167 squared prediction error achieved by a 10th order predictor. The SPER can be calculated for any
 168 predictor order but we choose $N = 10$, a common choice in narrowband speech codecs and the
 169 predictor order that, on average, captures most of the possible reduction in mean squared prediction
 170 error without including long term pitch prediction. For a normalized $MSE(X) = 1.0$, we see that for
 171 the speech frame in Fig. 4, the $SPER = 9.15$ dB.

172 Several other speech frames by the same speaker are analyzed in [5] and the results for these
 173 frames are tabulated in Table 5. From this table, it is evident that the mutual information between
 174 the input speech and a 10th order linear predictor can change quite dramatically across frames, even
 175 with the same speaker. We observe that the change in mutual information in some sense mirrors a
 176 change in SPER, with a larger SPER implying a larger mutual information. However, the explicit

Table 1. Change in Mutual Information from Eq. (16) as the Predictor Order is Increased: Frame 3237, $SPER = 9.15dB$

N	$MMSPE(X, X_N)$	$I(X; X_N) - I(X; X_{N-1})$
0	1.0	0 bits/letter
0-1	0.402	0.656 bits/letter
1-2	0.328	0.147 bits/letter
2-3	0.294	0.0795 bits/letter
3-4	0.2465	0.125 bits/letter
4-5	0.239	0.0234 bits/letter
5-6	0.2117	0.0869 bits/letter
6-7	0.212	0.0 bits/letter
7-8	0.125	0.381 bits/letter
8-9	0.1216	0.0206 bits/letter
9-10	0.1216	0.0 bits/letter
0-10 Total	0.1216	1.52 bits/letter

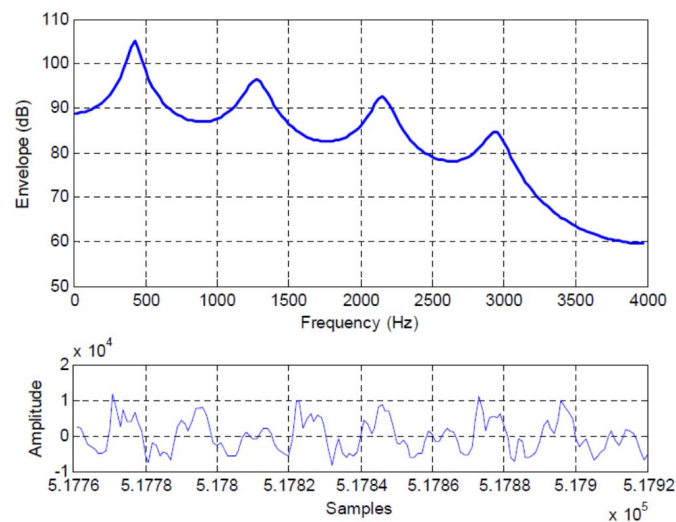


Figure 4. Frame 3237 Time Domain Waveform (Bottom) and Spectral Envelope, $SPER = 9.15$ dB

177 characterization in terms of mutual information is new and allows new analyses of the performance
 178 of CELP codecs in terms of bits/sample reduction in rate that can be associated with the short term
 179 predictor. We provide this analysis in Sec. 9.

180 8. Mutual Information in the Fixed and Adaptive Codebooks

181 How do we find the change in mutual information provided by the Adaptive Codebook and
 182 the Fixed Codebook? The Adaptive Codebook relies on long term prediction to model the vocal
 183 tract excitation due to the speaker pitch. As such, an approach similar to that used for the short
 184 term or linear predictor is possible. The Fixed Codebook contribution in terms of bits/sample is less
 185 direct. We could attempt to estimate the codebook complexity or we could simply use the number
 186 of bits/sample used to transmit the Fixed Codebook excitation. We elaborate on each of these in the
 187 following subsections.

