



INTERNATIONAL TELECOMMUNICATION UNION

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**G.723.1**

(03/96)

**GENERAL ASPECTS OF DIGITAL  
TRANSMISSION SYSTEMS**

---

**DUAL RATE SPEECH CODER  
FOR MULTIMEDIA COMMUNICATIONS  
TRANSMITTING AT 5.3 AND 6.3 kbit/s**

**ITU-T Recommendation G.723.1**

---

## FOREWORD

The ITU-T (Telecommunication Standardization Sector) is a permanent organ of the International Telecommunication Union (ITU). The ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Conference (WTSC), which meets every four years, establishes the topics for study by the ITU-T Study Groups which, in their turn, produce Recommendations on these topics.

The approval of Recommendations by the Members of the ITU-T is covered by the procedure laid down in WTSC Resolution No. 1 (Helsinki, March 1-12, 1993).

ITU-T Recommendation G.723.1 was prepared by ITU-T Study Group 15 (1993-1996) and was approved under the WTSC Resolution No. 1 procedure on the 19th of March 1996.

---

## NOTE

In this Recommendation, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

© ITU 1996

All rights reserved. No part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from the ITU.

## CONTENTS

		<i>Page</i>
1	Introduction.....	1
	1.1 Scope .....	1
	1.2 Bit rates.....	1
	1.3 Possible input signals.....	1
	1.4 Delay.....	1
	1.5 Speech coder description .....	1
2	Encoder principles.....	2
	2.1 General description .....	2
	2.2 Framer.....	2
	2.3 High pass filter.....	3
	2.4 LPC analysis .....	4
	2.5 LSP quantizer.....	4
	2.6 LSP decoder.....	5
	2.7 LSP interpolation .....	6
	2.8 Formant perceptual weighting filter.....	6
	2.9 Pitch estimation.....	7
	2.10 Subframe processing.....	7
	2.11 Harmonic noise shaping .....	7
	2.12 Impulse response calculator .....	8
	2.13 Zero input response and ringing subtraction.....	9
	2.14 Pitch predictor.....	9
	2.15 High rate excitation (MP-MLQ).....	9
	2.16 Low rate excitation (ACELP).....	11
	2.17 Excitation decoder .....	13
	2.18 Decoding of the pitch information.....	14
	2.19 Memory update.....	14
	2.20 Bit allocation.....	15
	2.21 Coder initialization .....	16
3	Decoder principles .....	17
	3.1 General description .....	17
	3.2 LSP decoder.....	17
	3.3 LSP interpolator.....	17
	3.4 Decoding of the pitch information.....	18
	3.5 Excitation decoder .....	18
	3.6 Pitch postfilter.....	18
	3.7 LPC synthesis filter.....	20
	3.8 Formant postfilter .....	20
	3.9 Gain scaling unit.....	21
	3.10 Frame interpolation handling.....	21
	3.11 Decoder initialization.....	22
4	Bitstream packing.....	22
5	ANSI C code .....	22
6	Glossary.....	22

## Summary

This Recommendation specifies a coded representation that can be used for compressing the speech or other audio signal component of multimedia services at a very low bit rate as part of the overall H.324 family of standards. This coder has two-bit rates associated with it, 5.3 and 6.3 kbit/s. The higher bit rate has greater quality. The lower bit rate gives good quality and provides system designers with additional flexibility. Both rates are a mandatory part of the encoder and decoder. It is possible to switch between the two rates at any frame boundary. An option for variable rate operation using discontinuous transmission and noise fill during non-speech intervals is also possible.

This coder was optimized to represent speech with a high quality at the above rates using a limited amount of complexity. It encodes speech or other audio signals in frames using linear predictive analysis-by-synthesis coding. The excitation signal for the high rate coder is Multipulse Maximum Likelihood Quantization (MP-MLQ) and for the low rate coder is Algebraic-Code-Excited Linear-Prediction (ACELP). The frame size is 30 ms and there is an additional look ahead of 7.5 msec, resulting in a total algorithmic delay of 37.5 msec. All additional delays in this coder are due to processing delays of the implementation, transmission delays in the communication link and buffering delays of the multiplexing protocol.

The description of this Recommendation is made in terms of bit-exact, fixed-point mathematical operations. The ANSI C code indicated in clause 5 constitutes an integral part of this Recommendation and shall take precedence over the mathematical descriptions in this text if discrepancies are found. A non-exhaustive set of test sequences which can be used in conjunction with the C code are available from the ITU.

## Recommendation G.723.1

**DUAL RATE SPEECH CODER FOR MULTIMEDIA COMMUNICATIONS  
TRANSMITTING AT 5.3 AND 6.3 kbit/s***(Geneva, 1996)***1 Introduction****1.1 Scope**

This Recommendation specifies a coded representation that can be used for compressing the speech or other audio signal component of multimedia services at a very low bit rate. In the design of this coder, the principal application considered was very low bit rate visual telephony as part of the overall H.324 family of standards.

**1.2 Bit rates**

This coder has two bit rates associated with it. These are 5.3 and 6.3 kbit/s. The higher bit rate has greater quality. The lower bit rate gives good quality and provides system designers with additional flexibility. Both rates are a mandatory part of the encoder and decoder. It is possible to switch between the two rates at any 30 ms frame boundary. An option for variable rate operation using discontinuous transmission and noise fill during non-speech intervals is also possible.

**1.3 Possible input signals**

This coder was optimized to represent speech with a high quality at the above rates using a limited amount of complexity. Music and other audio signals are not represented as faithfully as speech, but can be compressed and decompressed using this coder.

**1.4 Delay**

This coder encodes speech or other audio signals in 30 msec frames. In addition, there is a look ahead of 7.5 msec, resulting in a total algorithmic delay of 37.5 msec. All additional delays in the implementation and operation of this coder are due to:

- i) actual time spent processing the data in the encoder and decoder;
- ii) transmission time on the communication link;
- iii) additional buffering delay for the multiplexing protocol.

**1.5 Speech coder description**

The description of the speech coding algorithm of this Recommendation is made in terms of bit-exact, fixed-point mathematical operations. The ANSI C code indicated in clause 5, which constitutes an integral part of this Recommendation, reflects this bit-exact, fixed-point description approach. The mathematical descriptions of the encoder and decoder, given respectively in clauses 2 and 3, can be implemented in several other fashions, possibly leading to a codec implementation not complying with this Recommendation. Therefore, the algorithm description of the C code of clause 5 shall take precedence over the mathematical descriptions of clauses 2 and 3 whenever discrepancies are found. A non-exhaustive set of test sequences which can be used in conjunction with the C code are available from the ITU.

## 2 Encoder principles

### 2.1 General description

This coder is designed to operate with a digital signal obtained by first performing telephone bandwidth filtering (Recommendation G.712) of the analogue input, then sampling at 8000 Hz and then converting to 16-bit linear PCM for the input to the encoder. The output of the decoder should be converted back to analogue by similar means. Other input/output characteristics, such as those specified by Recommendation G.711 for 64 kbit/s PCM data, should be converted to 16-bit linear PCM before encoding or from 16-bit linear PCM to the appropriate format after decoding. The bitstream from the encoder to the decoder is defined within this Recommendation.

The coder is based on the principles of linear prediction analysis-by-synthesis coding and attempts to minimize a perceptually weighted error signal. The encoder operates on blocks (frames) of 240 samples each. That is equal to 30 msec at an 8 kHz sampling rate. Each block is first high pass filtered to remove the DC component and then divided into four subframes of 60 samples each. For every subframe, a 10th order Linear Prediction Coder (LPC) filter is computed using the unprocessed input signal. The LPC filter for the last subframe is quantized using a Predictive Split Vector Quantizer (PSVQ). The unquantized LPC coefficients are used to construct the short-term perceptual weighting filter, which is used to filter the entire frame and to obtain the perceptually weighted speech signal.

For every two subframes (120 samples), the open loop pitch period,  $L_{OL}$ , is computed using the weighted speech signal. This pitch estimation is performed on blocks of 120 samples. The pitch period is searched in the range from 18 to 142 samples.

From this point the speech is processed on a 60 samples per subframe basis.

Using the estimated pitch period computed previously, a harmonic noise shaping filter is constructed. The combination of the LPC synthesis filter, the formant perceptual weighting filter, and the harmonic noise shaping filter is used to create an impulse response. The impulse response is then used for further computations.

Using the pitch period estimation,  $L_{OL}$ , and the impulse response, a closed loop pitch predictor is computed. A fifth order pitch predictor is used. The pitch period is computed as a small differential value around the open loop pitch estimate. The contribution of the pitch predictor is then subtracted from the initial target vector. Both the pitch period and the differential value are transmitted to the decoder.

Finally the non-periodic component of the excitation is approximated. For the high bit rate, Multi-pulse Maximum Likelihood Quantization (MP-MLQ) excitation is used, and for the low bit rate, an algebraic-code-excitation (ACELP) is used.

The block diagram of the encoder is shown in Figure 1.

### 2.2 Framer

File : LBCCODEC.C	Procedure : main()	Reads 240 samples input frames
File : CODER.C	Procedure : Coder()	Performs subframe division

The coder processes the speech by buffering consecutive speech samples,  $y[n]$ , into frames of 240 samples,  $s[n]$ . Each frame is divided into two parts of 120 samples for pitch estimation computation. Each part is divided by two again, so that each frame is finally divided into four subframes of 60 samples each.

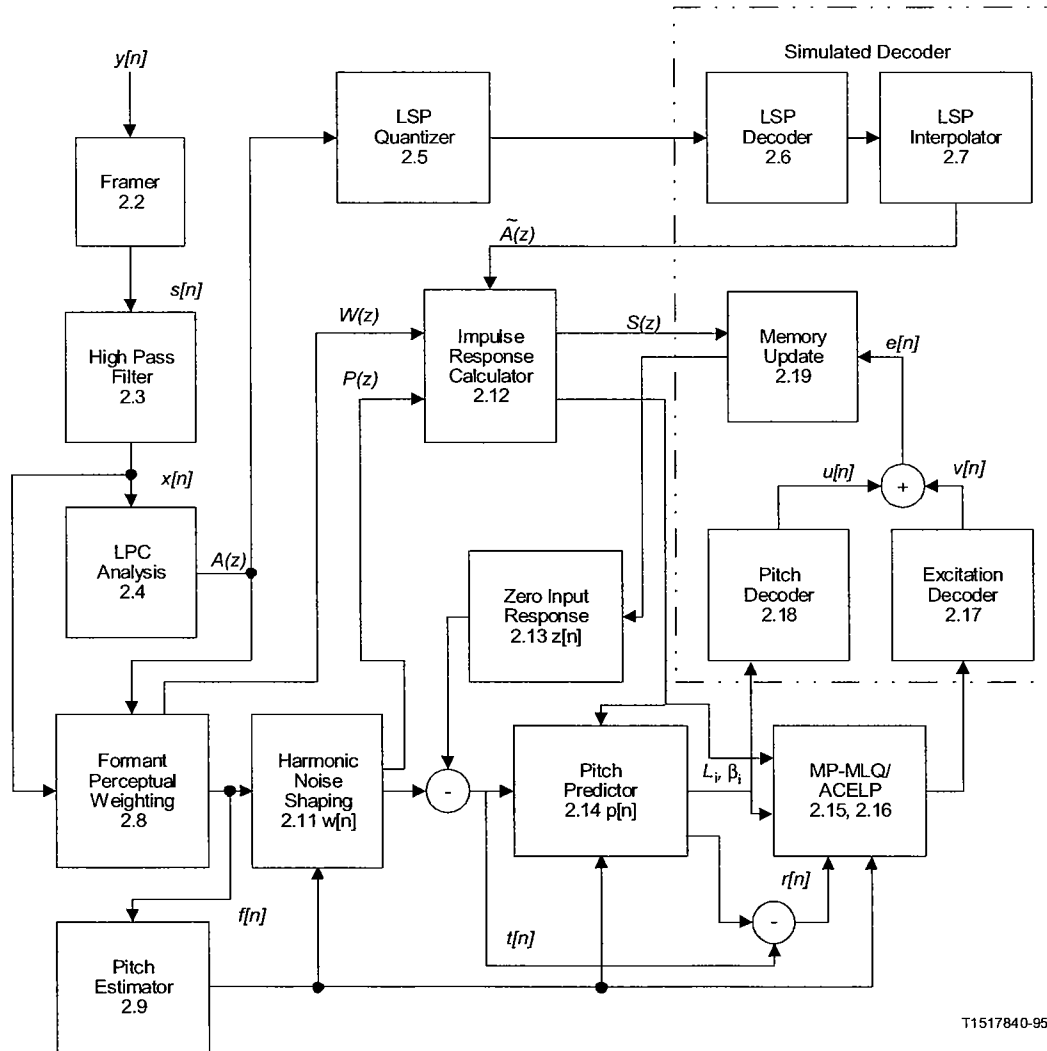
2.3 High pass filter

File : UTIL_LBC.C	Procedure : Rem_Dc()	Performs high pass filter
-------------------	----------------------	---------------------------

This block removes the DC element from the input speech,  $s[n]$ . The filter transfer function is:

$$H(z) = \frac{1 - z^{-1}}{1 - \frac{127}{128}z^{-1}} \tag{1}$$

The output of this filter is:  $x[n]_{n=0..239}$ .



T1517840-95/d01

FIGURE 1/G.723.1

Block diagram of the speech coder – For each block the corresponding reference number is indicated

## 2.4 LPC analysis

File : LPC.C	Procedure : Comp_Lpc()	Performs LPC coefficients calculation
File : LPC.C	Procedure : Durbin()	Levinson-Durbin recursion

The LPC analysis is performed on signal  $x[n]$  in the following way. Tenth order Linear Predictive (LP) analysis is used. For each subframe, a window of 180 samples is centered on the subframe. A Hamming window is applied to these samples. Eleven autocorrelation coefficients are computed from the windowed signal. A white noise correction factor of (1025/1024) is applied by using the formula  $R[0] = R[0](1 + 1/1024)$ . The other 10 autocorrelation coefficients are multiplied by the binomial window coefficients table. (The values for this table and all others are given in the C code.) The Linear Predictive Coefficients (LPC) are computed using the conventional Levinson-Durbin recursion. For every input frame, four LPC sets are computed, one for every subframe. These LPC sets are used to construct the short-term perceptual weighting filter. The LPC synthesis filter is defined as:

$$A_i(z) = \frac{1}{10 - \sum_{j=1}^{10} a_{ij}z^{-j}}, \quad 0 \leq i \leq 3 \quad (2)$$

where  $i$  is the subframe index and is defined to be between 0 and 3.

