



ARIB STD-T64-C.S0030-0 v3.0

**Selectable Mode Vocoder (SMV)
Service Option for Wideband
Spread Spectrum
Communication Systems**

Refer to "Industrial Property Rights (IPR)" in the preface of ARIB STD-T64 for Related Industrial Property Rights. Refer to "Notice" in the preface of ARIB STD-T64 for Copyrights

1 **Original Specification**

2 This standard, ARIB-T64-C.S0030-0 v3.0, was prepared by 3GPP2-WG of Association of Radio
3 Industries and Businesses (ARIB) based upon the 3GPP2 specification, C.S0030-0 v3.0.

4

5 **Modification to the original specification**

6 None.

7

8 **Notes**

9 None.

10



1

2 **Selectable Mode Vocoder (SMV) Service Option for**
3 **Wideband Spread Spectrum Communication**
4 **Systems**

5

6

COPYRIGHT

3GPP2 and its Organizational Partners claim copyright in this document and individual Organizational Partners may copyright and issue documents or standards publications in individual Organizational Partner's name based on this document. Requests for reproduction of this document should be directed to the 3GPP2 Secretariat at secretariat@3gpp2.org. Requests to reproduce individual Organizational Partner's documents should be directed to that Organizational Partner. See www.3gpp2.org for more information.

7

Intentionally left blank.

1
2
3
4
5
6
7

FOREWORD

These technical requirements form a standard for Service Option 56, a variable-rate two-way speech service option. The maximum speech-coding rate of the service option is 8.55 kbps.

This standard does not address the quality or reliability of Service Option 56, nor does it cover equipment performance or measurement procedures.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

NOTES

1. Accompanying “Recommended Minimum Performance Standard for the Selectable Mode Vocoder, Service Option 56,” provides specifications and measurement methods.
2. “Base station” refers to the functions performed on the land-line side, which are typically distributed among a cell, a sector of a cell, a mobile switching center, and a personal communications switching center.
3. This document uses the following verbal forms: “Shall” and “shall not” identify requirements to be followed strictly to conform to the standard and from which no deviation is permitted. “Should” and “should not” indicate that one of several possibilities is recommended as particularly suitable, without mentioning or excluding others; that a certain course of action is preferred but not necessarily required; or that (in the negative form) a certain possibility or course of action is discouraged but not prohibited. “May” and “need not” indicate a course of action permissible within the limits of the standard. “Can” and “cannot” are used for statements of possibility and capability, whether material, physical, or causal.
4. Footnotes appear at various points in this specification to elaborate and further clarify items discussed in the body of the specification.
5. Unless indicated otherwise, this document presents numbers in decimal form.

Binary numbers are distinguished in the text by the use of single quotation marks. In some tables, binary values may appear without single quotation marks if table notation clearly specifies that values are binary. The character ‘x’ is used to represent a binary bit of unspecified value. For example ‘xxx00010’ represents any 8-bit binary value such that the least significant five bits equal ‘00010’.

Hexadecimal numbers (base 16) are distinguished in the text by use of the form 0xh...h where h...h represents a string of hexadecimal digits. For example, 0x2fa1 represents a number whose binary value is ‘10111110100001’ and whose decimal value is 913.

1

NOTES

2

6. The following conventions apply to mathematical expressions in this standard:

3

- $\lfloor x \rfloor$ indicates the largest integer less than or equal to x : $\lfloor 1.1 \rfloor = 1$, $\lfloor 1.0 \rfloor = 1$.

4

- $\lceil x \rceil$ indicates the smallest integer greater than or equal to x : $\lceil 1.1 \rceil = 2$, $\lceil 2.0 \rceil = 2$.

5

- $|x|$ indicates the absolute value of x : $|-17| = 17$, $|17| = 17$.

6

- \oplus indicates exclusive OR.

7

- $\min(x, y)$ indicates the minimum of x and y .

8

- $\max(x, y)$ indicates the maximum of x and y .

9

- In figures, \otimes indicates multiplication. In formulas within the text, multiplication is implicit. For example, if $h(n)$ and $p_L(n)$ are functions, then $h(n) p_L(n) = h(n) \otimes p_L(n)$.

10

11

- $x \bmod y$ indicates the remainder after dividing x by y : $x \bmod y = x - (y \lfloor x/y \rfloor)$.

12

- $\text{round}(x)$ is traditional rounding: $\text{round}(x) = \lfloor x + 0.5 \rfloor$.

13

$$\text{sign}(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases} .$$

14

- \sum indicates summation. If the summation symbol specifies initial and terminal values, and the initial value is greater than the terminal value, then the value of the summation is 0. For example, if $N=0$, and if $f(n)$ represents an arbitrary function, then

15

16

17

18

$$\sum_{n=1}^N f(n) = 0.$$

19

- The bracket operator, $[]$, isolates individual bits of a binary value. $\text{VAR}[n]$ refers to bit n of the binary representation of the value of the variable VAR , such that $\text{VAR}[0]$ is the least significant bit of VAR . The value of $\text{VAR}[n]$ is either 0 or 1.

20

21

22

REFERENCES

The following standards contain provisions, through reference in this text constitute provisions of this Standard. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this Standard are encouraged to investigate the possibility of applying the most recent editions of the standards indicated below. ANSI and TIA maintain registers of currently valid national standards published by them.

—*American National Standards:*

1. ANSI/EIA/TIA-579-A-98, *Telecommunications Telephone Terminal Equipment Transmission Requirements for Digital Wireline Telephones*

—*Other Standards:*

2. ITU-T Recommendation G.711, *Pulse Code Modulation (PCM) of Voice Frequencies*, Vol. III, Geneva 1972.
3. ITU-T Recommendation G.712, *Transmission Performance Characteristics of Pulse Code Modulation*, November 1996.
4. IEEE Standard 269-2002, *IEEE Standard Methods for Measuring Transmission Performance of Analog and Digital Telephone Sets*, 2002.
5. IEEE Standard 661-1979, *Method for Determining Objective Loudness Ratings of Telephone Connections*, 1979.
6. 3GPP2 C.S0003-A, *Medium Access Control (MAC) Standard for cdma2000 Spread Spectrum Systems*, March, 2000.
7. 3GPP2 C.S0005-A, *Upper Layer (Layer 3) Signalling Standard for cdma2000 Spread Spectrum Systems*, March, 2000.
8. 3GPP2 C.S0034-0 Ver 1.0, *Recommended Minimum Performance Standard for the Selectable Mode Vocoder, Service Option 56*.
9. 3GPP2 C.S0014-0 Ver 1.0, *Enhanced Variable Rate Codec, Speech Service Option 3 for Wideband Spread Spectrum Digital Systems*, September, 1996.

1	Table of Contents		
2	1	Service Option 56: Selectable Mode Vocoder Two-Way Voice Communication.....	1
3	1.1	General Description.....	1
4	1.2	Overview of SMV Documentation.....	1
5	1.3	SMV Algorithmic Options	1
6	1.3.1	Standard Options for Noise Suppression.....	2
7	1.3.2	Standard Options for Voice Activity Detection	2
8	1.4	SMV Implementation Options	2
9	1.5	Service Option Number.....	2
10	1.6	Allowable Delays	2
11	1.6.1	Allowable Transmitting Speech Codec Encoding Delay	2
12	1.6.2	Allowable Receiving Speech Codec Decoding Delay	2
13	1.7	Special Cases.....	3
14	1.7.1	Blanked Packets	3
15	1.7.2	Null Traffic Channel Data.....	3
16	1.8	Terms and Numeric Information.....	3
17	2	Required Multiplex Option Support.....	2
18	2.1	Interface to Multiplex Option 1	2
19	2.1.1	Transmitted Packets	2
20	2.1.2	Received Packets.....	2
21	2.2	Negotiation for Service Option 56	3
22	2.2.1	Procedures Using Service Negotiation.....	3
23	2.2.1.1	Initialization and Connection	3
24	2.2.1.1.1	Mobile Station Requirements	3
25	2.2.1.1.2	Base Station Requirements	4
26	2.2.1.2	Service Option Control Messages	4
27	2.2.1.2.1	Mobile Station Requirements	4
28	2.2.1.2.2	Base Station Requirements	5
29	3	Audio Interfaces	9
30	3.1	Input Audio Interface	9
31	3.1.1	Input Audio Interface in the Mobile Station	9
32	3.1.1.1	Conversion and Scaling	9
33	3.1.1.2	Digital Audio Input	9
34	3.1.1.3	Analog Audio Input.....	9
35	3.1.1.3.1	Transmit Level Adjustment	9
36	3.1.1.3.2	Band Pass Filtering.....	9
37	3.1.1.3.3	Echo Return Loss.....	9
38	3.1.2	Input Audio Interface in the Base Station	10
39	3.1.2.1	Sampling and Format Conversion.....	10
40	3.1.2.2	Transmit Level Adjust.....	10
41	3.1.2.3	Line Echo Canceling	10
42	3.2	Output Audio Interface.....	10
43	3.2.1	Output Audio Interface in the Mobile Station.....	10
44	3.2.2	Band Pass Filtering.....	10
45	3.2.2.1	Receive Level Adjustment.....	10
46	3.2.3	Output Audio Interface in the Base Station.....	11
47	3.2.3.1	Receive Level Adjustment.....	11
48	4	The Selectable Mode Vocoder (SMV) Speech Coding – Introduction and Tables .	13

1	4.1	Introduction to the SMV Speech Coding Algorithm.....	13
2	4.2	Bit Allocation Tables.....	15
3	4.3	SMV Bit Streams.....	18
4	4.4	SMV Symbol Table.....	38
5	5	SMV Encoder.....	43
6	5.1	SMV Encoder Look Ahead Buffers and Delay.....	43
7	5.2	SMV Encoder Pre-Processing.....	44
8	5.2.1	Silence Enhancement.....	45
9	5.2.1.1	Identifying Minimum and Maximum Values.....	45
10	5.2.1.2	Setting Histogram Thresholds.....	45
11	5.2.1.3	Calculating the 3-Value Histogram.....	46
12	5.2.2	High-Pass Filtering.....	48
13	5.2.3	Noise Suppression.....	49
14	5.2.3.1	Noise Suppression Initialization.....	49
15	5.2.3.2	Frequency Domain Conversion.....	51
16	5.2.3.2.1	Windowing and Buffering - Option A.....	51
17	5.2.3.2.2	Windowing and Buffering - Option B.....	51
18	5.2.3.3	Channel Energy Estimator.....	53
19	5.2.3.4	Channel SNR Estimator.....	53
20	5.2.3.5	Voice Metric Calculation.....	54
21	5.2.3.6	Spectral Deviation Estimator.....	54
22	5.2.3.7	Background Noise Update Decision.....	55
23	5.2.3.8	SNR Estimate Modification.....	56
24	5.2.3.9	Channel Gain Computation.....	57
25	5.2.3.10	Frequency Domain Filtering.....	58
26	5.2.3.11	Background Noise Estimate Update.....	58
27	5.2.3.12	Time Domain Signal Reconstruction.....	58
28	5.2.3.12.1	Overlap-and-Add – Option A.....	58
29	5.2.3.12.2	Overlap-and-Add – Option B.....	59
30	5.2.4	Adaptive Tilt Filtering.....	60
31	5.3	SMV Encoder Frame Level Processing.....	61
32	5.3.1	Linear Prediction (LPC) Analysis.....	63
33	5.3.1.1	The LPC Analysis Weighting Windows.....	63
34	5.3.1.2	The LP Analysis Procedure.....	63
35	5.3.1.2.1	Energy Calculation.....	64
36	5.3.1.2.2	Calculation of the Autocorrelation Function.....	64
37	5.3.1.2.3	Levinson-Durbin Algorithm.....	64
38	5.3.1.2.4	Conversion to LSFs.....	65
39	5.3.1.2.5	LPC Prediction Error and LPC Prediction Gain.....	66
40	5.3.1.3	Interpolation of LPC Prediction Gains, Reflection Coefficients, and Prediction	
41		Coefficients.....	67
42	5.3.2	Voice Activity Detection – Option A.....	68
43	5.3.2.1	Computation of VAD Parameters.....	68
44	5.3.2.2	VAD Decision.....	70
45	5.3.2.3	Hangover.....	70
46	5.3.2.4	Parameter Update During Silence/Noise.....	70
47	5.3.2.5	Reset of Noise Floor Energy.....	70
48	5.3.3	Voice Activity Detection – Option B.....	71
49	5.3.3.1	Estimating the Data Rate and VAD Based on Current Signal Parameters...	71
50	5.3.3.1.1	Computing Band Energy.....	71

1	5.3.3.1.2	Calculating Rate Determination Thresholds.....	72
2	5.3.3.1.3	Comparing Thresholds.....	74
3	5.3.3.1.4	Performing Hangover.....	74
4	5.3.3.1.5	Constraining VAD Rate Selection.....	76
5	5.3.3.2	Updating RDA Parameters.....	76
6	5.3.3.2.1	Updating the Smoothed Band Energy.....	76
7	5.3.3.2.2	Updating the Background Noise Estimate.....	77
8	5.3.3.2.3	Updating the Signal Energy Estimate.....	78
9	5.3.4	Music Detection.....	79
10	5.3.4.1	Input to Music Detector.....	79
11	5.3.4.2	Computation of Music Detector Parameters.....	79
12	5.3.4.3	Music Detection Decision.....	80
13	5.3.4.4	Reset of tracking counters.....	80
14	5.3.4.5	Update of pitch lag, pitch correlation and pitch gain buffers.....	80
15	5.3.5	Calculation of Weighting Coefficient for 4-Subframe Structure.....	82
16	5.3.6	Generating the Weighted Residual.....	83
17	5.3.7	Initial Setting of Voiced/Unvoiced Level.....	84
18	5.3.7.1	Calculation of Decision Parameters.....	84
19	5.3.7.2	Initial Setting Voiced/Unvoiced Level.....	85
20	5.3.8	Generating the Weighted Speech.....	86
21	5.3.9	Open-Loop Pitch Detection.....	87
22	5.3.9.1	Open-Loop Pitch Detection on the Look Ahead Buffer.....	87
23	5.3.9.1.1	Selection of 4 Open-Loop Pitch Candidates.....	87
24	5.3.9.1.2	Correction of Open-Loop Pitch Multiples for the Lookahead.....	87
25	5.3.9.2	Open-Loop Pitch Detection on the Second Half of the Encoding Frame.....	88
26	5.3.9.3	Final Open-Loop Pitch Correction.....	88
27	5.3.9.3.1	Smoothing of the Open-Loop Pitch Candidates.....	88
28	5.3.9.3.2	Further Correction of Pitch Sub-Multiples.....	89
29	5.3.10	Initial Frame Classification.....	90
30	5.3.10.1	Average and Variance of Open-Loop Pitch Values.....	90
31	5.3.10.2	Noise-Suppressed Classification Parameters.....	90
32	5.3.10.2.1	First Order Reflection Coefficients for 4 Subframes.....	90
33	5.3.10.2.2	Noise-Suppressed Spectral Tilt, Absolute Maximum, and Open-Loop Pitch Correlation.....	91
34	5.3.10.2.3	Spectral Tilt Slope, Absolute Maximum Slope, and Slope Statistics.....	92
35	5.3.10.2.4	Statistics of Open-Loop Pitch Correlation.....	93
36	5.3.10.2.5	Classification Decision.....	93
37	5.3.11	Noise-to-Signal Ratio (NSR) Estimation.....	94
38	5.3.11.1	Signal Input to the NSR Estimation Function.....	94
39	5.3.11.2	NSR counters.....	94
40	5.3.11.3	NSR Calculation.....	94
41	5.3.12	Classification Refinement.....	96
42	5.3.12.1	Operation of the Classification Refinement.....	96
43	5.3.13	Rate Selection.....	98
44	5.3.13.1	Rate-Selection Algorithm.....	98
45	5.3.14	Open-Loop Pitch Quantization, Pitch Adjustment and Interpolation.....	99
46	5.3.14.1	Open-Loop Pitch Refinement and Quantization.....	99
47	5.3.14.2	Pitch Interpolation.....	100
48	5.3.15	Generating Modified Weighted Speech Signal.....	103
49			

1	5.3.15.1	Frame based processing for $C_{SM} \leq 0$	103
2	5.3.15.1.1	Region of Mapping of Weighted Speech to Modified Weighted Speech	
3		103	
4	5.3.15.1.2	Weighted Speech Mapping to Modified Weighted Speech for One Frame	
5		104	
6	5.3.15.1.3	Update of Modified Weighted Speech Buffer	104
7	5.3.15.2	Variable Subframe Processing $C_{SM} > 0$	104
8	5.3.15.2.1	Determining Parameters.....	105
9	5.3.15.2.2	Target Signal Determination for Current Local Delay Search	106
10	5.3.15.2.3	Searching Range Determination for Local Delay.....	107
11	5.3.15.2.4	Local Delay Search.....	109
12	5.3.15.2.4.1	Modification of Target Signal	109
13	5.3.15.2.4.2	Integer Local Delay Search	109
14	5.3.15.2.4.3	Fractional Local Delay Determination.....	110
15	5.3.15.2.4.4	Warping Current Weighted Speech With Determined Delay .	110
16	5.3.15.2.5	Further Modification of Weighted Speech for Half_Rate_Max.....	110
17	5.3.15.2.5.1	Generating a Backward Waveform Vector	111
18	5.3.15.2.5.2	Generating a Candidate Vector by Waveform Interpolation...	111
19	5.3.15.2.5.3	Making a Decision to Use Interpolated Modified Weighted	
20		Speech	112
21	5.3.15.2.6	Final Determination of Current Local Delay and Accumulated Delay	113
22	5.3.15.2.7	Updating Modified Weighted Speech Buffer	113
23	5.3.15.2.8	Pitch Pre-Enhancement and Periodicity Detection	113
24	5.3.15.3	Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch Gain	
25		for each Fixed Subframe	115
26	5.3.15.3.1	Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch	
27		Gains for Rate 1	115
28	5.3.15.3.2	Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch	
29		Gains for Rate 1/2.....	116
30	5.3.16	Final Frame Class and Type Decisions, Rate Selection, and Pitch Lag	
31		Refinement.....	118
32	5.3.16.1	Final Classification, Type and Rate Selection	118
33	5.3.16.2	Pitch Lags Refinement	119
34	5.3.17	Adaptive-codebook Gain De-emphasis for Type-1 Frame (Rate 1 and Rate 1/2)	
35		120	
36	5.3.18	Adaptive-codebook Gain Quantization for Type-1 Frame (Rate 1 and Rate 1/2)	
37		121	
38	5.3.19	Generating Modified Speech from Modified Weighted Speech.....	122
39	5.3.20	LSFs Smoothing Parameter.....	123
40	5.3.21	LSFs Smoothing.....	125
41	5.3.22	LSFs Quantization.....	126
42	5.3.22.1	Calculation of LSFs Quantization Weights	126
43	5.3.22.2	LSFs Mean Removal	127
44	5.3.22.3	LSFs Prediction	127
45	5.3.22.4	Quantization of the LSFs Prediction Error.....	127
46	5.3.22.4.1	LSFs Quantization for Rate 1	129
47	5.3.22.4.2	LSFs Quantization for Rate 1/2	129

1	5.3.22.4.3	LSFs Quantization for Rate 1/4	129
2	5.3.22.4.4	LSFs Quantization for Rate 1/8	130
3	5.3.23	Interpolation of Quantized and Unquantized LSFs.....	131
4	5.3.23.1	Interpolation of Quantized LSFs for Rate 1 Type-0 Frames.....	131
5	5.3.23.1.1	Calculating of Distance Weights	131
6	5.3.23.1.2	Determination of Best Interpolation Path.....	131
7	5.3.23.2	Interpolation of Quantized LSFs for Rate 1 Type-1 frames, Rate 1/4 frames, and Rate 1/8 Frames.....	132
9	5.3.23.3	Interpolation of Quantized and Unquantized LSFs for Rate 1/2 Type-0 Frames	132
11	5.3.23.4	Interpolation of Quantized and Unquantized LSFs for Rate 1/2 Type-1 Frames	132
13	5.3.24	Recalculating of Weighting Filter Coefficients for Rate 1/2	134
14	5.3.25	Identifying Long-Term Spectral Characteristic	135
15	5.4	Excitation Coding at Rate 1/8	137
16	5.4.1	Excitation and Unquantized Gain Determination	137
17	5.4.2	Gain quantization	137
18	5.5	Excitation Coding at Rate 1/4	139
19	5.5.1	Computation and Quantization of Gains	139
20	5.5.2	Random number generation	140
21	5.5.3	Creation of sparse non-zero excitation	141
22	5.5.4	Shaping the excitation	141
23	5.6	Excitation Coding at Rate 1 and Rate 1/2	143
24	5.6.1	Introduction to Excitation Coding for Rate 1 and Rate 1/2.....	143
25	5.6.2	Calculation of Impulse Response of the Combined Synthesis and Weighting Filter	146
27	5.6.3	Generating the Target Signal.....	147
28	5.6.4	Generating the Ideal Excitation.....	148
29	5.6.5	Adaptive-codebook Contribution for Type-1 Frames	149
30	5.6.5.1	The Adaptive-codebook Vector Contribution	149
31	5.6.5.2	The Filtered Adaptive Vector Contribution and Normalized Pitch Correlation 150	
33	5.6.6	Adaptive-codebook Contribution For Rate 1 Type-0 Frames.....	151
34	5.6.6.1	Normalized Cross Correlation Between the Target and the Filtered Excitation 152	
36	5.6.6.2	Fractional Pitch Search and Quantization.....	152
37	5.6.6.3	Adaptive-codebook Gain (Closed-Loop Pitch Gain) and Normalized Correlation.....	152
39	5.6.7	Adaptive-codebook Contribution For Rate 1/2 Type-0 Frames.....	154
40	5.6.7.1	Normalized Cross Correlation Between the Target and the Filtered Excitation 154	
42	5.6.7.2	Integer Pitch Search and Quantization.....	155
43	5.6.7.3	Adaptive-codebook Gain (Closed-Loop Pitch Gain) and Normalized Correlation.....	155
45	5.6.8	Scaling of Adaptive-Codebook Gain	156
46	5.6.9	Update of Target and Ideal Excitation	156
47	5.6.10	Calculation of Fixed-codebook Smoothing Parameters.....	157
48	5.6.11	Principles and Algorithms for Pulse Codebook Search	159
49	5.6.11.1	Analysis-by-Synthesis Approach for Fixed-Codebook Contribution	159
50	5.6.11.2	Principles of Iterative Search Procedure for Pulse Sub-Codebooks	159

1	5.6.11.3	Updates of Pulse Contributions	160
2	5.6.11.4	Concept of Pitch Enhancement for Pulse Codebooks.....	160
3	5.6.11.5	Pre-Calculation Procedure	161
4	5.6.11.6	An Iterative Algorithm Based on Position Search for a Single Pulse.....	162
5	5.6.11.6.1	Initialization of the Search Parameters	162
6	5.6.11.6.2	The Search Procedure	163
7	5.6.11.6.3	Building the Pulse Excitation.....	163
8	5.6.11.7	An Iterative Algorithm Based on Joint Position Search for Two Pulses	164
9	5.6.11.7.1	Initialization of the Search Parameters	164
10	5.6.11.7.2	Pre-Search for 2-Pulse Sub-Codebook	164
11	5.6.11.7.3	Pre-Search for 3-Pulse Sub-Codebook	164
12	5.6.11.7.4	The Search Procedure	165
13	5.6.11.7.5	Building the Pulse Excitation.....	165
14	5.6.11.8	Joint Position Search for Five Pulses	166
15	5.6.12	Gaussian Codebook Search.....	167
16	5.6.13	Fixed-codebook Contribution for Rate 1.....	168
17	5.6.13.1	Fixed-codebook Structure for Rate 1.....	168
18	5.6.13.2	Fixed-codebook Search for Rate 1	169
19	5.6.13.2.1	Pitch Enhancement Gain and Pre-Calculation of Search Parameters...	170
20	5.6.13.2.2	Fixed-codebook Search for Type-0 Frames.....	170
21	5.6.13.2.3	Fixed-codebook Search for Type-1 Frames.....	171
22	5.6.14	Fixed-codebook Contribution for Rate 1/2	172
23	5.6.14.1	Fixed-codebook Structure for Rate 1/2.....	172
24	5.6.14.2	Fixed-codebook Search for Rate 1/2	173
25	5.6.14.2.1	Parameters of Formant Correction Filter and Stability Index.....	174
26	5.6.14.2.2	Pulse Addition Based on Adaptive-codebook Correlations.....	175
27	5.6.14.2.3	Generating the Modified Impulse Response.....	176
28	5.6.14.2.4	Pitch Enhancement Gain and Pre-Calculation of Search Parameters...	176
29	5.6.14.2.5	Fixed-codebook Search for Type-0 Frames.....	177
30	5.6.14.2.6	Fixed-codebook Search for Type-1 Frames.....	177
31	5.6.15	Filtering of Fixed-Codebook Contribution.....	178
32	5.6.16	Gain Re-optimization and Normalization for Rate 1/2 and Rate 1	179
33	5.6.16.1	Normalization for $NSR > 0.125$	181
34	5.6.16.2	Normalization for $NSR \leq 0.125$	181
35	5.6.16.3	Fixed-Codebook Gain Normalization for Type-1 Frames and $NSR > 0.25$.	182
36	5.6.17	Adaptive-Codebook and Fixed-Codebook Gain Quantization for Type-0	
37		Frames	185
38	5.6.18	Fixed-Codebook Gain Quantization for Type-1 Frames.....	186
39	5.6.19	Random Number Generator	188
40	6	SMV Decoder	189
41	6.1	Rate and Frame Type Retrieval and Frame Error Detection.....	189
42	6.1.1	Handling NULL Rate 1/8 Frame.....	189
43	6.1.2	Handling All Zero Frames.....	189
44	6.1.3	Handling 'Rate 1 with Bit Errors' Frame.....	189
45	6.1.4	Channel Rate-Determination Algorithm Error Detection	189
46	6.2	General Structure of the Decoder	191
47	6.3	Decoding of the LSFs.....	191
48	6.4	Interpolation of Decoded LSFs	191

1	6.5	Excitation Decoding for Rate 1/8	192
2	6.6	Decoding of Rate 1/4.....	193
3	6.7	Decoding of Rate 1 and Rate 1/2.....	194
4	6.7.1	Identifying Long –Term Spectral Characteristics	195
5	6.7.2	Pitch Decoding (Adaptive-codebook Contribution).....	195
6	6.7.3	Adaptive-codebook Gain Decoding for Rate 1 and Rate 1/2 (Type-1).....	195
7	6.7.4	Fixed-codebook Excitation Decoding for Rate 1 and Rate 1/2.....	195
8	6.7.5	Fixed-codebook Excitation Contribution for Rate 1/8 and Rate 1/4.....	195
9	6.7.6	Gain Decoding.....	195
10	6.7.7	Signal-to-Noise Ratio (SNR) at the Decoder	196
11	6.8	Frame Erasure Concealment	197
12	6.8.1	LSF Extrapolation	197
13	6.8.2	Pitch Extrapolation	197
14	6.8.3	Gain Extrapolation	198
15	6.8.4	Excitation Generation.....	200
16	6.9	Post-Processing	200
17	6.9.1	Weighted residual signal	200
18	6.9.2	Pitch postfilter	200
19	6.9.3	Tilt Compensation	202
20	6.9.4	Short-Term Postfilter.....	202
21	6.9.5	Adaptive Gain Control	202
22	6.9.6	High-pass filtering and up-scaling	203
23	7	SMV Data Capabilities	205
24	7.1	TTY/TDD Payload Format	207
25	7.1.1	Baudot Code Payload	207
26	7.2	DTMF Payload Format	207
27	7.3	TTY/TDD Operation.....	208
28	7.3.1	TTY Header and Character Format.....	209
29	7.3.2	TTY Encoder Processing.....	210
30	7.3.2.1	TTY Encoder Inputs	210
31	7.3.2.2	Dit Classification	210
32	7.3.2.3	Dits to Bits	211
33	7.3.2.4	TTY Character Classification.....	212
34	7.3.2.5	TTY Baud Rate Determination	212
35	7.3.2.6	TTY State Machine.....	212
36	7.3.3	TTY/TDD Decoder Processing.....	212
37	7.3.3.1	TTY Decoder Inputs	212
38	7.3.3.2	Decoding the TTY/TDD Information.....	213
39	7.3.3.3	Baudot Generator	213
40	7.3.3.4	Tone Generator.....	214
41	7.4	DTMF	214
42	7.4.1	DTMF Detector	214
43	7.4.1.1	Description of the DTMF Detector.....	215
44	7.4.1.2	INPUT.....	215
45	7.4.1.3	OUTPUT	215
46			

1

2 **1 SERVICE OPTION 56: SELECTABLE MODE VOCODER TWO-WAY** 3 **VOICE COMMUNICATION**

4

5 **1.1 General Description**

6 Service Option 56, the Selectable Mode Vocoder (SMV), provides two-way voice
7 communication between the base station and the mobile station using the dynamically
8 variable data rate speech codec algorithm described in this standard. The transmitting
9 speech codec receives voice samples and generates an encoded speech packet for every
10 Fundamental Channel Traffic Channel frame.† The receiving station generates a speech
11 packet from every Fundamental Channel Traffic Channel frame and supplies it to the
12 speech codec for decoding into voice samples.

13 SMV communicates at one of four rates: 9600 bps, 4800 bps, 2400 bps and 1200 bps.

14 **1.2 Overview of SMV Documentation**

15 The SMV specification family consists of two standards. This standard provides the
16 algorithmic description of the SMV as well as the floating-point C-simulation of the codec.
17 The companion minimum performance standard, “*Recommended Minimum Performance*
18 *Standard for the Selectable Mode Vocoder, Service Option 56*”, consists of the master fixed-
19 point C-simulation of the SMV as well as a minimum performance specification. The
20 minimum performance specification consists of a set of objective and subjective tests used
21 to verify the quality of any non bit-exact SMV implementation.

22 Section 2 of this standard provides information on the interfaces between SMV and IS-2000
23 systems. Section 3 provides information on the audio interfaces to SMV. The specifications
24 given by Sections 4, 5, and 6 of this standard provide the algorithmic description of SMV at
25 a high level. Section 7 provides a description of the in-band data capabilities of SMV. The
26 floating point C-code accompanying this document provides a more detailed complementary
27 description of SMV. In the case of a discrepancy between the floating-point C-simulation
28 and the algorithmic description, the floating-point C-simulation will prevail.

29 **1.3 SMV Algorithmic Options**

30 The specifications defined in SMV Encoder Pre-Processing (Section 5.2), Voice Activity
31 Detection (Sections 5.3.2 and 5.3.3), Music Detection (Section 5.3.4), Frame Erasure
32 Concealment (Section 6.8), and Post-Processing (Section 6.9), are optional for non bit-exact
33 implementations intended for varying operational environments, such as base-station,
34 mobile station handsets and in-vehicle hands-free.

35 There are two algorithmic options for both the Noise Suppression and the Voice Activity
36 Detection modules as described below.

37

† IS-2000-A uses the term “frame” to represent a 20 ms grouping of data on the Fundamental Channel. Common speech codec terminology also uses the term “frame” to represent a quantum of processing. For Service Option 56, the speech codec frame corresponds to speech sampled over 20 ms. The speech samples are processed into a packet. This packet is transmitted in a Traffic Channel frame.

1 **1.3.1 Standard Options for Noise Suppression**

2 This standard defines two options (A and B) for Noise Suppression as described in Section
3 5.2.3. These two options are selected in the SMV simulation code by using the flags NS_A,
4 and NS_B. If the option NS_A is chosen, the total fixed algorithmic delay of the codec is 33
5 ms. If the option NS_B is chosen, the total fixed algorithmic delay at the codec is 30 ms.

7 **1.3.2 Standard Options for Voice Activity Detection**

8 This document defines two options (A and B) for Voice Activity Detection as described in
9 Sections 5.3.2 and 5.3.3. These two options are selected in the SMV simulation code by
10 using the flags VAD_A, and VAD_B respectively. VAD_A was used in the original SMV
11 candidate in the selection phase. VAD_B is based on the Rate Determination algorithm in
12 IS-127 (EVRC), requires less program memory, and reuses modules from earlier
13 implementations of IS-127.

14 **1.4 SMV Implementation Options**

15 There are two options available for implementation of the SMV:

16
17 If a bit-exact implementation approach is chosen, the implementation that is developed
18 shall be end-to-end bit-exact to the reference fixed-point C-simulation set to one of the four
19 possible combinations of noise suppression and voice activity detection. Bit-exactness is
20 verified via application of the associated SMV minimum performance standard [8], which
21 includes the reference fixed-point codec simulation.

22
23 Alternatively, an implementation that deviates from the end-to-end bit-exact SMV
24 specification described above may be developed. Such an implementation shall pass
25 objective and subjective tests defined by the SMV minimum performance standard [8]. The
26 master floating-point SMV in such a subjective test shall use Noise Suppression option B
27 and Voice Activity Detection option B since the SMV standard codec Characterization Test
28 was performed using Noise Suppression option B and Voice Activity Detection option B.

30 **1.5 Service Option Number**

31 The variable data rate two-way voice service option, using the speech codec algorithm
32 described by this standard, shall use service option number 56 and shall be called Service
33 Option 56.

35 **1.6 Allowable Delays**

36 **1.6.1 Allowable Transmitting Speech Codec Encoding Delay**

37 The transmitting speech codec shall supply a packet to the multiplex sublayer no later than
38 20 ms after it has obtained the last input sample for the current speech frame.

40 **1.6.2 Allowable Receiving Speech Codec Decoding Delay**

41 The receiving decoder shall generate the first sample of speech using parameters from a
42 packet supplied to it by the multiplex sublayer no later than 3 ms after being supplied the
43 packet.

1 **1.7 Special Cases**

2 **1.7.1 Blanked Packets**

3 A blanked frame occurs when the transmitting station uses the entire frame for either
4 signaling traffic or secondary traffic. The SMV does no special encode processing during
5 the generation of a blank packet; i.e., the generated voice packet is simply not used. The
6 decoder, in turn, treats a blank packet in the same manner as a frame erasure.

8 **1.7.2 Null Traffic Channel Data**

9 A Rate 1/8 packet with all bits set to '1' is considered as null Traffic Channel data. This
10 packet is declared an erased packet and handled as described in Section 6.8.

11

12 **1.8 Terms and Numeric Information**

13 **Autocorrelation Function.** A function showing the relationship of a signal with a time-
14 shifted version of itself.

15 **Base Station.** A station in the Public Radio Telecommunications Service, other than a
16 personal/mobile station, used for radio communications with personal/mobile stations.

17 **Codec.** The combination of an encoder and decoder in series (encoder/decoder).

18 **Code Excited Linear Predictive Coding (CELP).** A speech coding algorithm. CELP coders
19 use codebook excitation, a long-term pitch prediction filter, and a short-term formant
20 prediction filter.

21 **Codebook.** A set of vectors used by the speech codec. For each speech codec codebook
22 subframe, one particular vector is chosen and is used to excite the speech codec's filters.
23 The codebook vector is chosen to minimize the weighted error between the original and
24 synthesized speech after the pitch and formant synthesis filter coefficients have been
25 determined.

26 **Coder.** Same as "encoder."

27 **Decoder.** Generally, a device for the translation of a signal from a digital representation
28 into an analog format. For this standard, a device that converts speech encoded in the
29 format specified in this standard to analog or an equivalent PCM representation.

30 **Encoder.** Generally, a device for the translation of a signal into a digital representation.
31 For this standard, a device that converts speech from an analog, or from its equivalent PCM
32 representation, to the digital representation described in this standard.

33 **Formant.** A resonant frequency of the human vocal tract causing a peak in the short term
34 spectrum of speech.

35 **IIR Filter.** (Infinite-duration impulse response filter) A filter for which the output, in
36 response to an impulse input, never totally converges to zero. This term is usually used in
37 reference to digital filters.

38 **Linear Predictive Coding (LPC).** A method of predicting future samples of a sequence by a
39 linear combination of the previous samples of the same sequence. Linear Predictive Coding
40 is frequently used in reference to a class of speech codecs.

41 **Line Spectral Frequencies (LSFs).** A representation of digital filter coefficients in a
42 pseudo-frequency domain. This representation has good quantization and interpolation
43 properties.

- 1 **LSB.** Least significant bit.
- 2 **MSB.** Most significant bit.
- 3 **Normalized Autocorrelation Function.** A measure used to determine the pitch period
4 and the degree of periodicity of the input speech. This measure is useful in distinguishing
5 voiced from unvoiced speech.
- 6 **Packet.** The unit of information exchanged between service option applications in the base
7 station and the personal/mobile station.
- 8 **Personal/Mobile Station.** A station in the Public Radio Telecommunications Service
9 intended to be used while in motion or during halts at unspecified points.
- 10 **Pitch.** The fundamental frequency in speech caused by the periodic vibration of the
11 human vocal cords.
- 12 **RDA.** Rate Determination Algorithm.
- 13 **Receive Objective Loudness Rating (ROLR).** A measure of receive audio sensitivity.
14 ROLR is a frequency-weighted ratio of the line voltage input signal to a reference encoder to
15 the acoustic output of the receiver. IEEE 269 defines the measurement of sensitivity and
16 IEEE 661 defines the calculation of objective loudness rating.
- 17 **SPL.** Sound Pressure Level.
- 18 **Transmit Objective Loudness Rating (TOLR).** A measure of transmit audio sensitivity.
19 TOLR is a frequency-weighted ratio of the acoustic input signal at the transmitter to the
20 line voltage output of the reference decoder. IEEE 269 defines the measurement of
21 sensitivity and IEEE 661 defines the calculation of objective loudness rating.
- 22 **Voiced Speech.** Speech generated when the vocal cords are vibrating at a fundamental
23 frequency. Characterized by high energy, periodicity, and a large ratio of energy below
24 2 kHz to energy above 2 kHz.
- 25 **Unvoiced Speech.** Speech generated by forcing air through constrictions in the vocal tract
26 without vibration of the vocal cords. Characterized by a lack of periodicity, and a near-
27 unity ratio of energy below 2 kHz to energy above 2 kHz.
- 28 **WAEPL.** Weighted Acoustic Echo Path Loss. A measure of the echo performance under
29 normal conversation. ANSI/EIA/TIA-579 defines the measurement of WAEPL.
- 30 **Zero Input Response (ZIR).** The filter output caused by the non-zero initial state of the
31 filter when no input is present.
- 32 **Zero State Response (ZSR).** The filter output caused by an input when the initial state of
33 the filter is zero.
- 34 **ZIR.** See Zero Input Response.
- 35 **ZSR.** See Zero State Response.
- 36

2 REQUIRED MULTIPLEX OPTION SUPPORT

Service Option 56 shall support an interface with Multiplex Option 1. Speech packets for Service Option 56 shall only be transported as primary traffic.