Table 2. Change in Mutual Information from Eq. (16) for 10th Order Predictors and Corresponding SPERs for Several Speech Frames [5]

Speech Frame No.	SPER in dB	$I(X; X_{10}) - I(X; X_0)$
45	16	2.647 bits/letter
23	11.85	1.968 bits/letter
3314	7.74	1.29 bits/letter
87	5	0.83 bits/letter

188 8.1. Adaptive Codebook: Long Term Prediction

The long term predictor that makes up the adaptive codebook in CELP may have 1 or 3 taps and is of the form

$$P(z) = \beta_{-1}z^{-(P-1)} + \beta_0z^{-P} + \beta_1z^{-(P+1)} \quad (18)$$

where $\beta_i, i = -1, 0, 1$, are the predictor coefficients that are updated on a frame by frame basis and P is the lag of the periodic correlation that captures the pitch excitation of the vocal tract. A 3 tap predictor as shown in Eq. (18) requires a stability check to guarantee that this component does not cause divergence [22]. The single tap form given by

$$P(z) = \beta z^{-P} \quad (19)$$

189 only needs the stability check that $|\beta| < 1$, which should hold for any normalized autocorrelation. The
 190 3 tap form can often provide improved performance over the single tap predictor, but at increased
 191 complexity.

192 We denote the long term predicted value as X_p^P to distinguish it from the short term predictor of
 193 order N , X_N , which contains all terms up to and including N . Thus, we can write the MMSPE between
 194 X and X_p^P as $MMSPE(X, X_p^P)$.

195 The prediction gain in dB is often used as a performance indicator for a pitch predictor. The long
 196 term prediction gain has the form of Eq. (17) but with $MMSPE(X, X_{10})$ replaced by $MMSPE(X, X_p^P)$,
 197 where X_p^P is the long term predicted value for a pitch lag of P samples, where usually $P = 20$ up to
 198 $P = 140$ samples at 8000 samples/sec. Rather than calculate the prediction gain for selected speech
 199 frames, we consult the literature to get a range of values that might be expected.

200 Cuperman and Pettigrew [23] indicate a pitch prediction gain of 1 to 5 dB, but with around 2 dB
 201 on the average. Other sources indicate that the prediction gains can be 3-4 dB [24] or up to 5-6 dB [25].
 202 Table 3 shows the mutual information in bits/letter that can be associated with prediction gains of 1, 2,
 203 3, 4, and 6 dB.

204 It is interesting to consult the standardized speech codecs to see how many bits/sample are
 205 allocated to coding the pitch gain (coefficient or coefficients) and the pitch delay (lag). For the
 206 narrowband form of the AMR codec, we see that for 7.95 bits/sec rate, the pitch gain and pitch delay
 207 are allocated 16 and 28 bits, respectively, per 20 msec frame, or 44 bits/160 samples = 0.275 bits/sample
 208 [3]. For the highest rate of 12.2 kbits/sec, pitch gain and delay are allocated 46 bits/160 samples =
 209 0.2875 bits/sample. Consulting Table 3, we see that this corresponds to an SPER of between 1 and 2
 210 dB.

211 8.2. Fixed Codebook

212 The most often used fixed codebooks in popular speech coders, are sparse and have several
 213 shifted phase tracks, where each track consists of equally spaced pulses with 4 or 5 zero values in
 214 between. The best codeword or excitation from the fixed codebook is then selected by searching these

Table 3. Change in Mutual Information from Eq. (16) for Various Pitch Predictors and Corresponding SPER

SPER in dB	$I(X; X_p^p) - I(X; X_0)$
1	0.1667 bits/letter
2	0.33 bits/letter
3	0.5 bits/letter
4	0.65 bits/letter
6	1.0 bits/letter

215 tracks in some specific order. Getting an estimate of the number of bits/sample to be associated with a
 216 fixed codebook is therefore challenging. However, we know a few specific things.

217 The initial studies of analysis-by-synthesis codecs used Gaussian random sequences as the
 218 excitation. In particular, Atal and Schroeder used 1024 Gaussian random sequences that were 40
 219 samples long. Thus, this codebook used 10 bits/40 samples or 0.25 bits/sample [26]. However, this
 220 does not include the fixed codebook gain. The U. S. Federal Standard FS1016 CELP at 4.8 kbits/sec has
 221 allocations of 56 bits per frame of 240 samples or 0.233 bits/sample to the fixed codebook [10]. For a
 222 CELP codec that operates at 6 kbits/sec and above and a sparse codebook with five tracks and 8 pulses
 223 per track, an estimate of 3 bits plus a sign bit per track gives a fixed codebook rate of 0.5 bits/sample
 224 [24].