## 2.5 LSP quantizer

File : LSP.C	Procedure : AtoLsp()	Converts LPC to LSP coefficients
File : LSP.C	Procedure : LspQnt()	LSP vector quantization
File : LSP.C	Procedure : Lsp_Svq()	LSP sub-vectors quantization

First, a small additional bandwidth expansion (7.5 Hz) is performed. Then the resulting  $A_3(z)$  LP filter is quantized using a predictive split vector quantizer. The quantization is performed in the following way:

- 1) The LP coefficients,  $\{a_j\}_{j=1..10}$ , are converted to LSP coefficients,  $\{p'_j\}_{j=1..10}$ , by searching along the unit circle and interpolating for zero crossings.
- 2) The long term DC component,  $p_{DC}$ , is removed from the LSP coefficients,  $p'$ , and a new DC removed LSP vector,  $p$ , is obtained.
- 3) A first order fixed predictor,  $b = (12/32)$ , is applied to the previously decoded LSP vector  $\bar{p}_{n-1}$ , to obtain the DC removed predicted LSP vector,  $\bar{p}_n$ , and the residual LSP error vector,  $e_n$  at time (frame)  $n$ .

$$p_n^T = [p_{1,n} p_{2,n} \dots p_{10,n}] \quad (3.1)$$

$$\bar{p}_n^T = [\bar{p}_{1,n} \bar{p}_{2,n} \dots \bar{p}_{10,n}] \quad (3.2)$$



$$\bar{p}_n = b[\bar{p}_{n-1} - p_{DC}] \quad (3.3)$$

$$e_n = p_n - \bar{p}_n \quad (3.4)$$

- 4) The unquantized LSP vector,  $p'_n$ , the quantized LSP vector,  $\tilde{p}_n$ , the residual LSP error vector,  $e_n$ , are divided into 3 sub-vectors with dimension 3, 3 and 4 respectively. Each  $m$ th sub-vector is vector quantized using an 8 bit codebook. The index,  $l$ , of the appropriate sub-vector codebook entry that minimizes the error criterion  $E_{l,m}$  is selected.

$$p'_m{}^T = [p'_{1+3m} \ p'_{2+3m} \ \dots \ p'_{K_m+3m}], \quad K_m = \begin{cases} 3, & m=0 \\ 3, & m=1 \\ 4, & m=2 \end{cases} \quad (4.1)$$

$$\tilde{p}_{l,m}{}^T = [\tilde{p}_{1,l,m} \ \tilde{p}_{2,l,m} \ \dots \ \tilde{p}_{K_m,l,m}], \quad \begin{matrix} 0 \leq m \leq 2 \\ 1 \leq l \leq 256 \end{matrix} \quad (4.2)$$

$$p' = p + p_{DC} \quad (4.3)$$

$$\tilde{p}_{l,m} = \bar{p}_m + p_{DC_m} + e_{l,m}, \quad \begin{matrix} 0 \leq m \leq 2 \\ 1 \leq l \leq 256 \end{matrix} \quad (4.4)$$

$$E_{l,m} = (p'_m - \tilde{p}_{l,m})^T W_m (p'_m - \tilde{p}_{l,m}), \quad \begin{matrix} 0 \leq m \leq 2 \\ 1 \leq l \leq 256 \end{matrix} \quad (4.5)$$

where  $e_{l,m}$  is the  $l$ th entry of the  $m$ th split residual LSP codebook and  $W_m$  is a diagonal weighting matrix, determined from the unquantized LSP coefficients vector  $p'$ , with weights defined by:

$$\begin{aligned} w_{j,j} &= \frac{1}{\min \{p'_j - p'_{j-1}, p'_{j+1} - p'_j\}}, \quad 2 \leq j \leq 9 \\ w_{1,1} &= \frac{1}{p'_2 - p'_1} \\ w_{10,10} &= \frac{1}{p'_{10} - p'_9} \end{aligned} \quad (5)$$

- 5) The selected indices are transmitted to the channel.

## 2.6 LSP decoder

File : LSP.C	Procedure: Lsp_Inq()	Inverse quantization of LSP
--------------	----------------------	-----------------------------

The decoding of the LSP coefficients is performed in the following way:

- 1) First, the three sub-vectors,  $\{e_{m,n}\}_{m=0,2}$ , are decoded to form a tenth order vector,  $\tilde{e}_n$ .
- 2) The predicted vector,  $\bar{p}_n$ , is added to the decoded vector,  $\tilde{e}_n$ , and DC vector,  $p_{DC}$ , to form the decoded LSP vector,  $\tilde{p}_n$ .

- 3) A stability check is performed on the decoded LSP vector,  $\tilde{p}_n$ , to ensure that the decoded LSP vector is ordered according to the following condition:

$$\tilde{p}_{j+1,n} - \tilde{p}_{j,n} \geq \Delta_{\min}, 1 \leq j \leq 9 \quad (6)$$

$\Delta_{\min}$  is equal to 31.25 Hz. If this stability check (6) fails for  $\tilde{p}_i$  and  $\tilde{p}_{i+1}$ , then  $\tilde{p}_j$  and  $\tilde{p}_{j+1}$  are modified in the following way:

$$\tilde{p}_{avg} = (\tilde{p}_j + \tilde{p}_{j+1}) / 2 \quad (7.1)$$

$$\tilde{p}_j = \tilde{p}_{avg} - \Delta_{\min} / 2 \quad (7.2)$$

$$\tilde{p}_{j+1} = \tilde{p}_{avg} + \Delta_{\min} / 2 \quad (7.3)$$

The modification is performed until condition (6) is met. If after 10 iterations the condition of stability is not met, the previous LSP vector is used.

## 2.7 LSP interpolation

File : LSP.C	Procedure : Lsp_Int()	LSP interpolator
File : LSP.C	Procedure : LsptoA()	Converts LSP to LPC coefficients

Linear interpolation is performed between the decoded LSP vector,  $\tilde{p}_n$ , and the previous LSP vector,  $\tilde{p}_{n-1}$ , for each subframe. Four interpolated LSP vectors,  $\{\tilde{p}_i\}_{i=0..3}$ , are converted to LPC vectors,  $\{\tilde{a}_i\}_{i=0..3}$ .

$$\tilde{p}_{ni} = \begin{cases} 0.75\tilde{p}_{n-1} + 0.25\tilde{p}_n, & i = 0 \\ 0.5\tilde{p}_{n-1} + 0.5\tilde{p}_n, & i = 1 \\ 0.25\tilde{p}_{n-1} + 0.75\tilde{p}_n, & i = 2 \\ \tilde{p}_n, & i = 3 \end{cases} \quad (8)$$

$$\tilde{a}_i^T = [\tilde{a}_{i1}\tilde{a}_{i2} \dots \tilde{a}_{i10}]^T, 0 \leq i \leq 3 \quad (9)$$

The quantized LPC synthesis filter,  $\tilde{A}_i(z)$ , is used for generating the decoded speech signal and is defined as:

$$\tilde{A}_i(z) = \frac{1}{1 - \sum_{j=1}^{10} \tilde{a}_{ij} z^{-j}}, 0 \leq i \leq 3 \quad (10)$$

## 2.8 Formant perceptual weighting filter

File : LPC.C	Procedure : Wght_Lpc()	Computes perceptual filter coefficients
File : LPC.C	Procedure : Error_Wght()	Applies perceptual weighting filter

For each subframe a formant perceptual weighting filter is constructed, using the unquantized LPC coefficients  $\{a_{ij}\}_{j=1,\dots,10}$ . The filter has a transfer function:

$$W_i(z) = \frac{1 - \sum_{j=1}^{10} a_{ij} z^{-j} \gamma_1^j}{1 - \sum_{j=1}^{10} a_{ij} z^{-j} \gamma_2^j}, \quad 0 \leq i \leq 3 \quad (11)$$

where  $\gamma_1 = 0.9$  and  $\gamma_2 = 0.5$ . The input speech frame,  $\{x[n]\}_{n=0..239}$ , is then divided to four subframes and each subframe is filtered using the  $W_i(z)$  filter, and the weighted output speech signal,  $\{f[n]\}_{n=0..239}$  is obtained.

## 2.9 Pitch estimation

File : EXC_LBC.C	Procedure : Estim_Pitch()	Open loop pitch estimation
------------------	---------------------------	----------------------------

Two pitch estimates are computed for every frame, one for the first two subframes and one for the last two. The open loop pitch period estimate,  $L_{OL}$ , is computed using the perceptually weighted speech  $f[n]$ . A crosscorrelation criterion,  $C_{OL}(j)$ , maximization method is used to determine the pitch period, using the following expression:

$$C_{OL}(j) = \frac{\left( \sum_{n=0}^{119} f[n] \cdot f[n-j] \right)^2}{\sum_{n=0}^{119} f[n-j] \cdot f[n-j]}, \quad 18 \leq j \leq 142 \quad (12)$$

The index  $j$  which maximizes the crosscorrelation,  $C_{OL}(j)$ , is selected as the open loop pitch estimation for the appropriate two subframes. While searching for the best index, some preference is given to smaller pitch periods to avoid choosing pitch multiples. Maximums of  $C_{OL}(j)$  are searched for beginning with  $j = 18$ . For every maximum  $C_{OL}(j)$  found, its value is compared to the best previous maximum found,  $C_{OL}(j')$ . If the difference between indices  $j$  and  $j'$  is less than 18 and  $C_{OL}(j) > C_{OL}(j')$ , the new maximum is selected. If the difference between the indices is greater than or equal to 18, the new maximum is selected only if  $C_{OL}(j)$  is greater than  $C_{OL}(j')$  by 1.25 dB.

## 2.10 Subframe processing

From this point on, all the computational blocks are performed on a once per subframe basis.

## 2.11 Harmonic noise shaping

File : EXC_LBC.C	Procedure : Comp_Pw()	Computes harmonic noise filter coefficients
File : EXC_LBC.C	Procedure : Filt_Pw()	Applies harmonic noise filter

In order to improve the quality of the encoded speech, a harmonic noise shaping filter is constructed. The filter is:

$$P_i(z) = 1 - \beta z^{-L} \quad (13)$$

The optimal lag,  $L$ , for this filter is the lag which maximizes the criterion,  $C_{PW}(j)$ , while considering only positive correlation values for the numerator,  $N(j)$ , before squaring:

$$N(j) = \sum_{n=0}^{59} f[n] \cdot f[n-j] \quad (14.1)$$

$$C_{PW}(j) = \frac{(N(j))^2}{\sum_{n=0}^{59} f[n-j] \cdot f[n-j]}, L_1 \leq j \leq L_2 \quad (14.2)$$

where  $L_1 = L_{OL} - 3$  and  $L_2 = L_{OL} + 3$ . The maximum value will be defined as  $C_L$ . The optimal filter gain,  $G_{opt}$ , is:

$$G_{opt} = \frac{\sum_{n=0}^{59} f[n]f[n-L]}{\sum_{n=0}^{59} f[n-L]f[n-L]} \quad (15)$$

$G_{opt}$  is limited to the range [0,1]. The energy,  $E$ , of the weighted speech signal,  $\{f[n]\}_{n=0..59}$  is given by:

$$E = \sum_{n=0}^{59} f^2[n] \quad (16)$$

Then the coefficient  $\beta$  of the harmonic noise shaping filter,  $P(z)$ , is given by:

$$\beta = \begin{cases} 0.3125 G_{opt}, & \text{if } -10 \log_{10} \left( 1 - \frac{C_L}{E} \right) \geq 2.0 \\ 0.0, & \text{otherwise} \end{cases} \quad (17)$$

After computing the harmonic noise filter coefficients, the formant perceptually weighted speech,  $f[n]$ , is filtered using  $P(z)$  to obtain the target vector,  $w[n]$ .

$$w[n] = f[n] - \beta f[n-L], \quad 0 \leq n \leq 59 \quad (18)$$

## 2.12 Impulse response calculator

File : LPC.C	Procedure : Comp_Ir()	Impulse response computation
--------------	-----------------------	------------------------------

For closed loop analysis the following combined filter,  $S_i(z)$ , is used:

$$S_i(z) = \tilde{A}_i(z) \cdot W_i(z) \cdot P_i(z), \quad 0 \leq i \leq 3 \quad (19)$$

where the components of  $S_i(z)$  are defined in the formulas 10, 11 and 13. The impulse response of this filter is computed and will be referred as  $\{h_i[n]\}_{n=0..59, i=0..3}$ .

### 2.13 Zero input response and ringing subtraction

File : LPC.C	Procedure : Sub_Ring ()	Performs ringing subtraction
--------------	-------------------------	------------------------------

The zero input response of the combined filter,  $S_i(z)$ , is obtained by computing the output of that filter when the input signal is all zero-valued samples. The zero input response is denoted  $\{z[n]\}_{n=0..59}$ . The ringing subtraction is performed by subtracting the zero input response from the harmonic weighted speech vector,  $\{w[n]\}_{n=0..59}$ . The resulting vector is defined as  $t[n] = w[n] - z[n]$ .

### 2.14 Pitch predictor

File : EXC_LBC.C	Procedure : Find_Acbk()	Adaptive codebook contribution. Calls Get_Rez() and Decod_Acbk()
File : EXC_LBC.C	Procedure : Get_Rez()	Gets residual from the excitation buffer
File : EXC_LBC.C	Procedure : Decod_Acbk()	Decodes the adaptive codebook contribution

The pitch prediction contribution is treated as a conventional adaptive codebook contribution. The pitch predictor is a fifth order pitch predictor (see equation 41.2). For subframes 0 and 2 the closed loop pitch lag is selected from around the appropriate open loop pitch lag in the range  $\pm 1$  and coded using 7 bits. (Note that the open loop pitch lag is never transmitted.) For subframes 1 and 3 the closed loop pitch lag is coded differentially using 2 bits and may differ from the previous subframe lag only by  $-1, 0, +1$  or  $+2$ . The quantized and decoded pitch lag values will be referred to as  $L_i$  from this point on. The pitch predictor gains are vector quantized using two codebooks with 85 or 170 entries for the high bit rate and 170 entries for the low bit rate. The 170 entry codebook is the same for both rates. For the high rate if  $L_0$  is less than 58 for subframes 0 and 1 or if  $L_2$  is less than 58 for subframes 2 and 3, then the 85 entry codebook is used for the pitch gain quantization. Otherwise the pitch gain is quantized using the 170 entry codebook. The contribution of the pitch predictor,  $\{p[n]\}_{n=0..59}$ , is subtracted from the target vector  $\{t[n]\}_{n=0..59}$ , to obtain the residual signal  $\{r[n]\}_{n=0..59}$ .

$$r[n] = t[n] - p[n] \quad (20)$$

### 2.15 High rate excitation (MP-MLQ)

File : EXC_LBC.C	Procedure : Find_Fcbk()	Fixed codebook contribution
File : EXC_LBC.C	Procedure : Find_Best()	Residual signal quantization
File : EXC_LBC.C	Procedure : Gen_Trn()	Generates a train of Dirac functions
File : EXC_LBC.C	Procedure : Fcbk_Pack()	Combinatorial coding of pulse positions

The residual signal  $\{r[n]\}_{n=0..59}$ , is transferred as a new target vector to the MP-MLQ block. This block performs the quantization of this vector. The quantization process is approximating the target vector  $r[n]$  by  $r'[n]$ :

$$r'[n] = \sum_{j=0}^n h[j] \cdot v[n-j], \quad 0 \leq n \leq 59 \quad (21)$$

where  $v[n]$  is the excitation to the combined filter  $S(z)$  with impulse response  $h[n]$  and defined as:

$$v[n] = G \sum_{k=0}^{M-1} \alpha_k d[n-m_k], \quad 0 \leq n \leq 59 \quad (22)$$

where  $G$  is the gain factor,  $\delta[n]$  is a Dirac function,  $\{\alpha_k\}_{k=0..M-1}$  and  $\{m_k\}_{k=0..M-1}$  are the signs ( $\pm 1$ ) and the positions of Dirac functions respectively and  $M$  is the number of pulses, which is 6 for even subframes and is 5 for odd subframes. There is a restriction on pulses positions. The positions can be either all odd or all even. This will be indicated by a grid bit. So, the problem is to estimate the unknown parameters,  $G$ ,  $\{\alpha_k\}_{k=0..M-1}$ , and  $\{m_k\}_{k=0..M-1}$ , that minimize the mean square of the error signal  $err[n]$ :

$$err[n] = r[n] - r'[n] = r[n] - G \sum_{k=0}^{M-1} \alpha_k h[n-m_k] \quad (23)$$

The parameters estimation and quantization processes are based on an analysis-by-synthesis method. The  $G_{max}$  parameter is estimated and quantized as follows. First the crosscorrelation function  $d[j]$  between the impulse response,  $h[n]$ , and the new target vector,  $r[n]$ , is computed:

$$d[j] = \sum_{n=j}^{59} r[n] \cdot h[n-j], \quad 0 \leq j \leq 59 \quad (24)$$

The estimated gain is given by:

$$G_{max} = \frac{\max \{ |d[j]| \}_{j=0..59}}{\sum_{n=0}^{59} h[n] \cdot h[n]} \quad (25)$$

Then the estimated gain  $G_{max}$  is quantized by a logarithmic quantizer. This scalar gain quantizer is common to both rates and consists of 24 steps, of 3.2 dB each. Around this quantized value,  $\tilde{G}_{max}$ , additional gain values are selected within the range  $[\tilde{G}_{max} - 3.2, \tilde{G}_{max} + 6.4]$ . For each of these gain values the signs and locations of the pulses are sequentially optimized. This procedure is repeated for both the odd and even grids. Finally the combination of the quantized parameters that yields the minimum mean square of  $err[n]$  is selected. The optimal combination of pulse locations and gain is transmitted.  $\binom{30}{M}$  combinatorial coding is used to transmit the pulse locations. Furthermore, using the fact that the number of codewords in the fixed codebooks is not a power of 2, 3 additional bits are saved by combining the 4 most significant bits of the combinatorial codes for the 4 subframes to form a 13-bit index. The C code provides the details on how this information is packed.