2.1 Interface to Multiplex Option 1

2.1.1 Transmitted Packets

The speech codec shall generate and shall supply exactly one packet to the multiplex sublayer every 20 milliseconds. The packet contains the service option information bits that are transmitted as primary traffic. The packet shall be one of four types as shown in Table 2.1-1. The number of bits supplied to the multiplex sublayer for each type of packet shall also be as shown in Table 2.1-1. Unless otherwise commanded, the speech codec may supply a Rate 1, Rate 1/2, Rate 1/4, or Rate 1/8 packet. Upon command, the speech codec shall generate a Blank packet. Also upon command, the speech codec shall generate a non-blank packet with a maximum rate of Rate 1/2.

A Blank packet contains no bits and is used for blank-and-burst transmission of signaling traffic or secondary traffic (see IS-2000.3-A and IS-2000.5-A).

Table 2.1-1. Packet Types Supplied by Service Option 56 to the Multiplex Sublayer

Packet Type	Bits per Packet
Rate 1	171
Rate 1/2	80
Rate 1/4	40
Rate 1/8	16
Blank	0

2.1.2 Received Packets

The multiplex sublayer in the receiving station categorizes every received Traffic Channel frame, and supplies the packet type and accompanying bits, if any, to the speech codec as shown in ~~Table 2.1-1~~[Table 2.1-1](#). The speech codec processes the bits of the packet as described in Section 6. The received packet types shown in ~~Table 2.1-2~~[Table 2.1-2](#) correspond to the transmitted packet types shown in ~~Table 2.1-1~~[Table 2.1-1](#). The Blank packet type occurs when the receiving station determines that a blank-and-burst frame for signaling traffic or secondary traffic was transmitted. The Rate 1 with bit errors packet type occurs when the receiving station determines that the frame was transmitted at 9600 bps and the frame has one or more bit errors. The insufficient frame quality packet type occurs when the mobile station is unable to decide upon the data rate of the received frame or when the mobile station detects a frame in error that does not belong to the Rate 1 with bit errors packet type.

1 **Table 2.1-2. Packet Types Supplied by the Multiplex Sublayer to Service Option 56**

Packet Type	Bits per Packet
Rate 1	171
Rate 1/2	80
Rate 1/4	40
Rate 1/8	16
Blank	0
Rate 1 with bit errors	171
Insufficient frame quality (erasure)	0

2 **2.2 Negotiation for Service Option 56**

3 The mobile station and base station can negotiate for Service Option 56 service negotiation,
4 as described in IS-2000.5-A.

5 **2.2.1 Procedures Using Service Negotiation**

6 The mobile station and base station shall perform service negotiation for Service Option 56
7 as described in IS-2000.5-A, and the negotiated service configuration shall include only
8 valid attributes for the service option as specified in [Table 2.2-1](#)~~Table 2.2-1~~.
9

10 **Table 2.2-1. Valid Service Configuration Attributes for Service Option 56**

Service Configuration Attribute	Valid Selections
Forward Multiplex Option	Multiplex Option 1
Reverse Multiplex Option	Multiplex Option 1
Forward Transmission Rates	Rate Set 1 with all rates enabled
Reverse Transmission Rates	Rate Set 1 with all rates enabled
Forward Traffic Type	Primary Traffic
Reverse Traffic Type	Primary Traffic

11 **2.2.1.1 Initialization and Connection**

12 **2.2.1.1.1 Mobile Station Requirements**

13 If the mobile station accepts a service configuration, as specified in a *Service Connect*
14 *Message*, *General Handoff Direction Message*, or *Universal Handoff Direction Message* that
15 includes a service option connection using Service Option 56, the mobile station shall
16 perform the following:
17

- 18 • If the service option connection is new (that is, not part of the previous service
19 configuration), the mobile station shall perform speech codec initialization at the time
20 specified by the maximum of the action time associated with the message carrying
21

1 the Service Configuration Record, and the time that the corresponding Call Control
 2 Instance is instantiated. The mobile station shall initialize its SMV encoder mode of
 3 operation to a default value of 0. The mobile station shall complete the initialization
 4 within 40 ms.

- 5 • Beginning at the time specified by the maximum of the action time associated with
 6 the message carrying the Service Configuration Record, and the time that the
 7 corresponding Call Control Instance is instantiated, and continuing for as long as the
 8 service configuration includes the service option connection, Service Option 56 shall
 9 process received packets and generate and supply packets for transmission as
 10 follows:

- 11 - If the Call Control Instance is in the *Conversation Substate*, Service Option 56
 12 shall process the received packets and generate and supply packets for
 13 transmission in accordance with this standard.

- 14 - If the Call Control Instance is not in the *Conversation Substate*, Service Option
 15 56 shall process the received packets in accordance with this standard, and
 16 shall generate and supply Rate 1/8 Packets with all bits set to '1' for
 17 transmission, except when commanded to generate a Blank packet.

18 2.2.1.1.2 Base Station Requirements

19 If the base station establishes a service configuration, as specified in a *Service Connect*
 20 *Message*, *General Handoff Direction Message*, or *Universal Handoff Direction Message* that
 21 includes a service option connection using Service Option 56, the base station shall
 22 perform the following:

- 23 • If the service option connection is new (that is, not part of the previous service
 24 configuration), the base station shall perform speech codec initialization no later
 25 than the time specified by the maximum of the action time associated with the
 26 message carrying the Service Configuration Record, and the time that the
 27 corresponding Call Control Instance is Instantiated. The base station shall initialize
 28 its SMV encoder mode of operation to a default value of 0.
- 29 • Commencing at the time specified by the maximum of the action time associated with
 30 the message carrying the Service Configuration Record, and the time that the
 31 corresponding Call Control Instance is Instantiated, and continuing for as long as the
 32 service configuration includes the service option connection, Service Option 56 shall
 33 process received packets, and shall generate and supply packets for transmission in
 34 accordance with this standard. The base station may defer enabling the audio input
 35 and output.

37 2.2.1.2 Service Option Control Messages

38 2.2.1.2.1 Mobile Station Requirements

39 The mobile station shall support one pending *Service Option Control Message* for Service
 40 Option 56.

41 If the mobile station receives a *Service Option Control Message* for Service Option 56, then,
 42 at the action time associated with the message, the mobile station shall process the
 43 message as follows:

- 44 1. If the MOBILE_TO_MOBILE field is equal to '1', the mobile station should disable the
 45 audio output of the speech codec for 1 second after initialization.

1 If the MOBILE_TO_MOBILE field is equal to '0', the mobile station shall process each
2 received packet as described in Section 6.

3 2. If the INIT_CODEEC field is equal to '1', the mobile station shall perform speech codec
4 initialization. The mobile station shall complete the initialization within 40 ms.

5 3. SMV accepts as input, a mode of operation through the RATE_REDUCE field as defined
6 in Table 2.2-2. Using this mode of operation, SMV generates Rate 1, Rate 1/2, and
7 Rate 1/4 packets in a proportion that results in the average data rate given by Table
8 2.2-2.

9 Service Option 56 shall continue to use these fractions until either of the following events
10 occur:

- 11 • The mobile station receives a *Service Option Control Message* specifying a
12 different RATE_REDUCE, or
- 13 • Service Option 56 is initialized.

14 Service Option 56 was developed using reduced rate operation (encoding mode of
15 operation) as a network control criteria. Unlike the rate reduction mechanism in EVRC,
16 Service Option 56 does not use a deterministic pattern of Rate 1 and Rate 1/2 to reduce
17 the average encoding rate. Instead, the SMV codec selects the encoding rate based
18 upon the characteristics of the input speech: Voiced, unvoiced, transient, stationary, as
19 well as the encoding mode selected.

20

21

Table 2.2-2. SMV Encoding Rate Mode Control Parameters

RATE_REDUCE (binary)	SMV Encoding Mode of Operation	Estimated Average Encoding Rate for Active Speech (Source Encoding Rates)
'000'	0	7.95 kbps
'001'	1	5.82 kbps
'010'	2	4.50 kbps
'011'	3	3.95 kbps
'100'	4 (1/2 rate Maximum applied to Mode 0)	4.0 kbps
'101'	5 (1/2 rate Maximum applied to Mode 1)	3.67 kbps
All other RATE_REDUCE values are reserved.		

22 If the RATE_REDUCE field is not equal to a value defined in [Table 2.2-2](#), the
23 mobile station shall reject the message by sending a *Mobile Station Reject Order* with the
24 ORDQ field set equal to '00000100'.

25 2.2.1.2.2 Base Station Requirements

26 The base station may send a *Service Option Control Message* to the mobile station. If the
27 base station sends a *Service Option Control Message*, the base station shall include the
28 following type-specific fields for Service Option 56:
29

1

Table 2.2-3. Service Option Control Message Type-Specific Fields

Field	Length (bits)
RATE_REDUCE	3
RESERVED	3
MOBILE_TO_MOBILE	1
INIT_CODEC	1

2

3

RATE_REDUCE - SMV mode of operation.

4

The base station shall set this field to the RATE_REDUCE value from [Table 2.2-4](#) corresponding to the mode of operation that the mobile station is to operate in.

5

6

7

RESERVED - Reserved bits.

8

The base station shall set this field to '000'.

9

MOBILE_TO_MOBILE - Mobile-to-mobile processing.

10

If the mobile station is to perform mobile-to-mobile processing (Section 2.2.1.2.1), the base station shall set this field to '1'. In addition, if the mobile station is to disable the audio output of the speech codec for 1 second after initialization, the base station shall set the INIT_CODEC field and the MOBILE_TO_MOBILE field to '1'. If the mobile station is not to perform mobile-to-mobile processing, the base station shall set the MOBILE_TO_MOBILE field to '0'.

11

12

13

14

15

16

17

18

INIT_CODEC - Initialize speech codec.

19

If the mobile station is to initialize the speech codec, the base station shall set this field to '1'; otherwise, the base station shall set this field to '0'.

20

21

22

1

Table 2.2-4. SMV Encoding Rate Mode Control Parameters

RATE_REDUCE (binary)	SMV Encoding Mode of Operation	Estimated Average Encoding Rate for Active Speech (Source encoding rates)
'000'	0	7.95 kbps
'001'	1	5.82 kbps
'010'	2	4.50 kbps
'011'	3	3.95 kbps
'100'	4 (1/2 rate Maximum applied to Mode 0)	4.0 kbps
'101'	5 (1/2 rate Maximum applied to Mode 1)	3.67 kbps
All other RATE_REDUCE values are reserved.		

2
3
4

1 **3 AUDIO INTERFACES**

2 **3.1 Input Audio Interface**

3 **3.1.1 Input Audio Interface in the Mobile Station**

4 The input audio may be either an analog or digital signal.

5 **3.1.1.1 Conversion and Scaling**

6 Whether the input is analog or digital, the signal presented to the input of the speech codec
7 shall be sampled at a rate of 8000 samples per second and shall be quantized to a uniform
8 PCM format with at least 13 bits of dynamic range.

9 The quantities in this standard assume a 16-bit integer input normalization with a range
10 from -32,768 through +32,767. The following speech codec discussion assumes this 16-bit
11 integer normalization. If an input audio interface uses a different normalization scheme,
12 then appropriate scaling should be used.

13 **3.1.1.2 Digital Audio Input**

14 If the input audio is an 8-bit μ -Law PCM signal, it shall be converted to a uniform PCM
15 format according to Table 2 in ITU-T Recommendation G.711 "Pulse Code Modulation
16 (PCM) of Voice Frequencies.

17 **3.1.1.3 Analog Audio Input**

18 If the input is in analog form, the mobile station shall sample the analog speech and shall
19 convert the samples to a digital format for speech codec processing. This shall be
20 accomplished by either the following or an equivalent method. First, the input gain audio
21 level is adjusted. Then, the signal is bandpass filtered to prevent aliasing. Finally, the
22 filtered signal is sampled and quantized (Section 3.1.1.1).

23 3.1.1.3.1 Transmit Level Adjustment

24 The mobile station shall have a transmit objective loudness rating (TOLR) equal to -46 dB,
25 when transmitting to a reference base station. The loudness ratings are described in IEEE
26 Standard 661-1979 "IEEE Standard Method for Determining Objective Loudness Ratings of
27 Telephone Connections."

28 When a 97dBspl \pm 8dB C-message weighted white noise acoustic signal from an artificial
29 mouth is applied to the microphone of the subscriber unit, as specified by "IEEE Standard
30 Method for Measuring Transmission Performance of Telephone Sets," IEEE standard 269-
31 1983, the input samples shall be approximately three bits less than the maximum PCM
32 value. With PCM samples generated from μ -Law /uniform code, resulting in samples with a
33 maximum of 14 bits, only 11 bits of precision should be generated.

34 3.1.1.3.2 Band Pass Filtering

35 The analog-to-digital conversion shall conform to the mask shown in ITU-T G.712 (11/96)
36 Figure 4/G.712, "Attenuation/frequency distortion for 4-wire analogue-to-digital channels.
37 (E4in to Tout)". Additional anti-aliasing filtering may be provided by the manufacturer.

38 3.1.1.3.3 Echo Return Loss

1 Provision shall be made to ensure adequate isolation between receive and transmit audio
2 paths in all modes of operation, such that when no external transmit audio is present, the
3 speech codec does not generate packets at rates higher than Rate 1/8 due to acoustic
4 coupling of the receive audio into the transmit audio path (specifically with the receive
5 audio at full volume). Target levels of 45 dB WAEPL should be met. See ANSI/EIA/TIA
6 Standard 579 “Acoustic-to-Digital and Digital-to-Acoustic Transmission Requirements for
7 ISDN Terminals.”

8 **3.1.2 Input Audio Interface in the Base Station**

9 **3.1.2.1 Sampling and Format Conversion**

10 The base station converts the input speech (analog, μ -Law companded Pulse Code
11 Modulation, or other format) into a uniform quantized PCM format with at least 13 bits of
12 dynamic range. The sampling rate is 8000 samples per second. The sampling and
13 conversion process shall be as in Section 3.1.1.1.

14 **3.1.2.2 Transmit Level Adjust**

15 The base station shall set the transmit level so that a 1004 Hz tone at a level of 0 dBm0 at
16 the network interface produces a level 3.17 dB below maximum amplitude at the output of
17 the quantizer. Measurement techniques and tolerances are described in “*Recommended*
18 *Minimum Performance Standard for the Selectable Mode Vocoder, Service Option 56*”.

19 **3.1.2.3 Line Echo Canceling**

20 The base station shall provide a method to cancel echoes returned by the PSTN interface.†
21 The echo canceling function shall meet ITU-T G.168 requirement. The echo canceling
22 function should work over a range of PSTN echo return delays, from 0 to 48 ms; however,
23 the latter requirement is subject to local PSTN configuration demands, e.g., a 64 ms (or
24 greater) echo canceling capability may be required in cases where the PSTN does not
25 provide echo cancellation for long distance service.

26 **3.2 Output Audio Interface**

27 **3.2.1 Output Audio Interface in the Mobile Station**

28 The PCM samples shall be converted to μ -Law samples according to ITU-T Recommendation
29 G.711. This μ -Law signal shall be converted to analog as specified by ITU-T Red Book
30 G.711.

31 **3.2.2 Band Pass Filtering**

32 The digital-to-analog conversion shall conform to the mask shown in ITU-T G.712 (11/96)
33 Figure 4/G.712, “Attenuation/frequency distortion for 4-wire analogue-to-digital channels.
34 (E4in to Tout)”. Additional reconstruction filtering may be provided by the manufacturer.

35 **3.2.2.1 Receive Level Adjustment**

36 The mobile station shall have a nominal receive objective loudness rating (ROLR) equal to
37 51 dB when receiving from a reference base station. The loudness ratings are described in

† Because of the relatively long delays inherent in the speech coding and transmitting processes, echoes that are not sufficiently suppressed are noticeable to the mobile station user.

1 IEEE Standard 661-1979 "IEEE Standard Method for Determining Objective Loudness
2 Ratings of Telephone Connections."

3 When a 97dBspl \pm 8dB C-message weighted white noise acoustic signal from an artificial
4 mouth is applied to the microphone of the subscriber unit as specified by "IEEE Standard
5 Method for Measuring Transmission Performance of Telephone Sets," IEEE standard 269-
6 1983, the input samples shall be of approximately three bits less than the maximum PCM
7 value. With PCM samples generated from μ -Law /uniform code resulting in samples with a
8 maximum of 14 bits, only 11 bits of precision should be generated.

9 **3.2.3 Output Audio Interface in the Base Station**

10 Details of the digital and analog interfaces to the network are outside the scope of this
11 document.

12 **3.2.3.1 Receive Level Adjustment**

13 The base station shall set the audio level so that a received 1004 Hz tone 3.17 dB below
14 maximum amplitude produces a level of 0 dBm0 at the network interface. Measurement
15 techniques and tolerances are described in "*Recommended Minimum Performance Standard*
16 *for the Selectable Mode Vocoder, Service Option 56*".
17

- 1
- 2 No Text
- 3

1

1 4 THE SELECTABLE MODE VOCODER (SMV) SPEECH CODING – 2 INTRODUCTION AND TABLES

3 4.1 Introduction to the SMV Speech Coding Algorithm

4 The SMV speech-coding algorithm is based on four codecs (encoder/decoder) operating at
5 the rates of 8.55 kbps, 4.0 kbps, 2.0 kbps and 0.8 kbps. These codecs are called Rate 1
6 (full-rate), Rate 1/2 (half-rate), Rate 1/4 (quarter-rate) and Rate 1/8 (eighth-rate),
7 respectively.

8 The SMV algorithm has 4 network-controlled operating modes: Mode 0 (premium mode),
9 Mode 1 (standard mode), Mode 2 (economy mode), and Mode 3 (capacity-saving mode). The
10 different modes allow a tradeoff between average data rate (ADR) and speech quality. In
11 addition, Mode 0 and Mode 1 can be configured to operate in a half-rate max (HRM) mode,
12 where the maximum rate allowed is the Rate 1/2.

13

14

	Mode 0	Mode 1	Mode 2	Mode 3
Rate 1	68.90%	38.14%	15.43%	7.49%
Rate 1/2	6.03%	15.82%	38.34%	46.28%
Rate 1/4	0.00%	17.37%	16.38%	16.38%
Rate 1/8	25.07%	28.67%	29.85%	29.85%
ADR	7205 bps	5182 bps	4073 bps	3692 bps

15

16 **Figure 4.1-1: Rate Percentages and Average Data Rates for Clean Speech at Nominal**
17 **Level (-22 dBov)**

18

19

	Mode 0	Mode 1	Mode 2	Mode 3
Rate 1	55.13%	29.13%	16.60%	6.86%
Rate 1/2	10.65%	19.49%	32.39%	42.13%
Rate 1/4	0.00%	13.17%	14.93%	14.93%
Rate 1/8	34.21%	38.20%	36.08%	36.08%
ADR	6214 bps	4507 bps	3940 bps	3472 bps

20

21 **Figure 4.1-2: Rate Percentages and Average Data Rates for Speech in the Presence of**
22 **-15 dB Street Noise**

23 Figure 4.1-1 gives the percentages of the selected rates and the ADR for a clean speech
24 database at nominal level (-22 dBov) that was used for speech quality tests during the SMV
25 selection. The voice activity for this database is about 70%; hence, that database does not
26 represent a typical conversational speech traffic.

27 Figure 4.1-2 gives the percentages of the selected rates and the ADR for speech that was
28 generated from the same database, with the addition of street noise at the level of -15
29 signal-to-noise level.

30

1 A 20 ms frame of signal sampled at the rate of 8000 Hz (160 samples) is processed by one
2 of the four codecs, according to a rate-determination algorithm (RDA). The rate selection is
3 based upon the frame characteristics (voiced speech, unvoiced speech, background noise,
4 stationarity, etc.), and is controlled by the SMV mode. For each frame, the best rate for the
5 frame is selected under the ADR constraint dictated by the SMV operating mode.

6 The four codecs are based upon the excitation-filter approach, where a Linear-Prediction
7 (LPC) filter is excited by an excitation source. The LPC filter parameters and the excitation
8 parameters are transmitted from the encoder to the decoder.

9 The Rate 1 codec and the Rate 1/2 codec are based on the code-excited linear prediction
10 (CELP) approach, where the excitation of the LPC filter is represented by entries in a fixed
11 codebook and in an adaptive codebook, multiplied by a fixed-codebook gain and an
12 adaptive-codebook gain, respectively. The SMV system also selects a frame type for each
13 Rate 1 or Rate 1/2 frame. The type can be 0 or 1. Type-1 frames are frames of stationary
14 voiced speech, while frames of type 0 are all other types of speech. The bit allocation for the
15 excitation parameters and the LPC filter parameters differs between type-1 frames and
16 type-0 frames. For Rate 1 and Rate 1/2, the number of subframes in each frame depends
17 on the rate and on the frame type. Rate 1 frames are divided into 4 subframes of 40
18 samples, or 5 ms each. Rate 1/2 type-0 frames are divided into 2 subframes of 80 samples
19 or 10 ms each. Rate 1/2 type-1 frames are divided into 3 subframes of 53, 53, and 54
20 samples, or 6.625, 6.625, and 6.750 ms, respectively.

21 The excitation of the LPC filter for Rate 1/4 is generated by a random number generator.
22 The frame is divided 10 subframes of 2 ms, and each subframe is multiplied by a gain
23 factor. The LPC excitation is then filtered by a selected frequency-shaping filter.

24 The excitation of the LPC filter for Rate 1/8 is generated by a random number generator,
25 and multiplied by a single gain for the 20 ms frame.

26

27

28

1 **4.2 Bit Allocation Tables**

2 The bit allocation for Rate 1 is given in Table 4.2-1, the bit allocation for Rate 1/2 is given
 3 in Table 4.2-2, the bit allocation for Rate 1/4 is given in Table 4.2-3, and the bit allocation
 4 for Rate 1/8 is given in Table 4.2-4.

5
 6

Parameter	Bits per 20 ms frame			
	Type 0 (4 subframes)		Type 1 (4 subframes)	
LSFs	Interpolation	2 bit		
	Predictor switch	1 bit	Predictor switch	1 bit
	1st stage	7 bits	1st stage	7 bits
	2nd stage	7 bits	2nd stage	7 bits
	3rd stage	6 bits	3rd stage	6 bits
	4th stage	4 bits	4th stage	4 bits
		27 bits		25 bits
Type	1 bit			
Adaptive codebook	8,5,8,5 bits/subframe	26 bits	8 bits/frame	8 bits
Fixed codebook	5-pulse codebook	21 bits/subframe	8-pulse codebook	30 bits/subframe
	5-pulse codebook	20 bits/subframe		
	5-pulse codebook	20 bits/subframe		
		22 bits/subframe		
	22 bits/subframe	88 bits	30 bits/subframe	120 bits
Adaptive-codebook gain	2D VQ/subframe	7 bits/subframe	4D preVQ/frame	6 bits
Fixed-codebook gain		28 bits	4D delayed VQ/frame	10 bits
Rate sanity flag	1 bit		1 bit	
TOTAL	171 bits		171 bits	

7

8

Table 4.2-1: Bit Allocation Table for the SMV Rate 1 Codec

9

1

2

Parameter	Bits per 20 ms frame			
	Type 0 (2 subframes)		Type 1 (3 subframes)	
LSFs	Predictor switch		1 bit	
	1st stage		7 bits	
	2nd stage		7 bits	
	3rd stage		6 bits	
	21 bits			
Type	1 bit			
Adaptive codebook	7 bits/subframe	14 bits	7 bits/frame	7 bits
Fixed codebook	2-pulse codebook	14 bits/subframe	2-pulse codebook	12 bits/subframe
	3-pulse codebook	13 bits/subframe	3-pulse codebook	11 bits/subframe
	Gaussian codebook	13 bits/subframe	5-pulse codebook	11 bits/subframe
		15 bits/subframe		13 bits/subframe
	15 bits/subframe	30 bits	13 bits/subframe	39 bits
Adaptive-codebook gain	2D VQ/subframe	7 bits/subframe	3D preVQ/frame	4 bits
Fixed codebook gain	14 bits		3D delayed VQ/frame	8 bits
TOTAL	80 bits		80 bits	

3

4

Table 4.2-2: Bit Allocation Table for the SMV Rate 1/2 Codec

5

6

7

8

9

10

Note: since the entries of 2-pulse codebook for Rate 1/2 type-0 frames occupy only 12800 entries in the $2^{14}=16384$ possible entry locations, 3472 additional entries of the Gaussian codebook occupy the higher entry locations of the 2-pulse codebook.

1

Parameter	Bits per 20 ms frame	
LSFs	1st stage	7 bits
	2nd stage	7 bits
	3rd stage	6 bits
	20 bits	
Gains	17 bits	
Shape	2 bits	
Rate sanity flag	1 bit	
TOTAL	40 bits	

2

3

4

Table 4.2-3: Bit Allocation Table for the SMV Rate 1/4 codec

5

6

Parameter	Bits per 20 ms frame	
LSFs	1st stage	4 bits
	2nd stage	4 bits
	3rd stage	3 bits
	11 bits	
Energy	5 bits/subframe	5 bits
TOTAL	16 bits	

7

8

Table 4.2-4: Bit Allocation Table for the SMV Rate 1/8 Codec

9

10

11

12

13

14

1 4.3 SMV Bit Streams

2 After encoding, the encoded speech packets shall be formatted as described in Table 4.3-1
 3 through Table 4.3-7. For each parameter, bit index 0 corresponds to the most significant
 4 bit. The packet buffer bit index 1 indicates the bit position corresponding to the beginning
 5 of the packet buffer.

6

Bit Index	Rate 1 Frame Type 0 Packet Bits	Rate 1 Frame Type 1 Packet Bits
1	SVS_DECI=0	SVS_DECI=1
2	LSF_STAGE1[0]	LSF_STAGE1[0]
3	LSF_STAGE1[1]	LSF_STAGE1[1]
4	LSF_STAGE1[2]	LSF_STAGE1[2]
5	LSF_STAGE1[3]	LSF_STAGE1[3]
6	LSF_STAGE1[4]	LSF_STAGE1[4]
7	LSF_STAGE1[5]	LSF_STAGE1[5]
8	LSF_STAGE1[6]	LSF_STAGE1[6]
9	LSF_STAGE2[0]	LSF_STAGE2[0]
10	LSF_STAGE2[1]	LSF_STAGE2[1]
11	LSF_STAGE2[2]	LSF_STAGE2[2]
12	LSF_STAGE2[3]	LSF_STAGE2[3]
13	LSF_STAGE2[4]	LSF_STAGE2[4]
14	LSF_STAGE2[5]	LSF_STAGE2[5]
15	LSF_STAGE2[6]	LSF_STAGE2[6]
16	LSF_STAGE3[0]	LSF_STAGE3[0]
17	LSF_STAGE3[1]	LSF_STAGE3[1]
18	LSF_STAGE3[2]	LSF_STAGE3[2]
19	LSF_STAGE3[3]	LSF_STAGE3[3]
20	LSF_STAGE3[4]	LSF_STAGE3[4]
21	LSF_STAGE3[5]	LSF_STAGE3[5]
22	LSF_STAGE4[0]	LSF_STAGE4[0]
23	LSF_STAGE4[1]	LSF_STAGE4[1]
24	LSF_STAGE4[2]	LSF_STAGE4[2]
25	LSF_STAGE4[3]	LSF_STAGE4[3]
26	LSF_SWTCHPRD	LSF_SWTCHPRD
27	LSF_INTRPATH[0]	PITCH_FRM[0]
28	LSF_INTRPATH[1]	PITCH_FRM[1]
29	PITC_SBFPM1[0]	PITCH_FRM[2]

30	PITC_SBF1[1]	PITCH_FRM[3]
31	PITC_SBF1[2]	PITCH_FRM[4]
32	PITC_SBF1[3]	PITCH_FRM[5]
33	PITC_SBF1[4]	PITCH_FRM[6]
34	PITC_SBF1[5]	PITCH_FRM[7]
35	PITC_SBF1[6]	PITCH_VQ_GAIN[0]
36	PITC_SBF1[7]	PITCH_VQ_GAIN[1]
37	PITC_SBF2[0]	PITCH_VQ_GAIN[2]
38	PITC_SBF2[1]	PITCH_VQ_GAIN[3]
39	PITC_SBF2[2]	PITCH_VQ_GAIN[4]
40	PITC_SBF2[3]	PITCH_VQ_GAIN[5]
41	PITC_SBF2[4]	PLCB_VQ_GAIN[0]
42	PITC_SBF3[0]	PLCB_VQ_GAIN[1]
43	PITC_SBF3[1]	PLCB_VQ_GAIN[2]
44	PITC_SBF3[2]	PLCB_VQ_GAIN[3]
45	PITC_SBF3[3]	PLCB_VQ_GAIN[4]
46	PITC_SBF3[4]	PLCB_VQ_GAIN[5]
47	PITC_SBF3[5]	PLCB_VQ_GAIN[6]
48	PITC_SBF3[6]	PLCB_VQ_GAIN[7]
49	PITC_SBF3[7]	PLCB_VQ_GAIN[8]
50	PITC_SBF4[0]	PLCB_VQ_GAIN[9]
51	PITC_SBF4[1]	PLCB_PSIGN1_5_SBF1
52	PITC_SBF4[2]	PLCB_PSIGN2_6_SBF1
53	PITC_SBF4[3]	PLCB_PSIGN3_7_SBF1
54	PITC_SBF4[4]	PLCB_PSIGN4_8_SBF1
55	GAIN_SBF1[0]	PLCB_PLOC1_SBF1[0]
56	GAIN_SBF1[1]	PLCB_PLOC1_SBF1[1]
57	GAIN_SBF1[2]	PLCB_PLOC1_SBF1[2]
58	GAIN_SBF1[3]	PLCB_PLOC1_SBF1[3]
59	GAIN_SBF1[4]	PLCB_PLOC2_SBF1[0]
60	GAIN_SBF1[5]	PLCB_PLOC2_SBF1[1]
61	GAIN_SBF1[6]	PLCB_PLOC2_SBF1[2]
62	GAIN_SBF2[0]	PLCB_PLOC3_SBF1[0]
63	GAIN_SBF2[1]	PLCB_PLOC3_SBF1[1]
64	GAIN_SBF2[2]	PLCB_PLOC3_SBF1[2]
65	GAIN_SBF2[3]	PLCB_PLOC4_SBF1[0]
66	GAIN_SBF2[4]	PLCB_PLOC4_SBF1[1]

67	GAIN_SBF2[5]		PLCB_PLOC4_SBF1[2]
68	GAIN_SBF2[6]		PLCB_PLOC5_SBF1[0]
69	GAIN_SBF3[0]		PLCB_PLOC5_SBF1[1]
70	GAIN_SBF3[1]		PLCB_PLOC5_SBF1[2]
71	GAIN_SBF3[2]		PLCB_PLOC5_SBF1[3]
72	GAIN_SBF3[3]		PLCB_PLOC6_SBF1[0]
73	GAIN_SBF3[4]		PLCB_PLOC6_SBF1[1]
74	GAIN_SBF3[5]		PLCB_PLOC6_SBF1[2]
75	GAIN_SBF3[6]		PLCB_PLOC7_SBF1[0]
76	GAIN_SBF4[0]		PLCB_PLOC7_SBF1[1]
77	GAIN_SBF4[1]		PLCB_PLOC7_SBF1[2]
78	GAIN_SBF4[2]		PLCB_PLOC8_SBF1[0]
79	GAIN_SBF4[3]		PLCB_PLOC8_SBF1[1]
80	GAIN_SBF4[4]		PLCB_PLOC8_SBF1[2]
81	GAIN_SBF4[5]		PLCB_PSIGN1_5_SBF2
82	GAIN_SBF4[6]		PLCB_PSIGN2_6_SBF2
83	PLCB_SLCT1_SBF1=1	PLCB_SLCT1_SBF1=0	PLCB_PSIGN3_7_SBF2
84	PLCB_PSIGN1_3_SBF1	PLCB_SLCT2_SBF1	PLCB_PSIGN4_8_SBF2
85	PLCB_PSIGN2_4_SBF1	PLCB_PSIGN1_SBF1	PLCB_PLOC1_SBF2[0]
86	PLCB_PSIGN5_SBF1	PLCB_PSIGN2_SBF1	PLCB_PLOC1_SBF2[1]
87	PLCB_PLOC1_SBF1[0]	PLCB_PSIGN3_SBF1	PLCB_PLOC1_SBF2[2]
88	PLCB_PLOC1_SBF1[1]	PLCB_PSIGN4_SBF1	PLCB_PLOC1_SBF2[3]
89	PLCB_PLOC1_SBF1[2]	PLCB_PSIGN5_SBF1	PLCB_PLOC2_SBF2[0]
90	PLCB_PLOC1_SBF1[3]	PLCB_PLOC1_SBF1[0]	PLCB_PLOC2_SBF2[1]
91	PLCB_PLOC2_SBF1[0]	PLCB_PLOC1_SBF1[1]	PLCB_PLOC2_SBF2[2]
92	PLCB_PLOC2_SBF1[1]	PLCB_PLOC1_SBF1[2]	PLCB_PLOC3_SBF2[0]
93	PLCB_PLOC2_SBF1[2]	PLCB_PLOC2_SBF1[0]	PLCB_PLOC3_SBF2[1]
94	PLCB_PLOC3_SBF1[0]	PLCB_PLOC2_SBF1[1]	PLCB_PLOC3_SBF2[2]
95	PLCB_PLOC3_SBF1[1]	PLCB_PLOC2_SBF1[2]	PLCB_PLOC4_SBF2[0]
96	PLCB_PLOC3_SBF1[2]	PLCB_PLOC3_SBF1[0]	PLCB_PLOC4_SBF2[1]
97	PLCB_PLOC3_SBF1[3]	PLCB_PLOC3_SBF1[1]	PLCB_PLOC4_SBF2[2]
98	PLCB_PLOC4_SBF1[0]	PLCB_PLOC3_SBF1[2]	PLCB_PLOC5_SBF2[0]
99	PLCB_PLOC4_SBF1[1]	PLCB_PLOC4_SBF1[0]	PLCB_PLOC5_SBF2[1]
100	PLCB_PLOC4_SBF1[2]	PLCB_PLOC4_SBF1[1]	PLCB_PLOC5_SBF2[2]
101	PLCB_PLOC5_SBF1[0]	PLCB_PLOC4_SBF1[2]	PLCB_PLOC5_SBF2[3]
102	PLCB_PLOC5_SBF1[1]	PLCB_PLOC5_SBF1[0]	PLCB_PLOC6_SBF2[0]
103	PLCB_PLOC5_SBF1[2]	PLCB_PLOC5_SBF1[1]	PLCB_PLOC6_SBF2[1]

104	PLCB_PLOC5_SBF1[3]	PLCB_PLOC5_SBF1[2]	PLCB_PLOC6_SBF2[2]
105	PLCB_SLCT1_SBF2=1	PLCB_SLCT1_SBF2=0	PLCB_PLOC7_SBF2[0]
106	PLCB_PSIGN13_SBF2	PLCB_SLCT2_SBF2	PLCB_PLOC7_SBF2[1]
107	PLCB_PSIGN24_SBF2	PLCB_PSIGN1_SBF2	PLCB_PLOC7_SBF2[2]
108	PLCB_PSIGN5_SBF2	PLCB_PSIGN2_SBF2	PLCB_PLOC8_SBF2[0]
109	PLCB_PLOC1_SBF2[0]	PLCB_PSIGN3_SBF2	PLCB_PLOC8_SBF2[1]
110	PLCB_PLOC1_SBF2[1]	PLCB_PSIGN4_SBF2	PLCB_PLOC8_SBF2[2]
111	PLCB_PLOC1_SBF2[2]	PLCB_PSIGN5_SBF2	PLCB_PSIGN1_5_SBF3
112	PLCB_PLOC1_SBF2[3]	PLCB_PLOC1_SBF2[0]	PLCB_PSIGN2_6_SBF3
113	PLCB_PLOC2_SBF2[0]	PLCB_PLOC1_SBF2[1]	PLCB_PSIGN3_7_SBF3
114	PLCB_PLOC2_SBF2[1]	PLCB_PLOC1_SBF2[2]	PLCB_PSIGN4_8_SBF3
115	PLCB_PLOC2_SBF2[2]	PLCB_PLOC2_SBF2[0]	PLCB_PLOC1_SBF3[0]
116	PLCB_PLOC3_SBF2[0]	PLCB_PLOC2_SBF2[1]	PLCB_PLOC1_SBF3[1]
117	PLCB_PLOC3_SBF2[1]	PLCB_PLOC2_SBF2[2]	PLCB_PLOC1_SBF3[2]
118	PLCB_PLOC3_SBF2[2]	PLCB_PLOC3_SBF2[0]	PLCB_PLOC1_SBF3[3]
119	PLCB_PLOC3_SBF2[3]	PLCB_PLOC3_SBF2[1]	PLCB_PLOC2_SBF3[0]
120	PLCB_PLOC4_SBF2[0]	PLCB_PLOC3_SBF2[2]	PLCB_PLOC2_SBF3[1]
121	PLCB_PLOC4_SBF2[1]	PLCB_PLOC4_SBF2[0]	PLCB_PLOC2_SBF3[2]
122	PLCB_PLOC4_SBF2[2]	PLCB_PLOC4_SBF2[1]	PLCB_PLOC3_SBF3[0]
123	PLCB_PLOC5_SBF2[0]	PLCB_PLOC4_SBF2[2]	PLCB_PLOC3_SBF3[1]
124	PLCB_PLOC5_SBF2[1]	PLCB_PLOC5_SBF2[0]	PLCB_PLOC3_SBF3[2]
125	PLCB_PLOC5_SBF2[2]	PLCB_PLOC5_SBF2[1]	PLCB_PLOC4_SBF3[0]
126	PLCB_PLOC5_SBF2[3]	PLCB_PLOC5_SBF2[2]	PLCB_PLOC4_SBF3[1]
127	PLCB_SLCT1_SBF3=1	PLCB_SLCT1_SBF3=0	PLCB_PLOC4_SBF3[2]
128	PLCB_PSIGN13_SBF3	PLCB_SLCT2_SBF3	PLCB_PLOC5_SBF3[0]
129	PLCB_PSIGN24_SBF3	PLCB_PSIGN1_SBF3	PLCB_PLOC5_SBF3[1]
130	PLCB_PSIGN5_SBF3	PLCB_PSIGN2_SBF3	PLCB_PLOC5_SBF3[2]
131	PLCB_PLOC1_SBF3[0]	PLCB_PSIGN3_SBF3	PLCB_PLOC5_SBF3[3]
132	PLCB_PLOC1_SBF3[1]	PLCB_PSIGN4_SBF3	PLCB_PLOC6_SBF3[0]
133	PLCB_PLOC1_SBF3[2]	PLCB_PSIGN5_SBF3	PLCB_PLOC6_SBF3[1]
134	PLCB_PLOC1_SBF3[3]	PLCB_PLOC1_SBF3[0]	PLCB_PLOC6_SBF3[2]
135	PLCB_PLOC2_SBF3[0]	PLCB_PLOC1_SBF3[1]	PLCB_PLOC7_SBF3[0]
136	PLCB_PLOC2_SBF3[1]	PLCB_PLOC1_SBF3[2]	PLCB_PLOC7_SBF3[1]
137	PLCB_PLOC2_SBF3[2]	PLCB_PLOC2_SBF3[0]	PLCB_PLOC7_SBF3[2]
138	PLCB_PLOC3_SBF3[0]	PLCB_PLOC2_SBF3[1]	PLCB_PLOC8_SBF3[0]
139	PLCB_PLOC3_SBF3[1]	PLCB_PLOC2_SBF3[2]	PLCB_PLOC8_SBF3[1]
140	PLCB_PLOC3_SBF3[2]	PLCB_PLOC3_SBF3[0]	PLCB_PLOC8_SBF3[2]