225 The narrowband AMR codec at 7.95 kbits/sec allocates 88 bits/20 msec frame or 0.55 bits/sample
 226 to the fixed codebook, and at 12.2 kbits/sec, 1 bit/sample is devoted to the fixed codebook gain and
 227 delay [3]. Thus, we see that at around 4 kbits/sec the allocation is about 0.25 bits/sample, at 8 kbits/sec
 228 the fixed codebook gets about 0.5 bits/sample, and at 12.2 kbits/sec, 1 bit/sample.

229 We now have estimates of the bits/sample devoted to the short term predictor, adaptive codebook,
 230 and fixed codebook for a CELP codec operating at different bit rates. We show how to exploit these
 231 estimates in the following to predict the rate of CELP codecs for different speech sources.

232 9. Estimated versus Actual CELP Codec Performance

233 The analyses determining the mutual information in bits/sample between the input speech and
 234 the short term linear prediction, the adaptive codebook, and the fixed codebook individually, are
 235 entirely new and provide new ways to analyze CELP codec performance by only analyzing the input
 236 source. In this section we estimate the performance of a CELP codec by analyzing the input speech
 237 source to find the mutual information provided by each CELP component about the input speech, and
 238 then subtracting the three mutual informations from a reference codec rate in bits/sample for a chosen
 239 MOS value to get the final estimate of the rate required in bits/sample to achieve the target MOS.

240 For waveform speech coding, such as differential pulse code modulation (DPCM), for a particular
 241 utterance, we can study the rate in bits/sample versus the mean squared reconstruction error or SNR
 242 to obtain some idea of the codec performance for this input speech segment [14,15]. However, while
 243 SNR may order DPCM subjective codec performance correctly, it does not provide an indicator of the
 244 actual difference in subjective performance. Subjective performance is most accurately available by
 245 conducting subjective listening tests to obtain Mean Opinion Scores (MOS); alternatively, reasonable
 246 views of subjective performance can be obtained from software such as PESQ/MOS [12]. We use the
 247 latter. In either case, however, MOS cannot be generated on a per frame basis as listening tests and
 248 PESQ values are generated from longer utterances.

249 Therefore, we cannot use the per frame estimates of mutual information from Sec. 7 and need to
 250 calculate estimates of mutual information over longer sentences. Toward this end, Table 4 contains
 251 autocorrelations for two narrowband (200 to 3400 Hz) utterances sampled at 8,000 samples/sec, "We
 252 were away a year ago," spoken by a male and "A lathe is a big tool. Grab every dish of sugar", spoken

Table 4. Composite Source Models for Narrowband Speech Sentences[6]

Sequence	Mode	Autocorrelation coefficients for V, ON, H	Mean Square	Probability
	Average frame energy for UV	Prediction Error		
"lathe" (Female) (active speech level: -18.1 dBov) (sampling rate: 8 kHz)	V	[1 0.8217 0.5592 0.3435 0.1498 0.0200 -0.0517 -0.0732 -0.0912 -0.1471 -0.2340]	0.0656	0.5265
	ON	[1 0.8495 0.5962 0.3979 0.2518]	0.0432	0.0093
	H	[1 0.2709 0.2808 0.1576 0.1182]	0.7714	0.0186
	UV	0.1439	0.1439	0.0771
	S			0.3685
"we were away" (Male) (active speech level: -16.5 dBov) (sampling rate: 8 kHz)	V	[1 0.8014 0.5176 0.2647 0.0432 -0.1313 -0.2203 -0.3193 -0.3934 -0.4026 -0.3628]	0.0780	0.9842
	ON	[1 0.8591 0.7215 0.6128 0.5183]	0.0680	0.0053
	H			0
	UV			0
	S			0.0105

Table 5. Rate in Bits[5]

Speech Frame No.	SPER in dB	$I(X; X_{10}) - I(X; X_0)$
45	16	2.647 bits/letter
23	11.85	1.968 bits/letter
3314	7.74	1.29 bits/letter
87	5	0.83 bits/letter

by a female, including the decomposition of the sentences into five modes, namely, Voiced, Unvoiced, Onset, Hangover, and Silence, and their corresponding relative frequencies. These data are excerpted from tables in Gibson and Hu [6], and are not calculated on a perframe basis but averaged over all samples of the utterance falling in the specified subsourse model.