To improve the quality of speech with a short pitch period, the following additional procedure is used. If  $L_0$  is less than 58 for subframes 0 and 1 or if  $L_2$  is less than 58 for subframes 2 and 3, a train of Dirac functions with the period of the pitch index,  $L_0$  or  $L_2$ , is used for each location  $m_k$  instead of a single Dirac function in the above quantization

procedure. The choice between a train of Dirac functions or a single Dirac function to represent the residual signal is made based on the mean square error computation. The configuration which yields the lowest mean square error is selected and its parameter indices are transmitted.

## 2.16 Low rate excitation (ACELP)

File : EXC_LBC.C	Procedure : search_T0()	Pitch synchronous excitation
File : EXC_LBC.C	Procedure : ACELP_LBC_code()	Computes innovative vector
File : EXC_LBC.C	Procedure : Cor_h()	Correlations of impulse response
File : EXC_LBC.C	Procedure : Cor_h_X()	Correlation of target vector with impulse response
File : EXC_LBC.C	Procedure : D4i64_LBC()	Algebraic codebook search
File : EXC_LBC.C	Procedure : G_code()	Computes innovation vector gain

A 17-bit algebraic codebook is used for the fixed codebook excitation  $v[n]$ . Each fixed codevector contains, at most, four non-zero pulses. The 4 pulses can assume the signs and positions given in Table 1:

TABLE 1/G.723.1

### ACELP excitation codebook

Sign	Positions
$\pm 1$	0, 8, 16, 24, 32, 40, 48, 56
$\pm 1$	2, 10, 18, 26, 34, 42, 50, 58
$\pm 1$	4, 12, 20, 28, 36, 44, 52, (60)
$\pm 1$	6, 14, 22, 30, 38, 46, 54, (62)

The positions of all pulses can be simultaneously shifted by one (to occupy odd positions) which needs one extra bit. Note that the last position of each of the last two pulses falls outside the subframe boundary, which signifies that the pulse is not present.

Each pulse position is encoded with 3 bits and each pulse sign is encoded in 1 bit. This gives a total of 16 bits for the 4 pulses. Further, an extra bit is used to encode the shift resulting in a 17-bit codebook.

The codebook is searched by minimizing the mean square error between the weighted speech signal,  $r[n]$ , and the weighted synthesis speech given by:

$$E_{\xi} = \|r - \mathbf{G}\mathbf{H}\mathbf{v}_{\xi}\|^2 \quad (26)$$

where  $\mathbf{r}$  is the target vector consisting of the weighted speech after subtracting the zero-input response of the weighted synthesis filter and the pitch contribution,  $\mathbf{G}$  is the codebook gain;  $\mathbf{v}_{\xi}$  is the algebraic codeword at index  $\xi$ ; and  $\mathbf{H}$  is a lower triangular Toeplitz convolution matrix with diagonal  $h(0)$  and lower diagonals  $h(1), \dots, h(L-1)$ , with  $h(n)$  being the impulse response of the weighted synthesis filter  $S_f(z)$ .

It can be shown that the optimum codeword is the one which maximizes the term:

$$\tau_{\xi} = \frac{C_{\xi}^2}{\varepsilon_{\xi}} = \frac{(d^T v_{\xi})^2}{v_{\xi}^T \Phi v_{\xi}} \quad (27)$$

where  $d = H^T r$  is the correlation between the target vector signal,  $r[n]$ , and the impulse response,  $h(n)$ , and  $\Phi = H^T H$  is the covariance matrix of the impulse response. The vector  $d$  and the matrix  $\Phi$  are computed prior to the codebook search. The elements of the vector  $d$  are computed by:

$$d(j) = \sum_{n=j}^{59} r[n] \cdot h[n-j], \quad 0 \leq j \leq 59 \quad (28)$$

and the elements of the symmetric matrix  $\Phi(i, j)$  are computed by:

$$\Phi(i, j) = \sum_{n=j}^{59} h[n-i] \cdot h[n-j], \quad \begin{matrix} j \geq i \\ 0 \leq i \leq 59 \end{matrix} \quad (29)$$

NOTE – Only the elements actually needed are computed and an efficient storage has been designed that speeds the search procedure up.

The algebraic structure of the codebook allows for very fast search procedures since the excitation vector  $v_{\xi}$  contains only 4 non-zero pulses. The search is performed in 4 nested loops, corresponding to each pulse positions, where in each loop the contribution of a new pulse is added. The correlation in equation (27) is given by:

$$C = \alpha_0 d[m_0] + \alpha_1 d[m_1] + \alpha_2 d[m_2] + \alpha_3 d[m_3] \quad (30)$$

where  $m_k$  is the position of the  $k$ th pulse and  $\alpha_k$  is its sign ( $\pm 1$ ). The energy for even pulse position codevectors in equation (27) is given by:

$$\begin{aligned} \varepsilon &= \Phi(m_0, m_0) \\ &+ \Phi(m_1, m_1) + 2\alpha_0\alpha_1\Phi(m_0, m_1) \\ &+ \Phi(m_2, m_2) + 2[\alpha_0\alpha_2\Phi(m_0, m_2) + \alpha_1\alpha_2\Phi(m_1, m_2)] \\ &+ \Phi(m_3, m_3) + 2[\alpha_0\alpha_3\Phi(m_0, m_3) + \alpha_1\alpha_3\Phi(m_1, m_3) + \alpha_2\alpha_3\Phi(m_2, m_3)] \end{aligned} \quad (31)$$

For odd pulse position codevectors, the energy in equation (27) is approximated by the energy of the equivalent even pulse position codevector obtained by shifting the odd position pulses to one sample earlier in time. To simplify the search procedure, the functions  $d[j]$  and  $\Phi(m_1, m_2)$  are modified. The simplification is performed as follows (prior to the codebook search). First, the signal  $s[j]$  is defined and then the signal  $d'[j]$  is constructed.

$$\begin{aligned} s[2j] &= s[2j+1] = \text{sign}(d[2j]) && \text{if } |d[2j]| > |d[2j+1]| \\ s[2j] &= s[2j+1] = \text{sign}(d[2j+1]) && \text{otherwise} \end{aligned} \quad (32)$$

and the signal  $d'$  is given by  $d'[j] = d[j]s[j]$ . Second, the matrix  $\Phi$  is modified by including the signal information; that is,  $\Phi'(i, j) = s[i]s[j]\Phi(i, j)$ . The correlation in equation (30) is now given by:

$$C = d'[m_0] + d'[m_1] + d'[m_2] + d'[m_3] \quad (33)$$



and the energy in equation (31) is given by:

$$\begin{aligned}
 \epsilon &= \Phi'(m_0, m_0) \\
 &+ \Phi'(m_1, m_1) + 2\Phi'(m_0, m_1) \\
 &+ \Phi'(m_2, m_2) + 2[\Phi'(m_0, m_2) + \Phi'(m_1, m_2)] \\
 &+ \Phi'(m_3, m_3) + 2[\Phi'(m_0, m_3) + \Phi'(m_1, m_3) + \Phi'(m_2, m_3)]
 \end{aligned} \tag{34}$$

A focused search approach is used to further simplify the search procedure. In this approach a precomputed threshold is tested before entering the last loop, and the loop is entered only if this threshold is exceeded. The maximum number of times the loop can be entered is fixed so that a low percentage of the codebook is searched. The threshold is computed based on the correlation  $C$ . The maximum absolute correlation and the average correlation due to the contribution of the first three pulses,  $max_3$  and  $av_3$ , are found prior to the codebook search. The threshold is given by:

$$thr_3 = av_3 + (max_3 - av_3) / 2 \tag{35}$$

The fourth loop is entered only if the absolute correlation (due to three pulses) exceeds  $thr_3$ . Note that this results in a variable complexity search. To further control the search, the number of times the last loop is entered (for the 4 subframes) is not allowed to exceed 600. (The average worst case per subframe is 150 times. This can be viewed as searching only  $150 \times 8$  entries of the codebook, ignoring the overhead of the first three loops.)

A special feature of the codebook is that, for pitch delays less than 60, a pitch contribution depending on the index  $PGInd_i$  of the LTP pitch predictor gain vector is added to the code. That is, after the optimum algebraic code  $v[n]$  is determined, it is modified by  $v[n] \leftarrow v[n] + b(PGInd_i)v[n - L_i - e(PGInd_i)]$ , the values  $\beta(PGInd_i)$  and  $\epsilon(PGInd_i)$  are tabulated and  $L_i$  being the integer pitch period. Note that prior to the codebook search, the impulse response should be modified in a similar fashion if  $L_i < 60$ .

The last step, after getting the  $v[n]$  sequence, is quantizing the gain  $G$ . The gain is quantized in the same way as in the high rate excitation. This is done by stepping through the gain quantization table and selecting the index  $MGInd_i$  which minimizes the following expression:  $|G - \tilde{G}_j|, 0 \leq j \leq 23$ .

## 2.17 Excitation decoder

File : EXC_LBC.C	Procedure : Fcbk_Unpk()	Decode fixed codebook excitation
------------------	-------------------------	----------------------------------

The decoding of pulses is performed as follows:

- 1) First the maximum gain  $G_{max}$  index is derived using:

$$MGInd_i = GInd_i - PGInd_i \cdot GSize \tag{36}$$

where  $GSize = 24$  is the size of the  $\tilde{G}$  table and  $PGInd_i$  is obtained in 2.18.

- 2) The positions of the pulses are decoded using  $\binom{30}{M}$  combinatorial decoding where  $M$  is either 6 or 5, for the high rate. For the low rate, direct decoding of the position indices is performed.
- 3) The grid position (even/odd) is derived from the grid bit.
- 4) The pulse signs are derived from the sign bits.
- 5) For the high rate coder, the decoding of the pulse train bit is performed only if  $L_i < 58$ .
- 6) Then the  $v[n]$  vector is reconstructed using the decoded parameters.
- 7) Finally, the pitch contribution,  $u[n]$ , and the pulse contributions,  $v[n]$ , are summed together to form the excitation vector  $e[n]$ .

## 2.18 Decoding of the pitch information

File : EXC_LBC.C	Procedure : Get_Rez()	Gets residual from the excitation buffer
File : EXC_LBC.C	Procedure : Decod_Acbk()	Decodes the adaptive codebook contribution

The decoding of pitch information is performed as described below:

- 1) First, the lag of pitch predictor for even subframes is decoded:

$$L_i = PInd_i + 18, \quad i = 0,2 \quad (37)$$

- 2) The lag of pitch predictor for odd subframes is decoded as follows:

$$L_i = L_{i-1} + \Delta_i, \quad i = 1,3 \quad (38)$$

where  $\Delta_i \in \{-1,0,+1,+2\}$ .

- 3) The gain vector of the pitch predictor in  $i$ th subframe is derived from the gain index  $GInd_i$ . For the low rate this index contains the information about the pitch predictor gain vector and the index of the gain of the pulse sequence. In this case pitch gain index  $PGInd_i$  is derived as follows:

$$PGInd_i = \lfloor GInd_i / GSize \rfloor, \quad i = 0..3 \quad (39)$$

where  $\lfloor x \rfloor$  indicates the greatest integer  $\leq x$ . For the high rate, in the case that the condition  $L_i \geq 58$  is met, this index is derived in the same manner as described in (39). In the above cases  $PGInd_i$  is a pointer to 170 entries gain vector codebook. Otherwise, this index is a pointer to 85 entries gain vector codebook and contains an additional information about the impulse train bit. In this case pitch gain index is derived as follows:

$$PGInd_i = \lfloor GInd_i \&0x7FF / GSize \rfloor, \quad i = 0..3 \quad (40)$$

The pitch predictor lag and gain vector are decoded from these indices and utilized for the pitch contribution  $u[n]$  extraction as described below. First, a signal  $e'[n]$  is defined by:

$$\begin{aligned} e'[0] &= e[-L_i - 2] \\ e'[1] &= e[-L_i - 1] \\ e'[n] &= e[(n \bmod L_i) - L_i], \quad 2 \leq n \leq 63 \end{aligned} \quad (41.1)$$

where mod stands for the modulus operation. Then,

$$u[n] = \sum_{j=0}^{j=4} \beta_{ij} e'[n + j], \quad 0 \leq n \leq 59 \quad (41.2)$$

## 2.19 Memory update

File : LPC.C	Procedure : Upd_Ring ()	Memory update
--------------	-------------------------	---------------

The last task of the  $i$ th subframe before proceeding to encode the next subframe is to update the memories of the synthesis filter  $A_i(z)$ , the formant perceptual weighting filter  $W_i(z)$ , and the harmonic noise shaping filter  $P_i(z)$ . To accomplish this, the complete response of combined filter  $S_i(z)$ , is computed by passing the reconstructed excitation sequence through this filter. At the end of the excitation filtering, the memory of the combined filter is saved and will be used to compute the zero input response during the encoding of the next speech vector.

## 2.20 Bit allocation

File : UTIL_LBC.C	Procedure : Line_Pack()	Bitstream packing
-------------------	-------------------------	-------------------

This subclause presents the bit allocation tables for both high and low bit rates. The major differences between two rates are in the pulse positions and amplitudes coding. Also, at the lower rate 170 codebook entries are always used for the gain vector of the long term predictor. See Tables 2, 3 and 4.

TABLE 2/G.723.1

### Bit allocation of the 6.3 kbit/s coding algorithm

Parameters coded	Subframe 0	Subframe 1	Subframe 2	Subframe 3	Total
LPC indices					24
Adaptive codebook lags	7	2	7	2	18
All the gains combined	12	12	12	12	48
Pulse positions	20	18	20	18	73 (Note)
Pulse signs	6	5	6	5	22
Grid index	1	1	1	1	4
Total:					189

NOTE – By using the fact that the number of codewords in the fixed codebook is not a power of 2, 3 additional bits are saved by combining the 4 MSB of each pulse position index into a single 13-bit word.

TABLE 3/G.723.1

### Bit allocation of the 5.3 kbit/s coding algorithm

Parameters coded	Subframe 0	Subframe 1	Subframe 2	Subframe 3	Total
LPC indices					24
Adaptive codebook lags	7	2	7	2	18
All the gains combined	12	12	12	12	48
Pulse positions	12	12	12	12	48
Pulse signs	4	4	4	4	16
Grid index	1	1	1	1	4
Total:					158

TABLE 4/G.723.1

## List of transmitted parameters

Name	Transmitted parameters	High rate # bits	Low rate # bits
LPC	LSP VQ index	24	24
ACL0	Adaptive CodeBook Lag	7	7
ACL1	Differential Adaptive CodeBook Lag	2	2
ACL2	Adaptive CodeBook Lag	7	7
ACL3	Differential Adaptive CodeBook Lag	2	2
GAIN0	Combination of adaptive and fixed gains	12	12
GAIN1	Combination of adaptive and fixed gains	12	12
GAIN2	Combination of adaptive and fixed gains	12	12
GAIN3	Combination of adaptive and fixed gains	12	12
POS0	Pulse positions index	20 (Note)	12
POS1	Pulse positions index	18 (Note)	12
POS2	Pulse positions index	20 (Note)	12
POS3	Pulse positions index	18 (Note)	12
PSIG0	Pulse sign index	6	4
PSIG1	Pulse sign index	5	4
PSIG2	Pulse sign index	6	4
PSIG3	Pulse sign index	5	4
GRID0	Grid index	1	1
GRID1	Grid index	1	1
GRID2	Grid index	1	1
GRID3	Grid index	1	1

NOTE – The 4 MSB of these codewords are combined to form a 13-bit index, MSBPOS.