141	PLCB_PLOC3_SBF3[3]	PLCB_PLOC3_SBF3[1]	PLCB_PSIGN1_5_SBF4
142	PLCB_PLOC4_SBF3[0]	PLCB_PLOC3_SBF3[2]	PLCB_PSIGN2_6_SBF4
143	PLCB_PLOC4_SBF3[1]	PLCB_PLOC4_SBF3[0]	PLCB_PSIGN3_7_SBF4
144	PLCB_PLOC4_SBF3[2]	PLCB_PLOC4_SBF3[1]	PLCB_PSIGN4_8_SBF4
145	PLCB_PLOC5_SBF3[0]	PLCB_PLOC4_SBF3[2]	PLCB_PLOC1_SBF4[0]
146	PLCB_PLOC5_SBF3[1]	PLCB_PLOC5_SBF3[0]	PLCB_PLOC1_SBF4[1]
147	PLCB_PLOC5_SBF3[2]	PLCB_PLOC5_SBF3[1]	PLCB_PLOC1_SBF4[2]
148	PLCB_PLOC5_SBF3[3]	PLCB_PLOC5_SBF3[2]	PLCB_PLOC1_SBF4[3]
149	PLCB_SLCT1_SBF4=1	PLCB_SLCT1_SBF4=0	PLCB_PLOC2_SBF4[0]
150	PLCB_PSIGN13_SBF4	PLCB_SLCT2_SBF4	PLCB_PLOC2_SBF4[1]
151	PLCB_PSIGN24_SBF4	PLCB_PSIGN1_SBF4	PLCB_PLOC2_SBF4[2]
152	PLCB_PSIGN5_SBF4	PLCB_PSIGN2_SBF4	PLCB_PLOC3_SBF4[0]
153	PLCB_PLOC1_SBF4[0]	PLCB_PSIGN3_SBF4	PLCB_PLOC3_SBF4[1]
154	PLCB_PLOC1_SBF4[1]	PLCB_PSIGN4_SBF4	PLCB_PLOC3_SBF4[2]
155	PLCB_PLOC1_SBF4[2]	PLCB_PSIGN5_SBF4	PLCB_PLOC4_SBF4[0]
156	PLCB_PLOC1_SBF4[3]	PLCB_PLOC1_SBF4[0]	PLCB_PLOC4_SBF4[1]
157	PLCB_PLOC2_SBF4[0]	PLCB_PLOC1_SBF4[1]	PLCB_PLOC4_SBF4[2]
158	PLCB_PLOC2_SBF4[1]	PLCB_PLOC1_SBF4[2]	PLCB_PLOC5_SBF4[0]
159	PLCB_PLOC2_SBF4[2]	PLCB_PLOC2_SBF4[0]	PLCB_PLOC5_SBF4[1]
160	PLCB_PLOC3_SBF4[0]	PLCB_PLOC2_SBF4[1]	PLCB_PLOC5_SBF4[2]
161	PLCB_PLOC3_SBF4[1]	PLCB_PLOC2_SBF4[2]	PLCB_PLOC5_SBF4[3]
162	PLCB_PLOC3_SBF4[2]	PLCB_PLOC3_SBF4[0]	PLCB_PLOC6_SBF4[0]
163	PLCB_PLOC3_SBF4[3]	PLCB_PLOC3_SBF4[1]	PLCB_PLOC6_SBF4[1]
164	PLCB_PLOC4_SBF4[0]	PLCB_PLOC3_SBF4[2]	PLCB_PLOC6_SBF4[2]
165	PLCB_PLOC4_SBF4[1]	PLCB_PLOC4_SBF4[0]	PLCB_PLOC7_SBF4[0]
166	PLCB_PLOC4_SBF4[2]	PLCB_PLOC4_SBF4[1]	PLCB_PLOC7_SBF4[1]
167	PLCB_PLOC5_SBF4[0]	PLCB_PLOC4_SBF4[2]	PLCB_PLOC7_SBF4[2]
168	PLCB_PLOC5_SBF4[1]	PLCB_PLOC5_SBF4[0]	PLCB_PLOC8_SBF4[0]
169	PLCB_PLOC5_SBF4[2]	PLCB_PLOC5_SBF4[1]	PLCB_PLOC8_SBF4[1]
170	PLCB_PLOC5_SBF4[3]	PLCB_PLOC5_SBF4[2]	PLCB_PLOC8_SBF4[2]
171	Rate Sanity Flag = 0		

1

2

3

4

5

6

Table 4.3-1: Bit Stream for Rate 1

1

Bit Index	Rate 1/2 Frame Type 0 Packet Bits
1	SVS_DECI=0
2	LSF_STAGE1[0]
3	LSF_STAGE1[1]
4	LSF_STAGE1[2]
5	LSF_STAGE1[3]
6	LSF_STAGE1[4]
7	LSF_STAGE1[5]
8	LSF_STAGE1[6]
9	LSF_STAGE2[0]
10	LSF_STAGE2[1]
11	LSF_STAGE2[2]
12	LSF_STAGE2[3]
13	LSF_STAGE2[4]
14	LSF_STAGE2[5]
15	LSF_STAGE2[6]
16	LSF_STAGE3[0]
17	LSF_STAGE3[1]
18	LSF_STAGE3[2]
19	LSF_STAGE3[3]
20	LSF_STAGE3[4]
21	LSF_STAGE3[5]
22	LSF_SWTCHPRD
23	PITC_SBF1[0]
24	PITC_SBF1[1]
25	PITC_SBF1[2]
26	PITC_SBF1[3]
27	PITC_SBF1[4]
28	PITC_SBF1[5]
29	PITC_SBF1[6]
30	PITC_SBF2[0]
31	PITC_SBF2[1]
32	PITC_SBF2[2]
33	PITC_SBF2[3]

34	PITC_SBF2[4]		
35	PITC_SBF2[5]		
36	PITC_SBF2[6]		
37	GAIN_SBF1[0]		
38	GAIN_SBF1[1]		
39	GAIN_SBF1[2]		
40	GAIN_SBF1[3]		
41	GAIN_SBF1[4]		
42	GAIN_SBF1[5]		
43	GAIN_SBF1[6]		
44	GAIN_SBF2[0]		
45	GAIN_SBF2[1]		
46	GAIN_SBF2[2]		
47	GAIN_SBF2[3]		
48	GAIN_SBF2[4]		
49	GAIN_SBF2[5]		
50	GAIN_SBF2[6]		
51	PLCB_SLCT1_SBF1=1	PLCB_SLCT1_SBF1=0	
52	PLCB_PSIGN1_2_SBF1	PLCB_SLCT2_SBF1=1	PLCB_SLCT2_SBF1=0
53	PLCB_PLOC1_2_SBF1[0]	PLCB_PSIGN1_SBF1	GCB_VSIGN1_SBF1
54	PLCB_PLOC1_2_SBF1[1]	PLCB_PSIGN2_SBF1	GCB_VSIGN2_SBF1
55	PLCB_PLOC1_2_SBF1[2]	PLCB_PSIGN3_SBF1	GCB_VEC1_2_SBF1[0]
56	PLCB_PLOC1_2_SBF1[3]	PLCB_PLOC1_SBF1[0]	GCB_VEC1_2_SBF1[1]
57	PLCB_PLOC1_2_SBF1[4]	PLCB_PLOC1_SBF1[1]	GCB_VEC1_2_SBF1[2]
58	PLCB_PLOC1_2_SBF1[5]	PLCB_PLOC1_SBF1[2]	GCB_VEC1_2_SBF1[3]
59	PLCB_PLOC1_2_SBF1[6]	PLCB_PLOC1_SBF1[3]	GCB_VEC1_2_SBF1[4]
60	PLCB_PLOC1_2_SBF1[7]	PLCB_PLOC2_SBF1[0]	GCB_VEC1_2_SBF1[5]
61	PLCB_PLOC1_2_SBF1[8]	PLCB_PLOC2_SBF1[1]	GCB_VEC1_2_SBF1[6]
62	PLCB_PLOC1_2_SBF1[9]	PLCB_PLOC2_SBF1[2]	GCB_VEC1_2_SBF1[7]
63	PLCB_PLOC1_2_SBF1[10]	PLCB_PLOC3_SBF1[0]	GCB_VEC1_2_SBF1[8]
64	PLCB_PLOC1_2_SBF1[11]	PLCB_PLOC3_SBF1[1]	GCB_VEC1_2_SBF1[9]
65	PLCB_PLOC1_2_SBF1[12]	PLCB_PLOC3_SBF1[2]	GCB_VEC1_2_SBF1[10]
66	PLCB_SLCT1_SBF2=1	PLCB_SLCT1_SBF2=0	
67	PLCB_PSIGN1_2_SBF2	PLCB_SLCT2_SBF2=1	PLCB_SLCT2_SBF2=0
68	PLCB_PLOC1_2_SBF2[0]	PLCB_PSIGN1_SBF2	GCB_VSIGN1_SBF2
69	PLCB_PLOC1_2_SBF2[1]	PLCB_PSIGN2_SBF2	GCB_VSIGN2_SBF2
70	PLCB_PLOC1_2_SBF2[2]	PLCB_PSIGN3_SBF2	GCB_VEC1_2_SBF2[0]

71	PLCB_PLOC1_2_SBF2[3]	PLCB_PLOC1_SBF2[0]	GCB_VEC1_2_SBF2[1]
72	PLCB_PLOC1_2_SBF2[4]	PLCB_PLOC1_SBF2[1]	GCB_VEC1_2_SBF2[2]
73	PLCB_PLOC1_2_SBF2[5]	PLCB_PLOC1_SBF2[2]	GCB_VEC1_2_SBF2[3]
74	PLCB_PLOC1_2_SBF2[6]	PLCB_PLOC1_SBF2[3]	GCB_VEC1_2_SBF2[4]
75	PLCB_PLOC1_2_SBF2[7]	PLCB_PLOC2_SBF2[0]	GCB_VEC1_2_SBF2[5]
76	PLCB_PLOC1_2_SBF2[8]	PLCB_PLOC2_SBF2[1]	GCB_VEC1_2_SBF2[6]
77	PLCB_PLOC1_2_SBF2[9]	PLCB_PLOC2_SBF2[2]	GCB_VEC1_2_SBF2[7]
78	PLCB_PLOC1_2_SBF2[10]	PLCB_PLOC3_SBF2[0]	GCB_VEC1_2_SBF2[8]
79	PLCB_PLOC1_2_SBF2[11]	PLCB_PLOC3_SBF2[1]	GCB_VEC1_2_SBF2[9]
80	PLCB_PLOC1_2_SBF2[12]	PLCB_PLOC3_SBF2[2]	GCB_VEC1_2_SBF2[10]

1
2
3
4

Table 4.3-2: Bit Stream for Rate 1/2 Type-0 Frames

Bit Index	Rate 1/2 Frame Type 1 Packet Bits
1	SVS_DECI=1
2	LSF_STAGE1[0]
3	LSF_STAGE1[1]
4	LSF_STAGE1[2]
5	LSF_STAGE1[3]
6	LSF_STAGE1[4]
7	LSF_STAGE1[5]
8	LSF_STAGE1[6]
9	LSF_STAGE2[0]
10	LSF_STAGE2[1]
11	LSF_STAGE2[2]
12	LSF_STAGE2[3]
13	LSF_STAGE2[4]
14	LSF_STAGE2[5]
15	LSF_STAGE2[6]
16	LSF_STAGE3[0]
17	LSF_STAGE3[1]
18	LSF_STAGE3[2]
19	LSF_STAGE3[3]
20	LSF_STAGE3[4]

21	LSF_STAGE3[5]		
22	LSF_SWCHPRD		
23	PITCH_FRM[0]		
24	PITCH_FRM[1]		
25	PITCH_FRM[2]		
26	PITCH_FRM[3]		
27	PITCH_FRM[4]		
28	PITCH_FRM[5]		
29	PITCH_FRM[6]		
30	PITCH_VQ_GAIN[0]		
31	PITCH_VQ_GAIN[1]		
32	PITCH_VQ_GAIN[2]		
33	PITCH_VQ_GAIN[3]		
34	PLCB_VQ_GAIN[0]		
35	PLCB_VQ_GAIN[1]		
36	PLCB_VQ_GAIN[2]		
37	PLCB_VQ_GAIN[3]		
38	PLCB_VQ_GAIN[4]		
39	PLCB_VQ_GAIN[5]		
40	PLCB_VQ_GAIN[6]		
41	PLCB_VQ_GAIN[7]		
42	PLCB_SLCT1_SBF1=1	PLCB_SLCT1_SBF1=0	
43	PLCB_PSIGN1_SBF1	PLCB_SLCT2_SBF1=1	PLCB_SLCT2_SBF1=0
44	PLCB_PSIGN2_SBF1	PLCB_PSIGN1_SBF1	PLCB_PSIGN1_SBF1
45	PLCB_PLOC1_SBF1[0]	PLCB_PSIGN2_SBF1	PLCB_PSIGN2_SBF1
46	PLCB_PLOC1_SBF1[1]	PLCB_PSIGN3_SBF1	PLCB_PSIGN3_SBF1
47	PLCB_PLOC1_SBF1[2]	PLCB_PLOC1_SBF1[0]	PLCB_PSIGN4_SBF1
48	PLCB_PLOC1_SBF1[3]	PLCB_PLOC1_SBF1[1]	PLCB_PSIGN5_SBF1
49	PLCB_PLOC1_SBF1[4]	PLCB_PLOC1_SBF1[2]	PLCB_PLOC1_SBF1[0]
50	PLCB_PLOC2_SBF1[0]	PLCB_PLOC1_SBF1[3]	PLCB_PLOC1_SBF1[1]
51	PLCB_PLOC2_SBF1[1]	PLCB_PLOC2_SBF1[0]	PLCB_PLOC2_SBF1[0]
52	PLCB_PLOC2_SBF1[2]	PLCB_PLOC2_SBF1[1]	PLCB_PLOC3_SBF1[0]
53	PLCB_PLOC2_SBF1[3]	PLCB_PLOC3_SBF1[0]	PLCB_PLOC4_SBF1[0]
54	PLCB_PLOC2_SBF1[4]	PLCB_PLOC3_SBF1[1]	PLCB_PLOC5_SBF1[0]
55	PLCB_SLCT1_SBF2=1	PLCB_SLCT1_SBF2=0	
56	PLCB_PSIGN1_SBF2	PLCB_SLCT2_SBF2=1	PLCB_SLCT2_SBF2=0
57	PLCB_PSIGN2_SBF2	PLCB_PSIGN1_SBF2	PLCB_PSIGN1_SBF2

58	PLCB_PLOC1_SBF2[0]	PLCB_PSIGN2_SBF2	PLCB_PSIGN2_SBF2
59	PLCB_PLOC1_SBF2[1]	PLCB_PSIGN3_SBF2	PLCB_PSIGN3_SBF2
60	PLCB_PLOC1_SBF2[2]	PLCB_PLOC1_SBF2[0]	PLCB_PSIGN4_SBF2
61	PLCB_PLOC1_SBF2[3]	PLCB_PLOC1_SBF2[1]	PLCB_PSIGN5_SBF2
62	PLCB_PLOC1_SBF2[4]	PLCB_PLOC1_SBF2[2]	PLCB_PLOC1_SBF2[0]
63	PLCB_PLOC2_SBF2[0]	PLCB_PLOC1_SBF2[3]	PLCB_PLOC1_SBF2[1]
64	PLCB_PLOC2_SBF2[1]	PLCB_PLOC2_SBF2[0]	PLCB_PLOC2_SBF2[0]
65	PLCB_PLOC2_SBF2[2]	PLCB_PLOC2_SBF2[1]	PLCB_PLOC3_SBF2[0]
66	PLCB_PLOC2_SBF2[3]	PLCB_PLOC3_SBF2[0]	PLCB_PLOC4_SBF2[0]
67	PLCB_PLOC2_SBF2[4]	PLCB_PLOC3_SBF2[1]	PLCB_PLOC5_SBF2[0]
68	PLCB_SLCT1_SBF3=1	PLCB_SLCT1_SBF3=0	
69	PLCB_PSIGN1_SBF3	PLCB_SLCT2_SBF3=1	PLCB_SLCT2_SBF3=0
70	PLCB_PSIGN2_SBF3	PLCB_PSIGN1_SBF3	PLCB_PSIGN1_SBF3
71	PLCB_PLOC1_SBF3[0]	PLCB_PSIGN2_SBF3	PLCB_PSIGN2_SBF3
72	PLCB_PLOC1_SBF3[1]	PLCB_PSIGN3_SBF3	PLCB_PSIGN3_SBF3
73	PLCB_PLOC1_SBF3[2]	PLCB_PLOC1_SBF3[0]	PLCB_PSIGN4_SBF3
74	PLCB_PLOC1_SBF3[3]	PLCB_PLOC1_SBF3[1]	PLCB_PSIGN5_SBF3
75	PLCB_PLOC1_SBF3[4]	PLCB_PLOC1_SBF3[2]	PLCB_PLOC1_SBF3[0]
76	PLCB_PLOC2_SBF3[0]	PLCB_PLOC1_SBF3[3]	PLCB_PLOC1_SBF3[1]
77	PLCB_PLOC2_SBF3[1]	PLCB_PLOC2_SBF3[0]	PLCB_PLOC2_SBF3[0]
78	PLCB_PLOC2_SBF3[2]	PLCB_PLOC2_SBF3[1]	PLCB_PLOC3_SBF3[0]
79	PLCB_PLOC2_SBF3[3]	PLCB_PLOC3_SBF3[0]	PLCB_PLOC4_SBF3[0]
80	PLCB_PLOC2_SBF3[4]	PLCB_PLOC3_SBF3[1]	PLCB_PLOC5_SBF3[0]

1
2
3
4
5

Table 4.3-3: Bit Stream for Rate 1/2 Type-1 Frames

1

Bit Index	Rate 1/4 Packet Bits
1	LSF_STAGE1[0]
2	LSF_STAGE1[1]
3	LSF_STAGE1[2]
4	LSF_STAGE1[3]
5	LSF_STAGE1[4]
6	LSF_STAGE1[5]
7	LSF_STAGE1[6]
8	LSF_STAGE2[0]
9	LSF_STAGE2[1]
10	LSF_STAGE2[2]
11	LSF_STAGE2[3]
12	LSF_STAGE2[4]
13	LSF_STAGE2[5]
14	LSF_STAGE2[6]
15	LSF_STAGE3[0]
16	LSF_STAGE3[1]
17	LSF_STAGE3[2]
18	LSF_STAGE3[3]
19	LSF_STAGE3[4]
20	LSF_STAGE3[5]
21	GAIN_VQ[0]
22	GAIN_VQ[1]
23	GAIN_VQ[2]
24	GAIN_VQ[3]
25	GAIN_VQ[4]
26	GAIN_VQ_1[0]
27	GAIN_VQ_1[1]
28	GAIN_VQ_1[2]
29	GAIN_VQ_1[3]
30	GAIN_VQ_1[4]
31	GAIN_VQ_1[5]
32	GAIN_VQ_2[0]
33	GAIN_VQ_2[1]
34	GAIN_VQ_2[2]

35	GAIN_VQ_2[3]
36	GAIN_VQ_2[4]
37	GAIN_VQ_2[5]
38	SF_IDX[0]
39	SF_IDX[1]
40	Rate Sanity Flag = 0

1

2

3

4

Table 4.3-4: Bit Stream for Rate 1/4

1

Bit Index	Rate 1/8 Packet Bits
1	LSF_STAGE1[0]
2	LSF_STAGE1[1]
3	LSF_STAGE1[2]
4	LSF_STAGE1[3]
5	LSF_STAGE2[0]
6	LSF_STAGE2[1]
7	LSF_STAGE2[2]
8	LSF_STAGE2[3]
9	LSF_STAGE3[0]
10	LSF_STAGE3[1]
11	LSF_STAGE3[2]
12	GAIN_08K[0]
13	GAIN_08K[1]
14	GAIN_08K[2]
15	GAIN_08K[3]
16	GAIN_08K[4]

2

3

4

5

Table 4.3-5: Bit Stream for Rate 1/8

1

Bit Index	Rate 1 Data + 1/2 Rate Speech Packet Bits	Rate 1 Data + Auxiliary Data Packet Bits
1	SVS_DECI = 0	SVS_DECI = 0
2	HALF_RATE_SPEECH[0]	AUX_DATA[0]
3	HALF_RATE_SPEECH [1]	AUX_DATA [1]
4	HALF_RATE_SPEECH [2]	AUX_DATA [2]
5	HALF_RATE_SPEECH [3]	AUX_DATA [3]
6	HALF_RATE_SPEECH [4]	AUX_DATA [4]
7	HALF_RATE_SPEECH [5]	AUX_DATA [5]
8	HALF_RATE_SPEECH [6]	AUX_DATA [6]
9	HALF_RATE_SPEECH [7]	AUX_DATA [7]
10	HALF_RATE_SPEECH [8]	AUX_DATA [8]
11	HALF_RATE_SPEECH [9]	AUX_DATA [9]
12	HALF_RATE_SPEECH [10]	AUX_DATA [10]
13	HALF_RATE_SPEECH [11]	AUX_DATA [11]
14	HALF_RATE_SPEECH [12]	AUX_DATA [12]
15	HALF_RATE_SPEECH [13]	AUX_DATA [13]
16	HALF_RATE_SPEECH [14]	AUX_DATA [14]
17	HALF_RATE_SPEECH [15]	AUX_DATA [15]
18	HALF_RATE_SPEECH [16]	AUX_DATA [16]
19	HALF_RATE_SPEECH [17]	AUX_DATA [17]
20	HALF_RATE_SPEECH [18]	AUX_DATA [18]
21	HALF_RATE_SPEECH [19]	AUX_DATA [19]
22	HALF_RATE_SPEECH [20]	AUX_DATA [20]
23	HALF_RATE_SPEECH [21]	AUX_DATA [21]
24	HALF_RATE_SPEECH [22]	AUX_DATA [22]
25	HALF_RATE_SPEECH [23]	AUX_DATA [23]
26	HALF_RATE_SPEECH [24]	AUX_DATA [24]
27	HALF_RATE_SPEECH [25]	AUX_DATA [25]
28	HALF_RATE_SPEECH [26]	AUX_DATA [26]
29	HALF_RATE_SPEECH [27]	AUX_DATA [27]
30	HALF_RATE_SPEECH [28]	AUX_DATA [28]
31	HALF_RATE_SPEECH [29]	AUX_DATA [29]
32	HALF_RATE_SPEECH [30]	AUX_DATA [30]
33	HALF_RATE_SPEECH [31]	AUX_DATA [31]
34	HALF_RATE_SPEECH [32]	AUX_DATA [32]

35	HALF_RATE_SPEECH [33]	AUX_DATA [33]
36	HALF_RATE_SPEECH [34]	AUX_DATA [34]
37	HALF_RATE_SPEECH [35]	AUX_DATA [35]
38	HALF_RATE_SPEECH [36]	AUX_DATA [36]
39	HALF_RATE_SPEECH [37]	AUX_DATA [37]
40	HALF_RATE_SPEECH [38]	AUX_DATA [38]
41	HALF_RATE_SPEECH [39]	AUX_DATA [39]
42	HALF_RATE_SPEECH [40]	AUX_DATA [40]
43	HALF_RATE_SPEECH [41]	AUX_DATA [41]
44	HALF_RATE_SPEECH [42]	AUX_DATA [42]
45	HALF_RATE_SPEECH [43]	AUX_DATA [43]
46	HALF_RATE_SPEECH [44]	AUX_DATA [44]
47	HALF_RATE_SPEECH [45]	AUX_DATA [45]
48	HALF_RATE_SPEECH [46]	AUX_DATA [46]
49	HALF_RATE_SPEECH [47]	AUX_DATA [47]
50	HALF_RATE_SPEECH [48]	AUX_DATA [48]
51	HALF_RATE_SPEECH [49]	AUX_DATA [49]
52	HALF_RATE_SPEECH [50]	AUX_DATA [50]
53	HALF_RATE_SPEECH [51]	AUX_DATA [51]
54	HALF_RATE_SPEECH [52]	AUX_DATA [52]
55	HALF_RATE_SPEECH [53]	AUX_DATA [53]
56	HALF_RATE_SPEECH [54]	AUX_DATA [54]
57	HALF_RATE_SPEECH [55]	AUX_DATA [55]
58	HALF_RATE_SPEECH [56]	AUX_DATA [56]
59	HALF_RATE_SPEECH [57]	AUX_DATA [57]
60	HALF_RATE_SPEECH [58]	AUX_DATA [58]
61	HALF_RATE_SPEECH [59]	AUX_DATA [59]
62	HALF_RATE_SPEECH [60]	AUX_DATA [60]
63	HALF_RATE_SPEECH [61]	AUX_DATA [61]
64	HALF_RATE_SPEECH [62]	AUX_DATA [62]
65	HALF_RATE_SPEECH [63]	AUX_DATA [63]
66	HALF_RATE_SPEECH [64]	AUX_DATA [64]
67	HALF_RATE_SPEECH [65]	AUX_DATA [65]
68	HALF_RATE_SPEECH [66]	AUX_DATA [66]
69	HALF_RATE_SPEECH [67]	AUX_DATA [67]
70	HALF_RATE_SPEECH [68]	AUX_DATA [68]
71	HALF_RATE_SPEECH [69]	AUX_DATA [69]

72	HALF_RATE_SPEECH [70]	AUX_DATA [70]
73	HALF_RATE_SPEECH [71]	AUX_DATA [71]
74	HALF_RATE_SPEECH [72]	AUX_DATA [72]
75	HALF_RATE_SPEECH [73]	AUX_DATA [73]
76	HALF_RATE_SPEECH [74]	AUX_DATA [74]
77	HALF_RATE_SPEECH [75]	AUX_DATA [75]
78	HALF_RATE_SPEECH [76]	AUX_DATA [76]
79	HALF_RATE_SPEECH [77]	AUX_DATA [77]
80	HALF_RATE_SPEECH [78]	AUX_DATA [78]
81	HALF_RATE_SPEECH [79]	AUX_DATA [79]
82	DATA_HEADER[0]	DATA_HEADER[0]
83	DATA_HEADER[1]	DATA_HEADER[1]
84	DATA_HEADER[2]	DATA_HEADER[2]
85	DATA_BLOCK[0]	DATA_BLOCK[0]
86	DATA_BLOCK [1]	DATA_BLOCK [1]
87	DATA_BLOCK [2]	DATA_BLOCK [2]
88	DATA_BLOCK [3]	DATA_BLOCK [3]
89	DATA_BLOCK [4]	DATA_BLOCK [4]
90	DATA_BLOCK [5]	DATA_BLOCK [5]
91	DATA_BLOCK [6]	DATA_BLOCK [6]
92	DATA_BLOCK [7]	DATA_BLOCK [7]
93	DATA_BLOCK [8]	DATA_BLOCK [8]
94	DATA_BLOCK [9]	DATA_BLOCK [9]
95	DATA_BLOCK[10]	DATA_BLOCK[10]
96	DATA_BLOCK [11]	DATA_BLOCK [11]
97	DATA_BLOCK [12]	DATA_BLOCK [12]
98	DATA_BLOCK [13]	DATA_BLOCK [13]
99	DATA_BLOCK [14]	DATA_BLOCK [14]
100	DATA_BLOCK [15]	DATA_BLOCK [15]
101	DATA_BLOCK [16]	DATA_BLOCK [16]
102	DATA_BLOCK [17]	DATA_BLOCK [17]
103	DATA_BLOCK [18]	DATA_BLOCK [18]
104	DATA_BLOCK [19]	DATA_BLOCK [19]
105	DATA_BLOCK[20]	DATA_BLOCK[20]
106	DATA_BLOCK [21]	DATA_BLOCK [21]
107	DATA_BLOCK [22]	DATA_BLOCK [22]
108	DATA_BLOCK [23]	DATA_BLOCK [23]

109	DATA_BLOCK [24]	DATA_BLOCK [24]
110	DATA_BLOCK [25]	DATA_BLOCK [25]
111	DATA_BLOCK [26]	DATA_BLOCK [26]
112	DATA_BLOCK [27]	DATA_BLOCK [27]
113	DATA_BLOCK [28]	DATA_BLOCK [28]
114	DATA_BLOCK [29]	DATA_BLOCK [29]
115	DATA_BLOCK[30]	DATA_BLOCK[30]
116	DATA_BLOCK [31]	DATA_BLOCK [31]
117	DATA_BLOCK [32]	DATA_BLOCK [32]
118	DATA_BLOCK [33]	DATA_BLOCK [33]
119	DATA_BLOCK [34]	DATA_BLOCK [34]
120	DATA_BLOCK [35]	DATA_BLOCK [35]
121	DATA_BLOCK [36]	DATA_BLOCK [36]
122	DATA_BLOCK [37]	DATA_BLOCK [37]
123	DATA_BLOCK [38]	DATA_BLOCK [38]
124	DATA_BLOCK [39]	DATA_BLOCK [39]
125	DATA_BLOCK[40]	DATA_BLOCK[40]
126	DATA_BLOCK [41]	DATA_BLOCK [41]
127	DATA_BLOCK [42]	DATA_BLOCK [42]
128	DATA_BLOCK [43]	DATA_BLOCK [43]
129	DATA_BLOCK [44]	DATA_BLOCK [44]
130	DATA_BLOCK [45]	DATA_BLOCK [45]
131	DATA_BLOCK [46]	DATA_BLOCK [46]
132	DATA_BLOCK [47]	DATA_BLOCK [47]
133	DATA_BLOCK [48]	DATA_BLOCK [48]
134	DATA_BLOCK [49]	DATA_BLOCK [49]
135	DATA_BLOCK[50]	DATA_BLOCK[50]
136	DATA_BLOCK [51]	DATA_BLOCK [51]
137	DATA_BLOCK [52]	DATA_BLOCK [52]
138	DATA_BLOCK [53]	DATA_BLOCK [53]
139	DATA_BLOCK [54]	DATA_BLOCK [54]
140	DATA_BLOCK [55]	DATA_BLOCK [55]
141	DATA_BLOCK [56]	DATA_BLOCK [56]
142	DATA_BLOCK [57]	DATA_BLOCK [57]
143	DATA_BLOCK [58]	DATA_BLOCK [58]
144	DATA_BLOCK [59]	DATA_BLOCK [59]
145	DATA_BLOCK[60]	DATA_BLOCK[60]

146	DATA_BLOCK [61]	DATA_BLOCK [61]
147	DATA_BLOCK [62]	DATA_BLOCK [62]
148	DATA_BLOCK [63]	DATA_BLOCK [63]
149	ILLEGAL_CODEWORD[0] = 1	ILLEGAL_CODEWORD[0] = 1
150	ILLEGAL_CODEWORD[1] = 0	ILLEGAL_CODEWORD[1] = 0
151	ILLEGAL_CODEWORD[2] = 0	ILLEGAL_CODEWORD[2] = 0
152	ILLEGAL_CODEWORD[3] = 0	ILLEGAL_CODEWORD[3] = 0
153	ILLEGAL_CODEWORD[4] = 0	ILLEGAL_CODEWORD[4] = 0
154	ILLEGAL_CODEWORD[5] = 0	ILLEGAL_CODEWORD[5] = 0
155	ILLEGAL_CODEWORD[6] = 0	ILLEGAL_CODEWORD[6] = 0
156	ILLEGAL_CODEWORD[7] = 0	ILLEGAL_CODEWORD[7] = 0
157	ILLEGAL_CODEWORD[8] = 0	ILLEGAL_CODEWORD[8] = 0
158	ILLEGAL_CODEWORD[9] = 0	ILLEGAL_CODEWORD[9] = 0
159	ILLEGAL_CODEWORD[10] = 0	ILLEGAL_CODEWORD[10] = 0
160	ILLEGAL_CODEWORD[11] = 0	ILLEGAL_CODEWORD[11] = 0
161	ILLEGAL_CODEWORD[12] = 1	ILLEGAL_CODEWORD[12] = 1
162	ILLEGAL_CODEWORD[13] = 0	ILLEGAL_CODEWORD[13] = 0
163	ILLEGAL_CODEWORD[14] = 1	ILLEGAL_CODEWORD[14] = 1
164	ILLEGAL_CODEWORD[15] = 1	ILLEGAL_CODEWORD[15] = 1
165	ILLEGAL_CODEWORD[16] = 1	ILLEGAL_CODEWORD[16] = 1
166	ILLEGAL_CODEWORD[17] = 0	ILLEGAL_CODEWORD[17] = 0
167	ILLEGAL_CODEWORD[18] = 1	ILLEGAL_CODEWORD[18] = 1
168	ILLEGAL_CODEWORD[19] = 0	ILLEGAL_CODEWORD[19] = 0
169	ILLEGAL_CODEWORD[20] = 1	ILLEGAL_CODEWORD[20] = 1
170	ILLEGAL_CODEWORD[21] = 0	ILLEGAL_CODEWORD[21] = 0
171	Rate Sanity Flag = 0	Rate Sanity Flag = 0

1

2

3

Table 4.3-6: Data Bit Streams for Rate 1

1

Bit Index	Rate 1/2 Half Rate Data Packet Bits	Rate 1/2 Half Rate TTY Data Packet Bits	Rate 1/2 Half Rate DTMF Data Packet Bits
1	SVS_DECI = 0	SVS_DECI = 0	SVS_DECI = 0
2	DATA_HEADER[0]	DATA_HEADER[0] = 0	DATA_HEADER[0] = 0
3	DATA_HEADER[1]	DATA_HEADER[1] = 1	DATA_HEADER[1] = 1
4	DATA_HEADER[2]	DATA_HEADER[2] = 0	DATA_HEADER[2] = 1
5	DATA_BLOCK[0]	TTY_TYPE[0]	DTMF_DIGIT[0]
6	DATA_BLOCK [1]	TTY_TYPE [1]	DTMF_DIGIT [1]
7	DATA_BLOCK [2]	TTY_TYPE [2]	DTMF_DIGIT [2]
8	DATA_BLOCK [3]	TTY_TYPE [3]	DTMF_DIGIT [3]
9	DATA_BLOCK [4]	TTY_TYPE [4]	DTMF_DIGIT [4]
10	DATA_BLOCK [5]	TTY_TYPE [5]	DTMF_DIGIT [5]
11	DATA_BLOCK [6]	TTY_TYPE [6]	DTMF_DIGIT [6]
12	DATA_BLOCK [7]	TTY_TYPE [7]	DTMF_DIGIT [7]
13	DATA_BLOCK [8]	TTY_HEADER[0]	DTMF_DIGIT [8]
14	DATA_BLOCK [9]	TTY_HEADER [1]	DTMF_DIGIT [9]
15	DATA_BLOCK[10]	TTY_HEADER [2]	DTMF_DIGIT[10]
16	DATA_BLOCK [11]	TTY_HEADER [3]	DTMF_DIGIT [11]
17	DATA_BLOCK [12]	TTY_HEADER [4]	DTMF_DIGIT [12]
18	DATA_BLOCK [13]	TTY_HEADER [5]	DTMF_DIGIT [13]
19	DATA_BLOCK [14]	TTY_HEADER [6]	DTMF_DIGIT [14]
20	DATA_BLOCK [15]	TTY_HEADER [7]	DTMF_DIGIT [15]
21	DATA_BLOCK [16]	TTY_CHAR[0]	RESERVED
22	DATA_BLOCK [17]	TTY_CHAR [1]	RESERVED
23	DATA_BLOCK [18]	TTY_CHAR [2]	RESERVED
24	DATA_BLOCK [19]	TTY_CHAR [3]	RESERVED
25	DATA_BLOCK[20]	TTY_CHAR [4]	RESERVED
26	DATA_BLOCK [21]	TTY_CHAR [5]	RESERVED
27	DATA_BLOCK [22]	TTY_CHAR [6]	RESERVED
28	DATA_BLOCK [23]	TTY_CHAR [7]	RESERVED
29	DATA_BLOCK [24]	TTY_CHAR [8]	RESERVED
30	DATA_BLOCK [25]	TTY_CHAR [9]	RESERVED
31	DATA_BLOCK [26]	TTY_CHAR[10]	RESERVED
32	DATA_BLOCK [27]	TTY_CHAR [11]	RESERVED
33	DATA_BLOCK [28]	TTY_CHAR [12]	RESERVED