From Table 4, the Voiced subsourse models are set as $N = 10th$ order, with the 1 in the column vector representing the $a_0th = 1$ term. The Onset and Hangover modes are modeled as $N = 4th$ order autoregressive (AR). We see from this table that the sentence "We were away . . . ," is Voiced with a relative frequency of 0.98, and that the sentence "A lathe is a big . . . ," has a breakdown of Voiced (0.5265), Silence (0.3685), Unvoiced (0.0771), Onset (0.0093), and Hangover (0.0186). The *MMSPEs* for each mode are also shown in the table.

We focus on G.726 as our reference codec. Generally, G.726 ADPCM at 32 kbits/sec or 4 bits/sample for narrowband speech is considered to be "toll quality." G.726 is selected as the reference codec because ADPCM is a waveform-following codec and is the best performing speech codec which does not rely on a more sophisticated speech model. In particular, G.726 utilizes only two poles and six zeros and the parameter adaptation algorithms rely only on polarities. G.726 will track pure tones as well as other temporal waveforms in addition to speech. In G.726, no parameters are coded and transmitted, only the quantized and coded prediction error signal. Finally, both mean squared reconstruction error or SNR and MOS have meaning for G.726, which is useful since SPER plays a role in estimating the change in mutual information from the log ratio of entropy powers.

From Tables 6 and 7, G.726 achieves a PESQ/MOS of about 4.0 for 4 bits/sample and for both sentences, "We were away a year ago," spoken by a Male Speaker and "A lathe is a big tool. Grab every dish of sugar," spoken by a Female [6]. Therefore, we use 4 bits/sample as our reference point for toll quality speech for these two sentences. We then subtract from the rate of 4 bits/sample, the rates in bits/sample we associate with each of the CELP codec components as estimated from Secs. 7 and 8.

For "We were away . . . ," we see from Table 4 that the $10th$ order model of the Voiced mode has a *MMSPE* = 0.078, which corresponds to a mutual information reduction of 1.84 bits/sample. For this highly voiced sentence, we estimate the mutual information corresponding to the adaptive codebook as 0.5 bits/sample, and at a codec rate of 8,000 bits/sec, the fixed codebook mutual information would correspond to 0.5 bits/sample. Silence corresponds to about only 1 percent of the total utterance. If we sum these contributions up and subtract them from 4 bits/sample, we obtain $4 - 2.84 = 1.16$ bits/sample.

Table 6. Codec Rates to Achieve PESQ-MOS of 4.0 for We were away[6]

Codec	<i>R</i> bits/sample
G.726	4.0 bits/sample
G.729	1.1 bits/sample
AMR	1.1 bit/sample

Table 7. Codec Rates to Achieve PESQ-MOS of 4.0 for Lathe[6]

Codec	<i>R</i> bits/sample
G.726	4.0 bits/sample
G.726 w/CNG	2.8 bits/sample
G.729	1.0 bits/sample
AMR	1.0 bit/sample

283 Inspecting Table 6, the rates for the AMR and G.729 codecs at this MOS are 1.1 bits/sample, so there is
 284 surprisingly good agreement.