## 2.21 Coder initialization

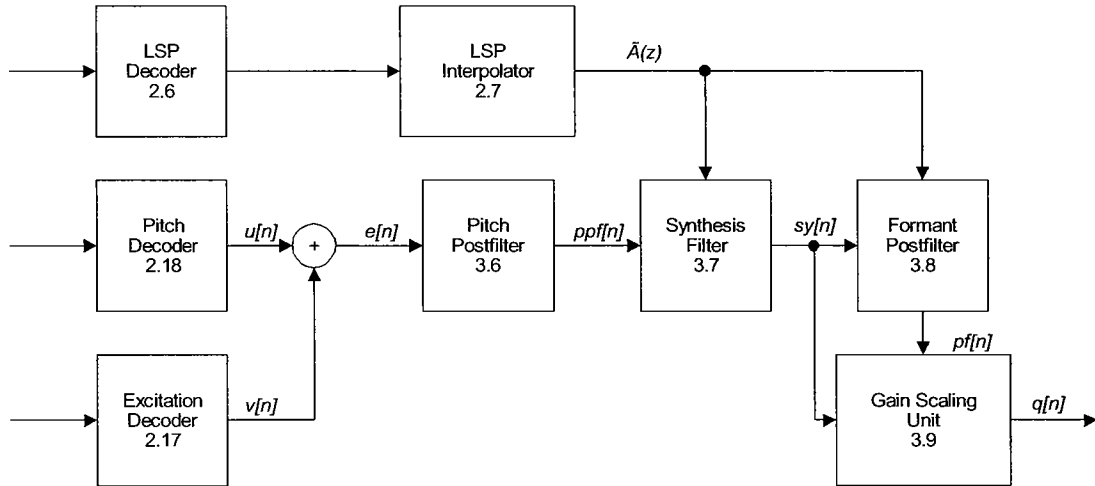
File : CODER.C	Procedure : Init_Coder()	Coder initialization
----------------	--------------------------	----------------------

All the static coder variables should be initialized to 0, with the exception of the previous LSP vector, which should be initialized to LSP DC vector,  $p_{DC}$ .

### 3 Decoder principles

#### 3.1 General description

The decoder operation is also performed on a frame-by-frame basis. First the quantized LPC indices are decoded, then the decoder constructs the LPC synthesis filter. For every subframe, both the adaptive codebook excitation and fixed codebook excitation are decoded and input to the synthesis filter. The adaptive postfilter consists of a formant and a forward-backward pitch postfilter. The excitation signal is input to the pitch postfilter, which in turn is input to the synthesis filter whose output is input to the formant postfilter. A gain scaling unit maintains the energy at the input level of the formant postfilter.



T1517850-95/d02

FIGURE 2/G.723.1

Block diagram of the speech decoder – For each block the corresponding reference number is indicated

#### 3.2 LSP decoder

File : LSP.C	Procedure : Lsp_Inq()	Inverse quantization of LSP
--------------	-----------------------	-----------------------------

This block is the same as in 2.6.

#### 3.3 LSP interpolator

File : LSP.C	Procedure : Lsp_Int()	LSP interpolator
File : LSP.C	Procedure : LsptoA()	Converts LSP to LPC coefficients

This block is the same as in 2.7.

### 3.4 Decoding of the pitch information

File : EXC_LBC.C	Procedure : Get_Rez()	Gets residual from the excitation buffer
File : EXC_LBC.C	Procedure : Decod_Acbk()	Decodes the adaptive codebook contribution

This block is the same as in 2.18.

### 3.5 Excitation decoder

File : EXC_LBC.C	Procedure : Fcbk_Unpk()	Decode fixed codebook excitation
------------------	-------------------------	----------------------------------

This block is the same as in 2.17.

### 3.6 Pitch postfilter

File : EXC_LBC.C	Procedure : Comp_Lpf()	Computes pitch postfilter parameters
File : EXC_LBC.C	Procedure : Find_F()	Forward crosscorrelation maximization
File : EXC_LBC.C	Procedure : Find_B()	Backward crosscorrelation maximization
File : EXC_LBC.C	Procedure : Get_Ind()	Gain computation
File : EXC_LBC.C	Procedure : Filt_Lpf()	Pitch postfiltering

A pitch postfilter is used to improve the quality of the synthesized signal. It is important to note that the pitch postfilter is performed for every subframe, and to implement it, it is required that the whole frame excitation signal  $\{e[n]\}_{n=0..239}$  is generated and saved. The quality improvement is obtained by increasing the SNR at multiples of the pitch period. This is done in the following way. The postfiltered signal  $\{ppf[n]\}_{n=0..59}$  is obtained from the decoded excitation signal  $\{e[n]\}_{n=0..59}$  as given by the following expression:

$$\begin{aligned} ppf[n] &= g_p \cdot \{e[n] + g_{lp}(w_f \cdot g_f \cdot e[n + M_f] + w_b \cdot g_b \cdot e[n - M_b])\} \\ &= g_p \cdot ppf'[n] \end{aligned} \quad (42)$$

where  $e[n]$  is the decoded excitation signal. The computation of the gains,  $g_p$ ,  $g_f$ ,  $g_b$ , and the delays,  $M_f$ ,  $M_b$ , is based on a forward and backward crosscorrelation analysis. The weights  $w_f$ ,  $w_b$  may have the following values: (0,0), (0,1), and (1,0). The delays are selected by maximizing the crosscorrelations. The crosscorrelation for the forward pitch lag is given by:

$$C_f = \sum_{n=0}^{59} e[n]e[n + M_f], M_1 \leq M_f \leq M_2 \quad (43.1)$$

and the crosscorrelation for the backward pitch lag is given by:

$$C_b = \sum_{n=0}^{59} e[n]e[n - M_b], M_1 \leq M_b \leq M_2 \quad (43.2)$$

where  $M_1 = L_i - 3$  and  $M_2 = L_i + 3$  and  $\{L_i\}_{i=0,2}$  are the received pitch lags for the first and third subframes.  $L_0$  is utilized for the first two subframes and  $L_2$  for the last two. Note that if one of the maximum correlations is negative or if for some  $n \in [0...59]$  there is no sample value  $e[n + M_f]$  available, then the corresponding weight and delay are set to 0. This makes four possible cases: (0) both maxima are negative and no pitch postfilter weights need to be computed, (1) only the forward maximum is positive, so it is selected, (2) only the backward maximum is positive, so it is selected, or (3) both maxima are positive and the one making the larger contribution is selected. This procedure is described below. For cases (1), (2) and (3), the relevant signal energies ( $T_{en}$ ,  $D_f$  and/or  $D_b$ ) for the optimum pitch lag ( $M_f$  or  $M_b$ ) are computed according to equations (44.1), (44.2) and (44.3):

$$D_f = \sum_{n=0}^{59} e[n + M_f] e[n + M_f] \quad (44.1)$$

$$D_b = \sum_{n=0}^{59} e[n - M_b] e[n - M_b] \quad (44.2)$$

$$T_{en} = \sum_{n=0}^{59} e[n] e[n] \quad (44.3)$$

The forward energy is given by:

$$E_f = \sum_{n=0}^{59} (e[n] - g_f e[n + M_f])^2 \quad (45.1)$$

and the backward energy is given by:

$$E_b = \sum_{n=0}^{59} (e[n] - g_b e[n - M_b])^2 \quad (45.2)$$

$E$  is minimized by maximizing  $\frac{C^2}{D}$ . In case (3), the selection between forward and backward is made by selecting the larger of  $\frac{C_f^2}{D_f}$  and  $\frac{C_b^2}{D_b}$ . The prediction gain is equal to  $-10 \log_{10} \left( 1 - \frac{C^2}{DT_{en}} \right)$ . If this gain is less than 1.25 dB, then the contribution is judged to be negligible and no pitch postfilter is used. If a pitch postfilter is used, then the optimum gain is given by:

$$g = \frac{C}{D} \quad (46)$$

According to the speech coder bit rate, the optimum gain is multiplied by a weighting factor,  $\gamma_{lp}$ , which equals 0.1875 for the high rate and 0.25 for the low rate. Finally, the scaling gain,  $g_p$ , is computed as:

$$g_p = \sqrt{\frac{\sum_{n=0}^{59} e^2[n]}{\sum_{n=0}^{59} (ppf'[n])^2}} \quad (47)$$

If the denominator in equation (47) is less than the numerator, the gain is set to 1.

### 3.7 LPC synthesis filter

File : LPC.C	Procedure : Synt()	Synthesizes reconstructed speech
--------------	--------------------	----------------------------------

The 10th order LPC synthesis filter  $\tilde{A}_i(z)$  is used to synthesize the speech signal  $sy[n]$  from the decoded pitch postfiltered residual  $ppf[n]$ .

$$sy[n] = ppf[n] + \sum_{j=1}^{10} \tilde{a}_{ij} sy[n - j] \quad (48)$$

### 3.8 Formant postfilter

File : LPC.C	Procedure : Spf()	Formant postfiltering
File : UTIL_LBC.C	Procedure : Comp_En()	Computes synthesized signal energy

A conventional ARMA postfilter is used. The transfer function of the short term postfilter is given by the following equations:

$$k = \frac{\sum_{n=1}^{59} sy[n]sy[n-1]}{\sum_{n=0}^{59} sy[n]sy[n]} \quad (49.1)$$

$$k_1 = \frac{3}{4} k_{old} + \frac{1}{4} k \quad (49.2)$$

$$F(z) = \frac{1 - \sum_{i=1}^{10} \tilde{a}_i \lambda_1^i z^{-i}}{1 - \sum_{i=1}^{10} \tilde{a}_i \lambda_2^i z^{-i}} (1 - 0.25k_1 z^{-1}) \quad (49.3)$$



where  $\lambda_1 = 0.65$  and  $\lambda_2 = 0.75$ ,  $k$  is the first autocorrelation coefficient that is estimated from the synthesized speech  $sy[n]$ , and  $k_{old}$  is the value of  $k_1$  from the previous subframe. The postfiltered signal  $pf[n]$  is obtained as the output of the formant postfilter with input signal  $sy[n]$ .

### 3.9 Gain scaling unit

File : UTIL_LBC.C	Procedure : Scale()	Gain adjustment of postfiltered signal
-------------------	---------------------	--

This unit receives two input vectors, the synthesized speech vector  $\{sy[n]\}_{n=0..59}$  and the postfiltered output vector  $\{pf[n]\}_{n=0..59}$ . First the amplitude ratio  $g_s$  is computed, using:

$$g_s = \sqrt{\frac{\sum_{n=0}^{59} sy^2[n]}{\sum_{n=0}^{59} pf^2[n]}} \quad (50)$$

If the denominator is equal to 0,  $g_s$  is set to 1.

Then the output vector  $q[n]$  is obtained by scaling the postfiltered signal  $pf[n]$  and the gain  $g[n]$  is updated using the following expressions respectively:

$$g[n] = (1 - \alpha)g[n - 1] + \alpha g_s \quad (51)$$

$$q[n] = pf[n] \cdot g[n] \cdot (1 + \alpha) \quad (52)$$

where  $\alpha$  is equal to 1/16.

### 3.10 Frame interpolation handling

File : EXC_LBC.C	Procedure : Comp_Info()	Computes interpolation index
File : EXC_LBC.C	Procedure : Regen()	Current frame regeneration

This coder has been designed to be robust for indicated frame erasures. An error concealment strategy for frame erasures has been included in the decoder. However, this strategy must be triggered by an external indication that the bitstream for the current frame has been erased. Because the coder was designed for burst errors, there is no error correction mechanism provided for random bit errors. If a frame erasure has occurred, the decoder switches from regular decoding to frame erasure concealment mode. The frame interpolation procedure is performed independently for the LSP coefficients and the residual signal.

#### 3.10.1 LSP interpolation

The decoding of the LSP coefficients in frame interpolation mode is performed in the following way:

- 1) The vector  $\tilde{e}_n$  is set to zero.
- 2) The predicted vector,  $\bar{p}_n$ , is added to the vector,  $\tilde{e}_n$ , and DC vector,  $p_{DC}$ , to form the decoded LSP vector,  $\bar{p}_n$ . For  $\bar{p}_n$  generation a different fixed predictor is used:  $b_c = 23/32$ .

From this point the decoding of the LSP continues as in 2.6, except that  $\Delta_{min} = 62.5$  Hz, rather than 31.25 Hz.

### 3.10.2 Residual interpolation

The residual interpolation is performed in two different ways, depending on the last previous good frame prior to the erased frame. The frame is checked with a voiced/unvoiced classifier.

The classifier is based on a cross-correlation maximization function. The last 120 samples of the frame are cross-correlated with  $L_2 \pm 3$ . The index which reaches the maximum correlation value, is chosen as the interpolation index candidate. Then the prediction gain of the best vector is tested. If this gain is more than 0.58 dB the frame is declared as voiced, otherwise the frame is declared as unvoiced.

The classifier returns 0 for the unvoiced case and the estimated pitch value for the voiced case. If the frame was declared unvoiced, the average of the gain indices for subframes 2 and 3 is saved.

If the current frame was marked as erased, and the previous frame was classified as unvoiced, the current frame excitation is generated using a uniform random number generator. The random number generator output is scaled using the previously computed gain value.

In the voiced case, the current frame is regenerated with periodic excitation having a period equal to the value provided by the classifier.

If the frame erasure state continues for the next two frames, the regenerated vector is attenuated by an additional 2.5 dB for each frame. After 3 interpolated frames, the output is muted completely.

### 3.11 Decoder initialization

File : DECOD.C	Procedure : Init_Decod()	Decoder initialization
----------------	--------------------------	------------------------

All the static decoder variables should be initialized to 0, with the exception of:

- 1) Previous LSP vector, which should be initialized to LSP DC vector,  $p_{DC}$ .
- 2) Postfilter gain  $g[-1]$ , which should be initialized to 1.

## 4 Bitstream packing

NOTE – Each bit of transmitted parameters is named PAR(x)\_By: where PAR is the name of the parameter and x indicates the subframe index if relevant and y stands for the bit position starting from 0 (LSB) to the MSB. The expression PAR<sub>x</sub>\_By..PAR<sub>x</sub>\_Bz stands for the range of transmitted bits from bit y to bit z. The unused bit is named UB (value = 0). RATEFLAG\_B0 tells whether the high rate (0) or the low rate (1) is used for the current frame. VADFLAG\_B0 tells whether the current frame is active speech (0) or non-speech (1). The combination of RATEFLAG and VADFLAG both being set to 1 is reserved for future use. Octets are transmitted in the order in which they are listed in Tables 5 and 6. Within each octet, the bits are with the MSB on the left and the LSB on the right.

## 5 ANSI C code

ANSI C code simulating the dual rate encoder/decoder in 16-bit fixed point arithmetic is available from ITU-T. Table 7 lists all the files included in this code.

## 6 Glossary

This is a glossary containing the mathematical symbols used in the text and a brief description of what they represent. See Table 8.