34	DATA_BLOCK [29]	TTY_CHAR [13]	RESERVED
35	DATA_BLOCK[30]	TTY_CHAR [14]	RESERVED
36	DATA_BLOCK [31]	TTY_CHAR [15]	RESERVED
37	DATA_BLOCK [32]	TTY_RATE[0]	RESERVED
38	DATA_BLOCK [33]	TTY_RATE [1]	RESERVED
39	DATA_BLOCK [34]	TTY_RATE [2]	RESERVED
40	DATA_BLOCK [35]	TTY_RATE [3]	RESERVED
41	DATA_BLOCK [36]	TTY_RATE [4]	RESERVED
42	DATA_BLOCK [37]	TTY_RATE [5]	RESERVED
43	DATA_BLOCK [38]	TTY_RATE [6]	RESERVED
44	DATA_BLOCK [39]	TTY_RATE [7]	RESERVED
45	DATA_BLOCK[40]	TTY_RATE [8]	RESERVED
46	DATA_BLOCK [41]	TTY_RATE [9]	RESERVED
47	DATA_BLOCK [42]	TTY_RATE[10]	RESERVED
48	DATA_BLOCK [43]	TTY_RATE [11]	RESERVED
49	DATA_BLOCK [44]	TTY_RATE [12]	RESERVED
50	DATA_BLOCK [45]	TTY_RATE [13]	RESERVED
51	DATA_BLOCK [46]	TTY_RATE [14]	RESERVED
52	DATA_BLOCK [47]	TTY_RATE [15]	RESERVED
53	DATA_BLOCK [48]	RESERVED	RESERVED
54	DATA_BLOCK [49]	RESERVED	RESERVED
55	DATA_BLOCK[50]	RESERVED	RESERVED
56	DATA_BLOCK [51]	RESERVED	RESERVED
57	DATA_BLOCK [52]	RESERVED	RESERVED
58	DATA_BLOCK [53]	RESERVED	RESERVED
59	DATA_BLOCK [54]	RESERVED	RESERVED
60	DATA_BLOCK [55]	RESERVED	RESERVED
61	DATA_BLOCK [56]	RESERVED	RESERVED
62	DATA_BLOCK [57]	RESERVED	RESERVED
63	DATA_BLOCK [58]	RESERVED	RESERVED
64	DATA_BLOCK [59]	RESERVED	RESERVED
65	DATA_BLOCK[60]	RESERVED	RESERVED
66	ILLEGAL_CODEWORD[0] = 1	ILLEGAL_CODEWORD[0] = 1	ILLEGAL_CODEWORD[0] = 1
67	ILLEGAL_CODEWORD[1] = 1	ILLEGAL_CODEWORD[1] = 1	ILLEGAL_CODEWORD[1] = 1
68	ILLEGAL_CODEWORD[2] = 1	ILLEGAL_CODEWORD[2] = 1	ILLEGAL_CODEWORD[2] = 1
69	ILLEGAL_CODEWORD[3] = 1	ILLEGAL_CODEWORD[3] = 1	ILLEGAL_CODEWORD[3] = 1
70	ILLEGAL_CODEWORD[4] = 1	ILLEGAL_CODEWORD[4] = 1	ILLEGAL_CODEWORD[4] = 1

71	ILLEGAL_CODEWORD[5] = 1	ILLEGAL_CODEWORD[5] = 1	ILLEGAL_CODEWORD[5] = 1
72	ILLEGAL_CODEWORD[6] = 1	ILLEGAL_CODEWORD[6] = 1	ILLEGAL_CODEWORD[6] = 1
73	ILLEGAL_CODEWORD[7] = 1	ILLEGAL_CODEWORD[7] = 1	ILLEGAL_CODEWORD[7] = 1
74	ILLEGAL_CODEWORD[8] = 1	ILLEGAL_CODEWORD[8] = 1	ILLEGAL_CODEWORD[8] = 1
75	ILLEGAL_CODEWORD[9] = 0	ILLEGAL_CODEWORD[9] = 0	ILLEGAL_CODEWORD[9] = 0
76	ILLEGAL_CODEWORD[10] = 0	ILLEGAL_CODEWORD[10] = 0	ILLEGAL_CODEWORD[10] = 0
77	ILLEGAL_CODEWORD[11] = 1	ILLEGAL_CODEWORD[11] = 1	ILLEGAL_CODEWORD[11] = 1
78	ILLEGAL_CODEWORD[12] = 0	ILLEGAL_CODEWORD[12] = 0	ILLEGAL_CODEWORD[12] = 0
79	ILLEGAL_CODEWORD[13] = 0	ILLEGAL_CODEWORD[13] = 0	ILLEGAL_CODEWORD[13] = 0
80	ILLEGAL_CODEWORD[14] = 0	ILLEGAL_CODEWORD[14] = 0	ILLEGAL_CODEWORD[14] = 0

1
2
3
4
5
6
7
8
9
10
11

Table 4.3-7: Data Bit Streams for Rate 1/2

4.4 SMV Symbol Table

Description	Symbol	Comments
Half-rate max flag	F_{HRM}	
Selected SMV mode	M_{SMV}	
The frame counter	frm	
Size of encoding frame	L_{frm}	$L_{frm} = 160$
Size of subframe	L_{SF}	Subframe size depends on the coding rate and the frame type
Input speech	$s(n)$	
Output of silence enhancement	$s'(n)$	
Output of high-pass filter	$s''(n)$	

Output of noise suppression	$s^m(n)$	
Pre-processed speech	$\tilde{s}(n)$	
Size of LPC window	L_m^{LPC}	$m = 0,1,2$
LPC window	$W_m^{LPC}(n)$	$m = 0,1,2$
LPC Autocorrelation function	$R_m(i)$	$m = 0,1,2$
Order of autocorrelation function	O_R	10 for VAD-A, 16 for VAD-B
Prediction coefficients	$a_m(i)$	$m = 0,1,2$
Line spectral frequencies - LSF	$lsf_m(i)$	$m = 0,1,2$
Reflection coefficients	$k_m(i)$	$m = 0,1,2$
LPC prediction error	PE_m^{LPC}	$m = 0,1,2$
LPC prediction gain	G_m^{LPC}	$m = 0,1,2$
Interpolated LPC prediction gain	\hat{G}_m^{LPC}	$m = 0,1,2,3,4$
Interpolated reflection coefficients	$\hat{k}_m(i)$	$m = 0,1,2,3,4$
Interpolated prediction coefficients	$\hat{a}_m(i)$	$m = 0,1,2,3$ for Rate 1, 1/4, and 1/8 $m = 0,1$ for Rate 1/2 type 0 $m = 0,1,2$ for Rate 1/2 type 1
Quantized and interpolated prediction coefficients	$\hat{a}_m^q(i)$	$m = 0,1,2,3$ for Rate 1, 1/4, and 1/8 $m = 0,1$ for Rate 1/2 type 0 $m = 0,1,2$ for Rate 1/2 type 1
Buffer of previous open-loop pitch lags	$l_p^B(i)$	
Buffer of previous unquantized pitch gains	$g_p^B(i)$	
Buffer of previous normalized pitch correlation	$corr_p^B(i)$	
Voice activity decision	VAD	
Previous frame voice activity decision	VAD^{prv}	
Music detection flag	F_{MUS}	

Perceptual filter numerator weight	w_m^{num}	$m = 0,1,2,3,4$
Perceptual filter denominator weight	w_m^{den}	$m = 0,1,2,3,4$
Weighing filter numerator coefficients	$\hat{a}_m^{num}(i)$	$m = 0,1,2,3,4$
Weighing filter denominator coefficients	$\hat{a}_m^{den}(i)$	$m = 0,1,2,3,4$
Adaptive low-pass filter parameters	μ_m	$m = 0,1,2,3,4$
Weighted residual	$r_w(n)$	
Voiced/Unvoiced level	C_{VUV}	
Weighted residual half-frame energy	$E_{wr}(m)$	$m = 0,1$
Sharpness of weighted residual	S_{rp}^{wr}	
Weighted speech	$s_w(n)$	
Open-loop pitch lag	$l_p^{ol}(m)$	$m = 0,1,2$
Open-loop normalized correlation	$R_p^{ol}(m)$	$m = 0,1,2$
Frame classifier	C_{FR}	
Previous frame classifier	C_{FR}^{prv}	
Noise to signal ratio	NSR	
Auxiliary class decision	C_{SM}	
Onset flag	F_{onset}	
LPC onset flag	F_{LPC_onset}	
Noisy-voiced flag	F_{noisyV}	
SMV rate for the frame	$Rate$	
SMV rate for the previous frame	$Rate^{prv}$	
Low boundary of open-loop pitch refinement at the end of the frame	l_p^{low}	
High boundary of open-loop pitch refinement at the end of the frame	l_p^{high}	
Signal modification accumulated delay	Δ_{SM}	
Quantized open-loop pitch at the end of the encoding frame	l_{pp}^q	

Integer part of quantized open-loop pitch at the end of the encoding frame	$l_{pp,INT}^q$	
Frame type	$Type$	
Pitch contour	$Pit(n)$	
Modified weighted speech	$s_{mw}(n)$	
Subframe pitch gains (unquantized or quantized)	$g_a(m)$	
Subframe integer open-loop pitch	$l_p^{SF}(m)$	
Modified speech	$s_{mod}(n)$	
LSFs smoothing factor	β_{SMO}^{LSF}	
Smoothed LSFs for the encoding frame	$\underline{lsf}_1(i)$	
Smoothed LSFs for the previous encoding frame	$\underline{lsf}_{1,prv}(i)$	
Quantized LSFs for the frame	$lsf_1^q(i)$	
Interpolated LSFs for each subframe	$\hat{lsf}_m(i)$	$m = 0,1,2,3$ for Rate 1, 1/4, and 1/8 $m = 0,1$ for Rate 1/2 type 0 $m = 0,1,2$ for Rate 1/2 type 1
Quantized interpolated LSFs for each subframe	$\hat{lsf}_m^q(i)$	$m = 0,1,2,3$ for Rate 1, 1/4, and 1/8 $m = 0,1$ for Rate 1/2 type 0 $m = 0,1,2$ for Rate 1/2 type 1
Quantized adaptive-codebook gain for the last subframe of previous frame	g_a^{prv}	
Spectral flatness flag	F_{flat}	
LPC prediction gain from quantized parameters of the first subframe	$G_{0,q}^{LPC}$	
Random number seed	$seed$	
Random excitation (Rate 1/8)	$r(n)$	
Shaped random excitation (Rate 1/4)	$r_q(n)$	

Impulse response of combined synthesis and weighting filter	$h_{imp}(n)$	
Modified residual	$r_{mod}(n)$	
Target for closed-loop pitch search / pitch interpolation	$T_{gs}(n)$	
Adaptive codebook	$C^{adp}(n)$	
Ideal excitation	$I_{ext}(n)$	
Adaptive-codebook contribution (vector)	$c_{unf}^{adp}(n)$	
Filtered adaptive-codebook contribution	$c_f^{adp}(n)$	
Fixed-codebook contribution	$c_{unf}^{fix}(n)$	
Filtered fixed-codebook contribution	$c_f^{fix}(n)$	
Normalized closed-loop pitch correlation	R_p^{wgt}	
Fractional pitch for subframe	l_{FRC}^p	
Fractional pitch for previous subframe	$l_{FRC}^{p,prv}$	
Integer pitch for subframe	l_{INT}^p	
Normalized closed-loop pitch correlation	R_p^{wgt}	
Subframe mode factor	C_{sub_mod}	
Gain normalization factor	β_E	
Modified pitch gain	\hat{g}_a	
(Possibly) modified impulse response of combined synthesis and weighting filter	$\hat{h}_{imp}(n)$	

1
2
3
4
5

1 5 SMV ENCODER

2

3 5.1 SMV Encoder Look Ahead Buffers and Delay

4 The SMV encoder buffers 160 samples (20 ms) in an encoding frame buffer, and 80 samples
 5 (10 ms) in a lookahead buffer. Two options of windowing and buffering schemes are
 6 implemented for the SMV noise-suppression module, and are selected in the SMV
 7 simulation code by option A and option B.

8 If the noise-suppression simulation code is compiled with option A, an additional 24
 9 samples are used for the noise suppression overlap-and-add procedure. In this case, the
 10 total fixed buffering delay at the encoder is 264 samples, or 33 ms.

11 If the noise-suppression simulation code is compiled with option B, no additional samples
 12 are used for the noise suppression overlap-and-add procedure. In this case, the total fixed
 13 buffering delay at the encoder is 240 samples, or 30 ms.

14 The encoder buffers are schematically represented in Figure 5.1-1.

15

16

17

18

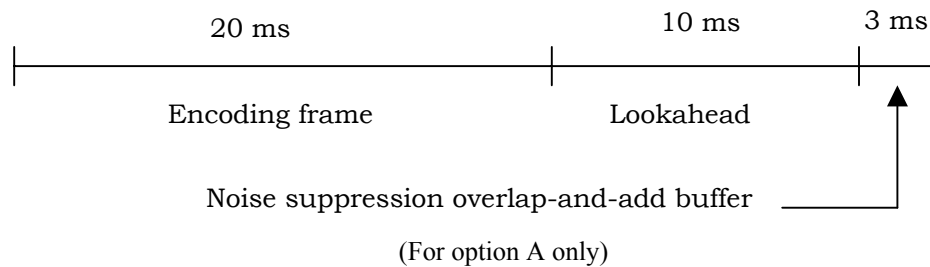
19

20

21

22

23



24

Figure 5.1-1: Selectable Mode Vocoder Encoder Buffering

25

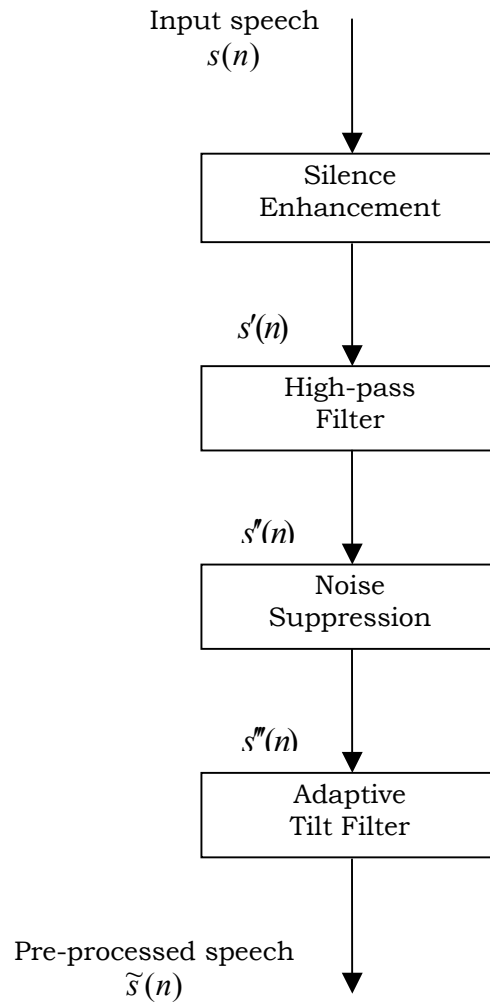
26 A speech modification procedure introduces a variable buffering delay, which is limited to
 27 20 samples (2.5 ms) (see Section 5.3.15).

28

28

1 5.2 SMV Encoder Pre-Processing

2 The speech samples input to the SMV encoder are processed by four pre-processing units.
3 The pre-processing units are silence enhancement, high-pass filtering, noise suppression,
4 and adaptive tilt filtering. The pre-processing flowchart is given in Figure 5.2-1.



28 **Figure 5.2-1: SMV Pre-Processing**

30
31

1 **5.2.1 Silence Enhancement**

2 Routine name: PPR_silence_enhan

3 **Input:**

- 4 • Input speech: $s(n)$

5 **Output:**

- 6 • Silence enhanced speech: $s'(n)$

7

8 The input signal is analyzed and modified to minimize low-level noise, typical of A-law and
9 μ -law PCM quantization process.

10

11 **5.2.1.1 Identifying Minimum and Maximum Values**

12 The two maximum values in the frame and two minimum values in the frame are found by:

$$13 \quad \begin{aligned} Max_0 &= \{s(n_0) : s(n_0) \geq s(n), n = 0, \dots, L_{frm} - 1\} \\ Max_1 &= \{s(n_1) : s(n_1) \geq s(n), n = 0, \dots, L_{frm} - 1, n_1 \neq n_0\}, \end{aligned} \quad (5.2.1.1-1)$$

$$14 \quad \begin{aligned} Min_0 &= \{s(n_0) : s(n_0) \leq s(n), n = 0, \dots, L_{frm} - 1\} \\ Min_1 &= \{s(n_1) : s(n_1) \leq s(n), n = 0, \dots, L_{frm} - 1, n_1 \neq n_0\}. \end{aligned} \quad (5.2.1.1-2)$$

15 Two smallest positive values in the frame and two largest negative values in the frame are
16 found by:

$$17 \quad MinPos_0 = \{s(n_0) : s(n_0) \leq s(n), n = 0, \dots, L_{frm} - 1, s(n_0) \geq 0\}, \quad (5.2.1.1-3)$$

$$18 \quad MinPos_1 = \{s(n_1) : s(n_1) \leq s(n), n = 0, \dots, L_{frm} - 1, s(n_1) \geq 0, n_1 \neq n_0\}, \quad (5.2.1.1-4)$$

$$19 \quad MaxNeg_0 = \{s(n_0) : s(n_0) \geq s(n), n = 0, \dots, L_{frm} - 1, s(n_0) < 0\}, \quad (5.2.1.1-5)$$

$$20 \quad MaxNeg_1 = \{s(n_1) : s(n_1) \geq s(n), n = 0, \dots, L_{frm} - 1, s(n_1) < 0, n_1 \neq n_0\}. \quad (5.2.1.1-6)$$

21 A level that represents the zero value is initialized to $V_{ZERO} = 8$. If $MinPos_0 < V_{ZERO}$, the
22 value of V_{ZERO} is set to $MinPos_0$.

23 **5.2.1.2 Setting Histogram Thresholds**

24 The four threshold values are initialized to $Th_{neg}^2 = -24$, $Th_{neg}^1 = -8$, $Th_{pos}^1 = 8$, and
25 $Th_{pos}^2 = 24$. The difference between the maximum value in the frame and the minimum
26 value in the frame is:

$$27 \quad \Delta_{lev} = Max_0 - Min_0. \quad (5.2.1.2-1)$$

28 If $\Delta_{lev} < \Delta_{lev}^{\min}$ and $\Delta_{lev} \neq 0$

$$1 \quad \Delta_{lev}^{\min} = \Delta_{lev} \quad (5.2.1.2-2)$$

2 If $\Delta_{lev} < \Delta_{lev}^{\min}$, $\Delta_{lev} \neq 0$ and $\Delta_{lev}^{\min} \leq 48$, the threshold values for the histogram are updated
3 according to:

$$\begin{aligned}
 & \text{if } (Max_1 \geq 0 \text{ and } Max_0 > 0) \{ \\
 & \quad Th_{pos}^1 = Max_1 \\
 & \quad Th_{pos}^2 = Max_0 \\
 & \quad \} \\
 4 \quad & \text{else} \{ \quad (5.2.1.2-3) \\
 & \quad \text{if } (MinPos_0 < 32767) \\
 & \quad \quad Th_{pos}^1 = MinPos_0 \\
 & \quad \text{if } (MinPos_1 < 32767) \\
 & \quad \quad Th_{pos}^2 = MinPos_1 \\
 & \quad \}
 \end{aligned}$$

$$\begin{aligned}
 5 \quad & \text{if } (Min_1 < 0 \text{ and } Min_0 < 0) \{ \\
 & \quad Th_{neg}^1 = Min_1 \\
 & \quad Th_{neg}^2 = Min_0 \\
 & \quad \} \\
 6 \quad & \text{else} \{ \quad (5.2.1.2-4) \\
 & \quad \text{if } (MaxNeg_0 > -32766) \cdot \\
 & \quad \quad Th_{neg}^1 = MaxNeg_0 \\
 & \quad \text{if } (MaxNeg_1 > -32766) \\
 & \quad \quad Th_{neg}^2 = MaxNeg_1 \\
 & \quad \}
 \end{aligned}$$

7 **5.2.1.3 Calculating the 3-Value Histogram**

8 A 3-value histogram is calculated, where H_1^{SE} counts the number of samples in the frame
9 on the interval $[Th_{neg}^1, 0)$, H_2^{SE} counts the number of samples in the frame on the interval
10 $[0, Th_{pos}^1]$, and H_3^{SE} counts the number of samples in the frame on the combined interval
11 $[-32768, Th_{neg}^2) \cup (Th_{pos}^1, 32767]$. Three population ratios for the frame are calculated
12 according to:

$$13 \quad PR_{ZERO}^0 = \frac{H_1^{SE}}{L_{frm}}, \quad (5.2.1.3-1)$$

$$1 \quad PR_{LOW}^0 = \frac{H_1^{SE} + H_2^{SE}}{L_{frm}} , \quad (5.2.1.3-2)$$

$$2 \quad PR_{HIGH}^0 = \frac{H_3^{SE}}{L_{frm}} . \quad (5.2.1.3-3)$$

3 The population ratios for the current frame and for the last 3 frames, as well as the
 4 previous frame ‘silencing’ decision, are used to generate the current frame ‘silencing’
 5 decision. If the current frame is ‘silenced’, but the previous frame was not ‘silenced’, the
 6 samples in the current frame are given by:

$$7 \quad s'(n) = \begin{cases} \frac{(39-n) \cdot s(n) + n \cdot V_{ZERO}}{40} & ; 0 \leq n < 40 \\ V_{ZERO} & ; 40 < n < L_{frm} \end{cases} . \quad (5.2.1.3-4)$$

8 If the current frame is ‘silenced’ and the previous frame was also ‘silenced’, all of the values
 9 in the frame are set to V_{ZERO} . The samples of the current frame are not changed if the
 10 current frame is not “silenced” and the previous frame was not “silenced”. If the current
 11 frame is not ‘silenced’ but previous frame was ‘silenced’, the samples in the current frame
 12 are given by:

$$13 \quad s'(n) = \begin{cases} \frac{(39-n) \cdot V_{ZERO} + n \cdot s(n)}{40} & ; 0 \leq n < 40 \\ s(n) & ; 40 < n < L_{frm} \end{cases} . \quad (5.2.1.3-5)$$

14

1 **5.2.2 High-Pass Filtering**

2 Routine name: FLT_filterPZ (generic pole-zero filter)

3 **Input:**

- 4 • Silence enhanced speech, $s'(n)$

5 **Output:**

- 6 • High-pass filtered speech, $s''(n)$

7

8 The speech is filtered by a high-pass filter with the transfer function:

9
$$H(z) = 0.5 \cdot \frac{0.92727435 - 1.8544941 \cdot z^{-1} + 0.92727435 \cdot z^{-2}}{1 - 1.9059465 \cdot z^{-1} + 0.9114024 \cdot z^{-2}}. \quad (5.2.2-1)$$

10 The input is scaled down 3 dB by the multiplication factor of 0.5.

11

1 **5.2.3 Noise Suppression**

2 Routine name: SNS_modified_noise_suprs

3 **Input:**

- 4 • High-pass filtered speech, $s''(n)$
- 5 • Previous frame class decision, C_{FR}^{Prv}

6 **Output:**

- 7 • Noise suppressed speech, $s'''(n)$

8 Noise Suppression is used to improve the signal quality that is presented to the Model
 9 Parameter Estimator. The procedures by which the noise suppression is implemented is as
 10 described in Section 5.2.3. The input to (and the output from) the core noise suppression
 11 algorithm can be processed through two different options that are described in the following
 12 section. One of the option flags must be set prior to compiling the SMV simulation C code.
 13 The difference between the two options generally corresponds to the use of different
 14 analysis/synthesis windowing and buffering methodology, resulting in zero algorithmic
 15 delay for the noise suppression module when option B is selected.

16 Processing: Although the frame size of the speech codec is 20 ms, the Noise Suppressor
 17 subframe size is 10 ms; therefore, the following procedures shall be executed two times per
 18 20 ms speech frame and the current 10 ms subframe shall be denoted m .

19

20 **5.2.3.1 Noise Suppression Initialization**

21 The following variables shall be set to zero at initialization ($m = 0$):

- 22 • The overlapped portion of the input subframe buffer, $\{d(m)\}$
- 23 • The pre-emphasis and de-emphasis memories
- 24 • The overlap-and-add buffer, $\{h(n)\}$
- 25 • The output buffer history, $s'''(n)$; $0 \leq n < L_{frm}$
- 26 • The onset counter C_{onset} .

27 The following shall be initialized to a startup value other than zero:

- 28 • The channel energy estimate, $E_{ch}(m)$, (see 5.2.3.3)
- 29 • The long-term power spectral estimate, $\bar{E}_{dB}(m)$, (see 5.2.3.6)
- 30 • The channel noise estimate, $E_n(m)$, (see 5.2.3.11)

31
 32

33 [Figure 5.2-2](#)

34

35 ~~Figure 5.2-2~~ depicts the Noise Suppression system that is described in the following
 36 sections.

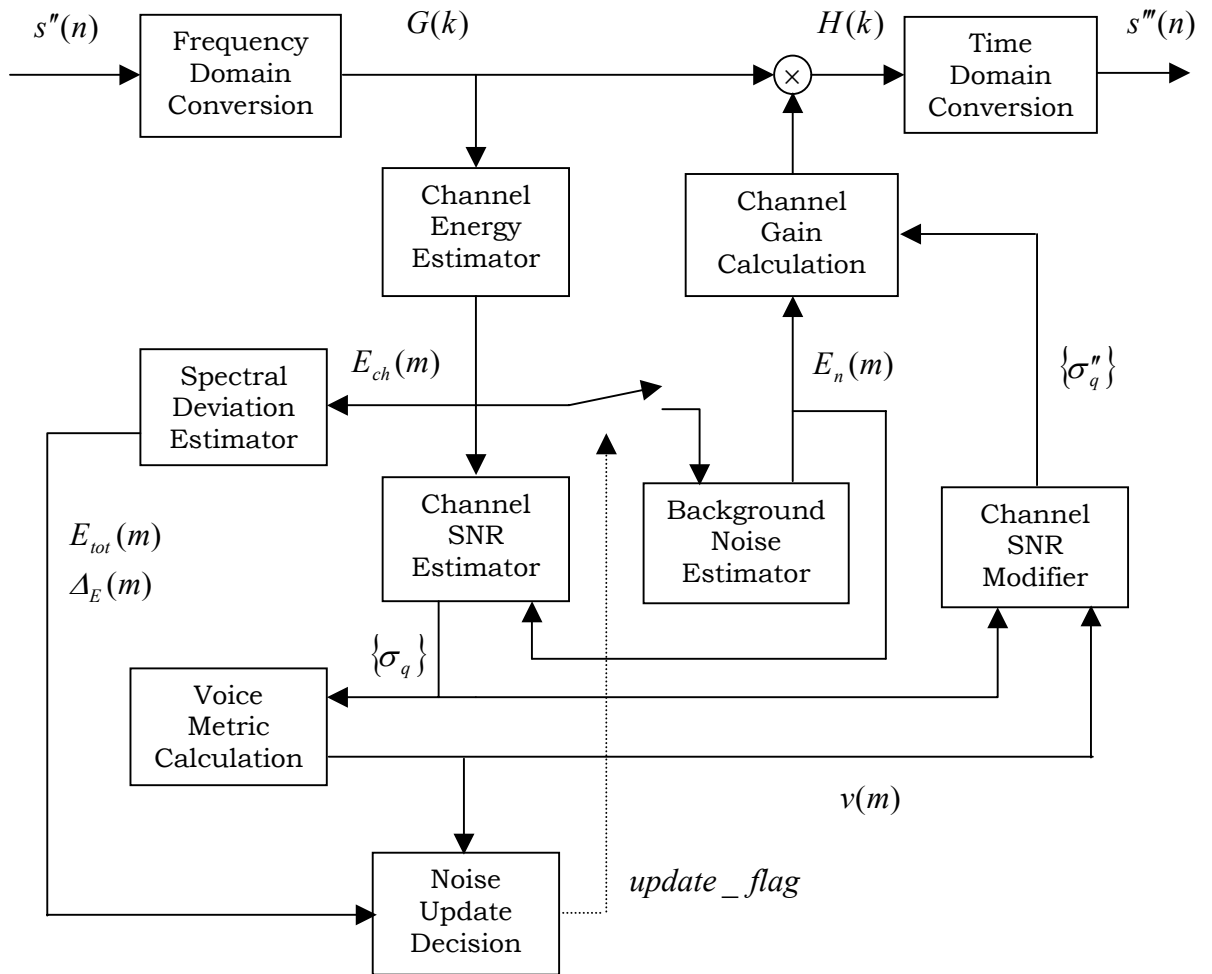


Figure 5.2-2: Noise Suppression Block Diagram

1 5.2.3.2 Frequency Domain Conversion

2 5.2.3.2.1 Windowing and Buffering - Option A

3 In this approach, the input signal is windowed using a smoothed trapezoid window wherein
4 the first D samples of the input subframe buffer $\{d(m)\}$ overlap the last D samples of the
5 previous subframe. This overlap is described as:

$$6 \quad d(m, n) = d(m-1, L+n) \quad 0 \leq n < D, \quad (5.2.3.2-1)$$

7 where m is the current subframe, n is the sample index to the buffer $\{d(m)\}$, $L = 80$ is
8 the subframe length, and $D = 24$ is the overlap (or delay) in samples.

9 The remaining samples of the input buffer are then pre-emphasized according to the
10 following:

$$11 \quad d(m, D+n) = s''(n) + \zeta_p \cdot s''(n-1); \quad 0 \leq n < L, \quad (5.2.3.2-2)$$

12 where $\zeta_p = -0.8$ is the pre-emphasis factor. This results in the input buffer containing
13 $L + D = 104$ samples of which the first D samples are the pre-emphasized overlap from
14 the previous subframe, and the following D samples are pre-emphasized input from the
15 current subframe.

16 Next, a smoothed trapezoidal window is applied to the input buffer to form the DFT data
17 buffer, $\{g(n)\}$ defined as:

$$18 \quad g(n) = \begin{cases} d(m, n) \sin^2(\pi(n+0.5)/2D) & ; 0 \leq n < D \\ d(m, n) & ; D \leq n < L \\ d(m, n) \sin^2(\pi(n-L+D+0.5)/2D) & ; L \leq n < D+L \\ 0 & ; D+L \leq n < M \end{cases}, \quad (5.2.3.2-3)$$

19 where $M = 128$ is the DFT sequence length and all other terms are previously defined.

20

21 5.2.3.2.2 Windowing and Buffering - Option B

22 A tapered trapezoidal window function as shown in Figure 5.2-3 is applied to the input
23 sequence. The analysis and synthesis windows have slopes of $D_{ana} = 12$ and $D_{syn} = 6$
24 samples in length, respectively.

25

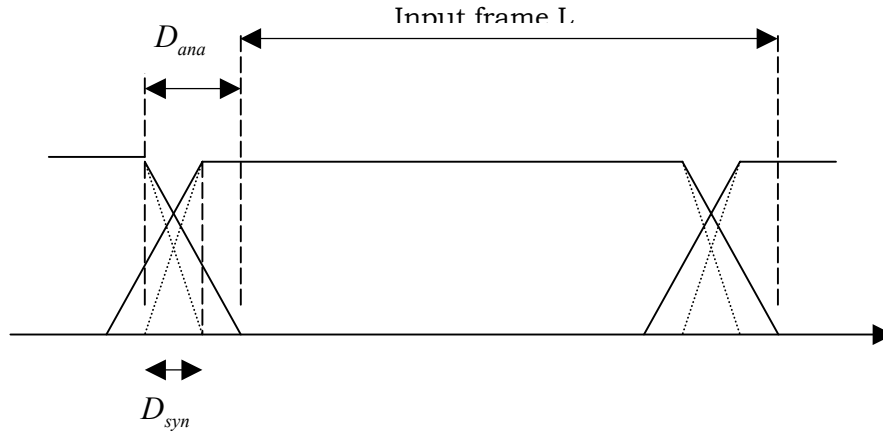


Figure 5.2-3: Two-Phase Trapezoidal Window with Overlap-and-Add. (The analysis window is illustrated using a solid line. The slopes of the synthesis window are shown with a dotted line.)

Similar to option A, the input is overlapped and pre-emphasized according to equations [5.2.3.2-1) and (5.2.3.2-2), respectively. Unlike option A, the constant D is now equal to $2D_{ana} - D_{syn}$. The slope of the analysis window is defined as follows:

$$f_{ana}(n) = (n+1)/(D_{ana} + 1), \quad 0 \leq n < D_{ana} \quad (5.2.3.2-4)$$

The trapezoidal analysis window is applied to the pre-emphasized signal.

$$g(n) = \begin{cases} d(m,n)f_{ana}(n) & ; 0 \leq n < D_{ana} \\ d(m,n) & ; D_{ana} \leq n < L + D_{ana} - D_{syn} \\ d(m,n)f_{ana}(L + 2D_{ana} - D_{syn} - n - 1) & ; L + D_{ana} - D_{syn} \leq n < L + 2D_{ana} - D_{syn} \\ 0 & ; L + 2D_{ana} - D_{syn} \leq n < M \end{cases} \quad (5.2.3.2-5)$$

where $M = 128$ is the length of the FFT buffer.

The steps in Sections 5.2.3.3 to 5.2.3.11 are common in both approaches and are described below.

The transformation of $g(n)$ to the frequency domain is performed using the Discrete Fourier Transform (DFT) defined† as:

† This atypical definition is used to exploit the efficiencies of the complex Fast Fourier Transform (FFT). The $2/M$ scale factor results from preconditioning the M point real sequence to form an $M/2$ point complex sequence that is transformed using an $M/2$ point complex FFT. Details on this technique can be found in Proakis, J. G. and Manolakis, D. G., *Introduction to Digital Signal Processing*, New York, Macmillan, 1988, pp. 721-722.

$$1 \quad G(k) = \frac{2}{M} \sum_{n=0}^{M-1} g(n) \cdot e^{j2\pi nk/M}; \quad 0 \leq k < M. \quad (5.2.3.2-6)$$

2 where $e^{j\omega}$ is a unit amplitude complex phasor with instantaneous radial position ω .

3

4 **5.2.3.3 Channel Energy Estimator**

5 Calculate the channel energy estimate $E_{ch}(m)$ for the current subframe, m , as:

6

$$7 \quad E_{ch}(m, i) = \max \left\{ E_{\min}, \alpha_{ch}(m) \cdot E_{ch}(m-1, i) + (1 - \alpha_{ch}(m)) \cdot \frac{\sum_{k=f_L(i)}^{f_H(i)} |G(k)|^2}{f_H(i) - f_L(i) + 1} \right\};$$

$$8 \quad 0 \leq i < N_c, \quad (5.2.3.3-1)$$

9 where $E_{\min} = 0.0625$ is the minimum allowable channel energy, $\alpha_{ch}(m)$ is the channel
 10 energy smoothing factor (defined below), $N_c = 16$ is the number of combined channels, and
 11 $f_L(i)$ and $f_H(i)$ are the i^{th} elements of the respective low and high channel combining tables,
 12 which are defined as:

$$13 \quad \mathbf{f}^L = \{2, 4, 6, 8, 10, 12, 14, 17, 20, 23, 27, 31, 36, 42, 49, 56\}, \quad (5.2.3.3-2)$$

$$14 \quad \mathbf{f}^H = \{3, 5, 7, 9, 11, 13, 16, 19, 22, 26, 30, 35, 41, 48, 55, 63\}. \quad (5.2.3.3-3)$$

15 The channel energy smoothing factor, $\alpha_{ch}(m)$ is defined as:

$$16 \quad \alpha_{ch}(m) = \begin{cases} 0 & ; m = 0 \\ 0.25 & ; m > 0, C_{onset} = 1 \\ 0.45 & ; m > 0, C_{onset} \neq 1 \end{cases} \quad (5.2.3.3-4)$$

17 Therefore, $\alpha_{ch}(m)$ assumes a value of zero for the first subframe ($m = 0$), the value of 0.25
 18 for onset frames, and 0.45 for all other frames. This allows the channel energy estimate to
 19 be initialized to the unfiltered channel energy of the first subframe, and to be adaptively
 20 updated for the subsequent subframes.

21

22 **5.2.3.4 Channel SNR Estimator**

23 Estimate the quantized channel SNR indices as:

$$24 \quad \sigma_q(i) = \max \left\{ 0, \min \left\{ 89, \text{round} \left(10 \log_{10} \frac{E_{ch}(m, i)}{E_n(m, i)} / 0.375 \right) \right\} \right\}; \quad 0 \leq i < N_c, \quad (5.2.3.4-1)$$

1 where $E_n(m,i)$ is the current channel noise energy estimate (Section 5.2.3.11), and the
 2 values of $\{\sigma_q\}$ are constrained to be between 0 and 89, inclusive.

3

4 **5.2.3.5 Voice Metric Calculation**

5 Next, calculate the sum of voice metrics as:

$$6 \quad v(m) = \sum_{i=0}^{N_c-1} V(\sigma_q(i)), \quad (5.2.3.5-1)$$

7 where $V(k)$ is the k^{th} value of the 90 element voice metric table \mathbf{V} , that is defined as:

$$8 \quad \mathbf{V} = \{ 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 4, 5, 5, 5, 6, 6, 7, 7, 7, 8, 8, 9, 9, \\ 9 \quad 10, 10, 11, 12, 12, 13, 13, 14, 15, 15, 16, 17, 17, 18, 19, 20, 20, 21, 22, 23, 24, 24, \\ 10 \quad 25, 26, 27, 28, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 37, 38, 39, 40, 41, 42, 43, 44, \\ 11 \quad 45, 46, 47, 48, 49, 50, 50, 50, 50, 50, 50, 50, 50, 50, 50 \}.$$

$$12 \quad (5.2.3.5-2)$$

14 **5.2.3.6 Spectral Deviation Estimator**

15 Calculate the estimated log power spectrum as:

$$16 \quad E_{dB}(m,i) = 10 \log_{10}(E_{ch}(m,i)); \quad 0 \leq i < N_c. \quad (5.2.3.6-1)$$

17 Then, calculate the estimated spectral deviation between the current power spectrum and
 18 the average long-term power spectral estimate:

$$19 \quad \Delta_E(m) = \sum_{i=0}^{N_c-1} |E_{dB}(m,i) - \bar{E}_{dB}(m,i)|, \quad (5.2.3.6-2)$$

20 where $\bar{E}_{dB}(m)$ is the average long-term power spectral estimate calculated during the
 21 previous subframe, as defined by Equation (5.2.3.6-7). The initial value of $\bar{E}_{dB}(m)$,
 22 however, is defined to be the estimated log power spectrum of subframe 0, or:

$$23 \quad \bar{E}_{dB}(m) = E_{dB}(0). \quad (5.2.3.6-3)$$

24 Calculate the total channel energy estimate, $E_{tot}(m)$ for the current subframe, m ,
 25 according to the following:

$$26 \quad E_{tot}(m) = 10 \log_{10} \left(\sum_{i=0}^{N_c-1} E_{ch}(m,i) \right). \quad (5.2.3.6-4)$$

27 Calculate the exponential windowing factor, $\alpha(m)$, as a function of total channel energy,
 28 $E_{tot}(m)$, as:

$$29 \quad \alpha(m) = \alpha_H - \left(\frac{\alpha_H - \alpha_L}{E_H - E_L} \right) (E_H - E_{tot}(m)), \quad (5.2.3.6-5)$$

1 and then limit the result to be between α_L and α_H by

$$2 \quad \alpha(m) = \max\{\alpha_L, \min\{\alpha_H, \alpha(m)\}\}, \quad (5.2.3.6-6)$$

3 where E_H and E_L are the energy endpoints (in dB) for the linear interpolation of $E_{tot}(m)$, that
 4 is transformed to $\alpha(m)$ that has the limits $\alpha_L \leq \alpha(m) \leq \alpha_H$. The values of these constants
 5 are defined as: $E_H = 50$, $E_L = 30$, $\alpha_H = 0.99$, $\alpha_L = 0.50$. As an example, a signal with relative
 6 energy of 40 dB would use an exponential windowing factor of $\alpha(m) = 0.745$ for the
 7 following calculation.