285 From Table 4, we see that for the utterance, "A lathe is . . . ," there is broader mix of speech
 286 modes, including significant Silence and Unvoiced components. Neither of these modes had to be
 287 dealt with for the sentence "We were away" Since Silence is almost never perfect silence and
 288 usually consists of low level background noise, we associate the bits/sample for Silence with a silence
 289 detection/comfort noise generation (CNG) method [24]. From Table 7, we see that G.726 w/CNG is
 290 about 1.2 bits/sample lower than G.726, even though it occurs for only 0.3685 portion of the utterance.
 291 For the short term prediction, the $MMSPE = 0.0656$ which corresponds to a mutual information of
 292 1.96 bits/sample. The adaptive codebook contribution at 8 kbits/sec, can again be estimated as 0.5
 293 bits/sample, and the fixed codebook component estimated at 0.5 bits/sample.

294 If we now combine all of these estimates with their associated relative frequencies of occurrence as
 295 indicated in Table 4, we obtain a total mutual information of $0.5265(1.96 + 0.5) + 1.2 + 0.5(0.0771) = 2.5$
 296 bits/sample. Subtracting this from 4 bits/sample, we estimate that the CELP codec rate in bits/sample
 297 for an MOS = 4.0 would be 1.5 bits/sample. We see from Table 7 that for AMR and G.729 their rate is
 298 1.0 bits/sample. This gap can be narrowed if the adaptive codebook contribution is toward the upper
 299 end of the expected *SPER* of say 6 dB. In this case the Voiced component has the mutual information of
 300 $0.5265(1.96 + 1.0) = 1.56$ so the total mutual information is 2.8, and then subtracting from 4 bits/sample
 301 we obtain a CELP codec rate estimate of 1.2 bits/sample to achieve an MOS = 4.0. The actual CELP rate
 302 needed by G.729 and AMR for this MOS is about 1.0 bits/sample, which constitutes good agreement.

303 10. Conclusions

304 We have introduced a new approach to estimating CELP codec performance on a particular
 305 speech utterance by analysis of the input speech source only. That is, a particular CELP codec does not
 306 need to be implemented and used to process the speech. We have presented examples that illustrate the
 307 steps in the process and the accuracy of the estimated performance. While the power of the approach
 308 is evident, it is clear that many more sentences need to be processed to gain experience in estimating
 309 the various components. However, this approach offers the possibility of conducting performance
 310 analyses prior to the implementation of new CELP codec architectures and perhaps other new speech
 311 codec designs. It is striking that estimates of codec performance are possible while only knowing the
 312 source model and the distortion measure. Thus, in one sense, this new method parallels rate distortion
 313 theory.

314 **Funding:** This research received no external funding.

315 **Conflicts of Interest:** The authors declare no conflict of interest.

316 Abbreviations

317 The following abbreviations are used in this manuscript:

318	AR	Autoregressive
	ADPCM	Adaptive differential pulse code modulation
	AMR	Adaptive multirate
	AbS	Analysis-by-Synthesis
	CELP	Code-excited linear prediction
	CNG	Comfort noise generation
	DPCM	Differential pulse code modulation
319	EVS	Enhanced voice services
	MOS	Mean Opinion Score
	MSPE	Mean squared prediction error
	MMSPE	Minimum mean squared prediction error
	MMSE	Minimum mean squared error
	SNR	Signal to quantization noise ratio
	SPER	Signal to prediction error ratio
	VoIP	Voice over Internet Protocol