TABLE 5/G.723.1

## Octet bit packing for the high bit rate codec

High rate

Transmitted octets	PARx By, ...
1	LPC_B5...LPC_B0, VADFLAG_B0, RATEFLAG_B0
2	LPC_B13...LPC_B6
3	LPC_B21...LPC_B14
4	ACL0_B5...ACL0_B0, LPC_B23, LPC_B22
5	ACL2_B4...ACL2_B0, ACL1_B1, ACL1_B0, ACL0_B6
6	GAIN0_B3...GAIN0_B0, ACL3_B1, ACL3_B0, ACL2_B6, ACL2_B5
7	GAIN0_B11...GAIN0_B4
8	GAIN1_B7...GAIN1_B0
9	GAIN2_B3...GAIN2_B0, GAIN1_B11...GAIN1_B8
10	GAIN2_B11...GAIN2_B4
11	GAIN3_B7...GAIN3_B0
12	GRID3_B0, GRID2_B0, GRID1_B0, GRID0_B0, GAIN3_B11...GAIN3_B8
13	MSBPOS_B6...MSBPOS_B0, UB
14	POS0_B1, POS0_B0, MSBPOS_B12...MSBPOS_B7
15	POS0_B9...POS0_B2
16	POS1_B2, POS1_B0, POS0_B15...POS0_B10
17	POS1_B10...POS1_B3
18	POS2_B3...POS2_B0, POS1_B13...POS1_B11
19	POS2_B11...POS2_B4
20	POS3_B3...POS3_B0, POS2_B15...POS2_B12
21	POS3_B11...POS3_B4
22	PSIG0_B5...PSIG0_B0, POS3_B13, POS3_B12
23	PSIG2_B2...PSIG2_B0, PSIG1_B4...PSIG1_B0
24	PSIG3_B4...PSIG3_B0, PSIG2_B5...PSIG2_B3

TABLE 6/G.723.1

## Octet bit packing for the low bit rate codec

Low rate

Transmitted octets	PARx By, ...
1	LPC_B5...LPC_B0, VADFLAG_B0, RATEFLAG_B0
2	LPC_B13...LPC_B6
3	LPC_B21...LPC_B14
4	ACL0_B5...ACL0_B0, LPC_B23, LPC_B22
5	ACL2_B4...ACL2_B0, ACL1_B1, ACL1_B0, ACL0_B6
6	GAIN0_B3...GAIN0_B0, ACL3_B1, ACL3_B0, ACL2_B6, ACL2_B5
7	GAIN0_B11...GAIN0_B4
8	GAIN1_B7...GAIN1_B0
9	GAIN2_B3...GAIN2_B0, GAIN1_B11...GAIN1_B8
10	GAIN2_B11...GAIN2_B4
11	GAIN3_B7...GAIN3_B0
12	GRID3_B0, GRID2_B0, GRID1_B0, GRID0_B0, GAIN3_B11...GAIN3_B8
13	POS0_B7...POS0_B0
14	POS1_B3...POS1_B0, POS0_B11...POS0_B8
15	POS1_B11...POS1_B4
16	POS2_B7...POS2_B0
17	POS3_B3...POS3_B0, POS2_B11...POS2_B8
18	POS3_B11...POS3_B4
19	PSIG1_B3...PSIG1_B0, PSIG0_B3...PSIG0_B0
20	PSIG3_B3...PSIG3_B0, PSIG2_B3...PSIG2_B0

TABLE 7/G.723.1

## List of software filenames

File name	Description
TYPEDEF.H	Data type definition is machine dependent
CST_LBC.H	Definition of constants for G.723
LBCCODEC.C	Main program for G.723 speech codecs
LBCCODEC.H	Functions prototypes
CODER.C	G.723 speech encoder for the two-bit rates
CODER.H	Functions prototypes
DECOD.C	G.723 speech decoder for the two-bit rates
DECOD.H	Functions prototypes
LPC.C	Linear predictive analysis
LPC.H	Functions prototypes
LSP.C	Line spectral pair related functions, quantizer
LSP.H	Functions prototypes
EXC_LBC.C	Adaptive and fixed (MP-MLQ, ACELP) excitation
EXC_LBC.H	Functions prototypes
UTIL_LBC.C	Miscellaneous functions (HPF, pack, unpack, I/O...)
UTIL_LBC.H	Functions prototypes
TAB_LBC.C	Tables of constants
TAB_LBC.H	External declaration for constant tables
BASOP.C	Fixed point arithmetic and logical operation
BASOP.H	Functions prototypes

TABLE 8/G.723.1

## Glossary of symbols in the text

Symbol	Description
$y[j]$	Input speech samples
$s[j]$	Input speech frame of 240 samples
$x[j]$	High pass filtered speech frame
$R[j]$	Autocorrelation function, $n = 0, 1, \dots, 10$
$a_i$	LPC coefficient vector of subframe $i$
$\tilde{a}_i$	Quantized LPC coefficient vector of subframe $i$
$p'$	Unquantized LSP vector
$p$	DC-removed LSP vector
$P_{DC}$	Long-term DC vector of LSP values
$\tilde{p}_n$	Decoded LSP vector for frame $n$
$\tilde{p}_n$	DC-removed predicted LSP vector
$e_n$	Residual LSP error vector for frame $n$
$\tilde{e}_n$	Quantized value of $e_n$
$W_n$	Diagonal weighting matrix for LSP quantization
$\gamma_1, \gamma_2$	Weights for perceptual weighting filter, 0.9, 0.5
$W_i$	Formant perceptual weighting filter for subframe $i$
$f[n]$	Formant perceptually weighted speech
$L_{OL}$	Open loop pitch period estimate
$C_{OL}(j)$	Open loop pitch estimate crosscorrelation criterion function
$C_{PW}(j)$	Harmonic noise shaping pitch estimate crosscorrelation criterion function
$\beta$	Harmonic noise shaping filter gain
$L$	Optimal lag for harmonic noise shaping filter
$G_{opt}$	Optimal gain for harmonic noise shaping filter
$E$	Energy of weighted speech signal
$P_i$	Harmonic noise shaping filter for subframe $i$
$S_i$	Combined harmonic and formant weighting and synthesis filters for subframe $i$
$h[n]$	Impulse response of combined filter
$z[n]$	Zero input response of combined filter
$w[n]$	Harmonic noise weighted speech
$t[n]$	Target vector
$p[n]$	Pitch predictor contribution vector
$r[n]$	Residual signal vector
$r'[n]$	Filtered excitation vector
$v[n]$	Fixed codebook excitation vector
$M$	Number of pulses

TABLE 8/G.723.1 (end)

## Glossary of symbols in the text

$\alpha_k$	Sign of pulse k
$m_k$	Position of pulse k
$d[n]$	Crosscorrelation function of $h[n]$ and $r[n]$
$G_{max}$	Estimated gain of pulses for high rate
$L_i$	Pitch lag for subframe $i$
$H$	Lower triangular Toeplitz convolution matrix with diagonals $h[n]$
$\Phi$	Covariance matrix formed by $H^T H$
$max_3$	Maximum correlation of 1st 3 pulses for low rate
$av_3$	Average correlation of 1st 3 pulses for low rate
$thr_3$	Threshold for correlation of 1st 3 pulses for low rate
$MGInd_i$	Maximum gain index of subframe $i$
$GInd_i$	Gain index of subframe $i$
$PInd_i$	Pitch index of subframe $i$
$PGInd_i$	Pitch lag index of subframe $i$
$GSize$	Size of excitation gain codebook, 24
$\tilde{G}$	Quantized gain
$\beta_{ij}$	Pitch predictor gain vector
$u[n]$	Adaptive codebook excitation vector
$e[n]$	Decoded combined excitation vector
$ppf[n]$	Pitch postfiltered excitation signal
$M_f, M_b$	Optimal forward and backward pitch postfilter lags
$ppf'[n]$	Unnormalized pitch postfiltered excitation signal
$\gamma_{tip}$	Gain weighting factor, 0.1875 or 0.25
$g_p$	Pitch postfilter scaling gain
$g_p, g_b$	Optimal forward and backward gains for pitch postfilter
$C_f, C_b$	Forward and backward excitation crosscorrelations
$D_f, D_b$	Forward and backward excitation energies
$E_f, E_b$	Forward and backward energies for pitch postfilter
$T_{en}$	Energy of excitation signal
$sy[n]$	LPC synthesized speech signal
$pf[n]$	Formant postfiltered signal
$q[n]$	Output speech signal
$g[n]$	Formant postfilter gain signal
$k_1$	Interpolated reflection coefficient for tilt compensation filter
$\lambda_1, \lambda_2$	Weights for formant postfilter, 0.65, 0.75
$g_s$	Amplitude ratio of sy and pf vectors
$\alpha$	0.0625 in gain scaling unit
$b_e$	Fixed predictor for LSP interpolation during frame erasure concealment, 23/32



INTERNATIONAL TELECOMMUNICATION UNION

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**G.723.1**

**Annex A**  
(11/96)

SERIES G: TRANSMISSION SYSTEMS AND MEDIA

Digital transmission systems – Terminal equipments –  
Coding of analogue signals by methods other than PCM

---

Dual rate speech coder for multimedia  
communications transmitting at 5.3 and 6.3 kbit/s

**Annex A: Silence compression scheme**

ITU-T Recommendation G.723.1 – Annex A

(Previously CCITT Recommendation)

---



ITU-T G-SERIES RECOMMENDATIONS  
TRANSMISSION SYSTEMS AND MEDIA

INTERNATIONAL TELEPHONE CONNECTIONS AND CIRCUITS	G.100–G.199
<b>INTERNATIONAL ANALOGUE CARRIER SYSTEM</b>	
GENERAL CHARACTERISTICS COMMON TO ALL ANALOGUE CARRIER-TRANSMISSION SYSTEMS	G.200–G.299
INDIVIDUAL CHARACTERISTICS OF INTERNATIONAL CARRIER TELEPHONE SYSTEMS ON METALLIC LINES	G.300–G.399
GENERAL CHARACTERISTICS OF INTERNATIONAL CARRIER TELEPHONE SYSTEMS ON RADIO-RELAY OR SATELLITE LINKS AND INTERCONNECTION WITH METALLIC LINES	G.400–G.449
COORDINATION OF RADIOTELEPHONY AND LINE TELEPHONY	G.450–G.499
<b>TRANSMISSION MEDIA CHARACTERISTICS</b>	G.600–G.699
<b>DIGITAL TRANSMISSION SYSTEMS</b>	
TERMINAL EQUIPMENTS	G.700–G.799
General	G.700–G.709
Coding of analogue signals by pulse code modulation	G.710–G.719
<b>Coding of analogue signals by methods other than PCM</b>	<b>G.720–G.729</b>
Principal characteristics of primary multiplex equipment	G.730–G.739
Principal characteristics of second order multiplex equipment	G.740–G.749
Principal characteristics of higher order multiplex equipment	G.750–G.759
Principal characteristics of transcoder and digital multiplication equipment	G.760–G.769
Operations, administration and maintenance features of transmission equipment	G.770–G.779
Principal characteristics of multiplexing equipment for the synchronous digital hierarchy	G.780–G.789
Other terminal equipment	G.790–G.799
DIGITAL NETWORKS	G.800–G.899
General aspects	G.800–G.809
Design objectives for digital networks	G.810–G.819
Quality and availability targets	G.820–G.829
Network capabilities and functions	G.830–G.839
SDH network characteristics	G.840–G.899
DIGITAL SECTIONS AND DIGITAL LINE SYSTEM	G.900–G.999
General	G.900–G.909
Parameters for optical fibre cable systems	G.910–G.919
Digital sections at hierarchical bit rates based on a bit rate of 2048 kbit/s	G.920–G.929
Digital line transmission systems on cable at non-hierarchical bit rates	G.930–G.939
Digital line systems provided by FDM transmission bearers	G.940–G.949
Digital line systems	G.950–G.959
Digital section and digital transmission systems for customer access to ISDN	G.960–G.969
Optical fibre submarine cable systems	G.970–G.979
Optical line systems for local and access networks	G.980–G.999

*For further details, please refer to ITU-T List of Recommendations.*

**ITU-T RECOMMENDATION G.723.1 – Annex A**

**SILENCE COMPRESSION SCHEME**

**Source**

Annex A to ITU-T Recommendation G.723.1, was prepared by ITU-T Study Group 15 (1993-1996) and was approved under the WTSC Resolution No. 1 procedure on the 8th of November 1996.

## FOREWORD

ITU (International Telecommunication Union) is the United Nations Specialized Agency in the field of telecommunications. The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of the ITU. The ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Conference (WTSC), which meets every four years, establishes the topics for study by the ITU-T Study Groups which, in their turn, produce Recommendations on these topics.

The approval of Recommendations by the Members of the ITU-T is covered by the procedure laid down in WTSC Resolution No. 1 (Helsinki, March 1-12, 1993).

In some areas of information technology which fall within ITU-T's purview, the necessary standards are prepared on a collaborative basis with ISO and IEC.

## NOTE

In this Recommendation, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

© ITU 1997

All rights reserved. No part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from the ITU.

## CONTENTS

	<b>Page</b>
A.1 Introduction .....	1
A.2 Description of the VAD.....	2
A.2.1 Adaptation enable flag computation .....	2
A.2.2 Inverse filtering .....	3
A.2.3 Filtered energy computation.....	3
A.2.4 Noise level computation.....	3
A.2.5 Threshold computation.....	4
A.2.6 The VAD decision.....	4
A.2.7 VAD hangover addition .....	4
A.2.8 VAD initialization .....	4
A.3 General description of the CNG .....	4
A.4 Description of the CNG encoder part.....	6
A.4.1 Computation of the frame autocorrelation function.....	6
A.4.2 Computation of the current frame type $F_{typ_t}$ .....	7
A.4.3 Quantization of the average energy.....	8
A.4.4 Computation and coding of SID parameters.....	9
A.4.5 Computation of the CNG excitation.....	10
A.4.6 Interpolation of LSPs and update.....	11
A.4.7 COD-CNG initialization .....	11
A.5 Description of the decoder part .....	11
A.5.1 Description of DEC-CNG .....	12
A.5.2 Frame erasure concealment with regards to the CNG.....	13
A.5.3 DEC-CNG initialization.....	14
A.6 Bit stream packing .....	14
A.7 Glossary.....	14
A.8 Bit-exact, fixed-point C source code.....	15



**Recommendation G.723.1 – Annex A****SILENCE COMPRESSION SCHEME***(Geneva, 1996)***A.1 Introduction**

This Annex describes the silence compression system that has been designed for the G.723.1 speech coder. Silence compression techniques are used to reduce the transmitted bit rate during silent intervals of speech. Systems allowing discontinuous transmission are based on a Voice Activity Detection (VAD) algorithm and a Comfort Noise Generator (CNG) algorithm that allows the insertion of an artificial noise during silence periods. This feature is necessary to avoid noise modulation introduced when the transmission is switched off: if the background acoustic noise that was present during active periods abruptly disappears, this very unpleasant noise modulation may even reduce the intelligibility of the speech.

The purpose of the VAD is to reliably detect the presence or absence of speech and to convey this information to the CNG algorithm. Typically, VAD algorithms base their decisions on several successive frames of information in order to make them more reliable and to avoid producing intermittent decisions. The VAD is constrained to operate on the same 30 ms speech frames which will subsequently either be encoded by the speech coder or filled with comfort noise by the comfort noise generator. The output of the VAD algorithm is passed to the CNG algorithm.

The largest difficulty in the detection of speech is the presence of any of a diverse range of background noise conditions. The VAD must be able to detect speech even in very low signal-to-noise ratio conditions. It is impossible to distinguish between speech and noise using simple level detection techniques when parts of the speech utterance are buried below the noise. The distinction between these conditions can only be made by taking into consideration the spectral characteristics of the input signal. In order to do this, the VAD incorporates an inverse filter, the coefficients of which are derived during noise-only periods by the CNG. All further details of the VAD are included in A.2.

The purpose of the CNG algorithm is to create a noise that matches the actual background noise with a global transmission cost as low as possible. At the transmitting end, the CNG algorithm uses the activity information given by the VAD for each frame, then computes the encoded parameters needed to synthesize the artificial noise at the receiving end. These encoded parameters compose the Silence Insertion Descriptor (SID) frames, which require less bits than the active speech frames and are transmitted during inactive periods.

The main feature of this CNG algorithm is that the transmission of SID frames is not periodic: for each inactive frame, the algorithm makes the decision of sending a SID frame or not, based on a comparison between the current inactive frame and the preceding SID frame. In this way, the transmission of the SID frames is limited to the frames where the power spectrum of the noise has changed.

During inactive frames, the comfort noise is synthesized at the decoder by introducing a pseudo-white excitation into the short-term synthesis filter. The parameters used to characterize the comfort noise are the LPC synthesis filter coefficients and the energy of the excitation signal. At the encoder, for each SID frame the algorithm computes a set of LPC parameters and quantizes the corresponding LSPs using the coder LSP quantizer on 24 bits. It also evaluates the excitation energy and quantizes it with 6 bits. This yields encoded SID frames of 4 bytes including the 2 bits for bit rate and DTX information.

A notable feature of this CNG algorithm is the method used to evaluate the spectrum of the ambient noise for each SID frame. It takes into account the local stationarity or non-stationarity of the input signal.

Finally, the excitation corresponds to the higher bit rate excitation of the G.723.1 codec. Since the fixed excitation has a rather poor spectrum, the long-term excitation is also used in order to obtain a better white-noise-type of excitation. The algorithm randomly chooses the codes of the long-term parameters (delays and gains) and the fixed codebook parameters (grid, pulse positions and signs). For every two subframes, it computes the gain of the fixed excitation to achieve a global energy derived from the transmitted SID energy.

The computation of the excitation needs to be performed both at the encoder and at the decoder to keep both parts synchronized.

At the receiver, to simplify the procedure, the harmonic postfilter is switched off during comfort noise generation since the generated noise is not a voiced signal.

The results of the tests on the VAD/DTX/CNG scheme as described in this Annex will be published at a later date as an appendix to Annex A to Recommendation G.723.1.<sup>1</sup>

## A.2 Description of the VAD

This subclause describes the Voice Activity Detector (VAD) used in the G.723.1 speech coder. The function of the VAD is to indicate whether each 30 msec frame produced by the speech encoder contains speech or not. The VAD decision at frame  $t$  is labelled as  $Vad_t$  and is the input to the COD-CNG block that computes  $Ftyp_t$ , as described in A.3 and Figure A.1. The performance of the VAD algorithm is characterized by the amount of audible speech clipping and the percentage of speech activity it indicates.