8 Update the average long-term power spectral estimate for the next subframe by:

$$9 \quad \bar{E}_{dB}(m+1, i) = \alpha(m) \cdot \bar{E}_{dB}(m, i) + (1 - \alpha(m)) \cdot E_{dB}(m, i);$$

$$10 \quad 0 \leq i < N_c, \quad (5.2.3.6-7)$$

11 where all the variables are previously defined.

12

13 **5.2.3.7 Background Noise Update Decision**

14 The following logic, shown in pseudo-code, demonstrates how the noise estimate update
 15 decision is ultimately made:

16

```

17     /* Normal update logic */
18     update_flag = FALSE
19     if ( v(m) ≤ UPDATE_THLD ) {
20         update_flag = TRUE
21         update_cnt = 0
22     }
23     /* Forced update logic */
24     else if (( E_tot(m) > NOISE_FLOOR_DB ) and ( Δ_E(m) < DEV_THLD )) {
25         update_cnt = update_cnt + 1
26         if ( update_cnt ≥ UPDATE_CNT_THLD )
27             update_flag = TRUE
28     }
29     /* “Hysteresis” logic to prevent long-term creeping of update_cnt */
30     if ( update_cnt == last_update_cnt )
31         hyster_cnt = hyster_cnt + 1
32     else
33         hyster_cnt = 0
34     last_update_cnt = update_cnt
35     if ( hyster_cnt > HYSTER_CNT_THLD )
36         update_cnt = 0

```

37 The values of the previously used constants are UPDATE_THLD = 35,
 38 NOISE_FLOOR_DB = $10\log_{10}(E_{floor})$ (see Section 5.2.3.9), DEV_THLD = 28,
 39 UPDATE_CNT_THLD = 50, HYSTER_CNT_THLD = 6.

1

2 **5.2.3.8 SNR Estimate Modification**3 Next, determine whether the channel SNR modification should take place, then proceed to
4 modify the appropriate SNR indices:

```

5       /* Set or reset modify flag */
6       index_cnt = 0
7       for ( i = NM to Nc - 1 step 1 ) {
8           if ( σq(i) ≥ INDEX_THLD )
9               index_cnt = index_cnt + 1
10       }
11       if ( index_cnt < INDEX_CNT_THLD )
12           modify_flag = TRUE
13       else
14           modify_flag = FALSE
15
16       /* Modify the SNR indices to get {σ'q} */
17       if ( modify_flag == TRUE ) {
18           for ( i = 0 to Nc - 1 step 1 )
19               if ( ( v(m) ≤ METRIC_THLD ) or ( σq(i) ≤ SETBACK_THLD ) )
20                   σ'q(i) = 1
21               else
22                   σ'q(i) = σq(i)
23           }
24       else {
25           {σ'q} = {σq}
26       }

```

27

28

29 /* Limit {σ''_q} to SNR threshold σ_{th} */

```

30       for ( i = 0 to Nc - 1 step 1 )
31           if ( σ'q(i) < σth )
32               σ''q(i) = σth
33           else
34               σ''q(i) = σ'q(i)

```

35

36 The previous constants and thresholds are given to be: N_M = 5, INDEX_THLD = 12,
37 INDEX_CNT_THLD = 5, METRIC_THLD = 45, SETBACK_THLD = 12, σ_{th} = 6.

38

1 **5.2.3.9 Channel Gain Computation**

2 Calculate the sum of the estimated noise spectrum $E_n(m)$ (see Section 5.2.3.7). Reduce the
3 sum by a constant factor, and limit it to a positive $\varepsilon = 0.00001$. The result, T_{ne} , is given by:

$$4 \quad 5 \quad T_{ne} = \max\left(\sum_{i=0}^{N_c-1} E_n(m, i) - 12.0412, 0.00001\right). \quad (5.2.3.9-1)$$

6 Compute the overall gain factor for the current subframe as:

$$7 \quad 8 \quad \gamma_n = \max\left(\gamma_{\min}, -10 \cdot \log\left(\frac{T_{ne}}{E_{\text{floor}}}\right)\right), \quad (5.2.3.9-2)$$

9 where $E_{\text{floor}} = 1$ is the noise floor energy, and the minimum overall gain $\gamma_{\min} = -13$ dB if
10 $F_{\text{flat}} = 0$, and $\gamma_{\min} = -14$ dB if $F_{\text{flat}} = 1$. Next, calculate channel gains (in dB) as:

$$11 \quad \gamma_{dB}^i(m) = \min\left(\left(\sigma_q''(i) - \sigma_{th}\right) \cdot \mu_g \cdot \left(1 - \frac{i}{4N_c}\right) + \gamma_n, 0\right); 0 \leq i < N_c, \quad (5.2.3.9-3)$$

12 where $\mu_g = 0.39$ is the gain slope, i is the channel index, and m is the subframe number.
13 Calculate the average gain for the low spectrum and for the high spectrum by:

$$14 \quad \gamma_L(m) = \frac{1}{N_M + 1} \sum_{i=0}^{N_M} \gamma_{dB}(i) ; \gamma_H(m) = \frac{1}{N_c - N_M} \sum_{i=N_M+1}^{N_c} \gamma_{dB}(i), \quad (5.2.3.9-4)$$

15 where m is the subframe number. C_{onset} is calculated as a function of previous frame speech
16 class C_{FR}^{Prv} , and is set to 1 if the previous frame was not an onset or voiced (stationary or
17 non-stationary), and if one of the following conditions is met: i) $\gamma_L(m) - \gamma_L(m-1) > 4$, ii)
18 $\gamma_L(m) - \gamma_L(m-1) > 2$ and $\gamma_L(m) - \gamma_H(m) > -3$.

19 If $C_{\text{onset}} \leq 3$ it is incremented, and if $C_{\text{onset}} > 3$ it is set to 0. Based on the *modify_flag*, the
20 previous frame class, $v(m)$, and C_{onset} , a voicing decision is made. If the voicing decision is
21 TRUE, correct the channel gains for the channels indexed lower than 9 according to:

$$22 \quad y_{dB}^i(m) = \max\left\{\left(0.25 + 0.75 \cdot i/8\right) \cdot y_{dB}^i(m) + \left(0.75 - 0.75 \cdot i/8\right) \cdot \max\{\gamma_L(m), \gamma_H(m)\}, y_{dB}^i(m)\right\}. \quad (5.2.3.9-5)$$

24 Further, if $C_{\text{onset}} \neq 1$, smooth the channel gains for the channels indexed lower than 9
25 according to:

$$26 \quad \gamma_{dB}^i(m) = \min\left(0.75 \cdot \gamma_{dB}^i(m) + 0.25 \cdot \gamma_{dB}^i(m-1), \gamma_{dB}^i(m)\right), \quad (5.2.3.9-6)$$

27 and then convert the channel gains in dB to linear channel gains:

$$28 \quad \gamma_{ch}(i) = 10^{\gamma_{dB}^i(m)/20} ; 0 \leq i < N_c. \quad (5.2.3.9-7)$$

29

1 **5.2.3.10 Frequency Domain Filtering**

2 Now, apply the channel gains to the transformed input signal $G(k)$:

$$3 \quad H(k) = \begin{cases} \gamma_{ch}(i) \cdot G(k) & ; f_L(i) \leq k \leq f_H(i) , \quad 0 \leq i < N_c \\ G(k) & ; 0 \leq k < f_L(0) , \quad f_H(N_c - 1) < k \leq M/2 \end{cases} . \quad (5.2.3.10-1)$$

4 where the bottom part of the equation represents the frequencies that are not altered by the
5 channel gains. Since it is also required that the magnitude of $H(k)$ be even, and the phase
6 be odd, the following condition is also imposed:

$$7 \quad H(M - k) = H^*(k) , \quad 0 < k < M/2 , \quad (5.2.3.10-2)$$

8 where $*$ denotes complex conjugate. This guarantees that the imaginary part of the inverse
9 DFT of $H(k)$ will be zero (in Equation (5.2.3.12-1)).

10

11 **5.2.3.11 Background Noise Estimate Update**

12 If (and only if) the update flag is set ($update_flag = \text{TRUE}$), then update the channel noise
13 estimate for the next subframe by:

$$14 \quad E_n(m + 1, i) = \max\{E_{\min}, \alpha_n \cdot E_n(m, i)\} + (1 - \alpha_n) \cdot E_{ch}(m, i); \quad 0 \leq i < N_c , \quad (5.2.3.11-1)$$

16 where $E_{\min} = 0.0625$ is the minimum allowable channel energy, and $\alpha_n = 0.9$ is the channel
17 noise smoothing factor. The channel noise estimate shall be initialized for each of the first
18 four subframes to the estimated channel energy, i.e.:

$$19 \quad E_n(m, i) = \max\{E_{\text{init}}, E_{ch}(m, i)\} ; \quad 0 \leq m \leq 3 , \quad 0 \leq i < N_c , \quad (5.2.3.11-2)$$

20 where $E_{\text{init}} = 16$ is the minimum allowable channel noise initialization energy.

21

22 **5.2.3.12 Time Domain Signal Reconstruction**

23 Convert the filtered signal to the time domain using the inverse DFT:

$$24 \quad h(m, n) = \frac{1}{2} \sum_{k=0}^{M-1} H(k) \cdot e^{j2\pi nk / M} ; \quad 0 \leq n < M . \quad (5.2.3.12-1)$$

25

26 5.2.3.12.1 Overlap-and-Add – Option A

27 Complete the frequency domain filtering process by applying overlap-and-add:

$$28 \quad h'(n) = \begin{cases} h(m, n) + h(m - 1, n + L) & ; 0 \leq n < M - L \\ h(m, n) & ; M - L \leq n < L \end{cases} . \quad (5.2.3.12-2)$$

29 Finally, apply signal de-emphasis by:

$$30 \quad s'''(n) = h'(n) + \zeta_d \cdot s'''(n - 1); \quad 0 \leq n < L , \quad (5.2.3.12-3)$$

1 where $\zeta_d = 0.8$ is the de-emphasis factor and $\{s^m(n)\}$ is the output buffer.

2

3 5.2.3.12.2 Overlap-and-Add – Option B

4 The slope of the synthesis window is defined as

$$5 \quad f_{syn}(n) = (n+1)/(D_{syn} + 1) \quad ; 0 \leq n < D_{syn} \quad (5.2.3.12-4)$$

6 The output overlap-add is calculated as follows:

$$7 \quad h'(m, n) = \begin{cases} h'(m-1, n+L+D_{ana}-D_{syn})f_{syn}(D_{syn}-n-1) + \\ h(m, n+D_{ana}-D_{syn})f_{syn}(n)/f_{ana}(n+D_{ana}-D_{syn}) & ; 0 \leq n < D_{syn} \\ h(m, n+D_{ana}-D_{syn}) & ; D_{syn} \leq n < L \\ h(m, n+D_{ana}-D_{syn})/f_{ana}(L+D_{ana}-n-1) & ; L \leq n < L+D_{ana} \end{cases} .$$

$$8 \quad (5.2.3.12-5)$$

9 Finally, the output signal is de-emphasized as follows:

$$10 \quad s^m(n) = h'(m, n+D_{ana}) + \zeta_d s^m(n-1); \quad -D_{ana} \leq n < L, \quad (5.2.3.12-6)$$

11 where $\zeta_d = 0.8$ is the de-emphasis factor and $\{s^m(n)\}$ is the output buffer.

12

13 Since a 10 ms subframe Noise Suppression is performed twice per 20 ms speech frame, the
14 output presented to the adaptive tilt filter is a 20 ms frame $s^m(n)$.

15

1 5.2.4 Adaptive Tilt Filtering

2 Routine name: PPR_lowpass

3 **Input:**

- 4 • Noise suppressed speech, $s'''(n)$
- 5 • Spectral flatness flag, F_{flat}
- 6 • SMV mode of operation, M_{SMV}
- 7 • Half-rate-max flag, F_{HRM}

8 **Output:**

- 9 • Pre-processed speech, $\tilde{s}(n)$

10

11 The speech is filtered by an adaptive tilt filter with the transfer function:

$$12 \quad T(z) = 1 + \lambda \cdot z^{-1}. \quad (5.2.4-1)$$

13 The parameter λ depends on the SMV operating mode, the spectral flatness flag, and on
 14 the half-rate-max flag, as given by Table 5.2-1.

15

	SMV Mode 0	SMV Mode 1	SMV Mode 2	SMV Mode 3
$F_{flat} = 1, F_{HRM} = 1$	0.0	0.0	0.0	0.0
$F_{flat} = 1, F_{HRM} = 0$	-0.1	0.0	0.0	0.0
$F_{flat} = 0, F_{HRM} = 1$	0.1	0.1	0.1	0.1
$F_{flat} = 0, F_{HRM} = 0$	0.0	0.05	0.075	0.1

16

17

Table 5.2-1: Input Spectral Tilt Coefficient Table

18

19

20

5.3 SMV Encoder Frame Level Processing

The frame processing in the SMV encoder includes linear prediction (LP) analysis, LSFs smoothing and quantization. The speech parameters are used for voice activity detection (VAD), frame classification, type selection, and rate determination. Using the perceptually weighted speech, open-loop pitch detection and signal modification is performed. Figure 5.3-1 is a block diagram of the SMV frame level processing.

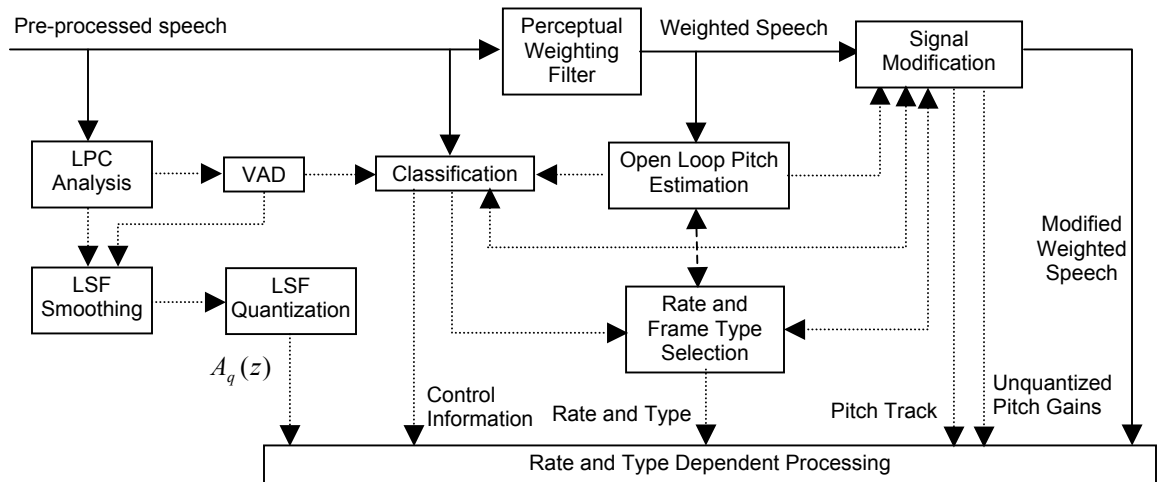


Figure 5.3-1: Signal Flow for Frame Processing in SMV Encoder

Depending on the encoding mode of operation, SMV classifies the input speech frame, and chooses an optimal Rate and Type for frame processing. Figure 5.3-2 illustrates the various stages of classification in the SMV encoder, as described in the following sections. Voiced activity detection is described in Sections 5.3.2 and 5.3.3. Music detection is described in Section 5.3.4. Voiced/unvoiced level detection is described in Section 5.3.7. Active speech classification is described in Section 5.3.10. Class correction is described in Section 5.3.12. The encoding mode dependent rate selection is described in Section 5.3.13. Voiced speech classification is described in Section 5.3.15.3. The final rate and class corrections are described in Section 5.3.16.

1

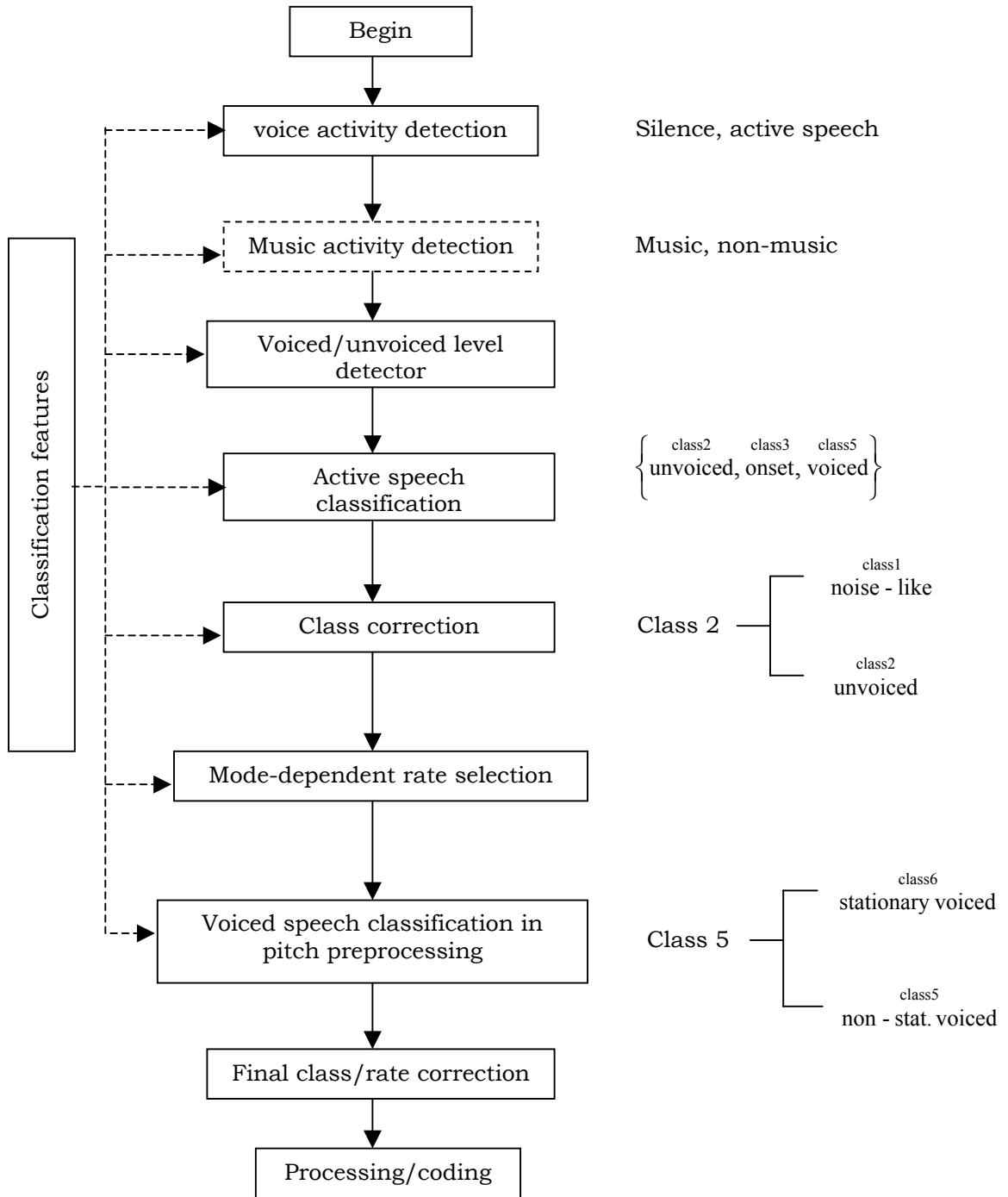


Figure 5.3-2: Stages of Classification in SMV Encoder

2
3
4
5
6

1

2

3 **5.3.1 Linear Prediction (LPC) Analysis**

4 The SMV encoder performs the LPC analysis three times for each frame. Each LPC analysis
5 is performed with a different weighting window, which is centered on a different section of
6 the frame or the lookahead.

7 The first LPC analysis is performed with the window centered on the second quarter of the
8 encoding frame. The LPC parameters from this LPC analysis are used for generating the
9 weighted speech, for the frame classification, and for determining the LSFs interpolation
10 path for Rate 1 Type-0 frames. The second LPC analysis is performed with the window
11 centered on the last quarter of the encoding frame. The LPC parameters from the second
12 LPC analysis are smoothed (Section 5.3.21) and quantized (Section 5.3.22.3). The third LPC
13 analysis is performed with the window centered on the lookahead buffer, and the LPC
14 parameters are used for generating the weighted speech for the lookahead buffer.

15 **5.3.1.1 The LPC Analysis Weighting Windows**

16 Each of the three LPC analyses uses a different weighting window, but all weighting
17 windows have the same size: $L_{LPC} = 240$. The first LPC weighting window, $W_0^{LPC}(n)$, is a
18 symmetric Hamming window of size 240. The second LPC weighting window, $W_1^{LPC}(n)$, is
19 formed by the concatenation of the rising half of a Hamming window of size 150 with the
20 falling half of a Hamming window of size 90. The third LPC weighting window, $W_2^{LPC}(n)$, is
21 formed by the concatenation of the rising half of a Hamming window of size 200 with the
22 falling half of a raised-cosine window of size 40.

23

24 **5.3.1.2 The LP Analysis Procedure**

25 Routine name: LPC_analysis

26 **Input** (for the m^{th} LPC analysis, $m = 0, 1, 2$):

- 27 • Pre-processed speech, $\tilde{s}(n)$
- 28 • Spectral flatness flag, F_{flat}
- 29 • LPC weighting window, $W_m^{LPC}(n)$
- 30 • The autocorrelation order, O_R , (10 for VAD-A, 16 for VAD-B)

31 **Output** (for the m^{th} LPC analysis, $m = 0, 1, 2$):

- 32 • Autocorrelation function, $R_m(i)$
- 33 • Reflection coefficients, $k_m(i)$
- 34 • Prediction coefficient, $a_m(i)$
- 35 • LSFs, $lsf_m(i)$

- 1 • LPC prediction error, PE_m^{LPC}
- 2 • LPC prediction gain, G_m^{LPC}

3

4 5.3.1.2.1 Energy Calculation

5 The segment energy is calculated as

$$6 \quad E_{LPC} = \sum_{n=0}^{L_{LPC}-1} \tilde{s}^2[n]. \quad (5.3.1.2-1)$$

7 If $E_{LPC} = 0$ the LPC analysis procedure returns zero values for the autocorrelation, the
 8 reflection coefficients, the prediction coefficients, and the LP prediction error. The LP
 9 prediction gain (the prediction error in dB) is set to $-DBL_MAX$. The LSFs are set to equally
 10 spaced values. If $E_{LPC} \neq 0$, the LPC analysis proceeds to the next steps, described in
 11 Sections 5.3.1.2.2 to 5.3.1.2.5.

12

13 5.3.1.2.2 Calculation of the Autocorrelation Function

14 The signal is multiplied by the LPC window (m indicates the index of the LPC analysis):

$$15 \quad s_w^{LPC}(n) = \tilde{s}(n) \cdot W_m^{LPC}(n) \quad ; \quad 0 \leq n \leq L_{LPC} - 1. \quad (5.3.1.2-2)$$

16 The autocorrelation function is calculated according to:

$$17 \quad R_m(i) = \sum_{n=0}^{L_{LPC}-i-1} s_w^{LPC}(n) \cdot s_w^{LPC}(n+i) \quad ; \quad 0 \leq i \leq O_R. \quad (5.3.1.2-3)$$

18 The autocorrelation function is multiplied by a bandwidth expansion factor of 60 Hz:

$$19 \quad R_m(i) \leftarrow R_m(i) \cdot \exp\left(-0.5 \cdot \left(\frac{2\pi \cdot i \cdot 60}{F_s}\right)^2\right) \quad ; \quad 1 \leq i \leq 10, \quad (5.3.1.2-4)$$

20 where $F_s = 8000$ is the sampling frequency in Hz. High frequency compensation is
 21 performed by multiplying the first autocorrelation coefficient by 1.0001, which is equivalent
 22 to adding a noise floor of -40 dB.

23

24 5.3.1.2.3 Levinson-Durbin Algorithm

25 The reflection coefficients, the prediction coefficients, and the LP prediction error are
 26 calculated by the standard Levinson-Durbin algorithm.¹ The algorithm is given by (the
 27 index m was omitted for the local variables of the algorithm):

28

1 See: Rabiner, L. R. and Schafer, R. W., *Digital Processing of Speech Signals*, (New Jersey: Prentice-Hall Inc, 1978), pp. 411-412.

1
2
3
4

$$E^{(0)} = R_m(0)$$

$$i = 1$$

while ($i < 10$)

{

$$k_m(i) = - \frac{R_m(i) + \sum_{j=1}^{i-1} \alpha_j^{(i-1)} \cdot R_m(i-j)}{E^{(i-1)}}$$

$$\alpha_i^{(i)} = k_m(i)$$

$$j = 1$$

while ($j < i - 1$)

{

$$\alpha_j^{(i)} = \alpha_j^{(i-1)} + k_m(i) \cdot \alpha_{i-j}^{i-1}$$

$$j = j + 1$$

}

$$E^{(i)} = (1 - k_m^2(i)) E^{(i-1)}$$

$$i = i + 1$$

5

}

(5.3.1.2-5)

6 The LPC prediction coefficients for the m^{th} LPC analysis are given by:

7

$$a_m(i) = \alpha_i^{(10)} \quad ; \quad 1 \leq i \leq 10$$

8

(5.3.1.2-6)

9 5.3.1.2.4 Conversion to LSFs

10 The prediction error filter transfer function, $A(z)$, is given by (the index m omitted for
11 clarity):

$$12 \quad A(z) = 1 + \sum_{i=1}^{10} a(i) \cdot z^{-i} \quad . \quad (5.3.1.2-7)$$

13 Define two new transfer functions $P_A(z)$ and $Q_A(z)$ as

$$14 \quad P_A(z) = A(z) - z^{-11} A(z^{-1}) = 1 + p_1 z^{-1} + \dots + p_5 z^{-5} - p_5 z^{-6} - \dots - p_1 z^{-10} - z^{-11},$$

15 (5.3.1.2-8)

16 and

$$Q_A(z) = A(z) + z^{-11}A(z^{-1}) = 1 + q_1z^{-1} + \dots + q_5z^{-5} + q_5z^{-6} + \dots + q_1z^{-10}z^{-11}, \quad (5.3.1.2-9)$$

3 where

$$4 \quad p_i = a_i - a_{11-i} \quad ; 1 \leq i \leq 5, \quad (5.3.1.2-10)$$

5 and

$$6 \quad q_i = a_i + a_{11-i} \quad ; 1 \leq i \leq 5. \quad (5.3.1.2-11)$$

7 The LSF frequencies are the ten roots that exist between $\omega = 0.0$ and $\omega = 1.0$ in the
8 following two equations:

$$9 \quad P'(\omega) = \cos(5\pi \cdot \omega) + p'_1 \cdot \cos(4\pi \cdot \omega) + \dots + p'_4 \cdot \cos(\pi \cdot \omega) + \frac{p'_5}{2}, \quad (5.3.1.2-12)$$

$$10 \quad Q'(\omega) = \cos(5\pi \cdot \omega) + q'_1 \cdot \cos(4\pi \cdot \omega) + \dots + q'_4 \cdot \cos(\pi \cdot \omega) + \frac{q'_5}{2}, \quad (5.3.1.2-13)$$

11 where the parameters p'_i and q'_i are computed recursively from the parameters p_i and q_i
12 as

$$13 \quad p'_0 = q'_0 = 1, \quad (5.3.1.2-14)$$

$$14 \quad p'_i = p_i + p'_{i-1} \quad ; 1 \leq i \leq 5, \quad (5.3.1.2-15)$$

$$15 \quad q'_i = p_i - q'_{i-1} \quad ; 1 \leq i \leq 5. \quad (5.3.1.2-16)$$

16 Since the formant synthesis (LPC) filter is stable, the roots of the two functions alternate in
17 the range from 0.0 to 1.0. If these ten roots are denoted as $[lsf(1), lsf(2), \dots, lsf(10)]$ in the
18 increasing order of magnitude, then $lsf(i)$ for $i = 1, 3, 5, 7, 9$ are roots of $P'(\omega)$ and $lsf(i)$
19 for $i = 2, 4, 6, 8, 10$ are those of $Q'(\omega)$.

20 The LSFs calculated with the window centered at the second quarter of the encoding frame
21 are denoted $lsf_0(i)$, the LSFs calculated with the window centered at the fourth quarter of
22 the encoding frame are denoted $lsf_1(i)$, and the LSFs calculated with the window centered
23 at the lookahead are denoted $lsf_2(i)$.

24

25 5.3.1.2.5 LPC Prediction Error and LPC Prediction Gain

26 The LPC prediction error is calculated by:

$$27 \quad PE_m^{LPC} = R_m(0) \cdot \prod_{i=1}^{10} (1 - k_m^2(i)). \quad (5.3.1.2-17)$$

28 Where $m = 0, 1, 2$ indicates the LPC analysis performed on the second quarter of the
29 encoding frame, the last quarter of the encoding frame, and the lookahead, respectively.
30 The LPC prediction gain is obtained from LPC prediction error by:

$$G_m^{LPC} = -10 \cdot \log_{10} \left(\frac{PE_m^{LPC}}{R_m(0)} \right) \quad ; m = 0, 1, 2 \quad (5.3.1.2-18)$$

2

3 5.3.1.3 Interpolation of LPC Prediction Gains, Reflection Coefficients, and Prediction 4 Coefficients

5 The encoding frame is divided into 4 subframes. The interpolated LPC prediction gain is
6 obtained for each of the four subframes and the lookahead, indexed $m = 0, 1, 2, 3, 4$, by:

$$\begin{aligned} \hat{G}_0^{LPC} &= 0.5 \cdot (G_{1,prv}^{LPC} + G_0^{LPC}) \\ \hat{G}_1^{LPC} &= G_0^{LPC} \\ \hat{G}_2^{LPC} &= 0.5 \cdot (G_0^{LPC} + G_1^{LPC}) \\ \hat{G}_3^{LPC} &= G_1^{LPC} \\ \hat{G}_4^{LPC} &= G_2^{LPC} \end{aligned} \quad (5.3.1.3-1)$$

8 Where $G_{1,prv}^{LPC}$ is G_1^{LPC} of the previous frame.

9 Five vectors of interpolated reflection coefficients are generated from the 3 vectors of
10 reflection coefficients obtained during LPC analysis, by:

$$\begin{aligned} \hat{k}_0(i) &= 0.5 \cdot (k_{1,prv}(i) + k_0(i)) \\ \hat{k}_1(i) &= k_0(i) \\ \hat{k}_2(i) &= 0.5 \cdot (k_0(i) + k_1(i)) \quad ; i = 1, \dots, 10 \\ \hat{k}_3(i) &= k_1(i) \\ \hat{k}_4(i) &= k_2(i) \end{aligned} \quad (5.3.1.3-2)$$

12 where $k_{1,prv}(i)$ is $k_1(i)$ of the previous frame.

13 Five vectors of prediction coefficients, denoted by $\hat{a}_m(i)$, are generated from the 3 vectors of
14 prediction coefficients obtained during LPC analysis and the 2 vectors of LSFs generated for
15 the second and the fourth quarter of the encoding frame. Similar to the interpolation of the
16 reflection coefficients, $\hat{a}_1(i)$ is the vector of prediction coefficients obtained from the LPC
17 analysis on the second quarter of the encoding frame, $\hat{a}_3(i)$ is the vector of prediction
18 coefficients obtained from the LPC analysis on the fourth quarter of the encoding frame,
19 and $\hat{a}_4(i)$ is the vector of prediction coefficients obtained from the LPC analysis on the
20 lookahead. The interpolated vector $\hat{a}_0(i)$ is obtained by interpolating $lsf_{1,prv}(i)$ and $lsf_0(i)$,
21 where $lsf_{1,prv}(i)$ is $lsf_1(i)$ of the previous frame, and converting the interpolated LSF vector
22 to the vector of prediction coefficients. The interpolated vectors $\hat{a}_2(i)$ obtained by
23 interpolating $lsf_0(i)$ and $lsf_1(i)$, and converting the interpolated LSF vector to the vector of
24 prediction coefficients.

25

26

1 **5.3.2 Voice Activity Detection – Option A**

2 Routine name: VAD_voice_detection

3 A voice activity detector (VAD) is used to classify the input signal as either active voice
4 (speech) or inactive voice (silence or background noise). The VAD extracts several
5 parameters from the input signal and makes a voicing decision using a set of thresholds.

6 **Input:**

- 7 • Pre-processed speech, $\tilde{s}(n)$
- 8 • Reflection coefficients from the LPC analysis on the last quarter of the coding
9 frame, $k_1(i)$
- 10 • PE_1^{LPC} , the LP prediction error from the LPC analysis on the last quarter of the
11 coding frame
- 12 • Buffer of 5 previous open-loop pitch lags, $l_p^B(i) \quad i = 1, \dots, 5$
- 13 • Buffer of 5 previous pitch gains (unquantized adaptive-codebook gains),
14 $g_p^B(i) \quad i = 1, \dots, 5$
- 15 • Vector of LSFs from the LPC analysis on the last quarter of the coding frame,
16 $lsf_1(i)$
- 17 • Autocorrelation function from the LPC analysis on the last quarter of the coding
18 frame, $R_1(i)$.

19

20 **Output**

- 21 • Voice activity decision (initial value), VAD

22

23

24 **5.3.2.1 Computation of VAD Parameters**

25 This section describes the computation of the VAD parameters. The frame energy, E , is
26 calculated:

$$27 \quad E = \max \left(10, 10 \cdot \log_{10} \left(\frac{R_1(0)}{L_{LPC}} \right) \right). \quad (5.3.2.1-1)$$

28 A partial residual energy, E^{res} , is calculated:

$$29 \quad E^{res} = \max \left(0, 10 \cdot \log_{10} \left(\frac{R_1(0) \cdot \prod_{i=1}^4 (1 - k_1^2(i))}{L_{LPC}} \right) \right). \quad (5.3.2.1-2)$$

30 The frame LPC gain, G_{LPC} , is given by

$$1 \quad G_{LPC} = -10 \cdot \log_{10} \left(\prod_{i=1}^{10} (1 - k_1^2(i)) \right). \quad (5.3.2.1-3)$$

2 Pitch lags standard deviation is calculated by:

$$3 \quad \sigma_p = \sqrt{\frac{1}{4} \sum_{i=1}^5 \left(l_p^B(i) - \frac{1}{5} \sum_{k=1}^5 l_p^B(k) \right)^2}. \quad (5.3.2.1-4)$$

4 A running mean of the pitch gain is given by:

$$5 \quad \bar{g}_p = 0.8 \cdot \bar{g}_p + 0.2 \cdot \left(\frac{1}{5} \cdot \sum_{i=1}^5 g_p^B(i) \right). \quad (5.3.2.1-5)$$

6 A periodicity flag, F_p , is set, using σ_p and \bar{g}_p . Two average vectors of the LSFs are given
7 by (the negative index j denotes the saved previous frame LSFs, and the index 0
8 corresponds to the LSFs of the current frame):

$$9 \quad \overline{lsf}_1(i) = \frac{1}{4} \sum_{j=-3}^{j=0} lsf_j^B(i) \quad \text{and} \quad \overline{lsf}_2(i) = \frac{1}{3} \sum_{j=-3}^{j=-1} lsf_j^B(i) \quad \text{for } i=1, \dots, 10 \quad (5.3.2.1-6)$$

10 A normalized maximum measure is calculated

$$11 \quad M_N = \frac{\max_{i=1}^{160} |\tilde{s}(i)|}{\bar{M}_N}, \quad (5.3.2.1-7)$$

12 where \bar{M}_N is a running mean of M_N , updated during silence segments. Three spectral
13 distortion measures are obtained. The first measure is given by:

$$14 \quad SD_1 = \frac{\bar{a}^T \cdot \mathbf{R} \cdot \bar{a}}{PE_1^{LPC}}, \quad (5.3.2.1-8)$$

15 where \bar{a} is the LP prediction coefficients vector obtained from \overline{lsf}_2 , and \mathbf{R} is a symmetric
16 Toeplitz matrix obtained for R . The second measure is given by

$$17 \quad SD_2 = \sum_{i=1}^{10} \left(lsf_1(i) - \overline{lsf}_N(i) \right)^2, \quad (5.3.2.1-9)$$

18 where $\overline{lsf}_N(i)$ is a running mean of the LSFs during silence segments (Section 5.3.2.4).
19 The third spectral distortion measure is given by:

$$20 \quad SD_3 = \sum_{i=1}^{10} \left(lsf_1(i) - \overline{lsf}(i) \right)^2, \quad (5.3.2.1-10)$$

21 where \overline{lsf} is a continuous running mean of the LSFs. A minimum residual energy, E_{\min}^{res} , is
22 tracked over a sliding window of the last 64 frames. The minimum is tracked on 8
23 contiguous sub-segment for fast implementation. An active-speech mean energy is updated
24 if $E^{res} > E_{\min}^{res} + 12$ according to:

$$\bar{E}_{speech}^{res} = 0.9 \cdot \bar{E}_{speech}^{res} + 0.1 \cdot E^{res} \quad (5.3.2.1-11)$$

A running mean of E^{res} that tracks the noise floor, \bar{E}_{noise}^{res} , is updated during inactive frames (Section 5.3.2.4). A signal-to-noise ratio measure is estimated by:

$$SNR = \bar{E}_{speech}^{res} - \bar{E}_{noise}^{res} \quad (5.3.2.1-12)$$

An onset flag F_{onset} is set according to energy increases and spectral changes. The running mean and history parameters are either initialized to default values or averages of their updating values over the first 32 frames.