320

- 321 1. Chen, J.H.; Thyssen, J. Analysis-by-Synthesis Speech Coding. In *Springer Handbook of Speech Processing*;
322 Springer, 2008.
- 323 2. Gibson, J. Speech Compression. *Information* **2016**, *32*.
- 324 3. ETSI. 3GPP AMR Speech Codec; Transcoding Functions. Technical report, ETSI, 2002.
- 325 4. Dietz, M.; Multrus, M.; Eksler, V.; Malenovsky, V.; Norvell, E.; Pobloth, H.; Miao, L.; Wang, Z.; Laaksonen,
326 L.; Vasilache, A.; Kamamoto, Y.; Kikuri, K.; Ragot, S.; Faure, J.; Ehara, H.; Rajendran, V.; Atti, V.; Sung, H.;
327 Oh, E.; Yuan, H.; Zhu, C. Overview of the EVS codec architecture. Proc. Speech and Signal Processing
328 (ICASSP) 2015 IEEE Int. Conf. Acoustics, 2015, pp. 5698–5702. doi:10.1109/ICASSP.2015.7179063.
- 329 5. Gibson, J.D. Entropy Power, Autoregressive Models, and Mutual Information. *Entropy* **2018**.
- 330 6. Gibson, J.D.; Hu, J. Rate distortion bounds for voice and video. *Foundations and Trends® in Communications
331 and Information Theory* **2014**, *10*, 379–514.
- 332 7. Shannon, C.E. A mathematical theory of communication. *Bell Sys. Tech. Journal* **1948**, *27*, 379–423.
- 333 8. Gibson, J.D. Log Ratio of Entropy Powers. UCSD Information Theory and Application Workshop, 2018.
- 334 9. Gibson, J.; Oh, H. Analysis of Cascaded Signal Processing Operations Using Entropy Rate Power. Asilomar
335 Conference on Signals, Systems and Computers, 2018.
- 336 10. Chu, W.C. *Speech Coding Algorithms*; Wiley, 2003.
- 337 11. Grancharov, V.; Kleijn, W.B. Speech Quality Assessment. In *Springer Handbook of Speech Processing*; Springer,
338 2008.
- 339 12. ITU-T Recommendation P.862. Perceptual Evaluation of Speech Quality (PESQ), an objective method for
340 end-to-end Speech Quality Assessment of Narrow-band telephone networks and Speech Codecs. 2001.
- 341 13. Kleijn, W.B. Principles of Speech Coding. In *Springer Handbook of Speech Processing*; Springer, 2008.
- 342 14. Gibson, J.D.; Berger, T.; Lookabaugh, T.; Lindbergh, D.; Baker, R.L. *Digital Compression for Multimedia:
343 Principles and Standards*; Morgan Kaufmann Publishers Inc.: San Francisco, CA, 1998.
- 344 15. Rabiner, L.R.; Schafer, R.W. *Digital processing of speech signals*; Pearson, 2011.
- 345 16. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*; Wiley-Interscience, 2006.
- 346 17. Darbellay, G.A.; Vajda, I. Estimation of the information by an adaptive partitioning of the observation
347 space. *IEEE Transactions on Information Theory* **1999**, *45*, 1315–1321. doi:10.1109/18.761290.
- 348 18. Messerschmitt, D.G. Accumulation of Distortion in tandem Communication links. *IEEE Trans. on
349 Information Theory* **1979**, *IT-25*, 692–698.
- 350 19. Gray, R.M. *Linear Predictive Coding and the Internet Protocol*; NOW: Hanover, MA, USA, 2010.
- 351 20. Shynk, J.J. *Probability, random variables, and random processes: theory and signal processing applications*; John
352 Wiley & Sons: Hoboken, NJ, USA, 2012.
- 353 21. Sayood, K. *Introduction to data compression*; Morgan Kaufmann: Waltham, MA, USA, 2017.

- 354 22. Ramachandran, R.P.; Kabal, P. Pitch prediction filters in speech coding. *and Signal Processing IEEE*
355 *Transactions on Acoustics, Speech* **1989**, *37*, 467–478. doi:10.1109/29.17527.
- 356 23. Cuperman, V.; Pettigrew, R. Robust low-complexity backward adaptive pitch predictor for
357 low-delay speech coding. *Speech and Vision IEE Proceedings I-Communications* **1991**, *138*, 338–344.
358 doi:10.1049/ip-i-2.1991.0044.
- 359 24. Kondo, A.M. *Digital Speech: Coding for Low Bit Rate Communication Systems*; Wiley, 2004.
- 360 25. Kleijn, W.B.; Paliwal, K.K. *Speech Coding and Synthesis*; Elsevier, 1995.
- 361 26. Schroeder, M.; Atal, B. Code-excited linear prediction(CELP): High-quality speech at very low bit rates.
362 Proc. and Signal Processing ICASSP '85. IEEE Int. Conf. Acoustics, Speech, 1985, Vol. 10, pp. 937–940.
363 doi:10.1109/ICASSP.1985.1168147.