The VAD is basically an energy detector. The energy of the inverse filtered signal is compared with a threshold. Speech is indicated whenever the threshold is exceeded. The threshold is computed by a two-step procedure. First, the noise level is updated based on its previous value and the energy of the filtered signal. Second, the threshold is computed from the noise level via a logarithmic approximation.

Hangover is a term describing the practice of declaring the first few frames of silence following a speech burst to still be speech. It is used to eliminate low level speech clipping. Hangover is only added to speech bursts which exceed a certain duration to avoid extending noise spikes.

### A.2.1 Adaptation enable flag computation

An adaptation enable flag, denoted  $Aen_t$  for the current frame  $t$ , is used to be sure that the VAD noise level is adapted only when speech is not present. It is based on the fact that the background noise or the silence is neither a voiced signal nor a sine wave:

– Voiced/Unvoiced detection:

The open loop pitch delays of the preceding and current frame are used to test voicing. Let us note  $L_{OL}^j, j=0,1,2,3$  those four values. The minimum delay  $L_{OL}^{\min} = \text{Min}(L_{OL}^j, j=0,1,2,3)$  is first computed. The counter  $pc \in [1,2,3,4]$  indicating how many delays  $L_{OL}^j$  lie in the neighborhood of a multiple of  $L_{OL}^{\min}$  ( $\pm 3$ ) is evaluated. If  $pc$  is equal to 4 the signal is considered as voiced.

<sup>1</sup> It should be noted that test conditions were not sufficiently severe for mobile conditions.

- Sine wave detection : (already present COM 15-255 Contribution)

The following sine wave detector is included in the LPC analysis of the G.723.1 encoder:

Let  $k'_i[2]$  be the second reflection coefficient computed by the Durbin recursion for each subframe  $i = 0, \dots, 3$  of frame  $t$ .

If  $k'_i[2] \geq 0.95$  for at least 14 of the 15 last values, then a sine wave is detected ( $SinD = 1$ ).

In the other case,  $SinD = 0$ .

- Compute the adaptation enable flag:

$$\begin{cases} Aen_t = Aen_{t-1} + 2 & \text{if } pc = 4 \text{ or } SinD = 1 \\ Aen_t = Aen_{t-1} - 1 & \text{otherwise} \end{cases}$$

$Aen_t$  is bounded into  $[0,6]$ .

### A.2.2 Inverse filtering

The input signal frame,  $\{s[n]\}_{n=60..239}$ , is inverse filtered by a FIR filter  $A_{no}(z)$  with coefficients  $\{a_{no}[j]\}_{j=1..10}$ . This filter is calculated by the CNG block and provides an estimation of the LPC filter associated to the current background noise.

$$e'_t[n] = s[n] + \sum_{j=1}^{10} a_{no}[j] \cdot s[n-j] \quad n = 60 \rightarrow 239 \quad (\text{A-1})$$

where  $e'_t[n]$  is the inverse filtered signal.

### A.2.3 Filtered energy computation

The energy,  $Enr_t$ , is computed from the inverse filtered signal of the current frame by:

$$Enr_t = \frac{1}{80} \sum_{n=60}^{239} e'^2_t[n] \quad (\text{A-2})$$

### A.2.4 Noise level computation

The noise level at frame  $t$ ,  $Nlev_t$ , is updated based on its previous value and on the previous energy,  $Enr_{t-1}$  and on the adaptation enable flag  $Aen_t$ . This update procedure is characterized by slow attack and fast decay. The dynamic range of the noise level at frame  $t$  is limited to the range  $[Nlev_{\min}, Nlev_{\max}]$ .

- 1) If  $Nlev_{t-1} > Enr_{t-1}$  then the noise level is first clipped:

$$Nlev_t = \begin{cases} 0.25 \cdot Nlev_{t-1} + 0.75 \cdot Enr_{t-1} & \text{if } Nlev_{t-1} > Enr_{t-1} \\ Nlev_{t-1} & \text{otherwise} \end{cases} \quad (\text{A-3})$$

- 2) Then  $Nlev_t$  is increased, if adaptation is enabled, otherwise it is decreased by a small amount:

$$Nlev_t = \begin{cases} 1.03125 \times Nlev_t & \text{if } Aen_t = 0 \\ 0.9995 \times Nlev_t & \text{otherwise} \end{cases} \quad (\text{A-4})$$

$$\text{with } \begin{cases} Nlev_{\min} = 128 \\ Nlev_{\max} = 131071 \end{cases}$$



### A.2.5 Threshold computation

The relationship between the noise level at frame  $t$ ,  $Nlev_t$ , and the threshold,  $Thr$ , is defined by logarithmic approximation and defined by the following formula:

$$Thr = \begin{cases} 5.012 & \text{if } Nlev = 128, \\ 10^{0.7-0.05\log_2 \frac{Nlev}{128}} & \text{if } 128 < Nlev < 16384 \\ 2.239 & \text{if } Nlev \geq 16384 \end{cases} \quad (\text{A-5})$$

### A.2.6 The VAD decision

The VAD decision is based on the comparison between the threshold,  $Thr$ , and the current energy,  $Enr_t$ .

$$Vad_t = \begin{cases} 1 & \text{if } Enr_t \geq Thr \\ 0 & \text{if } Enr_t < Thr \end{cases} \quad (\text{A-6})$$

### A.2.7 VAD hangover addition

A hangover of 6 frames is added only in the case of speech bursts ( $Vad_t = 1$ ) larger or equal than 2 frames.

### A.2.8 VAD initialization

All static variables of the VAD algorithm are initialized to zero, except the following variables:

$$\begin{aligned} Nlev_{-1} &= 1024 \\ Enr_{-1} &= 1024 \\ L_{OL}^j &= 1 & j = 0,1 \\ L_{OI}^j &= 60 & j = 2,3 \end{aligned} \quad (\text{A-7})$$

## A.3 General description of the CNG

The algorithm is divided into two blocks situated at the encoder and the decoder, that will be called respectively COD-CNG and DEC-CNG. At the encoder (see Figure A.1), the COD-CNG block uses the autocorrelation function of the speech signal computed for each 60 samples subframe, the past excitation samples and LSPs from the preceding frame.

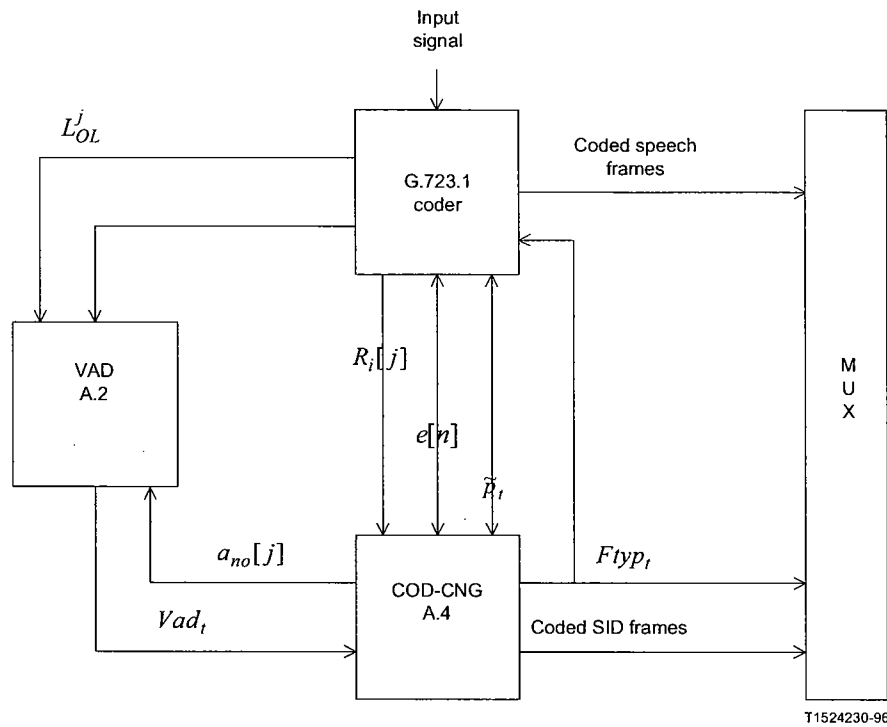


FIGURE A.1/G.723.1

**Block diagram of the encoder with VAD/CNG**

For inactive frames, COD-CNG computes the CNG excitation samples in order to synchronize the local decoder of the encoder with the distant decoder.

Because of the predictive coding of the LSPs in the G.723.1 scheme, a similar input/output with update is done for LSP parameters during inactive frames.

COD-CNG outputs the encoded SID frames and the final decision  $Ftyp_t$  (Frame type of frame  $t$ ) as one of the three values, 0, 1, or 2 corresponding to untransmitted frame, active speech frame or SID frame, respectively.

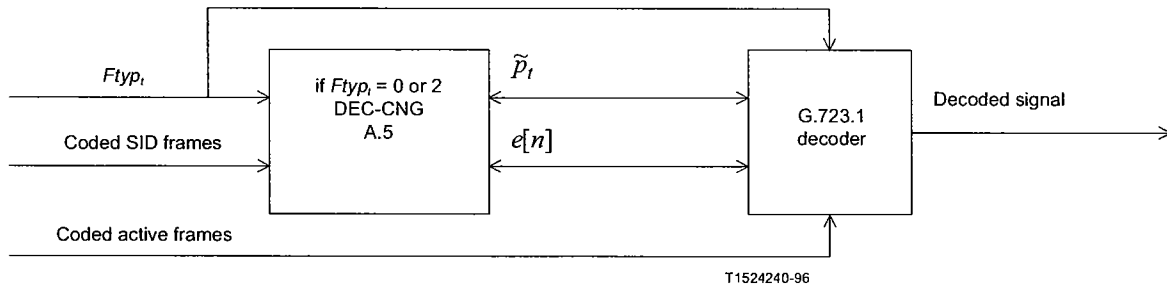


FIGURE A.2/G.723.1

**Block diagram of the decoder with VAD/DTX**

At the receiver (see Figure A.2), the DEC-CNG block processes only inactive speech frames, for which the input information  $Ftyp_t$  is equal to 0 or 2 (untransmitted/SID). DEC-CNG decodes the SID frames and both for SID and untransmitted frames, computes the current LSPs and excitation using the same method as COD-CNG.

Then the G.723.1 decoder synthesizes the comfort noise using the CNG excitation and LSPs.

**A.4 Description of the CNG encoder part**

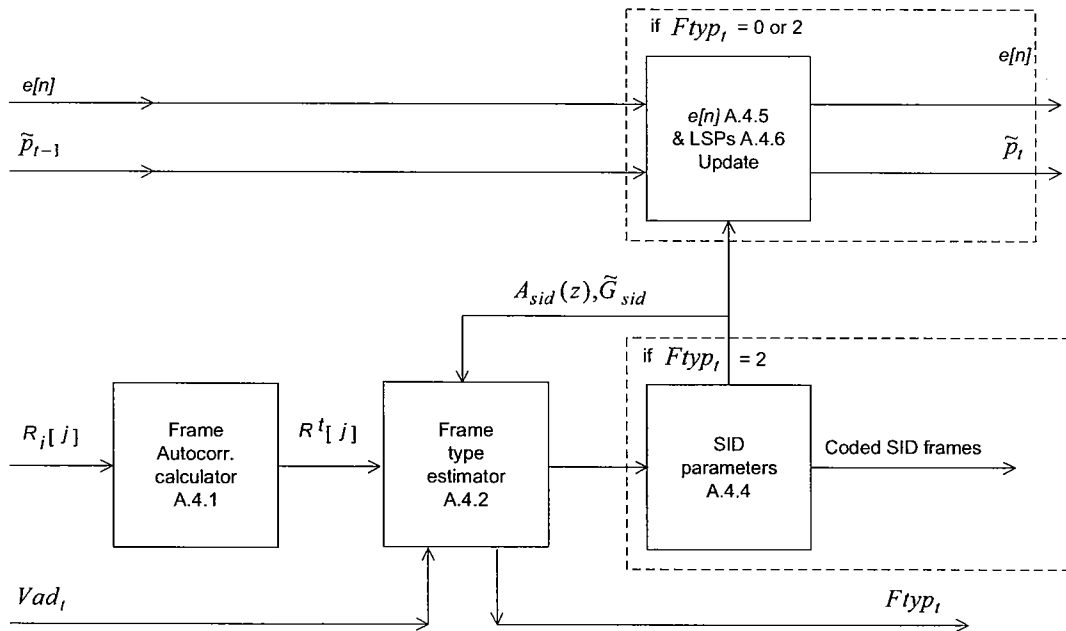
For each frame of 240 samples (active or inactive), the COD-CNG block processes the data coming from the VAD and the coder, produces the  $Ftyp_t$  information and the coded SID frame according to the procedure depicted by Figure A.3 and detailed in A.4.1 to A.4.7.

**A.4.1 Computation of the frame autocorrelation function**

File: COD_CNG.C	Procedure: Update_Acf()	Update autocorrelation function
-----------------	-------------------------	---------------------------------

For every frame  $t$  (active or inactive), the autocorrelation coefficients (calculated in the encoder as described in 2.4 of Recommendation G.723.1)  $R_i[j]$ ,  $j = 0$  to 10 of the four subframes indexed by  $i = 0$  to 3 are summed. The cumulated autocorrelation function of the current frame  $t$ , is given by:

$$R^t[j] = \sum_{i=0}^3 R_i[j], \text{ for } j = 0 \text{ to } 10 \tag{A-8}$$



T1524250-96

FIGURE A.3/G.723.1

Block diagram of the CNG at the encoder part

#### A.4.2 Computation of the current frame type $Ftyp_t$

File: COD_CNG.C	Procedure: Cod_Cng()	COD_CNG main body
File: COD_CNG.C	Procedure: LpcDiff()	Itakura distance comparison
File: LPC.C	Procedure: Durbin()	Levinson-Durbin recursion

If the current frame  $t$  is an active speech frame ( $Vad_t=1$ ), then  $Ftyp_t = 1$  and no other processing is performed.

In the other case, the decision SID/untransmitted frame is taken according to the following procedure:

The LPC filter  $A_t(z)$  of the current frame  $t$  is calculated by the Durbin procedure (see 2.4 of Recommendation G.723.1) using  $R'[j]$  as input. The coefficients of  $A_t(z)$  are noted  $a_t[j]$ ,  $j=1$  to 10. The Durbin procedure also provides the residual energy  $E_t$ , that will be used as an estimate of the frame excitation energy.

Then the current frame type  $Ftyp_t$  is determined in the following way:

- If the current frame is the first inactive frame of the inactive zone, the frame is selected as SID frame, the variable  $\bar{E}$  which reflects the energy sum is taken equal to  $E_t$ , and the number of frames involved in the summation,  $k_E$ , is initialized to 1:

$$(Vad_{t-1} = 1) \Rightarrow \begin{cases} Ftyp_t = 2 \\ \bar{E} = E \\ k_E = 1 \end{cases} \quad (\text{A-9})$$

- Else, if the current filter is significantly different from the preceding SID filter, or if the current excitation energy significantly differs from the preceding SID energy, then the frame is selected as SID ( $Ftyp_t = 2$ ).
- Otherwise, if the current frame is not the first of an inactive period, and if the current LPC filter and the excitation energy are similar to the SID ones, the frame is not transmitted ( $Ftyp_t = 0$ ).

The LPC filters and energies are compared according to the following methods:

##### Comparison of the LPC filters

The current LPC filter and SID filter are considered as significantly different if the Itakura distance between the two filters exceeds the given threshold, which is expressed by:

$$\sum_{j=0}^{10} R_a[j] \times R'[j] \geq E_t \times thr1 \quad (\text{A-10})$$

where  $R_a[j]$ ,  $j = 0$  to 10 is a function derived from the autocorrelation of the coefficients of the SID filter, given by :

$$\begin{cases} R_a[j] = 2 \sum_{k=0}^{10-j} a_{sid}[k] \times a_{sid}[k+j] & \text{if } j \neq 0 \\ R_a(0) = \sum_{k=0}^{10} a_{sid}[k]^2 \end{cases} \quad (\text{A-11})$$

with  $a_{sid}[0] = 1$

A value of 1.2136 is used for  $thr1$ .

### Comparison of the energies

$k_E$  being first incremented up to the maximum value 3, the sum the frame energies  $\bar{E} = \sum_{i=t-k_E+1}^t E_i$  is calculated.