8

9 **5.3.2.2 VAD Decision**

10 An initial voice activity decision is taken by comparing the VAD parameters to each other
11 and to a set of fixed thresholds, controlled by the various flags. For the first 32 frames the
12 thresholds are fixed, since the SNR convergence is not yet obtained, while for all
13 subsequent frames a combination of fixed and adaptive thresholds, controlled by the SNR ,
14 dictates the initial VAD decision.

15

16 **5.3.2.3 Hangover**

17 An active speech VAD segment is further extended to avoid speech clipping. The energy E ,
18 F_{onset} and SNR are used to control the amount of hangover used to extend the active speech
19 decision to the current frame if the initial decision were inactive.

20

21 **5.3.2.4 Parameter Update During Silence/Noise**

22 The running-mean parameters that track the silence are updated if either the extended
23 VAD decision is inactive speech or ($E^{res} > \max(\bar{E}_{noise}^{res}, E_{min}^{res}) + 3.5$) and F_p indicates no
24 periodicity. The running-mean parameters are updated using an auto-regressive scheme.

25

26 **5.3.2.5 Reset of Noise Floor Energy**

27 The noise floor energy, \bar{E}_{noise}^{res} is set for any frame subsequent to the 64th frame, according
28 to:

$$\bar{E}_{noise}^{res} = \begin{cases} \max(0.75 \cdot \max(0.5 \cdot (E_{min}^{res} + \bar{E}_{noise}^{res}), 0) + 0.25 \cdot E_{min}^{res}, 0) & ; \bar{E}_{noise}^{res} < E_{min}^{res} - 10 \text{ and } F_p = 0 \\ \max(E_{min}^{res}, 0) & ; \bar{E}_{noise}^{res} > E_{min}^{res} + 4 \\ \bar{E}_{noise}^{res} & ; \text{otherwise} \end{cases}$$

30

$$(5.3.2.5-1)$$

31

32

1 **5.3.3 Voice Activity Detection – Option B**

2 This scheme is used to select the voice activity decision (VAD) of the encoded frame prior to
3 signal classification and rate-selection. Active speech is encoded at Rate 1, Rate 1/2, or
4 Rate 1/4 and background noise is encoded at Rate 1/8. The actual data rate used for
5 encoding the input signal is determined later by the signal classifier module.

6 **Inputs:**

- 7 • Bandwidth expanded correlation coefficients, R_1
- 8 • Long-term prediction gain, $\beta = R_p^{ol}(1)$
- 9 • SMV Operating Mode, M_{SMV}

10 **Outputs:**

- 11 • Voice activity decision (initial value), VAD

12 **Initialization:**

- 13 • Local VAD Rate histories, $VAD_Rate(frm-1)$ and $VAD_Rate(frm-2)$, are set to Rate 1/8
14 frames.
- 15 • .
- 16 • Band energy, $E_{f(i)}^{sm}(frm)$, initialization is given in 5.3.3.2.1.
- 17 • Background noise estimate, $B_{f(i)}(frm)$, initialization is given in 5.3.3.2.2.
- 18 • Signal energy estimate, $S_{f(i)}(frm)$, initialization is given in 5.3.3.2.3.

19

20 **5.3.3.1 Estimating the Data Rate and VAD Based on Current Signal Parameters**

21 A data rate local to Voice Activity Detection Option B is computed to aid the determination
22 of the VAD. The local VAD data rate is determined by comparing the current frame energy
23 in each of two frequency bands, $f(1)$ and $f(2)$, to background noise estimates in these
24 respective bands. Thresholds above the background noise in each band are determined by
25 an estimated signal-to-noise ratio in each band. These thresholds are set for Rate 1, Rate
26 1/2, and Rate 1/8 encoding. The highest rate calculated from the two frequency bands is
27 then selected as the rate for the current frame. If Rate 1/8 encoding was selected VAD is
28 set to 0, and if otherwise VAD is set to 1.

29 5.3.3.1.1 Computing Band Energy

30 The rate determination algorithm uses energy thresholds to determine the rate for the
31 current frame. The input speech is divided into two distinct bands: band $f(1)$ spans 0.3-
32 2.0 kHz, band $f(2)$ spans 2.0-4.0 kHz.† The band energy for band $f(i)$, $B_{f(i)}$, is
33 calculated as

† Whenever a variable (or symbol or value) with a subscript $f(i)$ appears in any equation or pseudocode in 5.3.3.1.4, 5.3.3.2.2, and 5.3.3.2.3, it refers to a variable (or symbol or value) associated with either band $f(1)$ or band $f(2)$.

$$BE_{f(i)} = R_1(0)R_{f(i)}(0) + 2.0 \sum_{k=1}^{L_h-1} R_1(k)R_{f(i)}(k), \quad (5.3.3.1-1)$$

2 where:

$$R_{f(i)} = \sum_{n=0}^{L_h-1-k} h_i(n)h_i(n+k); \quad 0 \leq k < L_h, \quad (5.3.3.1-2)$$

4 and $h_i(n)$ is the impulse response of the band-pass filter i , where $i = 1, 2$. $R_1(k)$ is the
5 autocorrelation sequence defined in Section 5.3.1.2.2, and $L_h = 17$ is the length of the
6 impulse response of the band-pass filters.

7 The band-pass filters used for both frequency bands are defined in Table 5.3-1.

8

n	h₁(n) (lower band)	n	h₂(n) (upper band)
0	-0.0555725097656250	0	-0.0122985839843750
1	-0.0721588134765625	1	0.0437622070312500
2	-0.0103759765625000	2	0.0123901367187500
3	0.0234527587890625	3	-0.0624389648437500
4	-0.0607147216796875	4	-0.0124511718750000
5	-0.1398925781250000	5	0.1053619384765625
6	-0.0122528076171875	6	0.0124816894531250
7	0.2799224853515625	7	-0.3180694580078125
8	0.4375000000000000	8	0.4875030517578125
9	0.2799224853515625	9	-0.3180694580078125
10	-0.0122528076171875	10	0.0124816894531250
11	-0.1398925781250000	11	0.1053619384765625
12	-0.0607147216796875	12	-0.0124511718750000
13	0.0234527587890625	13	-0.0624389648437500
14	-0.0103759765625000	14	0.0123901367187500
15	-0.0721588134765625	15	0.0437622070312500
16	-0.0555725097656250	16	-0.0122985839843750

9

10 **Table 5.3-1 FIR Filter Coefficients Used for Band Energy Calculations**

11 5.3.3.1.2 Calculating Rate Determination Thresholds

12 The rate determination thresholds for each frequency band $f(i)$ are a function of both the
13 background noise estimate, $B_{f(i)}(frm-1)$, and the signal energy estimate, $B_{f(i)}(frm-1)$, of
14 the previous or $(frm-1)^{th}$ frame. Two thresholds for each band are computed as

$$T_1(B_{f(i)}(frm-1), SNR_{f(i)}(frm-1)) = k1(SNR_{f(i)}(frm-1))B_{f(i)}(frm-1), \quad (5.3.3.1-3)$$

$$T_2(B_{f(i)}(frm-1), SNR_{f(i)}(frm-1)) = k2(SNR_{f(i)}(frm-1))B_{f(i)}(frm-1), \quad (5.3.3.1-4)$$

where the integer $SNR_{f(i)}(frm-1)$ is:

$$SNR_{f(i)}(frm-1) = \begin{cases} 0 & ; QSNRU_{f(i)}(frm-1) < 0 \\ QSNRU_{f(i)}(frm-1) & ; 0 \leq QSNRU_{f(i)}(frm-1) \leq 7 \\ 7 & ; QSNRU_{f(i)}(frm-1) > 7 \end{cases}, \quad (5.3.3.1-5)$$

and where:

$$QSNRU_{f(i)}(frm-1) = \text{round} \left\{ \frac{(10 \log_{10}(S_{f(i)}(frm-1)B_{f(i)}(frm-1)) - 20)}{5} \right\}, \quad (5.3.3.1-6)$$

The functions $k1(\bullet)$ and $k2(\bullet)$ for the different SMV modes are defined in Table 5.3-2, Table 5.3-3 and Table 5.3-4, and $B_{f(i)}(frm-1)$ and $S_{f(i)}(frm-1)$ are defined in Sections 5.3.3.2.2 and 5.3.3.2.3, respectively.

SNR_{f(i)}(m-1)	k1(SNR_{f(i)}(m-1))	k2(SNR_{f(i)}(m-1))
0	7.0	9.0
1	7.0	12.5937
2	8.0	17.0
3	8.59375	18.5
4	8.90625	19.40625
5	9.40625	20.90625
6	22.0	51
7	63.1875	159.1875

Table 5.3-2 Threshold Scale Factors as a Function of SNR, $M_{SMV} = 0$

1

$SNR_{f(i)(m-1)}$	$k1(SNR_{f(i)(m-1)})$	$k2(SNR_{f(i)(m-1)})$
0	7.0	9.0
1	7.0	12.59375
2	8.0	17.0
3	8.59375	18.5
4	8.90625	19.40625
5	11.28125	25.09375
6	33.0	76.5
7	126.40625	318.40625

2

Table 5.3-3 Threshold Scale Factors as a Function of SNR, $M_{SMV} = 1$

3

$SNR_{f(i)(m-1)}$	$k1(SNR_{f(i)(m-1)})$	$k2(SNR_{f(i)(m-1)})$
0	7.0	9.0
1	8.6875	15.8125
2	9.0	18.5
3	10.0	20.0
4	10.59375	22.8125
5	13.8125	31.59375
6	50.1875	158.5
7	200.0	631.0

4

Table 5.3-4 Threshold Scale Factors as a Function of SNR, $M_{SMV} = 2, 3$

5 The threshold scale factors are identical for the low- and high-frequency bands.

6 5.3.3.1.3 Comparing Thresholds

7 Band energy, $BE_{f(i)}$, is compared with two thresholds: $T_1(B_{f(i)}(frm-1), SNR_{f(i)}(frm-1))$
8 and $T_2(B_{f(i)}(frm-1), SNR_{f(i)}(frm-1))$. If $BE_{f(i)}$ is greater than both thresholds, Rate 1 is
9 selected. If $BE_{f(i)}$ is greater than only one threshold, Rate 1/2 is selected. If $BE_{f(i)}$ is at or
10 below both thresholds, Rate 1/8 is selected. This procedure is performed for both
11 frequency bands and the higher of the two VAD rates selected from the individual bands is
12 chosen as the preliminary VAD rate, VAD_Rate , for the current frame.

13 5.3.3.1.4 Performing Hangover

14 If the preliminary VAD data rate decision (from Section 5.3.3.1.3) transitions from at least
15 two consecutive Rate 1 frames to a lower VAD rate frame, then the next M frames are
16 selected as Rate 1 before allowing the VAD rate to drop to Rate 1/2 and finally to Rate 1/8.
17 The number of hangover frames, M , is a function of the $SNR_{f(i)}(frm-1)$ (the SNR in the

1 lower frequency band) and is denoted as $Hangover(SNR_{f(l)}(frm-1))$ for the different SMV
 2 modes in Table 5.3-5, Table 5.3-6 and Table 5.3-7. $SNR_{f(l)}(frm-1)$ is calculated as
 3 defined in Equation (5.3.3.1-5). The hangover algorithm is defined by the following
 4 pseudocode:

```

5   {
6     if ( VAD_Rate(frm) == Rate 1 )
7       count = 0
8     if (VAD_Rate(frm-1) == Rate 1 and VAD_Rate(frm-2) == Rate 1 and VAD_Rate(frm) !=
9     Rate 1) {
10      if ( count == 0 )
11        M = Hangover(SNRf(1)(frm-1))
12      if ( count < M ) {
13        VAD_Rate(frm) = Rate 1
14        count = count + 1
15      }
16    }
17  }

```

18 where $VAD_Rate(frm)$ is the VAD rate of the current frame and $VAD_Rate(frm-1)$ and
 19 $VAD_Rate(frm-2)$ are the VAD rates of the previous two frames, respectively. Also, the VAD
 20 rates for the previous frames, as used here, are those generated by the RDA prior to
 21 override logic, and are, therefore, not subject to modification by external rate control as
 22 described in Section 5.3.3.1.5.

23

$SNR_{f(1)}(m-1)$	Hangover ($SNR_{f(1)}(m-1)$)
0	6
1	6
2	5
3	3
4	0
5	0
6	0
7	0

24

Table 5.3-5 Hangover Frames as a Function of SNR, $M_{SMV} = 0$

25

1

$SNR_{f(1)(m-1)}$	Hangover ($SNR_{f(1)(m-1)}$)
0	5
1	4
2	3
3	2
4	0
5	0
6	0
7	0

2

Table 5.3-6 Hangover Frames as a Function of SNR, $M_{SMV} = 1$

3

$SNR_{f(1)(m-1)}$	Hangover ($SNR_{f(1)(m-1)}$)
0	5
1	4
2	3
3	2
4	0
5	0
6	0
7	0

4

Table 5.3-7 Hangover Frames as a Function of SNR, $M_{SMV} = 2, 3$

5.3.3.1.5 Constraining VAD Rate Selection

The VAD rate selected by the procedures described in Sections 5.3.3.1.3 and 5.3.3.1.4 are used for the current frame except where it is modified by the following constraints.

If the previous frame was selected as Rate 1 and the current frame is selected as Rate 1/8, then the VAD rate of the current frame should be modified to Rate 1/2. There are no other restrictions on VAD rate transitions.

5.3.3.2 Updating RDA Parameters

After the RDA is complete, RDA parameters shall be updated as described in Sections 5.3.3.2.1 through 5.3.3.2.3.

5.3.3.2.1 Updating the Smoothed Band Energy

1 The band energy, $BE_{f(i)}$, calculated in Equation (5.3.3.1-1) is smoothed and used to
 2 estimate both the background noise energy (see Section 5.3.3.2.2) and signal energy (see
 3 Section 5.3.3.2.3) in each band. The smoothed band energy, $E_{f(i)}^{sm}(frm)$, is computed as:

$$4 \quad E_{f(i)}^{sm}(frm) = \alpha_E(frm)E_{f(i)}^{sm}(frm-1) + (1 - \alpha_E(frm))BE_{f(i)} \quad ; 1 \leq i \leq 2, \quad (5.3.3.2-1)$$

5 where frm refers to the current frame, and

$$6 \quad \alpha_E(frm) = \begin{cases} 0 & ; frm \leq 1 \\ 0.6 & ; frm > 1 \end{cases} \quad (5.3.3.2-2)$$

7 This allows the smoothed band energy to be initialized to the band energy of the first frame
 8 ($frm = 0$).

9

10 5.3.3.2.2 Updating the Background Noise Estimate

11 An estimate of the background noise level, $B_{f(i)}(frm)$, is computed for the current, or
 12 frm th, frame using $B_{f(i)}(frm-1)$, $E_{f(i)}^{sm}(frm)$ (see Section 5.3.3.2.1) and $SNR_{f(i)}(frm-1)$
 13 (see Section 5.3.3.1.2). Pseudocode describing the background noise update for band $f(i)$
 14 is given as

```

15 {
16     if (  $\beta < 0.30$  for 4 or more consecutive frames )
17          $B_{f(i)}(m) = \min\{ E_{f(i)}^{sm}(m), 80954304, \max\{ 1.06B_{f(i)}(m-1), B_{f(i)}(m-1)+ 1 \} \}$ 
18     else {
19         if (  $SNR_{f(i)}(m-1) > 1$  )
20              $B_{f(i)}(m) = \min\{ E_{f(i)}^{sm}(m), 80954304, \max\{ 1.00547B_{f(i)}(m-1), B_{f(i)}(m-$ 
21              $1)+ 1 \} \}$ 
22         else
23              $B_{f(i)}(m) = \min\{ E_{f(i)}^{sm}(m), 80954304, B_{f(i)}(m-1) \}$ 
24     }
25     if (  $B_{f(i)}(m) < lownoise(i)$  )
26          $B_{f(i)}(m) = lownoise(i)$ 
27 }
```

28 where $\beta = R_p^{ol}(1)$, is defined in Section 5.3.9.2, and $E_{f(i)}^{sm}(frm)$, as well as $SNR_{f(i)}(frm-1)$
 29 are defined in Section 5.3.3.2.1 and Section 5.3.3.1.2, respectively. $lownoise(1)$ equals
 30 160.0 and $lownoise(2)$ equals 80.0.

31 At initialization, the background noise estimate for the first frame, $B_{f(i)}(0)$, is set to
 32 80,954,304 for both frequency bands, and the consecutive frame counter for β is set to
 33 zero. Initialization also occurs if the audio input to the encoder is disabled and then
 34 enabled.†

35

† This prevents the silence before the audio is connected from being mistaken as unusually low background noise.

1 5.3.3.2.3 Updating the Signal Energy Estimate

2 The signal energy, $S_{f(i)}(frm)$, is computed as

3 {
 4 If ($\beta > 0.5$ for 5 or more consecutive frames)
 5 $S_{f(i)}(frm) = \max\{ E^{sm}_{f(i)}(frm), \min\{ 0.97S_{f(i)}(frm-1), S_{f(i)}(frm-1) - 1 \} \}$
 6 else
 7 $S_{f(i)}(frm) = \max\{ E^{sm}_{f(i)}(frm), S_{f(i)}(frm-1) \}$
 8 }

9 where $\beta = R_p^{ol}(1)$, is defined in Section 5.3.9.2, and $E_{f(i)}^{sm}(frm)$ is defined in Section
 10 5.3.3.2.1.

11 At initialization, the signal energy estimates for the first frame, $S_{f(1)}(frm)$ and $S_{f(2)}(frm)$,
 12 are set to 32000 and 3200, respectively, and the consecutive frame counter for β is set to
 13 zero. Initialization also occurs if the audio input to the encoder is disabled and then
 14 enabled.
 15

1

2 **5.3.4 Music Detection**3 Routine name: MUS_musdetect

4 A music detector is used to classify the input signal as either music or speech. The
 5 detection takes place after the VAD decision is made. The music detector reuses some of
 6 the parameters of the VAD and extracts additional parameters from the speech. It makes a
 7 music detection decision and also corrects the initial voice activity decision from the VAD
 8 module using a set of thresholds.

9 **5.3.4.1 Input to Music Detector**

10 In addition to several parameters that are obtained in the VAD option A, or that are
 11 calculated independently if VAD option B is used, the music detector uses the following
 12 input:

- 13 • a buffer of 5 previous normalized pitch correlations, $corr_p^B(i) \quad i = 1, \dots, 5$

14

15 **5.3.4.2 Computation of Music Detector Parameters**

16 This section describes the computation of the music detector parameters. A difference Δ_{lsf}
 17 between $lsf_1(2)$ and $lsf_1(1)$, obtained from the LPC analysis on the last quarter of the
 18 coding frame, is computed. A running mean of $lsf_1(1)$ is computed as:

$$19 \quad \overline{lsf}(1) = 0.75 \cdot \overline{lsf}(1) + 0.25 \cdot lsf_1(1) \quad (5.3.4.2-1)$$

20 A running mean energy, \overline{E} , is calculated as:

$$21 \quad \overline{E} = 0.75 \cdot \overline{E} + 0.25 \cdot E \quad (5.3.4.2-2)$$

22 where E is the frame energy. A spectral difference SD_4 is calculated as:

$$23 \quad SD_4 = \sum_{i=1}^{10} (k_1(i) - \overline{k}_N(i))^2 \quad (5.3.4.2-3)$$

24 where \overline{k}_N is the running mean reflection coefficients of noise/silence. The running mean of
 25 the partial residual energy \overline{E}_N^{res} is updated along \overline{k}_N when the input VAD is inactive as:

$$26 \quad \overline{E}_N^{res} = 0.9 \cdot \overline{E}_N^{res} + 0.1 \cdot E^{res} \quad (5.3.4.2-4)$$

27 and

$$28 \quad \overline{k}_N(i) = 0.75 \cdot \overline{k}_N(i) + 0.25 \cdot k_1(i) \quad i = 1, \dots, 10 \quad (5.3.4.2-5)$$

29 A running mean of the normalized pitch correlation is given by:

$$30 \quad \overline{corr}_p = 0.8 \cdot \overline{corr}_p + 0.2 \cdot \left(\frac{1}{5} \cdot \sum_{i=1}^{i=5} corr_p^B(i) \right) \quad (5.3.4.2-6)$$

1 A periodicity flag F_p , similar to the one in the VAD module, is calculated using
 2 $corr_p^B(\cdot)$ and different thresholds. A spectral continuity counter c_{sp} is incremented if
 3 $k_1(2) \geq 0.0$ and $\overline{corr}_p < 0.5$ and reset to 0 otherwise. A periodicity continuity counter c_{pr}
 4 is incremented each time F_p is set and reset to 0 every 32 frames. A running mean of the
 5 periodicity counter \bar{c}_{pr} is updated every 32 frames as:

$$6 \quad \bar{c}_{pr} = \alpha \cdot \bar{c}_{pr} + (1 - \alpha) \cdot c_{pr} \quad (5.3.4.2-7)$$

7 where

$$8 \quad \alpha = \begin{cases} 0.98 & c_{pr} > 12 \\ 0.95 & c_{pr} > 10 \\ 0.90 & \text{otherwise} \end{cases} \quad (5.3.4.2-8)$$

9 A very low frequency noise flag F_f is set if the initial VAD is inactive and either
 10 $lsf_1(1) < 110$ Hz or $\overline{lsf}(1) < 150$ Hz. The initial inactive VAD decision from the VAD module
 11 is corrected to an active VAD decision by comparing SD_4 , E^{res} , \bar{E}_N^{res} , E , and \bar{c}_{pr} to a set
 12 of thresholds. A noise continuity counter C_N is incremented each time the corrected VAD
 13 is inactive and is reset otherwise. A running mean of the normalized pitch correlation
 14 \overline{corr}_p^N is updated if either the corrected VAD is inactive or F_f is set:

$$15 \quad \overline{corr}_p^N = 0.8 \cdot \overline{corr}_p^N + 0.2 \cdot \left(\frac{1}{5} \cdot \sum_{i=1}^{i=5} corr_p^B(i) \right) \quad (5.3.4.2-9)$$

16 A music continuity counter c_M is adaptively incremented and decremented by comparing
 17 the MD parameters to each other and to a set of fixed thresholds, controlled by the various
 18 flags. A running mean of this counter \bar{c}_M is updated as:

$$19 \quad \bar{c}_M = 0.9 \cdot \bar{c}_M + 0.1 \cdot c_M \quad (5.3.4.2-10)$$

20 **5.3.4.3 Music Detection Decision**

21 The music detection flag F_{MUS} is set if either $\bar{c}_{pr} \geq 18$ or $\bar{c}_M > 200$. In this case, \bar{E}_N^{res} is
 22 reset to 0.

23 **5.3.4.4 Reset of tracking counters**

24 \bar{c}_{pr} , c_{pr} , and c_{sp} are reset to 0 if either $E < 13$ dB or F_f is set or $c_{pr} > 50$, or $c_{sp} > 20$.

25 **5.3.4.5 Update of pitch lag, pitch correlation and pitch gain buffers**

26 The pitch lag buffer $l_p^B(\cdot)$ is updated on a frame basis as:

$$27 \quad l_p^B(i) = l_p^B(i+2), \quad i = 0,1,2 \quad (5.3.4.5-1)$$

1 and

$$2 \quad l_p^B(3) = l_p^{ol}(0), l_p^B(4) = l_p^{ol}(1) \quad (5.3.4.5-2)$$

3 Similarly, the pitch correlation buffer $corr_p(\cdot)$ is updated on a frame basis as:

$$4 \quad corr_p^B(i) = corr_p^B(i+2), i = 0,1,2 \quad (5.3.4.5-3)$$

5 and

$$6 \quad corr_p^B(3) = R_p^{ol}(0), corr_p^B(4) = R_p^{ol}(1) \quad (5.3.4.5-4)$$

7 On the other hand, the pitch gain buffer $g_p^B(\cdot)$ is updated as on a subframe basis during
8 Rate 1 and Rate 1/2 as:

$$9 \quad g_p^B(i) = g_p^B(i+1), i = 0,1,2,3, \quad (5.3.4.5-5)$$

10 and

$$11 \quad g_p^B(4) = g_a, \quad (5.3.4.5-6)$$

12 where g_a is the unquantized adaptive-codebook gain.

13

14

15

1 5.3.5 Calculation of Weighting Coefficient for 4-Subframe Structure

2 Routine name: PWF_WeighingFilters

3

4 **Input:**

- 5 • First element of the interpolated reflection coefficient vector, $\hat{k}_m(0)$, for
- 6 $m = 0,1,2,3,4$
- 7 • Interpolated LPC prediction gain, \hat{G}_m^{LPC} , for $m = 0,1,2,3,4$
- 8 • Interpolated prediction coefficients, $\hat{a}_m(i)$, for $m = 0,1,2,3,4$

9 **Output:**

- 10 • Weighing filter numerator coefficients, $\hat{a}_m^{num}(i)$ for $m = 0,1,2,3,4$
- 11 • Weighing filter denominator coefficients, $\hat{a}_m^{den}(i)$, for $m = 0,1,2,3,4$

12

13 The weighting filter coefficients for the four subframes and the lookahead, indexed by
 14 $m = 0,1,2,3,4$, are calculated by the following procedure. A weighting coefficient for the
 15 numerator is given by:

$$16 \quad w_m^{num} = 0.9 + \min\left(0.03 \cdot \max\left(\hat{G}_m^{LPC} - 5, 0\right) \cdot \max\left(-\hat{k}_m(0) - 0.5, 0\right), 0.04\right) \quad (5.3.5-1)$$

17 and a weighting coefficient for the denominator is given by

$$18 \quad w_m^{den} = 0.6 + 1.25 \cdot \min\left(\hat{k}_m(0) + 0.9, 0\right) \quad (5.3.5-2)$$

19 The filter coefficient for the numerator are calculated by:

$$20 \quad \hat{a}_m^{num}(i) = \hat{a}_m(i) \cdot \left(w_m^{num}\right)^i \quad i = 1, \dots, 10 \quad (5.3.5-3)$$

21 The filter coefficients for the denominator are calculated by:

$$22 \quad \hat{a}_m^{den}(i) = \hat{a}_m(i) \cdot \left(w_m^{den}\right)^i \quad i = 1, \dots, 10 \quad (5.3.5-4)$$

23

1 **5.3.6 Generating the Weighted Residual**

2 Routine name: PWF_speech_to_residu

3 **Input:**

- 4 • Weighing filter numerator coefficients, $\hat{a}_m^{num}(i)$ for $m = 0,1,2,3,4$
- 5 • Pre-processed speech, $\tilde{s}(n)$

6 **Output:**

- 7 • Weighted residual, $r_w(n)$
- 8

9 For each subframe, indexed by $m = 0,1,2,3,4$, the weighted residual signal, $r_w(n)$, is
 10 generated from the speech by the filter given by

11
$$H_m(z) = 1 + \sum_{i=1}^{10} \hat{a}_m^{num}(i)z^{-i} \quad (5.3.6-1)$$

12 The filter parameters, $\hat{a}_m^{num}(i)$, are the filter coefficients for each of the four 5 ms subframes
 13 and the lookahead, calculated at Section 5.3.5.

14
 15

1 **5.3.7 Initial Setting of Voiced/Unvoiced Level**

2 Function name: CLA_NoiseUnvoicedDetect

3 **Input:**

- 4 • Pre-processed speech, $\tilde{s}(n)$
- 5 • Weighted residual, $r_w(n)$

6 **Output:**

- 7 • Voiced/Unvoiced level (initial value), C_{VUV}
- 8 • Energies of weighted residual for two halves of the frame, $E_{rw}(0)$ and $E_{rw}(1)$
- 9 • Sharpness of weighted residual, S_{rp}^{wr}

10

11 The initial setting of voiced/unvoiced level of speech, C_{VUV} , is to the values 0, 1, or 2. The
12 value of the parameter is further refined, and extended to include level 3, in Sections 5.3.12
13 and 5.3.15.

14

15 **5.3.7.1 Calculation of Decision Parameters**

16 The sharpness of the weighted residual is calculated by:

$$17 \quad S_{rp}^{wr} = \frac{\frac{1}{L_F} \sum_{n=0}^{L_{frm}-1} |r_w(n)|}{\max_{n=0, \dots, L_{frm}-1} |r_w(n)|}. \quad (5.3.7.1-1)$$

18 The first order prediction coefficient of the speech on the second half of the frame is given
19 by:

$$20 \quad \check{p} = \frac{\sum_{n=L_{frm}/2}^{L_{frm}-2} \tilde{s}(n) \cdot \tilde{s}(n+1)}{\sum_{n=L_{frm}/2}^{L_{frm}-1} \tilde{s}^2(n)} \quad (5.3.7.1-2)$$

21 The zero-crossing rate of the speech signal on the second half of the frame, C_{zc} , is
22 calculated by evaluating the number of sign changes of the speech signal over the second
23 half of the frame, and normalizing it with the half-frame size.

24 The energies of the weighted residual on the first half and on the second half of the frame
25 are given by:

$$26 \quad E_{rw}(0) = \sum_{n=0}^{L_{frm}/2-1} r_w^2(n) \quad E_{rw}(1) = \sum_{n=L_{frm}/2}^{L_{frm}-1} r_w^2(n) \quad (5.3.7.1-3)$$

1 The energy of the signal of the first $\frac{3}{4}$ of the frame is calculated by:

$$2 \quad E_{\frac{3}{4}} = \frac{1}{\frac{3}{4} L_{frm}} \sum_{n=0}^{\frac{3}{4} L_{frm} - 1} \tilde{s}^2(n) \quad (5.3.7.1-4)$$

3 The number of speech samples in the frame that has an absolute value less than 0.1 is
 4 counted, and denoted $C_{|0.1|}$.

5 **5.3.7.2 Initial Setting Voiced/Unvoiced Level**

6 Six of the 7 parameters (E_2 excluded) are compared to a set of thresholds, for an initial
 7 setting of the voiced/unvoiced level, C_{VUV} , for the frame. The initial setting is 0, 1, or 2.

8

1 **5.3.8 Generating the Weighted Speech**

2 Routine name: PWF_residu_to_wspeech

3 **Input:**

- 4 • Weighted residual, $r_w(n)$
- 5 • Weighing filter denominator coefficients, $\hat{a}_m^{den}(i)$, for $m = 0,1,2,3,4$
- 6 • First element of the interpolated reflection coefficient vector, $\hat{k}_m(0)$, for
- 7 $m = 0,1,2,3,4$
- 8 • Spectral flatness flag, F_{flat}

9 **Output:**

- 10 • Weighted speech, $s_w(n)$
- 11 • Adaptive low-pass filter parameters, μ_m , for $m = 0,1,2,3,4$

12 The weighted speech, s_w , is generated from the weighted residual signal with the
 13 application of a fixed weighted synthesis filter and an adaptive low-pass filter. The transfer
 14 function of the weighted synthesis filter for each subframe, indexed by $m = 0,1,2,3,4$, is
 15 given by:

$$16 \quad H_m(z) = \frac{1}{1 + \sum_{i=1}^{10} a_m^{den}(i) \cdot z^{-i}} \quad (5.3.7.2-1)$$

17 The weighting filter parameters, $a_m^{den}(i)$, are the coefficients for each of the four 5 ms
 18 subframes and the lookahead calculated in Section 5.3.5. The adaptive low-pass filter is
 19 given by:

$$20 \quad L(z) = 1 + \mu_m \cdot z^{-1} \quad (5.3.7.2-2)$$

21 Where μ_m is given by:

$$22 \quad \mu_m = \begin{cases} \max\{\min\{\hat{k}_m(0) + 0.5, 0.3\}, 0\} & F_{flat} = 0 \\ 0 & F_{flat} = 1 \end{cases} \quad (5.3.7.2-3)$$

23

24

1 **5.3.9 Open-Loop Pitch Detection**

2 Routine name: PIT_ol_2pitches

3 **Input:**

- 4 • Weighted speech, $s_w(n)$
- 5 • Previous frame class, C_{FR}^{prv}
- 6 • Previous frame open-loop pitch lags
- 7 • Previous frame open-loop pitch normalized correlation

8 **Output:**

- 9 • Open-loop pitch lags, $l_p^{ol}(m)$, $m = 0,1,2$
- 10 • Open-loop pitch normalized correlation, $R_p^{ol}(m)$, $m = 0,1,2$

11

12 The open-loop pitch detection module calculates 3 open-loop pitch values, for the first half
 13 of the encoding frame, the second half of the encoding frame, and the lookahead. It
 14 searches and determines the open-loop pitch lags for the look ahead and for the second
 15 half of the encoding frame. It obtains the open-loop pitch lag for the first half of the
 16 encoding frame by refining the open-loop pitch lag of the previous frame lookahead open-
 17 loop pitch lag.

18 **5.3.9.1 Open-Loop Pitch Detection on the Look Ahead Buffer**

19 The open-loop pitch on the lookahead buffer is estimated by selecting 4 open-loop pitch
 20 candidates on 4 intervals, selecting an initial best candidate, and selecting a final best
 21 candidate using smoothing.

22 5.3.9.1.1 Selection of 4 Open-Loop Pitch Candidates

23 The normalized correlation of the weighted speech is calculated for the interval [17,148].
 24 The normalized correlation is given by:

$$25 \quad R_p(l) = \frac{\sum_{n=0}^{79} s_w(n) \cdot s_w(n-l)}{\sqrt{\sum_{n=0}^{79} s_w^2(n) \cdot \sum_{n=0}^{79} s_w^2(n-l)}} \quad l = 17, \dots, 148 \quad (5.3.9.1-1)$$

26 The maximum of the normalized correlation is found for each of the four intervals [17,33],
 27 [34,67], [68,135], and [136,148], resulting in four maximum values of $R_p(l)$. The locations
 28 of these four $R_p(l)$ are the initial estimate of the four open-loop pitch candidates. The non-
 29 uniform size of the intervals is to help the detection of potential pitch multiples, as will be
 30 discussed in the next section.

31 5.3.9.1.2 Correction of Open-Loop Pitch Multiples for the Lookahead

1 The location of the maximum of the normalized correlation $R_p(l)$ from the 4 possible
 2 candidates is selected as the initial candidate, l_{cand}^{ol} , for the open-loop pitch lag. In order to
 3 avoid pitch doubling, pitch correction is performed by inspecting the pitch candidates from
 4 the lower pitch ranges. The lookahead pitch value is corrected using the last two frames
 5 open-loop pitch lags that were calculated on the second half of the respective frame, and
 6 their 4-level reliability indicator. To avoid pitch multiples, the pitch is corrected by
 7 inspecting the pitch candidates selected from the intervals of lower pitch lag. A weighting
 8 parameter, D , is calculated for each lower-lag pitch candidate, based on its fitting to the
 9 previous frames' pitch contour. A flag signifies if the lower-lag pitch candidate is a possible
 10 $1/2$, $1/3$, or $1/4$ of the initial candidate, l_{cand}^{ol} .

11 If the lower-lag pitch fits into the pitch contour, if it is a possible sub-multiple of the initial
 12 candidate, and if its normalized correlation $R_p(l)$ is larger than 0.25, the lower-lag pitch
 13 value replaces the initial candidate l_{cand}^{ol} as the selected open-loop pitch, $l_p^{ol}(2)$.

14

15 **5.3.9.2 Open-Loop Pitch Detection on the Second Half of the Encoding Frame**

16 Before detecting the open-loop pitch detection for the second half of the frame, the past
 17 pitch values and their corresponding reliability indicators, used for the pitch contour
 18 matching, are modified. The open-loop pitch from the previous frame lookahead, or the
 19 open-loop pitch from the current frame lookahead (computed in the last section), replaces
 20 the open-loop pitch kept from 2 frames in the past, if the pitch estimates of the last two
 21 frames were not reliable.

22 Following this correction, the normalized correlation and the initial 4 pitch candidates for
 23 this segment are calculated in similar fashion to Section 5.3.9.1.1, and pitch multiples are
 24 eliminated, in similar fashion to Section 5.3.9.1.2, to select an updated open-loop pitch
 25 candidate.

26

27 **5.3.9.3 Final Open-Loop Pitch Correction**

28 The final open-loop pitch for the lookahead obtained according to Section 5.3.9.1 is selected
 29 as $l_p^{ol}(2)$, and the normalized correlation is returned as:

$$30 \quad R_p^{ol}(2) = R_p(l_p^{ol}(2)) \quad (5.3.9.3-1)$$

31 The open-loop pitch candidate for the second half of the encoding frame is selected
 32 according to Section 5.3.9.2. The open-loop pitch candidate for the first half of the encoding
 33 frame is given by the lookahead open-loop pitch value of the previous frame. The two later
 34 candidates are tested and modified to select the final open-loop pitch values, providing a
 35 smooth pitch contour and further correction of pitch sub-multiples.

36

37 **5.3.9.3.1 Smoothing of the Open-Loop Pitch Candidates**

38 The open-loop pitch of the previous two frames, the open-loop pitch candidate for the first
 39 and the second halves of the frame, the open-loop pitch for the lookahead, and their
 40 corresponding normalized correlations, are examined for smoothness and reliability.
 41 According to this examination, the open-loop pitch for the second half of the frame can be

1 modified to either the value of the open-loop pitch on the lookahead, or an average of the
 2 open-loop pitch on the first half of the frame and the open-loop pitch on the lookahead, or
 3 can be reselected from the initial pitch candidates that were computed according to Section
 4 5.3.9.1.1. The open-loop pitch for the first half of the frame can be modified, according to
 5 this examination, to either the value of the average of the open-loop pitch of the previous
 6 frame and the open-loop pitch on the second half of the frame, or can be reselected from
 7 the initial pitch candidates that were computed according to Section 5.3.9.1.1.