Then  $\bar{E}$  is quantized, using the 6-bit pseudo-logarithmic quantizer described in A.4.3. The coded gain index  $GInd_t$  is compared to the previous coded SID gain index  $GInd_{sid}$ . If the difference exceeds the threshold  $thr2=3$ , the two energies will be considered as significantly different.

### A.4.3 Quantization of the average energy

File: UTIL_CNG.C	Procedure: Qua_SidGain()	Quantize Sid Gain
File: UTIL_CNG.C	Procedure: Dec_SidGain()	Decode Sid Gain

The quantization procedure operates on the sum of the energies  $\bar{E}$ , and the decoding provides a gain, which corresponds to the decoded value of the average energy square root.

A scaling factor  $\alpha_w = 2.70375$  is introduced to take into account the effect of windowing and bandwidth expansions present in the subframes autocorrelation functions  $R_i[j]$

The value used at the input of the gain quantizer is:

$$G = \alpha_w \times \sqrt{\frac{1}{k_E \times 240} \bar{E}}, \text{ bounded in } [0, 352].$$

The quantizer is a pseudo-log one, that divides  $[0, 352]$  into three segments indexed  $isg = 0$  to 2 of length  $N[isg] = 16, 16, 32$  with the associated resolutions 2, 4 and 8.

Let  $G_{isg}[j], j = 0$  to  $N[isg]-1$  be the decoded values for segment  $isg$ . Those values are given by:

$$G_{isg}[j] = G_{isg}[0] + j \times 2^{(isg+1)} \quad (\text{A-12})$$

The procedure uses  $G^2$  to calculate the index  $isg$  of the segment which contains  $G$  and the index  $i_s$  of  $G_{isg}(i_s)$  the closer to  $G$ .

The current quantization index is given by:

$$GInd_t = 16 \times isg + i_s \quad (\text{A-13})$$

The decoding is performed using the following formula:

$$Q^{-1}(GInd_t) = G_{isg}[0] + (GInd_t - \lfloor isg / 16 \rfloor) \times 2^{isg+1} \quad (\text{A-14})$$

where  $\lfloor x \rfloor$  denotes the greatest integer  $\leq x$ .

#### A.4.4 Computation and coding of SID parameters

File: COD_CNG.C	Procedure: Cod_Cng()	COD_CNG main body
File: COD_CNG.C	Procedure: ComputePastAvFilter()	Computes past average filter
File: COD_CNG.C	Procedure: LpcDiff()	Itakura distance comparison
File: COD_CNG.C	Procedure: CalcRc()	Compute function RC from LPC
File: LPC.C	Procedure: Durbin()	Levinson-Durbin recursion
File: LSP.C	Procedure: AtoLsp()	Converts LPC coefficients into LSP
File: LSP.C	Procedure: Lsp_Qnt()	LSP quantization
File: LSP.C	Procedure: Lsp_Inq()	LSP inverse quantization

When the current frame is a SID frame, the SID parameters are calculated and quantized. Notice that those parameters will serve in making the SID decision for the next inactive frames up to the next SID frame.

#### Computation of SID LPC filter $[A_{sid}(z)]$ and update of VAD LPC filter $[A_{no}(z)]$

First, the past average LPC filter  $\bar{A}_p(z)$  built from the three frames preceding the current one is estimated, using the Durbin procedure with the following autocorrelation function as input:

$$\bar{R}_p[j] = \sum_{k=t-3}^{t-1} R^k[j], \text{ for } j = 0 \text{ to } 10 \quad (\text{A-15})$$

the autocorrelation functions  $R^k[j]$  being the cumulated ones calculated by (A-8).

The past average LPC filter coefficients are denoted  $\bar{a}_p[j], j = 1 \text{ to } 10$ .

The VAD noise LPC filter used in A-2 is then updated with  $\bar{a}_p[j]$  but only when the adaptation enable flag  $Aen_t$  allows it:

$$\text{if } Aen_t = 0 \text{ then } a_{no}[j] = \bar{a}_p[j], j = 1, 2, \dots, 10 \quad (\text{A-16})$$

Then  $A_{sid}(z) = \begin{cases} A_t(z) & \text{if the distance between } A_t(z) \text{ and } \bar{A}_p(z) \text{ is } \geq thr1 \\ \bar{A}_p(z) & \text{otherwise} \end{cases}$  see eq. (A-10)

The distance between the current LPC filter  $A_t(z)$  and the average past LPC filter  $\bar{A}_p(z)$  is computed in the same manner as in A.4.2.

The coefficients  $a_{sid}[j], j = 1 \rightarrow 10$  of the new SID LPC filter are LSP converted and the LSPs are quantized using the encoder LSP 24-bit quantization procedure (see 2.5 of Recommendation G.723.1). The decoded value will be called  $\tilde{p}_{sid}$ .

#### SID gain

The quantized value of the SID gain is given by:

$$GInd_{sid} = GInd_t \quad (\text{A-17})$$

and the decoded value is denoted  $\tilde{G}_{sid}$ .

#### A.4.5 Computation of the CNG excitation

File: UTIL_CNG.C	Procedure: Calc_Exc_Rand()	Computation of the excitation
File: UTIL_CNG.C	Procedure: random_number()	Random number generation
File: UTIL_CNG.C	Procedure: distG()	Used to select excitation Gain
File: UTIL_LBC.C	Procedure: Sqrt_lbc()	Square root
File: UTIL_LBC.C	Procedure: Rand_lbc()	Pseudo-random sequence

The update of the excitation signal is performed both for SID frames and for untransmitted frames.

First, let us define the target excitation gain  $\tilde{G}_t$  as the square root of the average energy that must be obtained for the current frame  $t$  synthetic excitation.  $\tilde{G}_t$  is calculated using the following smoothing procedure:

$$\tilde{G}_t = \begin{cases} \tilde{G}_{sid} & \text{if } Vad_{t-1} = 1 \\ \frac{7}{8}\tilde{G}_{t-1} + \frac{1}{8}\tilde{G}_{sid} & \text{otherwise} \end{cases} \quad (\text{A-18})$$

The 240 samples of the frame are divided into two blocks of 120 samples, each block comprising two subframes of 60 samples.

For each block, the CNG excitation samples are synthesized using the following algorithm:

First the LTP parameters of the two subframes are selected:

- The pitch lag for the first subframe is randomly chosen in the interval [123, 143].
- The two subframes gain vector indices are randomly chosen into [0, 49], which corresponds to the first 50 vectors of the 170 entries gain codebook.
- The second subframe lag offset is taken equal to 0 for the first block, and 3 for the second block.

Next, the fixed codebook vectors of the two subframes are built by random selection of the grid, the pulses signs and positions, corresponding to the higher rate fixed excitation pattern.

Then a unique fixed excitation gain is computed for the two subframes of the block.

The adaptive excitation vector on the current block is noted  $u[n], n = 0$  to 119 and the fixed excitation  $v[n], n = 0$  to 119.

The fixed excitation gain is obtained by calculating the value  $Gf$  that yields a block average energy the closest to the target energy  $\tilde{G}_t^2$ :

$$\text{select } Gf \text{ such that } \left| \frac{1}{20} \sum_{n=0}^{119} (u[n] + Gf \times v[n])^2 - \tilde{G}_t^2 \right| \text{ minimum} \quad (\text{A-19})$$

Notice that  $Gf$  can take a negative value.

Let us define  $C(X) = aX^2 + 2bX + c$  such that:

$$a = \left( \sum_{n=0}^{119} v[n]^2 \right), b = \left( \sum_{n=0}^{119} u[n]v[n] \right), c = \left( \sum_{n=0}^{119} u[n]^2 - 120\tilde{G}_t^2 \right)$$

The equation  $C(X) = 0$  is then studied:

If the discriminant is  $\leq 0$  then  $Gf = -\frac{b}{a}$  is selected, else the two roots are calculated and the one with the lowest absolute value is selected.

Then  $Gf$  is bounded:  $Gf \leq 5000$

Finally the block CNG excitation is built, using:

$$e[n] = u[n] + Gf \times v[n], n = 0 \text{ to } 119 \quad (\text{A-20})$$

#### A.4.6 Interpolation of LSPs and update

File: COD_CNG.C	Procedure: Cod_Cng()	COD_CNG main body
File: LSP.C	Procedure: Lsp_Int()	LSP Interpolator

Both for SID frames and for untransmitted frames, the interpolated sets of LPC coefficients are calculated using  $\tilde{p}_{sid}$  and the previous LSP vector  $\tilde{p}_{t-1}$  provided to COD-CNG.

The LSP update is also performed:  $\tilde{p}_t = \tilde{p}_{sid}$ .

#### A.4.7 COD-CNG initialization

The following initialization must be performed on the frame autocorrelation functions, the target excitation gain, VAD information, and seed of the random generator used to compute the CNG excitation:

$$\left\{ \begin{array}{l} R^k[j] = 0 \text{ for } j = 0, \dots, 10 \text{ and } k = -1, -2, -3 \\ \tilde{G}_{-1} = 0 \\ Vad_{-1} = 1 \\ rseed = 12345 \end{array} \right.$$

No initialization is needed for the other static variables of COD-CNG.

#### A.5 Description of the decoder part

At the receiving end, DEC-CNG processes SID frames and untransmitted frames to produce the synthesized comfort noise.

The procedures developed to deal with frame erasures are described next.



A.5.1 Description of DEC-CNG

File: DEC CNG.C	Procedure: Dec Cng()	DEC CNG main body
File: UTIL CNG.C	Procedure: Calc Exc Rand()	Computation of the excitation
File: UTIL CNG.C	Procedure: random number()	Random number generation
File: UTIL CNG.C	Procedure: distG()	Used to select excitation Gain
File: UTIL LBC.C	Procedure: Sqrt lbc()	Square root
File: UTIL LBC.C	Procedure: Rand lbc()	Pseudo-random sequence
File: LSP.C	Procedure: Lsp-Inq()	LSP inverse quantization
File: LSP.C	Procedure: Lsp Int()	LSP Interpolator
File: UTIL CNG.C	Procedure: Qua SidGain()	Quantize Sid Gain
File: UTIL CNG.C	Procedure: Dec SidGain()	Decode Sid Gain

Figure A4 provides a general description of the comfort noise generation at the decoder part.

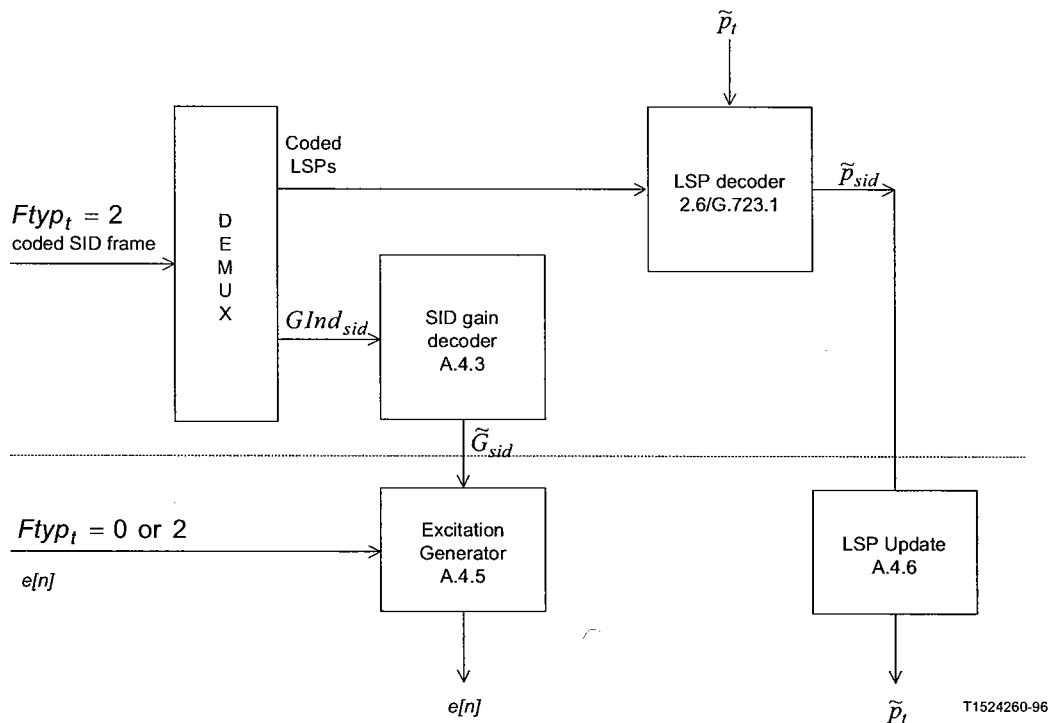


FIGURE A.4/G.723.1

Block diagram of the CNG at the decoder part

When the decoder receives an SID frame, DEC-CNG decodes the SID parameters.

Both in case of SID and untransmitted frames, the module DEC-CNG uses the decoded SID parameters to compute the LSPs and the excitation of the comfort noise that will be synthesized by the decoder synthesis module.

The CNG type of frame information  $Ftyp_t$ , (for frame  $t$ ) provided at the receiver is the same as the value computed by COD-CNG at the encoder.

- When  $Ftyp_t = 2$ , the parameters of the SID frame are decoded:  $\tilde{p}_{sid}$  for the LSPs and  $\tilde{G}_{sid}$  for the decoded gain.
- When  $Ftyp_t = 0$ ,  $Ftyp_{t-1}$  is tested to verify that the SID information has not been erased (see A.5.2). If  $Ftyp_{t-1} = 1$ , an energy term  $Enr$ , that has been calculated by the G.723.1 decoder during the processing of the last valid frame is quantized and decoded using the same procedure as the average energy in 3.1 of Recommendation G.723.1 except that there is no scaling factor  $\alpha_w$ . The decoded value is the restored  $\tilde{G}_{sid}$ .

Then, in both cases, the CNG excitation is calculated according to the procedure described in COD-CNG in A.4.5. The new LSP vector,  $\tilde{p}_{sid}$ , is used to compute the interpolated LPC coefficients, and the LSP updating is performed:  $\tilde{p}_t = \tilde{p}_{sid}$

### A.5.2 Frame erasure concealment with regards to the CNG

File: DECOD.C	Procedure: Decod()	Frame decoding
---------------	--------------------	----------------

When a frame erasure is detected by the decoder, the erased frame type depends on the preceding frame type:

- if the preceding frame was active, then the current erased frame is considered as active,
- else if the preceding frame was either an SID frame or an untransmitted frame, the current erased frame is considered as untransmitted:

$$\begin{cases} Ftyp_{t-1} = 1 & \Rightarrow & Ftyp_t = 1 \\ Ftyp_{t-1} = 0 \text{ or } 2 & \Rightarrow & Ftyp_t = 0 \end{cases} \quad (\text{A-21})$$

If an untransmitted frame has been erased, no error is then introduced.

If a SID frame is erased, there are two possibilities:

- If it is not the first SID frame of the current inactive period, then the previous SID parameters are kept.
- If it is the first SID frame of an inactive period, a special protection has been taken:

As stated in A5.1, this case is detected by the fact that  $Ftyp_{t-1} = 1$  and  $Ftyp_t = 0$ .

This combination of events does not imply that the preceding frame was a good active frame: several frames up to the preceding one may have been erased. What is certain is that the last good frame was an active frame, that the present frame was not erased, and that the SID frame supposed to provide information for the current untransmitted frame is lost.

To recover the SID information, DEC-CNG uses parameters provided by the G.723.1 decoder main part:

- The LSPs of the last valid active frame are used for  $\tilde{p}_{sid}$ .
- The energy term  $Enr$  calculated by the decoder during the the residual interpolation procedure (see 3.10.2 of Recommendation G.723.1) over the 120 last excitation samples of

the last valid active frame is used to recover  $\tilde{G}_{sid}$ , according to the method described in A.5.1.

Finally, to avoid de-synchronization of the random generator used to compute the excitation, the pseudo-random sequence reset is performed at each active frame, both at the encoder and coder part:  $rseed = 12345$ .

### A.5.3 DEC-CNG initialization

Only the following variables must be initialized:

$$\begin{cases} \tilde{G}_{sid} = 0 \\ \tilde{P}_{sid} = LSP \ DC \ vector \ P_{DC} \\ Vad_{-1} = 1 \\ rseed = 12345 \end{cases}$$

### A.6 Bit stream packing

Table A.1 shows the bit stream of the SID frames according to the notations used in clause 4 of Recommendation G.723.1.

TABLE A.1/G.723.1  
Bit packing for SID frames

Transmitted octets	PARx_By, ...
1	LPC_B5 ... LPC_B0, VADFLAG_B0, RATEFLAG_B0
2	LPC_B13 ... LPC_B6
3	LPC_B21 ... LPC_B14
4	GAIN_B5 ... GAIN_B0, LPC_B23, LPC_B22

### A.7 Glossary

$a_{no}[j]$	noise LPC filter coefficients
$L_{OL}^j$	preceding frame and current frame open loop pitch delays
$pc$	pitch delays counter for voicing estimation
$Aen_t$	adaptation enable flag
$e_t'[n]$	noise-inverse filtered input signal for frame $t$
$Enr_t$	noise-inverse filtered input signal energy for frame $t$
$Nlev_t$	noise level at frame $t$
$Nlev_{min}$	minimum bound on $Nlev_t$
$Nlev_{max}$	maximum bound on $Nlev_t$
$Thr$	adapted threshold for VAD decision
$k_i'[2]$	second reflection coefficient for subframe $i$ in frame $t$

$SinD$	sine wave detection flag (1: sine detected, 0: else)
$e[n]$	decoded combined excitation vector
$R_t[j]$	autocorrelation function for subframe $i$ , $j = 0, 1, \dots, 10$
$thr_2$	threshold for energies distance
$R_{a[j]}$	modified autocorrelation of LPC coefficients
$Gind_{sid}$	SID gain index
$\tilde{G}_{sid}$	decoded SID gain
$G$	excitation gain used at the SID quantizer input
$Gind_t$	gain index for frame $t$
$isg$	SID gain quantizer segment index
$N[isg]$	SID gain quantizer segment length
$G_{isg}[j]$	gain decoded values of segment $isg$ , $j=0, 1, \dots, N[isg]-1$
$i_s$	gain index relative to the segment
$\alpha_w$	energy scaling factor
$a_{sid}$	SID LPC coefficient vector
$\bar{a}_p$	past average LPC filter coefficients
$\bar{R}_p[j]$	sum of past autocorrelation functions
$\tilde{p}_{sid}$	decoded SID LSP vector
$u[n]$	adaptive codebook excitation vector
$v[n]$	fixed codebook excitation vector
$\tilde{G}_t$	target excitation gain for excitation synthesis
$a, b, c$	coefficients of energy minimization equation
$C(X)$	energy minimization equation
$Gf$	fixed codebook gain for CNG excitation synthesis
$rseed$	random generator seed

### A.8 Bit-exact, fixed-point C source code

All details of the silence compression algorithm are included as part of bit-exact, fixed-point ANSI C source code. In the event of any discrepancy between the above descriptions and the C source, the C source code is presumed to be correct. This C source code is a part of the code distributed by the ITU-T as Recommendation G.723.1.



**ITU-T RECOMMENDATIONS SERIES**

- Series A Organization of the work of the ITU-T
- Series B Means of expression
- Series C General telecommunication statistics
- Series D General tariff principles
- Series E Telephone network and ISDN
- Series F Non-telephone telecommunication services
- Series G Transmission systems and media**
- Series H Transmission of non-telephone signals
- Series I Integrated services digital network
- Series J Transmission of sound-programme and television signals
- Series K Protection against interference
- Series L Construction, installation and protection of cables and other elements of outside plant
- Series M Maintenance: international transmission systems, telephone circuits, telegraphy, facsimile and leased circuits
- Series N Maintenance: international sound-programme and television transmission circuits
- Series O Specifications of measuring equipment
- Series P Telephone transmission quality
- Series Q Switching and signalling
- Series R Telegraph transmission
- Series S Telegraph services terminal equipment
- Series T Terminal equipment and protocols for telematic services
- Series U Telegraph switching
- Series V Data communication over the telephone network
- Series X Data networks and open system communication
- Series Z Programming languages



INTERNATIONAL TELECOMMUNICATION UNION

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**G.723.1**

**Annex B**

(11/96)

SERIES G: TRANSMISSION SYSTEMS AND MEDIA

Digital transmission systems – Terminal equipments –  
Coding of analogue signals by methods other than PCM

---

Dual rate speech coder for multimedia  
communications transmitting at 5.3 and 6.3 kbit/s

**Annex B: Alternative specification based on  
floating point arithmetic**

ITU-T Recommendation G.723.1 – Annex B

(Previously CCITT Recommendation)

---

ITU-T G-SERIES RECOMMENDATIONS  
TRANSMISSION SYSTEMS AND MEDIA

INTERNATIONAL TELEPHONE CONNECTIONS AND CIRCUITS	G.100–G.199
<b>INTERNATIONAL ANALOGUE CARRIER SYSTEM</b>	
GENERAL CHARACTERISTICS COMMON TO ALL ANALOGUE CARRIER-TRANSMISSION SYSTEMS	G.200–G.299
INDIVIDUAL CHARACTERISTICS OF INTERNATIONAL CARRIER TELEPHONE SYSTEMS ON METALLIC LINES	G.300–G.399
GENERAL CHARACTERISTICS OF INTERNATIONAL CARRIER TELEPHONE SYSTEMS ON RADIO-RELAY OR SATELLITE LINKS AND INTERCONNECTION WITH METALLIC LINES	G.400–G.449
COORDINATION OF RADIOTELEPHONY AND LINE TELEPHONY	G.450–G.499
<b>TRANSMISSION MEDIA CHARACTERISTICS</b>	
<b>DIGITAL TRANSMISSION SYSTEMS</b>	
TERMINAL EQUIPMENTS	G.700–G.799
General	G.700–G.709
Coding of analogue signals by pulse code modulation	G.710–G.719
<b>Coding of analogue signals by methods other than PCM</b>	<b>G.720–G.729</b>
Principal characteristics of primary multiplex equipment	G.730–G.739
Principal characteristics of second order multiplex equipment	G.740–G.749
Principal characteristics of higher order multiplex equipment	G.750–G.759
Principal characteristics of transcoder and digital multiplication equipment	G.760–G.769
Operations, administration and maintenance features of transmission equipment	G.770–G.779
Principal characteristics of multiplexing equipment for the synchronous digital hierarchy	G.780–G.789
Other terminal equipment	G.790–G.799
DIGITAL NETWORKS	G.800–G.899
General aspects	G.800–G.809
Design objectives for digital networks	G.810–G.819
Quality and availability targets	G.820–G.829
Network capabilities and functions	G.830–G.839
SDH network characteristics	G.840–G.899
DIGITAL SECTIONS AND DIGITAL LINE SYSTEM	G.900–G.999
General	G.900–G.909
Parameters for optical fibre cable systems	G.910–G.919
Digital sections at hierarchical bit rates based on a bit rate of 2048 kbit/s	G.920–G.929
Digital line transmission systems on cable at non-hierarchical bit rates	G.930–G.939
Digital line systems provided by FDM transmission bearers	G.940–G.949
Digital line systems	G.950–G.959
Digital section and digital transmission systems for customer access to ISDN	G.960–G.969
Optical fibre submarine cable systems	G.970–G.979
Optical line systems for local and access networks	G.980–G.999

*For further details, please refer to ITU-T List of Recommendations.*



**ITU-T RECOMMENDATION G.723.1 – Annex B**

**ALTERNATIVE SPECIFICATION BASED ON FLOATING POINT ARITHMETIC**

**Source**

Annex B to ITU-T Recommendation G.723.1, was prepared by ITU-T Study Group 15 (1993-1996) and was approved under the WTSC Resolution No. 1 procedure on the 8th of November 1996.

### **FOREWORD**

ITU (International Telecommunication Union) is the United Nations Specialized Agency in the field of telecommunications. The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of the ITU. The ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Conference (WTSC), which meets every four years, establishes the topics for study by the ITU-T Study Groups which, in their turn, produce Recommendations on these topics.

The approval of Recommendations by the Members of the ITU-T is covered by the procedure laid down in WTSC Resolution No. 1 (Helsinki, March 1-12, 1993).

In some areas of information technology which fall within ITU-T's purview, the necessary standards are prepared on a collaborative basis with ISO and IEC.

### **NOTE**

In this Recommendation, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

© ITU 1997

All rights reserved. No part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from the ITU.

**CONTENTS**

	<b>Page</b>
B.1 Introduction .....	1
B.2 Algorithm description.....	1
B.3 ANSI C code.....	1
B.3 ANSI C code           1	



**Recommendation G.723.1 – Annex B****ALTERNATIVE SPECIFICATION BASED ON FLOATING POINT ARITHMETIC***(Geneva, 1996)***B.1 Introduction**

Recommendation G.723.1 provides a bit-exact, fixed-point specification of a dual rate 6.3 and 5.3 kbit/s speech coder for multimedia telecommunications applications. Exact details of this specification are given in bit-exact, fixed-point C code available from the ITU-T. This Annex describes an alternative implementation of G.723.1 contained in floating point C source code. A set of digital test vectors for this floating point specification is also available from the ITU-T in order to facilitate the implementation of Recommendation G.723.1. Note that passing the test vectors is a necessary but not a sufficient condition to comply with Recommendation G.723.1.

**B.2 Algorithm description**

The floating point version of G.723.1 has the same algorithmic steps as the fixed-point version. Similarly, the bit stream is identical to that of G.723.1. The reader is referred to clauses 2, 3, 4 and 6 of the main body of Recommendation G.723.1 for details.

**B.3 ANSI C code**

ANSI C code simulation the dual/rate encoder/decoder in floating point arithmetic is available from the ITU-T. Table B.1 lists all the files included in this code. Individual C function names are the same as in the main body text of Recommendation G.723.1.

TABLE B.1/G.723.1

**List of software filenames**

File name	Description
TYPEDEF2.H	data type definition is machine dependent
CST2.H	definition of constants for Recommendation G.723.1
LBCCODE2.C	main program for G.723.1 speech codecs
LBCCODE2.H	function prototypes
CODER2.C	G.723.1 speech encoder for the two bit rates
CODER2.H	function prototypes
DECOD2.C	G.723.1 speech decoder for the two bit rates
DECOD2.H	function prototypes
LPC2.C	linear predictive analysis
LPC2.H	function prototypes
LSP2.C	line spectral pair related functions, quantizer
LSP2.H	function prototypes
EXC2.C	adaptive and fixed (MP-MLQ, ACELP) excitation
EXC2.H	function prototypes
UTIL2.C	miscellaneous functions (HPF, pack, unpack, I/O ...)
UTIL2.H	function prototypes
TAB2.C	tables of constants
TAB2.H	external declaration for constant tables



INTERNATIONAL TELECOMMUNICATION UNION

# ITU-T G.723.1

## Implementers Guide

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

(25 October 2002)

SERIES G: TRANSMISSION SYSTEMS AND MEDIA,  
DIGITAL SYSTEMS AND NETWORKS

Digital terminal equipments – Coding of analogue signals  
by methods other than PCM

---

**Implementors' Guide for G.723.1**

***(Dual rate speech coder for multimedia  
communications transmitting at 5.3 & 6.3 kbit/s)***

---

Implementers Guide for Recommendation G.723.1

---

**Contact Information**

Rapporteur, ITU-T Study  
Group 16 / Question 10

Claude Lamblin  
France Telecom R&D /DIH  
Technopole Anticipa  
2 avenue Pierre Marzin  
22307 LANNION Cedex  
France

Tel: +33 2 96 05 13 03  
Fax: +33 2 96 05 35 30  
Email: [claude.lamblin@rd.francetelecom.com](mailto:claude.lamblin@rd.francetelecom.com)

**SUMMARY**

**Implementors' Guide for Recommendation G.723.1**

This document contains Implementers' Guide for the text of ITU-T Recommendation G.723.1 Annex A and for the software C-codes of G.723.1 Annexes A and B that corrects defects reported at SG 16's meeting on 15-25 October 2002.



**Table of Contents**

Implementors' Guide for Recommendation G.723.1 ..... ii

Implementers' Guide for G.723.1 Annexes A &B..... 1

1      Corrections to G.723.1 annexes A and B source codes..... 1

2      Editorial corrections to section 2.9/G.723.1 Annex A recommendation text ..... 2

## Implementers' Guide for G.723.1 Annexes A & B

### 1 Corrections to G.723.1 annexes A and B source codes

In March 1996, SG 15 approved Rec. G.723.1 Annex A. As described in COM16-D261, a problem in G.723.1 Annex A has been discovered. When the encoder input PCM file starts with absolute silence (i.e. all zeros input), the encoder never generates any silent frames, only speech frames and SID frames. For instance, the encoding of 30 seconds of absolute silence outputs 3 speech frames followed by 997 SID frames while only 1 SID frame followed by silent frames would have been expected.

This strange behaviour of G723.1 annex A (and also annex B) with silent input comes from the filters comparison (Itakura-Saito distance computation) in the section A4.2 (Comparison of the LPC filters) equation A-10.

The current LPC filter and SID filter are considered as significantly different if the Itakura distance between the two filters exceeds the given threshold, which is expressed by:

$$\sum_{j=0}^{10} R_a[j] \times R'[j] \geq E_t \times thr1 \tag{A-10}$$

where  $R_a[j]$ ,  $j = 0$  to  $10$  is a function derived from the autocorrelation of the coefficients of the SID filter.

With null input,  $E_t$ , all  $R_a[i]$  and  $R_s[i]$  (except  $R_s[0]$ ) are all equal to zero.

In the ANSI C-source code, however the following test is performed:  $\sum_{j=0}^{10} R_a[j] \times R'[j] < E_t \times thr1$

If the test is false, the filters are judged different otherwise they are not. However with zeros the test  $<$  is not true, so the filters are judged different and SID frames are sent. To fix this, it is sufficient to replace  $<$  by  $\leq$

In the fixed-point C source code (G.723.1 Annex A), it is proposed to modify the line 463 in file COD\_CNG.C, routine LpcDiff:

Replace line 463: `if (L_temp0 < L_temp1) diff=1` by `if (L_temp0 ≤ L_temp1) diff=1`

In the floating-point C source code (G.723.1 Annex B), it is proposed to modify the line 397 in file CODCNG2.C, routine LpcDiff

Replace line 397: `if (temp0 < temp1)` by `if (temp0 ≤ temp1)`

Table 1 summarizes these modifications.

**Table 1:**  
Modified lines in the comparison of the LPC filters routines

Annex	File name	Routine name	Modified line	Correction
A	COD_CNG.C	LpcDiff	463	if(L_temp0 ≤ L_temp1) diff = 1
B	CODCNG2.C	LpcDiff	397	if (temp0 ≤ temp1)

**2 Editorial corrections to section 2.9/G.723.1 Annex A recommendation text**

COM16-D261 also reported editorial errors in the ITU-T G.723.1 Annex A recommendation text in equation (A.10) and (A.11) in the sub-section "Comparison of the LPC filters" of the section A.4.2 "Computation of the current frame type Ftyp<sub>t</sub>"

In equation (A.10), the index i should be replaced by j. Furthermore, it was pointed out that the text should be aligned with the C-source code, modified as described in the previous section.

The equation (A.10) should be:  $\sum_{j=0}^{10} R_a[j] \times R'[j] > E_t \times thr1$  instead of:  $\sum_{j=0}^{10} R_a[i] \times R'[i] \geq E_t \times thr1$

In the first line of equation (A.11), the sum operand is wrongly set as an exponent. The equation (A.11) should be:

$$R_a[j] = 2 \sum_{k=0}^{10-j} a_{sid}[k] \times a_{sid}[k+j], \text{ if } j \neq 0 \text{ instead of } R_a[j] = 2 \sum_{k=0}^{10-j} a_{sid}[k] \times a_{sid}[k+j] \text{ if } j \neq 0$$

The exiting text should be corrected as shown below.

[Begin Correction]

**Comparison of the LPC filters**

The current LPC filter and SID filter are considered as significantly different if the Itakura distance between the two filters exceeds the given threshold, which is expressed by:

$$\sum_{j=0}^{10} R_a[j] \times R'[j] > E_t \times thr1 \tag{A-10}$$

where  $R_a[j], j = 0$  to 10 is a function derived from the autocorrelation of the coefficients of the SID filter, given by :

$$\begin{cases} R_a[j] = 2 \sum_{k=0}^{10-j} a_{sid}[k] \times a_{sid}[k+j] \text{ if } j \neq 0 \\ R_a(0) = \sum_{k=0}^{10} a_{sid}[k]^2 \end{cases} \tag{A-11}$$

[End Correction]