8

9 5.3.9.3.2 Further Correction of Pitch Sub-Multiples

10 For voiced speech frames, open-loop pitch lag candidates of less than 50 samples for the
 11 first and second halves of the frame are further tested against an average of the open-loop
 12 pitch lag. If they are significantly smaller than the average, they are modified to either one
 13 of initial pitch candidates that were computed according to Section 5.3.9.1.1, or to the
 14 open-loop pitch of the previous frame. The average pitch lag is updated on a frame-by-
 15 frame basis using the final open-loop pitch candidate from the second half of the frame.
 16 The update is done in a recursive manner with a leakage factor of 0.75 and it is done only
 17 for voiced speech with reliable and smooth pitch contour. The average pitch lag is reset to
 18 40 after 3 frames of inactive speech.

19 The refined pitch lag for the first half of the frame is returned as $l_p^{ol}(0)$ and the selected
 20 pitch lag for the second half of the frame is returned as $l_p^{ol}(1)$. The normalized correlations
 21 are returned as:

$$22 \quad R_p^{ol}(0) = R_p(l_p^{ol}(0)) \quad (5.3.9.3-2)$$

$$23 \quad R_p^{ol}(1) = R_p(l_p^{ol}(1)) \quad (5.3.9.3-3)$$

24

1 **5.3.10 Initial Frame Classification**

2 Routine name: CLA_signal_classifier

3 **Input:**

- 4 • Pre-processed speech, $\tilde{s}(n)$
- 5 • Open-loop pitch lag on the second half of the frame $l_p^{ol}(1)$
- 6 • Open-loop pitch normalized correlation on the second half of the frame, $R_p^{ol}(1)$

7 **Output:**

- 8 • Frame class decision (initial value), C_{FR}

9

10 The initial frame classification perform the initial classification of the frame into one of the
11 6 categories:

- 12 0 - Silence
- 13 1 - Noise-like Unvoiced
- 14 2 - Unvoiced
- 15 3 - Onset
- 16 5 - Non-stationary voiced
- 17 6 - Stationary voiced

18 Note that the category number 4 is not used. Also, class 1, of noise-like unvoiced, and class
19 6, of stationary voiced, are determined later.

20 **5.3.10.1 Average and Variance of Open-Loop Pitch Values**

21 An average of the open-loop pitch, $l_p^{ol}(1)$, on the last two frames and the current frame, is
22 calculated and denoted by l_p^{avg} . The variance of the same 3 open-loop pitch values is
23 calculated and denoted by l_p^{var} .

24 **5.3.10.2 Noise-Suppressed Classification Parameters**

25 5.3.10.2.1 First Order Reflection Coefficients for 4 Subframes

26 The frame is divided into 4 subframes of 40 samples each. For each of such subframe the
27 first order reflection coefficient is calculated, using a Hamming window of 80 samples
28 centered on the middle of the subframe. The first order reflection coefficient for the m^{th}
29 subframe is denoted k_m^{cl} and is given by:

$$30 \quad k_m^{cl} = \frac{\sum_{n=1}^{79} \tilde{s}_m(n) \cdot \tilde{s}_m(n-1)}{\sum_{n=0}^{79} \tilde{s}_m^2(n)} \quad m = 0, \dots, 3 \quad (5.3.10.2-1)$$

1 where

$$2 \quad \tilde{s}_m(n) = \tilde{s}(40 \cdot m - 20 + n) \cdot W_H(n) \quad n = 0, \dots, 79 \quad (5.3.10.2-2)$$

3 and $W_H(n)$ is a Hamming window of size 80.

4

5 5.3.10.2.2 Noise-Suppressed Spectral Tilt, Absolute Maximum, and Open-Loop Pitch 6 Correlation

7 At the next step, the frame is divided into 8 subframes of 20 samples each. For each of
8 such subframe, indexed by m , the spectral tilt for the m^{th} subframe is defined by

9 $T_m^{cl} = k_{\lfloor m/2 \rfloor}^{cl}$. For each of such subframe the absolute maximum value of the samples in an

10 interval is calculated and denoted M_m^{cl} . The interval for the calculation of the absolute
11 maximum value is related, but not identical, to the subframe, and is defined by its start
12 and end points. The start point is given by:

$$13 \quad \text{start}_m = \begin{cases} \max(-140, 100 - 1.5 \cdot l_p^{ol}(1)) + 20m & l_p^{ol}(1) > 20 \\ 80 + 20m & l_p^{ol}(1) \leq 20 \end{cases} \quad (5.3.10.2-3)$$

14 (Negative indices relate to past samples). The end point is given by $\text{end}_m = 100 + 20m - 1$.
15 The maximum of the absolute value is given by:

$$16 \quad M_m^{cl} = \max_i \{|s(i)|, \text{start}_m \leq i \leq \text{end}_m\} \quad (5.3.10.2-4)$$

17 For each such subframe the energy is calculated on an interval, which is related to the
18 subframe, and is denoted E_m^{cl} . Its start and end points define the interval. The start point
19 is given by:

$$20 \quad \text{start}_m = \begin{cases} \max(-140, 20 - 2 \cdot l_p^{ol}(1)) + 20m & l_p^{ol}(1) > 20 \\ 20m & l_p^{ol}(1) \leq 20 \end{cases} \quad (5.3.10.2-5)$$

21 The end point is given by $\text{end}_m = 20 + 20m - 1$. The energy of the segment is calculated by:

$$22 \quad E_m^{cl} = 1 + \frac{\sum_{i=\text{start}_m}^{\text{end}_m} s^2(i)}{\text{end}_m - \text{start}_m} \quad (5.3.10.2-6)$$

23 The running means of the energy, the spectral tilt, the absolute maximum and the open-
24 loop pitch normalized correlation ($\bar{E}_{noise}, \bar{T}_{noise}, \bar{M}_{noise}, \bar{R}_{noise}$) are updated for each on the
25 20-sample subframes. For the first 6 frames, the update is performed if

26 $T_n^{cl} < -0.4$ and $R_p^{ol}(1) < 0.5$ according to:

$$\begin{aligned}
\bar{E}_{noise} &= 0.80 \cdot \bar{E}_{noise} + 0.20 \cdot E_m^{cl} \\
\bar{T}_{noise} &= 0.75 \cdot \bar{T}_{noise} + 0.25 \cdot T_m^{cl} \\
\bar{M}_{noise} &= 0.75 \cdot \bar{M}_{noise} + 0.25 \cdot M_m^{cl} \\
\bar{R}_{noise} &= 0.75 \cdot \bar{R}_{noise} + 0.25 \cdot R_p^{ol}(1)
\end{aligned} \tag{5.3.10.2-7}$$

2 For frames later than the first six frames the update is controlled by the VAD decision, and
3 the update is performed if the $VAD = 0$ according to:

$$\begin{aligned}
\bar{E}_{noise} &= 0.999 \cdot \bar{E}_{noise} + 0.001 \cdot E_m^{cl} \\
\bar{T}_{noise} &= 0.990 \cdot \bar{T}_{noise} + 0.010 \cdot T_m^{cl} \\
\bar{M}_{noise} &= 0.990 \cdot \bar{M}_{noise} + 0.010 \cdot M_m^{cl} \\
\bar{R}_{noise} &= 0.990 \cdot \bar{R}_{noise} + 0.010 \cdot R_p^{ol}(1)
\end{aligned} \tag{5.3.10.2-8}$$

5 The noise-to-signal ratio is given by:

$$NSR_m^{cl} = \min \left(0.968, \sqrt{\frac{\bar{E}_{noise}}{E_m^{cl}}} \right) \tag{5.3.10.2-9}$$

7 A buffer of the last 7 values of the noise-suppressed spectral tilt and the noise-suppressed
8 absolute maximum is maintained. (The values in the buffer can be from the previous frame
9 or the current frame, as the subframe index progresses from the first to the last.) The
10 updated value of each noise-suppressed parameter is calculated according to:

$$\begin{aligned}
T_{ns}^7(m) &= T_m^{cl} - NSR_m^{cl} \cdot \bar{T}_{noise} \\
M_{ns}^7(m) &= T_m^{cl} - NSR_m^{cl} \cdot \bar{M}_{noise}
\end{aligned} \tag{5.3.10.2-10}$$

12 A noise-suppressed value of the open-loop pitch normalized correlation is calculated
13 according to:

$$R_{ns}(m) = R_p^{ol}(1) - NSR_m^{cl} \cdot \bar{R}_{noise} \tag{5.3.10.2-11}$$

15 5.3.10.2.3 Spectral Tilt Slope, Absolute Maximum Slope, and Slope Statistics

16 A first order derivative (slope) of the spectral tilt and the absolute maximum for the m^{th}
17 subframe is calculated by:

$$\begin{aligned}
T'_{ns}(m) &= \sum_{i=1}^7 i \cdot (T_{ns}^i(m) - T_{ns}^0(m)) / 140 \\
M'_{ns}(m) &= \sum_{i=1}^7 i \cdot (M_{ns}^i(m) - M_{ns}^0(m)) / 140
\end{aligned} \tag{5.3.10.2-12}$$

19 For each frame, using the noise-suppressed parameters, the minimum of the spectral tilt
20 slope for the frame, T'_{min} , the maximum of the slope of the absolute maximum, M'_{max} , the
21 sum of the slopes of the absolute maximums, \bar{M}' , the minimum of the spectral tilt, T'_{min} ,
22 and the sum of the noise suppressed spectral tilt, \bar{T} , are given by:

$$T'_{\min} = \min_{i=0}^7 T'_{ns}(i)$$

$$M'_{\max} = \max_{i=0}^7 M'_{ns}(i)$$

$$1 \quad \bar{M}' = \sum_{i=0}^7 M'_{ns}(i) \quad (5.3.10.2-13)$$

$$\bar{T} = \sum_{i=0}^7 T_{ns}^7(i)$$

$$T_{\min} = \min_{i=0}^7 T_{ns}^7(i)$$

2 An auto-regressive average of T_{\min} is generated by $T_{\min}^{avg} = 0.75 \cdot T_{\min}^{avg} + 0.25 \cdot T_{\min}$.

3 5.3.10.2.4 Statistics of Open-Loop Pitch Correlation

4 The maximum of the noise-suppressed open-loop pitch correlation, R_{\max} , and the sum of
5 noise-suppressed open-loop pitch correlation, \bar{R} , are given by:

$$6 \quad R_{\max} = \max_{i=1}^8 R_{ns}(i)$$

$$7 \quad \bar{R} = \sum_{i=1}^8 R_{ns}(i) \quad (5.3.10.2-14)$$

7 An auto-regressive average of \bar{R} is given by $\bar{R}^{avg} = 0.75 \cdot \bar{R}^{avg} + 0.25 \cdot \bar{R}$.

8 5.3.10.2.5 Classification Decision

9 The calculated parameters are compared to a set of thresholds to generate two Boolean
10 parameters, one indicating voiced speech, and the other indicates onset speech. The
11 Boolean parameters are further modified according to the stored Boolean parameters of the
12 previous frame. The voiced parameter, the onset parameter, the VAD decision, and the
13 previous frame initial classification decision, are used in a set of Boolean conditions to
14 provide the initial frame classification decision for the current frame, C_{FM} . The initial
15 frame classification is limited to the classes of Silence (0), Unvoiced (2), Onset (3), or Voiced
16 (5).

17

18

1 **5.3.11 Noise-to-Signal Ratio (NSR) Estimation**

2 Routine name: SNR_Calc_NSr_enc

3 **Input:**

- 4 • Pre-processed speech, $\tilde{s}(n)$
- 5 • LSFs calculated on the last quarter of the encoding frame, $lsf_1(i)$
- 6 • Voice activity detection flag, VAD
- 7 • Frame class (initial value), C_{FM}

8 **Output:**

- 9 • Noise-to-signal measure, NSR

10

11 The noise-to-signal ratio (NSR) is a measure of the level of background noise presented to
12 the SMV codec. The NSR is the reciprocal of the more familiar signal-to-noise ratio value,
13 but is bounded by 1.

14 **5.3.11.1 Signal Input to the NSR Estimation Function**

15 The input of $\tilde{s}(n)$ to the NSR calculation function is a segment of 160 samples that
16 includes the second half of the coding frame and the lookahead buffer.

17 **5.3.11.2 NSR counters**

18 Three counters are used in the NSR estimation. C_{noise} counts the number of consecutive
19 frames where the initial frame classification is below 3 (not onset or voiced frames). C_{vad}
20 counts the number of consecutive frames where the VAD indicates inactive voice. C_{NSR} is a
21 modulo 1000 frame counter, which helps in resetting the NSR estimate.

22 **5.3.11.3 NSR Calculation**

23 The energy of the current segment is calculated by:

$$24 \quad E^{NSR} = \sum_{i=0}^{L_{fm}-1} \tilde{s}^2(i) \quad (5.3.11.3-1)$$

25 If the current frame is zero in modulo 1000, the returned NSR is 0, and the average energy
26 memory, \bar{E}_p^{NSR} is set to E^{NSR} . Otherwise, a spectral difference measure is calculated, as the
27 sum of the absolute values of the current frame LSFs, $lsf_1(i)$, and the previous frame LSFs,
28 $lsf_{1,prv}(i)$:

$$29 \quad D_{lsf} = \sum_{i=1}^{10} |lsf_1(i) - lsf_{1,prv}(i)| \quad (5.3.11.3-2)$$

30 An auto-regressive average of the spectral distance is calculated by

$$1 \quad \bar{D}_{lsf} = 0.75 \cdot \bar{D}_{lsf} + 0.25 \cdot D_{lsf} \quad (5.3.11.3-3)$$

2 An energy difference measure is calculated between the current frame energy, E^{NSR} , and
3 the previous frame energy E_p^{NSR} , by:

$$4 \quad D_E = \frac{|E^{NSR} - E_p^{NSR}|}{\max(E^{NSR} + E_p^{NSR}, 8000)} \quad (5.3.11.3-4)$$

5 An auto-regressive average of the energy difference distance is calculated by

$$6 \quad \bar{D}_E = 0.75 \cdot \bar{D}_E + 0.25 \cdot D_E \quad (5.3.11.3-5)$$

7 An SNR measure is calculated as

$$8 \quad SNR = \frac{E^{NSR}}{\max\{\bar{E}_p^{NSR}, 0.1\}}, \quad (5.3.11.3-6)$$

9 and if the VAD indicates inactive voice, \bar{E}_p^{NSR} is updated according to:

$$10 \quad \bar{E}_p^{NSR} = \begin{cases} E^{NSR} & \bar{E}_p^{NSR} < 0.1 \\ 0.75 \cdot \bar{E}_p^{NSR} + 0.25 \cdot E^{NSR} & \bar{E}_p^{NSR} \geq 0.1 \end{cases} \quad (5.3.11.3-7)$$

11 The SNR , C_{vad} , C_{NSR} , D_{lsf} , \bar{D}_{lsf} , D_E , and \bar{D}_E , are compared to set of thresholds and a
12 locally refined VAD decision is obtained. This locally refined VAD is also controlled by the
13 average of the background noise energy \bar{E}_{noise} and the modulo counter C_{NSR} . If the locally
14 refined VAD decision indicates inactive speech, the average of the background noise is
15 updated according to:

$$16 \quad \bar{E}_{noise} = \begin{cases} E^{NSR} & \bar{E}_{noise} < 0 \\ \min(0.85 \cdot \bar{E}_{noise} + 0.15 \cdot E^{NSR}, E^{NSR}) & \bar{E}_{noise} \geq 0 \end{cases} \quad (5.3.11.3-8)$$

17 (\bar{E}_{noise} can be negative due to initialization.) The NSR is calculated as:

$$18 \quad NSR = \min\left(\sqrt{\frac{\max(\bar{E}_{noise} - 16000, 0)}{\max(E^{NSR}, 0.1)}}, 1.0\right) \quad (5.3.11.3-9)$$

19

1 **5.3.12 Classification Refinement**

2 Routine name: CLA_Class_Correct

3 **Input:**

- 4 • Music detection flag, F_{MUS}
- 5 • Frame class decision (initial value), C_{FR} , and previous frame class decision, C_{FR}^{prv}
- 6 • VAD decision (initial value), VAD
- 7 • Voiced/unvoiced level (initial value), C_{VUV}
- 8 • Open-loop normalized pitch correlations, $R_p^{ol}(l)$, $l = 0,1,2$
- 9 • Noise-to-signal ratio, NSR
- 10 • Interpolated reflection coefficients, $\hat{k}_m(i)$, $m = 0,1,2,3,4$
- 11 • Sharpness of weighted residual, S_{rp}^{wr} , (Section 5.3.7)
- 12 • SMV mode, M_{SMV}
- 13 • Energy of the lookahead buffer, given by $R_2(0)$.
- 14 • Energy of the weighted residual given by $E_{rw}(0) + E_{rw}(1)$, (Section 5.3.7)
- 15 • LPC gain on the last quarter of the encoding frame, G_1^{LPC}

16 **Output:**

- 17 • Frame class decision (refined), C_{FR}
- 18 • Auxiliary class decision (used mainly for signal modification), C_{SM}
- 19 • Onset flag, F_{onset}
- 20 • LPC onset flag, F_{LPC_onset}
- 21 • Noisy-voiced flag F_{noisyV}
- 22 • VAD decision (refined), VAD
- 23 • Voiced/unvoiced level (refined), C_{VUV}

24

25 **5.3.12.1 Operation of the Classification Refinement**

26 The class refinement uses previously calculated parameters, and saved values from
 27 previous frames, to refine the classification decision, the VAD decision, and the
 28 voiced/unvoiced level decision. It also sets three flags that are used later, F_{onset} , F_{LPC_onset} ,
 29 and $F_{noisyV} \cdot F_{onset}$ are used during the rate selection and the final classification correction,

1 F_{LPC_onset} is used to control the LSFs smoothing, and F_{noisyV} is used to control the gain
2 quantization, the perceptual weighting filter and the fixed-codebook search. The
3 classification refinement procedure also sets C_{SM} , an auxiliary classification decision, used
4 to control several aspects of the encoding algorithm, and in particular, the weighted speech
5 modification algorithm (Section 5.3.15). C_{SM} takes the values of $-1, \dots, 4$, where an
6 increasing number indicates an increasing level of voicing.

7 The decision parameters are compared to each other, and to a set of thresholds, for a
8 sequential refinement of initial decision parameters and for the setting of new decision
9 parameters and flags.

10

11

12

13

1 **5.3.13 Rate Selection**

2 Routine name: CLA_Rate_Select

3 **Input:**

- 4 • Music detection flag, F_{MUS}
- 5 • SMV mode, M_{SMV}
- 6 • Frame class decision (refined), C_{FR} , and previous frame class decision, C_{FR}^{prv}
- 7 • Open-loop normalized pitch correlations, $R_p^{ol}(l)$, $l = 0,1,2$
- 8 • Noise-to-signal ratio, NSR
- 9 • Sharpness of weighted residual, S_{rp}^{wr} , (Section 5.3.7)
- 10 • Energy of the weighted residual given by $E_{rw}(0) + E_{rw}(1)$, (Section 5.3.7)
- 11 • First element of third reflection coefficient vector, $\hat{k}_2(0)$
- 12 • Onset flag, F_{onset}

13 **Output:**

- 14 • SMV rate, $Rate$

16 **5.3.13.1 Rate-Selection Algorithm**

17 For each frame, the rate selection algorithm selects one of the 4 possible rates, Rate 1, Rate
 18 1/2, Rate 1/4, or Rate 1/8. The selection of the rate is controlled by the SMV operating
 19 mode, which can be Mode 0 (Premium), Mode 1 (Standard), Mode 2 (Economy), or Mode 3
 20 (Max-Capacity). The rate selection algorithm uses previously computed parameters and a
 21 set of fixed and adaptive thresholds for initial selection of the rate for each frame. The
 22 initial rate selection is later refined (Section 5.3.16).

23 The rate-selection algorithm sets the rate to Rate 1 if the music flag indicates a music
 24 signal. Otherwise, the algorithm selects the coding rate for the frame according to a set of
 25 fixed thresholds and conditions, where the set of thresholds and conditions depends on the
 26 dictated mode of operation for the SMV.

27 An initial selection of Rate 1 can be forced to Rate 1/2 by the half-rate max flag, F_{HRM} .

28

1 **5.3.14 Open-Loop Pitch Quantization, Pitch Adjustment and** 2 **Interpolation**

3 If the auxiliary classification decision $C_{SM} > 0$, the open loop $l_p^{ol}(1)$ has to be quantized,
4 adjusted and interpolated before the modified weighted speech can be generated as
5 described in the next section. Before the quantization begins, if $l_p^{ol}(1) > 120$ and the Rate is
6 not 1, C_{SM} will be set to zero and the frame will be classified into a type-0 frame and
7 hence, be excluded from this quantization process. The quantization of $l_p^{ol}(1)$ is performed
8 by searching through a pitch codebook of size 256 for Rate 1 and of size 128 for other
9 rates. To avoid an exhaustive search, the search range is limited to $[l_p^{low}, l_p^{high}]$ where

$$10 \quad l_p^{low} = \max\{l_p^{ol}(1) - 5, 17\} \quad (5.3.14-1)$$

$$11 \quad l_p^{high} = \begin{cases} \min\{l_p^{ol}(1) + 5, 148\} & \text{Rate is Rate 1} \\ \min\{l_p^{ol}(1) + 5, 120\} & \text{otherwise} \end{cases} \quad (5.3.14-2)$$

12 Since the pitch table is not necessarily expressed in integer values, the search requires to
13 upsample the normalized pitch correlation in a bandlimited fashion. As a result, the
14 quantized pitch is chosen to be the codeword that maximizes the upsampled correlations.
15 Based on this pitch, a pitch contour is constructed and will be used by the signal
16 modification algorithm (Section 5.3.15). The modification algorithm attempts to modify the
17 speech signal to match this pitch contour. If the matching is successful, the frame is
18 classified into type 1, and type 0 otherwise. For type-1 frames, the quantized pitch value
19 obtained by this procedure is not further modified.
20

21 **5.3.14.1 Open-Loop Pitch Refinement and Quantization**

22 Routine name: PIT_FinePitch

23 **Input:**

- 24 • SMV rate, $Rate$
- 25 • Pitch search range, $[l_p^{low}, \dots, l_p^{high}]$
- 26 • Weighted speech, $s_w(n)$
- 27 • Signal modification accumulated delay, Δ_{SM}

28 **Output:**

- 29 • Quantized open-loop pitch lag at the end of the encoding frame, l_{pp}^q
- 30 • Integer part of the quantized open-loop pitch lag at the end of the encoding frame,
31 $l_{pp,INT}^q$

32

1 The pitch refinement search starts with calculating the normalized pitch correlation on the
2 weighted speech according to

$$3 \quad R_p(l) = \frac{\sum_{n=0}^{79} s_w(n) \cdot s_w(n-l)}{\sqrt{\sum_{n=0}^{79} s_w^2(n) \cdot \sum_{n=0}^{79} s_w^2(n-l)}} \quad l = l_p^{low} - 6, \dots, l_p^{high} + 6 \quad (5.3.14.1-1)$$

4 and the optimal integer lag l_p^* is found by

$$5 \quad l_p^* = \arg \max_{l=l_p^{low}, \dots, l_p^{high}} R_p(l) \quad (5.3.14.1-2)$$

6 Note that the ± 6 samples are needed for the upsampling of $R_p(\cdot)$ during the fractional
7 search in the next step. After obtaining the optimal integer lag l_p^* , the fractional search
8 begins. The search interval is $[l_p^* - 1, \dots, l_p^* + 1]$ and $R_p(\cdot)$ is upsampled at all possible
9 fractional pitch values within that interval. The upsampling is bandlimited and uses a
10 series of pre-stored Sinc tables. The lag that maximizes the interpolated normalized
11 correlation is chosen. The corresponding codebook index is denoted by I_p .

12 Next, the quantized pitch lag is adjusted to reduce the risk of reaching the delay limits of
13 the signal modification algorithm, since the allowed level for the accumulated delay is
14 $\Delta_{SM} \leq 20$. The adjustment is carried out by modifying the index of the pitch quantization,
15 I_p , according to:

$$16 \quad I_p = \begin{cases} I_p + 1 & \text{Rate1 and } \Delta_{SM} > 10 \\ I_p + 2 & \text{Rate1 and } \Delta_{SM} > 15 \\ I_p - 1 & \text{Rate1 and } \Delta_{SM} < -10 \\ I_p - 2 & \text{Rate1 and } \Delta_{SM} < -15 \\ I_p + 1 & \text{Rate } \frac{1}{2} \text{ and } \Delta_{SM} > 10 \\ I_p - 1 & \text{Rate } \frac{1}{2} \text{ and } \Delta_{SM} < -10 \end{cases} \quad (5.3.14.1-3)$$

17 After the modification, attention is paid to ensure the modified I_p is still a valid codebook
18 index (i.e., I_p cannot be negative or exceed the codebook size). The final pitch value, l_{pp}^q is
19 obtained from the quantization table using the modified index. The $l_{pp,INT}^q$ is simply the
20 rounded version of l_{pp}^q .

21

22 **5.3.14.2 Pitch Interpolation**

23 If $C_{SM} > 0$, pitch interpolation path, $Pit(n)$, is constructed, using the open-loop quantized
24 pitch lags from previous two frames, denoted by $l_{pp}^q(frm-1)$ and $l_{pp}^q(frm-2)$ for the

1 previous frame and the one before, respectively. If the previous rate is Rate 1/4 or Rate
 2 1/8, these lags are set to $l_{pp}^q(frm-1) = l_{pp}^q(frm-2) = l_{pp}^q$. Otherwise, if the previous frame
 3 was of type 0, $l_{pp}^q(frm-1) = l_{pp}^q(frm-2)$, and both are set to the closed-loop pitch lag of
 4 the last subframe of the calculated in Sections 5.6.6 and 5.6.7.

5

6 Routine name: PIT_PitchInterpolat7 **Input:**

- 8 • Quantized open-loop pitch lag at the end of the encoding frame, l_{pp}^q
- 9 • Open-loop pitch lags from previous two frames, $l_{pp}^q(frm-1)$ and $l_{pp}^q(frm-2)$

10 **Output:**

- 11 • Quantized Interpolated pitch contour, $Pit(n)$

12

13 Before the pitch is interpolated, a mean pitch value is first obtained by:

$$14 \quad \bar{l}_{pp}^q = 0.5 \cdot (l_{pp}^q(frm-2) + l_{pp}^q) \quad (5.3.14.2-1)$$

15 The relationships between $l_{pp}^q(frm-1)$, l_{pp}^q , and \bar{l}_{pp}^q are tested, and if these values are
 16 close to each other, the pitch contour is set according to:

$$17 \quad Pit(n) = \begin{cases} \frac{n \cdot l_{pp}^q + (L_{frm} - n) \cdot l_{pp}^q(-1)}{L_{frm}} & n = 0, \dots, L_{frm} - 1 \\ l_{pp}^q & n = L_{frm}, \dots, L_{frm} + L_{lkhd} - 1 \end{cases} \quad (5.3.14.2-2)$$

18 Otherwise, the pitch contour is not interpolated and is set to:

$$19 \quad Pit(n) = l_{pp}^q \quad n = 0, \dots, L_{frm} + L_{lkhd} - 1 \quad (5.3.14.2-3)$$

1 **5.3.15 Generating Modified Weighted Speech Signal**

2 Generate a modified weighted speech signal, $s_{mw}(n)$, by warping the current weighted
3 speech signal, $s_w(n)$, to match the pitch contour, $Pit(n)$, for the current frame.

4

5 **Input:** The inputs to the weighted speech modification routine are:

- 6 • M_{SMV} : selected SMV mode
- 7 • $Rate$: fixed codec bit rate
- 8 • F_{flat} : spectral flatness flag
- 9 • NSR : background noise to the signal ratio
- 10 • $s_w(n)$: current input weighted speech signal
- 11 • $Pit(n)$: pitch track for the current frame
- 12 • C_{SM} : signal modification class

13

14 **Output:** The outputs from the weighted speech modification routine are:

- 15 • $s_{mw}(n)$: modified weighted speech signal
- 16 • Δ_{SM} : signal modification accumulated delay
- 17 • C_{FR} : frame class 6 determination that defines Type 1
- 18 • C_{VUV} : voiced/unvoiced level (updated)
- 19 • $g_a(m)$: pitch gain (adaptive-codebook gain) for each subframe (indexed by m)
20 of stationary-voiced frame (class 6)

21

22 **5.3.15.1 Frame based processing for $C_{SM} \leq 0$**

23 For these two classes of C_{SM} , ($C_{SM} \leq 0$), the accumulated delay will be set to zero or kept
24 unchanged.

25

26 5.3.15.1.1 Region of Mapping of Weighted Speech to Modified Weighted Speech

27 If $C_{SM} = -1$ reset the accumulated delay:

$$28 \quad \Delta_{SM} = 0 \quad (5.3.15.1-1)$$

29 Calculate the starting time point T_0 and the ending time point T_1 of the weighted speech
30 signal to be mapped to the current frame of the modified speech signal:

$$1 \quad T_0 = L_{olpit} - L_{lkhd} - L_{frm} + \Delta_{SM} \quad (5.3.15.1-2)$$

$$2 \quad T_1 = T_0 + L_{frm} \quad (5.3.15.1-3)$$

3

4 5.3.15.1.2 Weighted Speech Mapping to Modified Weighted Speech for One Frame

5 Since the delay in the region $[T_0, T_1]$ is 0 or constant, only one *Sinc* table based on the
6 fractional value of the delay Δ_{SM} is needed to do the signal mapping:

$$7 \quad s_{mw1}(n) = \sum_{i=-10}^{10} w_s(f(\Delta_{SM}), i) \cdot s_w(n + i + I(T_0)), \quad n = 0, \dots, L_{frm} - 1 \quad (5.3.15.1-4)$$

8 where $f(\circ)$ takes the fractional part, $I(\circ)$ takes the integer part, and $w_s(f, i)$ is the
9 Hamming weighted *Sinc* window.

10

11 5.3.15.1.3 Update of Modified Weighted Speech Buffer

12 Shift the old buffer:

$$13 \quad s_{mw_buff}(n) \leftarrow s_{mw_buff}(n + L_{frm}), \quad 0 \leq n < L_{mw_buff} - L_{frm} \quad (5.3.15.1-5)$$

14 where $s_{mw_buff}(n)$, $n = 0, \dots, L_{mw_buff} - 1$ is the buffer vector of the modified weighted speech
15 signal. Then, add the current modified weighted speech signal to the buffer:

$$16 \quad s_{mw_buff}(n + L_{mw_buff} - L_{frm}) = s_{mw}(n), \quad 0 \leq n < L_{frm} \quad (5.3.15.1-6)$$

17

18 **5.3.15.2 Variable Subframe Processing** $C_{SM} > 0$

19 For $C_{SM} > 0$ the analysis and the modification of the weighted voiced speech is conducted
20 on a variable subframe basis. The subframe size, L_{SF} , is related to pitch lag value and the
21 location of the subframe within a frame. For each subframe, a local delay, D_L , at the end of
22 the subframe is estimated in a fast way by maximizing the correlation between a current
23 subframe and a modified target signal. A continuous time warping of the signal as opposed
24 to simple shifting is then performed with a variable local delay from 0 to the estimated local
25 delay D_L . The absolute value of the accumulated delay is limited a maximum 20 samples
26 (2.5 ms). As a portion of the results, the estimated open-loop pitch gains and the
27 normalized pitch correlation values will be provided.

28

29 **Initialize parameters:**

- 30 • Set the starting point of the first subframe within a frame to zero, $i_0 = 0$
31 • Set the pitch gain vector to zero, $g_a(l) = 0$, $l = 0, 1, 2, 3$

- 1 • Set the pitch correlation vector to zero, $R_{SM}(l) = 0$, $l = 0,1,2,3$
 2 • Set the average pitch correlation value to zero, $R_{SM_avg} = 0$
 3 • Set the minimum pitch correlation value to 1.0, $R_{SM_min} = 1.0$

4

5 5.3.15.2.1 Determining Parameters

6 Before searching for a local delay for each subframe, some parameters need to be
 7 estimated.

8 Calculate the pulse length (pulse width):

$$9 \quad L_{puls} = \min\{\max\{0.2 \cdot Pit(i_0), 5\}, 30\} \quad (5.3.15.2-1)$$

10 Calculate the starting time, T_0 , of the current subframe on the current weighted speech
 11 signal, s_w :

$$12 \quad T_0 = L_{olpit} - L_{lkhd} - L_{frm} + i_0 + \Delta_{SM} \quad (5.3.15.2-2)$$

13 where Δ_{SM} is the accumulated delay for the last subframe.

14 Define the pulse peak location searching range $[T_0, T_2]$ for the current subframe, where the
 15 end time point, T_2 , is calculated as:

$$16 \quad T_2 = \min\{T_0 + Pit(i_0), L_{olpit} - 12\} \quad (5.3.15.2-3)$$

17 The pulse peak location, T_{peak} , is evaluated by maximizing the peak energy as a function of
 18 the location n :

$$19 \quad E_{peak}(n) = \sum_{k=-2}^2 s_w(n+k)^2, \quad T_0 \leq n \leq T_2, \quad (5.3.15.2-4)$$

20 That is,

$$21 \quad E_{peak}(T_{peak}) = \max\{E_{peak}(n) : n = T_0, \dots, T_2\} \quad (5.3.15.2-5)$$

22 From $n = T_0$ to $n = T_2$ the average peak energy is estimated as:

$$23 \quad E_{avg}(T_0) = E_{peak}(T_0) \quad (5.3.15.2-6)$$

$$24 \quad \begin{aligned} &\text{if } E_{peak}(n) < 4 \cdot E_{avg}(n-1) \\ &E_{avg}(n) = 0.875 \cdot E_{avg}(n-1) + 0.125 \cdot E_{peak}(n) \quad ; n = T_0 + 1, \dots, T_2 \end{aligned} \quad (5.3.15.2-7)$$

25 After obtaining the maximum peak energy and the average peak energy, an average energy
 26 to peak energy ratio is calculated by

$$27 \quad APR = \frac{E_{avg}(T_2)}{E_{peak}(T_{peak})} \quad (5.3.15.2-8)$$

28 For the first subframe, some protection for evaluating the peak location is set by

1 if $i_0 = 0$ and $T_{peak} < T_0 + L_{puls}$

$$T_{peak} = T_0 + L_{puls} \quad (5.3.15.2-9)$$

2 Based on the parameters obtained above, the current subframe size, L_{sub} , and the
3 extended subframe size, L_{ex_sub} , are calculated by

$$L_{sub} = \min\{T_{peak} - T_0 + L_{puls}, L_{frm} - i_0\} \quad (5.3.15.2-10)$$

$$L_{ex_sub} = T_{peak} - T_0 + 2 \cdot L_{puls} \quad (5.3.15.2-11)$$

6 If $T_{peak} + 1.25 \cdot Pit(i_0) > L_{olpit} - L_{lkh} + \Delta_{SM}$, the next possible peak location is evaluated in
7 order to judge whether the current subframe must include the remaining samples up to the
8 end of the current frame. The next peak is estimated by using the same approach as
9 described above with a new searching range of $[\underline{T}_0, \underline{T}_2]$:

$$\underline{T}_0 = T_{peak} + 0.5 \cdot Pit(i_0) \quad (5.3.15.2-12)$$

$$\underline{T}_2 = \min\{\underline{T}_0 + Pit(i_0), L_{olpit} - 12\} \quad (5.3.15.2-13)$$

12 If the newly found peak location is beyond the end of the current frame,

$$T'_{peak} > L_{olpit} - L_{lkh} + \Delta_{SM} \quad (5.3.15.2-14)$$

14 The current subframe size and the extended subframe size are modified according to:

$$L_{sub} = L_{frm} - i_0 \quad (5.3.15.2-15)$$

$$L_{ex_sub} = \max\{L_{sub}, T'_{peak} - \underline{T}_0 + 2 \cdot L_{puls}\} \quad (5.3.15.2-16)$$

$$L_{ex_sub} \leftarrow \min\{\min\{L_{ex_sub}, L_{frm} - i_0 + \max\{L_{lkh} - \Delta_{SM} - 15, 0\}\}, L_{max_ex_sub}\} \quad (5.3.15.2-17)$$

19 Each of the extended subframe sizes is always limited to $L_{max_ex_sub} = 180$.

20 Each of the defined subframes is split into two regions where the continuous time warping
21 will be applied on the first region and the time shifting will be performed on the second
22 region. The length of the first region is defined as

$$L_{center} = \min\left\{L_{sub}, \max\left\{T_{peak} - T_0 - L_{puls}, \frac{L_{puls}}{2}\right\}\right\} \quad (5.3.15.2-18)$$

25 5.3.15.2.2 Target Signal Determination for Current Local Delay Search

26 The target signal for the current local delay searching is determined according to the buffer
27 of the modified weighted speech signal and the current pitch track function.

28 Create the current pitch track function:

$$PitF(n) = Pit(i_0 + n) \quad n = 0, \dots, L_{imp} - 1 \quad (5.3.15.2-19)$$

$$1 \quad PitF(n) = PitF(L_{tmp} - 1) \quad n = L_{tmp}, \dots, L_{ex_sub} - 1 \quad (5.3.15.2-20)$$

2 where L_{tmp} is a temporal parameter:

$$3 \quad L_{tmp} = \min \left\{ L_{ex_sub}, L_{frm} - i_0 + \frac{L_{frm}}{2} \right\} \quad (5.3.15.2-21)$$

4 Calculate the target signal and memorize it in a temporal buffer extended from the past
5 modified weighted speech signal:

$$6 \quad s_{mw_buff}(L_{mw_buff} + n) = \sum_{i=-10}^{10} w_s(f(PitF(n)), i) \cdot s_{mw_buff}(L_{mw_buff} + n - I(PitF(n)) + i) \\ n = 0, \dots, L_{ex_sub} - 1 \quad (5.3.15.2-22)$$

8 In order to simplify later computation, copy the target signal into a target signal vector:

$$9 \quad T_{G_SM}(n) = s_{mw_buff}(L_{mw_buff} + n) \quad n = 0, \dots, L_{ex_sub} - 1 \quad (5.3.15.2-23)$$

10

11 5.3.15.2.3 Searching Range Determination for Local Delay

12 The local delay search range is limited by the positive boundary SR_1 and the negative
13 boundary SR_0 . Two APR values that could influence the searching range are calculated
14 first in the same way as described from by Equation (5.3.15.2-4) to (5.3.15.2-8). APR_1 is
15 based on the weighted speech in the region $[T_0, T_2]$ where

$$16 \quad T_0 = T_{peak} - 0.5 \cdot Pit(i_0) \quad (5.3.15.2-24)$$

$$17 \quad T_2 = \min \{ T_0 + Pit(i_0), L_{olpit} - 12 \} \quad (5.3.15.2-25)$$

18 APR_2 is measured on the target signal $T_{G_SM}(n)$ with $3 \leq n \leq L_{ex_sub} - 3$. The searching
19 range will be initialized first and then some modification on it could be performed.

20 Initialize the boundaries:

$$21 \quad SR_1 = I(\max \{ \min \{ 0.075 \cdot L_{center} + 0.5, 5 \}, 1 \}) \quad (5.3.15.2-26)$$

$$22 \quad SR_0 = SR_1 \quad (5.3.15.2-27)$$

23 Let $L_{tmp} = 4$ if $C_{SM} \geq 3$ and $L_{tmp} = 3$ otherwise. Check the influence from the two APR
24 values by:

1 if $APR_1 > \frac{1}{4}$

$$SR_1 = \min\{SR_1, 3\}$$

$$SR_0 = \max\{SR_0, -3\}$$

else if $APR_1 > \frac{1}{2}$ (5.3.15.2-28)

$$SR_1 = \min\{SR_1, 4\}$$

$$SR_0 = \max\{SR_0, -4\}$$

2 if $APR_2 > \frac{1}{3}$

$$SR_1 = \min\{SR_1, L_{imp}\}$$

$$SR_0 = \max\{SR_0, -L_{imp}\}$$

else if $APR_2 > \frac{1}{6}$ (5.3.15.2-29)

$$SR_1 = \min\{SR_1, L_{imp} + 1\}$$

$$SR_0 = \max\{SR_0, -L_{imp} - 1\}$$

3 Further consider the influence of the subframe size. If $L_{ex_sub} < T_{peak} - T_0 - 2$ or
 4 $Pit(i_0) < 25$ or $T_{peak} \leq T_{puls} + T_0$:

5 if $APR_2 > \frac{1}{4}$

$$SR_1 = \min\{SR_1, 1\}$$

$$SR_0 = \max\{SR_0, -1\}$$

else (5.3.15.2-30)

$$SR_1 = \min\{SR_1, 2\}$$

$$SR_0 = \max\{SR_0, -2\}$$

6 Further consider the influence from the accumulated delay:

7 if $\Delta_{SM} > 5$

$$SR_1 = \min\{SR_1, 4\}$$

if $\Delta_{SM} > 10$ (5.3.15.2-31)

$$SR_1 = \min\{SR_1, 2\}$$

if $\Delta_{SM} > 15$

$$SR_1 = \min\{SR_1, 0\}$$

8

9 if $\Delta_{SM} < -5$

$$SR_0 = \max\{SR_0, -4\}$$

if $\Delta_{SM} < -10$ (5.3.15.2-32)

$$SR_0 = \max\{SR_0, -2\}$$

if $\Delta_{SM} < -15$

$$SR_0 = 0$$

1

2 5.3.15.2.4 Local Delay Search

3 This section will be divided into (1) modification of the target signal, (2) searching for an
 4 integer part of the local delay, (3) searching for a fractional part of the local delay, and (4)
 5 generating a warped weighted-speech signal.

6 5.3.15.2.4.1 Modification of Target Signal

7 The target signal is modified in time domain before maximizing the correlation between the
 8 target signal and the current extended subframe by applying a weighting window:

if $C_{SM} \geq 3$

$$T_{G_SM_M}(n) = \left(0.25 + 0.75 \frac{n}{L_{center}} \right) \cdot T_{G_SM}(n)$$

9

else

(5.3.15.2-33)

$$T_{G_SM_M}(n) = \left(0.5 + 0.5 \frac{n}{L_{center}} \right) \cdot T_{G_SM}(n) \quad ; 0 \leq n < L_{center}$$

$$T_{G_SM_M}(n) = T_{G_SM}(n) \quad ; L_{center} \leq n < L_{ex_sub} - 1$$

11 where $T_{G_SM_M}$ is the modified target signal.

12 5.3.15.2.4.2 Integer Local Delay Search

13 In order to calculate the correlation values between the target signal and the current signal
 14 with different integer local delays, a signal keeping the latest accumulated delay with zero
 15 local delay is calculated first by mapping the weighted speech signal from the region
 16 $[TT_0, TT_1]$ to a signal vector $s_{w_shift}(n)$, $n = 0, \dots, TT_1 - TT_0 - 1$

$$s_{w_shift}(n) = \sum_{i=-10}^{10} w_s(f(TT_0), i) \cdot s_w(n + I(TT_0) + i), \quad ; n = 0, \dots, TT_1 - TT_0 - 1 \quad (5.3.15.2-34)$$

18 where

$$TT_0 = T_0 + SR_0 - 5 \quad (5.3.15.2-35)$$

$$TT_1 = T_0 + SR_1 + 5 + L_{ex_sub} + 1 \quad (5.3.15.2-36)$$

21 Then, the correlation value, $R_{P_SM_shift}(k - (SR_0 - 5))$, on the integer local delay of
 22 $SR_0 - 5 \leq k \leq SR_1 + 5$, is calculated by

$$R_{P_SM_shift}(k - (SR_0 - 5)) = \frac{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM_M}(n) s_{w_shift}(n + k - (SR_0 - 5))}{\sqrt{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM_M}^2(n) \sum_{n=0}^{L_{ex_sub}-1} s_{w_shift}^2(n + k - (SR_0 - 5))}} \quad (5.3.15.2-37)$$

24 The optimal integer local delay, k_{opt} , and the corresponding correlation value are obtained
 25 by maximizing the above correlation with a minimum value condition:

if $C_{SM} \geq 3$

$$R_{P_SM_shift}(k_{opt} - (SR_0 - 5)) = \max\{R_{P_SM_shift}(k - (SR_0 - 5)) > 0.3, SR_0 \leq k \leq SR_0\}$$

else

$$R_{P_SM_shift}(k_{opt} - (SR_0 - 5)) = \max\{R_{P_SM_shift}(k - (SR_0 - 5)) > 0.4, SR_0 \leq k \leq SR_0\}$$

(5.3.15.2-38)

5.3.15.2.4.3 Fractional Local Delay Determination

All the possible fractional local delays are generated and memorized in a table, $Tab(i)$,

$$Tab(i) = k_{opt} - 0.75 + 0.1 \cdot i \quad i = 0, \dots, 15$$

To find a best fractional delay, the correlation with integer delays is up-sampled for each fractional delay candidate:

$$R_{P_SM_F}(Tab(i)) = \sum_{j=-4}^4 w_s(f(Tab(i)), j) \cdot R_{P_SM_shift}(I(Tab(i)) + j - (SR_0 - 5))$$

(5.3.15.2-39)

The best fractional local delay, D_L , and the corresponding correlation value, $R_{P_SM_opt}$, will be found by maximizing the up-sampled correlation:

$$\begin{aligned} R_{P_SM_opt} &= R_{P_SM_F}(D_L) \\ &= \max\{R_{P_SM_F}(Tab(i)), i = 0, 1, \dots, 15\} \end{aligned}$$

(5.3.15.2-40)

5.3.15.2.4.4 Warping Current Weighted Speech With Determined Delay

After finding the local delay, the current subframe is warped from the original weighted speech signal to a temporal modified weighted speech signal vector:

$$\begin{aligned} s'_{mw}(n) &= \sum_{i=-10}^{10} w_s(f(T_0 + D_L \cdot n/L_{center}), i) \cdot s_w(n + I(T_0 + D_L \cdot n/L_{center}) + i), \\ &\quad ; n = 0, \dots, L_{center} - 1 \\ s'_{mw}(n) &= \sum_{i=-10}^{10} w_s(f(T_0 + D_L), i) \cdot s_w(n + I(T_0 + D_L) + i), \\ &\quad ; n = L_{center}, \dots, L_{ex_sub} - 1 \end{aligned}$$

(5.3.15.2-41)

5.3.15.2.5 Further Modification of Weighted Speech for Half_Rate_Max

If the maximum rate is limited to the Rate 1/2, further signal modification could be performed on some specific transition areas by smoothing the signal with waveform interpolation.

1 5.3.15.2.5.1 *Generating a Backward Waveform Vector*

2 Calculate two conditions for the possible further signal modification:

$$3 \quad Cond_1 = 0.25 < R_{P_SM_opt} < 0.7 \text{ and } Pit(i_0) > 25 \text{ and } L_{sub} > \max\{0.9 \cdot Pit(i_0), 17\}$$

4

(5.3.15.2-42)

$$5 \quad Cond_2 = I(T_0 + \Delta_{SM} + Pit(i_0) + L_{ex_sub} + 17) > L_{olpit}$$

(5.3.15.2-43)

6 If both conditions are met, calculate a following reference signal for estimating a backward
7 waveform:

$$8 \quad T'_{G_SM_M}(n) = 0.8 (0.5 + 0.5n / L_{center}) s'_{mw}(n) , \quad n = 0, \dots, L_{center} - 1$$

$$9 \quad T'_{G_SM_M}(n) = s'_{mw}(n) , \quad n = L_{center}, \dots, L_{ex_sub} - 1$$

10

(5.3.15.2-44)

11 Then, the current backward waveform will be generated by taking a future signal vector in
12 the interval $[TT_0 + D_{shift}, TT_0 + D_{shift} + L_{ex_sub}]$ of the weighted speech signal s_w , where

$$13 \quad TT_0 = T_0 + \Delta_{SM} + Pit(i_0)$$

(5.3.15.2-45)

14 The shifting parameter, $-3 \leq D_{shift} \leq 3$, is optimized by maximizing the correlation between
15 the future signal vector and the reference signal defined in Equation (5.3.15.2-44). The
16 same searching procedure as described from the Sections 5.3.15.2.4.2 to 5.3.15.2.4.4 is
17 used to obtain D_{shift} and the corresponding correlation value R_{P_shift} by setting $SR_0 = -3$,
18 $SR_1 = 3$, and replacing $T_{G_SM_M}$ with the reference signal $T'_{G_SM_M}$ defined by Equation
19 (5.3.15.2-44).

20 The future signal vector in the above determined area will be mapped to the backward
21 waveform vector, s_{w_bkwd} , if the following condition is satisfied:

$$22 \quad Cond_3 = (C_{SM} \geq 4 \text{ and } R_{P_shift} > 0.25) \text{ or } R_{P_shift} > 0.4$$

(5.3.15.2-46)

23 If $Cond_3$ is met,

$$24 \quad s_{w_bkwd}(n) = \sum_{i=-10}^{10} w_s(f(TT_0 + D_{shift}), i) \cdot s_w(n + I(TT_0 + D_{shift}) + i)$$

(5.3.15.2-47)

$n = 0, \dots, L_{ex_sub} - 1$

25 5.3.15.2.5.2 *Generating a Candidate Vector by Waveform Interpolation*

26 If $Cond_1$ and $Cond_3$ are met, or $R_{P_SM_opt} > 0.4$, the following waveform interpolation will
27 be performed to generate a candidate vector to possibly replace the modified weighted
28 speech signal for the current variable subframe.

29 Generate a temporal vector:

if $Cond_3$

$$s_{WI}(n) = T_{G_SM}(n) + s_{w_bkwd}(n) \quad n=0, \dots, L_{ex_sub} - 1$$

1 else (5.3.15.2-48)

$$s_{WI}(n) = T_{G_SM}(n) \quad , \quad n=0, \dots, L_{ex_sub} - 1$$

2 where the target signal T_{G_SM} acts as the forward waveform. The energy of the vector is
3 normalized by a gain factor G_{imp} :

$$G_{imp} = \sqrt{\frac{\sum_{n=0}^{L_{ex_sub}-1} s'_{mw}(n)^2}{\sum_{n=0}^{L_{ex_sub}-1} s_{WI}(n)^2}}$$

4 (5.3.15.2-49)

$$s_{WI}(n) \leftarrow s_{WI}(n) \cdot G_{imp} \quad , \quad n = 0, \dots, L_{ex_sub}$$

5 (5.3.15.2-50)

6 Set a length of an interpolation region:

if $C_{SM} \geq 3$

$$L_{imp} = L_{sub} / 2$$

7

else

$$L_{imp} = L_{sub} / 3$$

8 (5.3.15.2-51)

8 Perform the interpolation at the beginning of the subframe:

$$s_{WI}(n) \leftarrow (n \cdot s_{WI}(n) + (L_{imp} - n) \cdot s'_{mw}(n)) / L_{imp} \quad , \quad n = 0, 1, \dots, L_{imp} - 1$$

9 (5.3.15.2-52)

10 Perform the interpolation at the end of the subframe:

$$s_{WI}(L_{sub} - n) \leftarrow (n \cdot s_{WI}(L_{sub} - n) + (L_{imp} - n) \cdot s'_{mw}(L_{sub} - n)) / L_{imp} \quad , \quad n = 1, 2, \dots, L_{imp}$$

11 (5.3.15.2-53)

$$s_{WI}(n) = s'_{mw}(n) \quad , \quad n = L_{sub}, \dots, L_{ex_sub} - 1$$

12 (5.3.15.2-54)

14 5.3.15.2.5.3 Making a Decision to Use Interpolated Modified Weighted Speech

15 If $Cond_1$ and $Cond_3$ are met, or $R_{P_SM_opt} > 0.4$ calculate the normalized correlation
16 between the original modified target signal and the interpolated one:

$$R_{P_imp} = \frac{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM_M}(n) \cdot s_{WI}(n)}{\sqrt{\left(\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM_M}(n)^2 \right) \left(\sum_{n=0}^{L_{ex_sub}-1} s_{WI}(n)^2 \right)}}$$

17 (5.3.15.2-55)

18 The decision is made by:

$$\begin{aligned}
& \text{if } (R_{P_tmp} > 1.1 R_{P_SM_opt}) \{ \\
& \quad R_{P_SM_opt} = R_{P_tmp} \\
& \quad s'_{mw}(n) = s_{WT}(n), \quad n=0,1,\dots,L_{ex_sub} - 1 \\
& \}
\end{aligned} \tag{5.3.15.2-56}$$

5.3.15.2.6 Final Determination of Current Local Delay and Accumulated Delay

Correct the local delay:

$$\begin{aligned}
& \text{If } (L_{center} > L_{sub}) \\
& \quad D_L \leftarrow D_L \cdot L_{sub} / L_{center}
\end{aligned} \tag{5.3.15.2-57}$$

If the accumulated delay is larger than the maximum limit, the local delay will be set to zero and the accumulated delay will be kept the same:

$$\begin{aligned}
& \text{If } (|\Delta_{SM} + D_L| > 20), \\
& \quad s'_{mw}(n) = \sum_{i=-10}^{10} w_s(f(T_0), i) \cdot s_w(I(T_0) + n + i), \quad n = 0, 1, \dots, L_{ex_sub} - 1 \\
& \text{else} \\
& \quad \Delta_{SM} \leftarrow \Delta_{SM} + D_L
\end{aligned} \tag{5.3.15.2-58}$$

5.3.15.2.7 Updating Modified Weighted Speech Buffer

Shift the old buffer:

$$s_{mw_buff}(n) \leftarrow s_{mw_buff}(n + L_{sub}), \quad 0 \leq n < L_{mw_buff} - L_{sub}, \tag{5.3.15.2-59}$$

where $s_{mw_buff}(n)$, $n = 0, \dots, L_{mw_buff} - 1$, is the buffer vector of the modified weighted speech.

Then, add the current modified weighted speech signal to the buffer:

$$s_{mw_buff}(n + L_{mw_buff} - L_{sub}) = s'_{mw}(n), \quad 0 \leq n < L_{sub} \tag{5.3.15.2-60}$$

5.3.15.2.8 Pitch Pre-Enhancement and Periodicity Detection

The pitch periodicity of the modified original signal is pre-enhanced according to different conditions. The pitch gain and the normalized pitch correlation are also estimated.

Calculate the normalized correlation between the current modified weighted speech signal and the signal modification target:

$$R_{P_SM} = \frac{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM}(n) \cdot s'_{mw}(n)}{\sqrt{\left(\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM}(n)^2\right) \left(\sum_{n=0}^{L_{ex_sub}-1} s'_{mw}(n)^2\right)}} \quad (5.3.15.2-61)$$

2 Calculate the open-loop pitch gain:

$$G_{P_SM} = \frac{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM}(n) \cdot s'_{mw}(n)}{\sum_{n=0}^{L_{ex_sub}-1} T_{G_SM}(n)^2} \quad (5.3.15.2-62)$$

4 If $R_{P_SM} > 0.5$ the pitch pre-enhancement will be applied and the energy will be
5 normalized. The energy of the modified signal is saved first before adding the pre-
6 enhancement:

$$E_{tmp} = \sum_{n=0}^{L_{ex_sub}-1} s'_{mw}(n)^2 \quad (5.3.15.2-63)$$

9 Apply the pre-enhancement:

$$s'_{mw}(n) \Leftarrow s'_{mw}(n) + C_{tmp} \cdot T_{G_SM}(n), \quad n = 0, \dots, L_{ex_sub} - 1 \quad (5.3.15.2-64)$$

11 where the enhancement coefficient C_{tmp} is

$$\begin{aligned} & \text{if } M_{SMV} \geq 1 \\ & C_{tmp} = \min\{4NSR + 0.25, 0.4\} \cdot \min\{G_{P_SM}, 1.0\} \\ & \text{else} \\ & C_{tmp} = \min\{4NSR + 0.125, 0.25\} \cdot \min\{G_{P_SM}, 1.0\} \end{aligned} \quad (5.3.15.2-65)$$

13 Calculate the energy normalization gain:

$$G_{tmp} = \sqrt{\frac{E_{tmp}}{\sum_{n=0}^{L_{ex_sub}-1} s'_{mw}(n)^2}} \quad (5.3.15.2-66)$$

15 Apply the energy gain normalization:

$$s'_{mw}(n) \Leftarrow G_{tmp} \cdot s'_{mw}(n), \quad n = 0, 1, 2, \dots, L_{ex_sub} - 1 \quad (5.3.15.2-67)$$

17 Then, copy the temporal modified signal to the modified signal vector:

$$s_{mw}(i_0 + n) = s'_{mw}(n), \quad n = 0, \dots, L_{sub} - 1 \quad (5.3.15.2-68)$$

19 After obtaining the above information, the temporal class information C_{VUV} could be
20 corrected under some conditions:

if $((APR_1 < 1/30$ or $G_{P_SM} > 0.7$ or $R_{P_SM} > 0.6)$ and

$$\begin{aligned} & (L_{ex_sub} > T_{peak} - T_0 + 2) \text{ and } (C_{VUV} = 2) \\ & C_{VUV} = 3 \end{aligned} \quad (5.3.15.2-69)$$

2 Before finishing the subframe based processing for the current variable subframe, some
3 temporal pre-parameters are calculated and memorized as described below. These
4 parameters obtained from the current variable subframe will be later converted for the fixed
5 subframe.

6 Update the minimum pitch correlation:

$$\begin{aligned} & \text{if } (R_{P_SM} < R_{SM_min}), \\ & R_{SM_min} = R_{P_SM} \end{aligned} \quad (5.3.15.2-70)$$

8 Pre-calculate the average pitch correlation:

$$R_{SM_avg} \leftarrow R_{SM_avg} + R_{P_SM} \cdot L_{sub} \quad (5.3.15.2-71)$$

10 Pre-calculate the pitch gain and the pitch correlation for each fixed subframe:

$$\begin{aligned} & \text{if } (Rate = Full_Rate), \\ & L_{tmp} = \min\{i_0 + L_{sub}, (i+1) \cdot 40\} - \max\{i_0, i \cdot 40\} \\ & g_a(i) \leftarrow g_a(i) + G_{P_SM} \cdot \max\{L_{tmp}, 0\} \\ & R_{SM}(i) \leftarrow R_{SM}(i) + R_{P_SM} \cdot \max\{L_{tmp}, 0\} \\ & i = 0,1,2,3 \\ & \text{else} \end{aligned} \quad (5.3.15.2-72)$$

$$\begin{aligned} & L_{tmp} = \min\{i_0 + L_{sub}, (i+1) \cdot 53\} - \max\{i_0, i \cdot 53\} \\ & g_a(i) \leftarrow g_a(i) + G_{P_SM} \cdot \max\{L_{tmp}, 0\} \\ & R_{SM}(i) \leftarrow R_{SM}(i) + R_{P_SM} \cdot \max\{L_{tmp}, 0\} \\ & i = 0,1,2 \end{aligned}$$

12 Finally, the starting point for the next variable subframe is updated:

$$i_0 \leftarrow i_0 + L_{sub} \quad (5.3.15.2-73)$$

14 **5.3.15.3 Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch** 15 **Gain for each Fixed Subframe**

16 This section describes the detection of stationary voiced speech frame (type-1 frame) and
17 the calculation of pitch gain for each fixed subframe. The algorithm for Rate 1 is slightly
18 different from that for Rate 1/2. The type-1 frame detection will be helped by the average
19 pitch correlation, which is finalized by:

$$R_{SM_avg} \leftarrow R_{SM_avg} / L_{frm} \quad (5.3.15.3-1)$$

22 5.3.15.3.1 Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch Gains for 23 Rate 1

24 Calculate the pitch gain and the pitch correlation for each fixed subframe:

$$\begin{aligned}
 1 \quad R_{SM}(i) &\leftarrow R_{SM}(i)/40, \quad i = 0,1,2,3 \\
 g_a(i) &\leftarrow g_a(i)/40, \quad i = 0,1,2,3
 \end{aligned} \tag{5.3.15.3-2}$$

2 Modify pitch gain for each subframe:

$$\begin{aligned}
 & \text{if } (F_{flat} = 1), \\
 & \quad g_a(i) \leftarrow \max\{\min\{g_a(i) \cdot (0.75 + 0.25 R_{SM}(i)), 1.2\}, 0\} \\
 & \quad \quad \quad i = 0,1,2,3 \\
 3 \quad & \text{else} \\
 & \quad g_a(i) \leftarrow \max\{\min\{g_a(i) \cdot (0.7 + 0.3 R_{SM}(i)), 1.2\}, 0\} \\
 & \quad \quad \quad i = 0,1,2,3
 \end{aligned} \tag{5.3.15.3-3}$$

4 Detect stationary voiced speech:

$$\begin{aligned}
 & \text{if } (R_{SM_last} > 0.6 \text{ and } R_{SM_avg} > 0.6 \text{ and } C_{SM} > 0): \\
 5 \quad & \quad G_{P_avg} = \frac{g_a(0) + g_a(1) + g_a(2) + g_a(3)}{4}
 \end{aligned} \tag{5.3.15.3-4}$$

6 Calculate:

$$\begin{aligned}
 7 \quad C_1 &= |g_a(0) - G_{P_avg}| < 0.25 \text{ and } |g_a(1) - G_{P_avg}| < 0.25 \text{ and} \\
 & |g_a(2) - G_{P_avg}| < 0.25 \text{ and } |g_a(3) - G_{P_avg}| < 0.25
 \end{aligned} \tag{5.3.15.3-5}$$

$$8 \quad C_2 = g_a(0) > 0.5 \text{ and } g_a(1) > 0.5 \text{ and } g_a(2) > 0.5 \text{ and } g_a(3) > 0.5 \tag{5.3.15.3-6}$$

$$\begin{aligned}
 9 \quad C_3 &= R_{SM}(0) > 0.75 \text{ and } R_{SM}(1) > 0.75 \text{ and } R_{SM}(2) > 0.75 \text{ and } R_{SM}(3) > 0.75 \\
 10 &
 \end{aligned} \tag{5.3.15.3-7}$$

$$\begin{aligned}
 11 \quad & \text{If } (C_1 \text{ and } C_2 \text{ and } R_{SM_min} > 0.5) \text{ or } (C_3 \text{ and } G_{P_avg} > 0.5) \\
 & \quad C_{FR} = 6
 \end{aligned} \tag{5.3.15.3-8}$$

12 Update last pitch correlation:

$$13 \quad R_{SM_last} = R_{SM}(3) \tag{5.3.15.3-9}$$

14 5.3.15.3.2 Detecting Stationary Voiced Speech and Calculating Open-Loop Pitch Gains for
15 Rate 1/2

16 Calculate the pitch gain and the pitch correlation for each fixed subframe:

$$\begin{aligned}
 17 \quad R_{SM}(i) &\leftarrow R_{SM}(i)/53, \quad i = 0,1,2 \\
 g_a(i) &\leftarrow g_a(i)/53, \quad i = 0,1,2
 \end{aligned} \tag{5.3.15.3-10}$$

18 Modify pitch gain for each subframe:

$$\begin{aligned}
 19 \quad g_a(i) &\leftarrow \max\{\min\{g_a(i) \cdot (0.75 + 0.25 R_{SM}(i)), 1.2\}, 0\} \\
 & \quad \quad \quad i = 0,1,2
 \end{aligned} \tag{5.3.15.3-11}$$

20 Detect stationary voiced speech:

if ($R_{SM_avg} > 0.6$ and $C_{SM} > 0$):

$$G_{P_avg} = \frac{g_a(0) + g_a(1) + g_a(2)}{3} \quad (5.3.15.3-12)$$

Calculate:

$$C_1 = |g_a(0) - G_{P_avg}| < 0.25 \text{ and } |g_a(1) - G_{P_avg}| < 0.25 \text{ and } |g_a(2) - G_{P_avg}| < 0.25 \quad (5.3.15.3-13)$$

$$C_2 = G_p(0) > 0.4 \text{ and } G_p(1) > 0.4 \text{ and } G_p(2) > 0.4 \quad (5.3.15.3-14)$$

$$\text{If } (C_1 \text{ and } C_2 \text{ and } R_{SM_min} > 0.4) \\ C_{FR} = 6 \quad (5.3.15.3-15)$$

Update last pitch correlation:

$$R_{SM_last} = R_{SM}(2) \quad (5.3.15.3-16)$$

9
10

5.3.16 Final Frame Class and Type Decisions, Rate Selection, and Pitch Lag Refinement

5.3.16.1 Final Classification, Type and Rate Selection

The final decision on the frame class, the frame type, and the rate, is made according to the final values of several parameters, $Rate$, $Rate^{prv}$, C_{FR} , C_{FR}^{prv} , C_{VUV} , $l_p^{ol}(0)$, $l_p^{ol}(1)$, l_{pp}^q , $l_{pp}^q(frm-1)$, M_{SMV} , F_{MUS} , F_{HRM} . The parameters are compared to each other and to a set of fixed thresholds to obtain the final frame class, the final frame type, and the final coding rate for the frame.

The following frame rate/type transitions shall be disallowed at the SMV encoder

- a) Rate 1 type-1 frame following a Rate 1/8 frame
- b) Rate 1/2 type-1 frame following a Rate 1/8 frame
- c) Rate 1 type-1 frame following a Rate 1/4 frame
- d) Rate 1/2 type-1 frame following a Rate 1/4 frame
- e) Rate 1/8 frame following a Rate 1 type-1 frame
- f) Rate 1/8 frame following a Rate 1/2 type-1 frame
- g) Rate 1/4 frame following a Rate 1 type-1 frame
- h) Rate 1/4 frame following a Rate 1/2 type-1 frame

Transitions a, b, c, and d are disallowed by the rate-selection procedure described in Section 5.3.13. The transitions e, f, g, and h are disallowed here by using a Rate 1/2 type-0 frame in place of the first Rate 1/8 or Rate 1/4 frame that immediately follows a type-1 frame.

If the current frame is of type 1, and the previous frame is of type 1, and the pitch lag of the current frame is different from that of the previous frame by more than 15, the current frame shall be modified to a Rate 1, type-0 frame.

Table 5.3-8 and Table 5.3-9 show the different rates that are used to encode the different speech classes, C_{FM} , for the different encoding modes of SMV. The optimal rate selection for each of these classes is based on the various speech parameters described above.

C_{FM}	Rate 1/8	Rate 1/4	Rate 1/2	Rate 1
silence ($C_{FM} = 0$)	√		√	
noise-like ($C_{FM} = 1$)			√	√
unvoiced ($C_{FM} = 2$)			√	√
onset ($C_{FM} = 3$)				√
non-stat. voiced ($C_{FM} = 5$)				√
stat. voiced ($C_{FM} = 6$)				√

Table 5.3-8: SMV Class-Rate Map (Mode 0)

1
2
3
4
5
6

C_{FM}	Rate 1/8	Rate 1/4	Rate 1/2	Rate 1
silence ($C_{FM} = 0$)	√			
noise-like ($C_{FM} = 1$)		√	√	
unvoiced ($C_{FM} = 2$)		√	√	
onset ($C_{FM} = 3$)		√	√	√
non-stat. voiced ($C_{FM} = 5$)			√	√
stat. voiced ($C_{FM} = 6$)			√	√

7

8

Table 5.3-9: SMV Class-Rate Map (Modes 1, 2 and 3)

9

10

5.3.16.2 Pitch Lags Refinement

12 The open-loop pitch lags are modified if the rate is changed from Rate 1 to Rate 1/2 by the
13 F_{HRM} signal, and if $C_{SM} > 0$, according to:

$$14 \quad l_p^{ol}(0) = \begin{cases} l_p^{ol}(0) & 30 \leq l_p^{ol}(0) \\ 2 \cdot l_p^{ol}(0) & l_p^{ol}(0) < 30 \end{cases} \quad (5.3.16-1)$$

$$15 \quad l_p^{ol}(1) = \begin{cases} l_p^{ol}(1) & 30 \leq l_p^{ol}(1) \\ 2 \cdot l_p^{ol}(1) & l_p^{ol}(1) < 30 \end{cases} \quad (5.3.16-2)$$

16 The subframe lags to be used as initial search points for the closed-loop pitch search for
17 Rate 1 are set according to:

$$18 \quad l_p^{SF}(0) = l_p^{SF}(1) = l_p^{ol}(0), \quad l_p^{SF}(2) = l_p^{SF}(3) = l_p^{ol}(1) \quad (5.3.16-3)$$

19 For all other rates the subframe lags to be used as initial search points for the closed-loop
20 pitch search are set according to:

$$21 \quad l_p^{SF}(0) = l_p^{ol}(0), \quad l_p^{SF}(1) = l_p^{ol}(1) \quad (5.3.16-4)$$

22 If $C_{SM} > 0$, the last subframe lag to be used as an initial search point is further modified to
23 the integer part of the quantized open-loop pitch at the end of the encoding frame. For Rate
24 1 $l_p^{SF}(3) = l_{pp,INT}^q$, while for all other rates $l_p^{SF}(1) = l_{pp,INT}^q$.

25

26

1 **5.3.17 Adaptive-codebook Gain De-emphasis for Type-1 Frame (Rate 1**
 2 **and Rate 1/2)**

3 The adaptive-codebook gain vector $g_a(m)$ for type-1 frame will be slightly attenuated from
 4 their calculated values. The gains for the 4 subframes of Rate 1 or the 3 subframes of Rate
 5 1/2 is modified by the following steps:

6 If F_{noisyV} is 1

$$7 \quad g_a(m) = 0.95 \cdot g_a(m) \quad m = 1,2,3,[4] \quad (5.3.17-1)$$

8 else if F_{noisyV} is 2

$$9 \quad g_a(m) = 0.92 \cdot g_a(m) \quad m = 1,2,3,[4] \quad (5.3.17-2)$$

10 Furthermore, if F_{flat} is 1

$$11 \quad g_a(m) = \alpha \cdot g_a(m) \quad m = 1,2,3,[4] \quad (5.3.17-3)$$

12 where α is 0.95 if M_{SMV} is 0, otherwise, it is 0.97.

13

14

1 **5.3.18 Adaptive-codebook Gain Quantization for Type-1 Frame (Rate 1**
 2 **and Rate 1/2)**

3
 4 Routine name: GEQ_Quant_PitchGain

5
 6 **Input:**

- 7 • Unquantized adaptive-codebook gain vector, $g_a(m)$, $m = 1,2,3,[4]$
- 8 • Quantization table for the selected rate, $G_a^{Rate}(i, m)$, $m = 1,2,3,[4]$

9
 10 **Output:**

- 11 • Quantized adaptive-codebook gain vector, $g_a(m)$, $m = 1,2,3,[4]$

12
 13 The adaptive-codebook gain vector $g_a(m)$ for Type-1 frame is quantized to 6 and 4 bits for
 14 Rate 1 and Rate 1/2, respectively. The best vector for $g_a(m)$ is found as the vector that
 15 satisfies:

$$16 \quad \arg \min_{i=0, \dots, N_j-1} \left\{ \sum_{m=1}^{M_j} (g_a(m) - G_a^{Rate}(i, m))^2 \right\} \quad (5.3.18-1)$$

17 For Rate 1 $M_1 = 4$ and $N_1 = 64$, and for Rate 1/2 $M_2 = 3$ and $N_2 = 16$.

18 The quantization of the adaptive-codebook gains for type-1 frames is performed on a frame
 19 basis prior to any subframe processing.

20
 21

1 **5.3.19 Generating Modified Speech from Modified Weighted Speech**

2 Routine name: PWF_wspeech_to_speech

3

4 **Input:**

- 5 • Modified weighted speech, $s_{wm}(n)$
- 6 • Weighing filter numerator coefficients, $\hat{a}_m^{num}(i)$, for $m = 0,1,2,3,4$
- 7 • Weighing filter denominator coefficients, $\hat{a}_m^{den}(i)$, for $m = 0,1,2,3,4$
- 8 • Adaptive low-pass filter parameters, μ_m , for $m = 0,1,2,3,4$

9 **Output:**

- 10 • Modified speech, $s_{mod}(n)$

11

12 The modified speech, $s_{mod}(n)$, is generated from the modified weighted speech, $s_{wm}(n)$, by
 13 the inverse weighting filter and the inverse of the adaptive low-pass filtering. The transfer
 14 function of the inverse of the fixed weighted synthesis filter for each 5 ms subframe,
 15 indexed by $m = 0,1,2,3$, is given by:

$$16 \quad H_m(z) = \frac{1 + \sum_{i=1}^{10} a_m^{den}(i) \cdot z^{-i}}{1 + \sum_{i=1}^{10} a_m^{num}(i) \cdot z^{-i}} \quad (5.3.16.2-1)$$

17 The filter parameters, $a_m^{den}(i)$ and $a_m^{num}(i)$ for $m = 0,1,2,3$, are the weighted filter
 18 parameters for each of the four 5 ms subframes, calculated in Section 5.3.5. The inverse
 19 adaptive low-pass filter is given by:

$$20 \quad L_m(z) = \frac{1}{1 + \mu_m \cdot z^{-1}} \quad (5.3.16.2-2)$$

21 where μ_m , for $m = 0,1,2,3,4$, was calculated in Section 5.3.8.

22

23

24

1 **5.3.20 LSFs Smoothing Parameter**

2 Routine name: SMO_lsf_smooth_est

3

4 **Input:**

- 5 • Voice activity decision, VAD
- 6 • Previous frame voice activity decision, VAD^{prv}
- 7 • First element of reflection coefficient vector for the second quarter of the encoding
- 8 frame, $k_1(0)$
- 9 • LSFs vector calculated at the fourth quarter of the encoding frame, $lsf_1(i)$

10 **Output:**

- 11 • LSFs smooth factor, β_{SMO}^{LSF}

12

13 An LSFs smoothing factor is calculated, based on the VAD, the previous VAD, the first
 14 reflection coefficient corresponding to the second quarter of the encoding frame, and two
 15 spectral distortion measures that are calculated by the LSFs smoothing routine. The first
 16 spectral distortion measure is calculated between the current frame LSFs, $lsf_1(i)$, and the
 17 previous frame LSFs, $lsf_{1,prv}(i)$:

$$18 \quad SD_1^{SMO} = \sum_{i=1}^{10} \left(lsf_1(i) - lsf_{1,prv}(i) \right)^2 \quad (5.3.16.2-1)$$

19 The second spectral distortion is calculated by finding the difference between the current
 20 frame LSFs, $lsf_1(i)$, and an average of the LSFs, \overline{lsf}^{SMO} :

21

$$22 \quad sd^{SMO} = \sum_{i=1}^{10} \left(lsf_1(i) - \overline{lsf}^{SMO}(i) \right)^2 \quad (5.3.16.2-2)$$

23 The sum of sd^{SMO} over the current frame and the last 3 frames gives the second spectral
 24 distortion measure, SD_2^{SMO} . The input parameters and the two spectral distortion
 25 measures are tested and compared to a set of thresholds, to generate a noise stationarity
 26 index, I_{nst} , which has a value at the range $[0, \dots, 5]$. The LSFs smoothing factor is given by:

$$27 \quad \beta_{SMO}^{LSF} = SMO_{tab}[I_{nst}] \quad (5.3.16.2-3)$$

28 where the values in the LSF smoothing table, SMO_{tab} , are:

$$1 \quad SMO_{tab} = \left[0, 0, \frac{0.9}{16}, \frac{0.9 \cdot 4}{16}, \frac{0.9 \cdot 9}{16}, 0.9 \right].$$

(5.3.16.2-4)

2

3

1 **5.3.21 LSFs Smoothing**

2

3 The average of the LSFs is updated according to:

$$4 \quad \overline{lsf}^{SMO}(i) = \beta_{LSF}^{SMO} \cdot \overline{lsf}^{SMO}(i) + (1 - \beta_{LSF}^{SMO}) \cdot lsf_1(i) \quad i = 1, \dots, 10 \quad (5.3.16.2-1)$$

5 For the SMV operation Mode 3, or if the SMV in a Half-Rate-Max Mode, or if the flag
6 F_{LPC_onset} is on, the LSFs that are quantized and transmitted, $\underline{lsf}_1(i)$ are generated
7 according to:

$$8 \quad \underline{lsf}_1(i) = \beta_{LSF}^{SMO} \cdot \underline{lsf}_{1,prv}(i) + (1 - \beta_{LSF}^{SMO}) \cdot (0.75 \cdot lsf_2(i) + 0.25 \cdot lsf_1(i)) \quad i = 1, \dots, 10$$

9 (5.3.16.2-2)

10 where $\underline{lsf}_{1,prv}$ are the previous frame smoothed LSFs.

11 For all other cases, the LSFs that are quantized and transmitted, $\underline{lsf}_1(i)$ are generated
12 according to:

$$13 \quad \underline{lsf}_1(i) = \beta_{LSF}^{SMO} \cdot \underline{lsf}_{1,prv}(i) + (1 - \beta_{LSF}^{SMO}) \cdot lsf_1(i) \quad i = 1, \dots, 10 \quad (5.3.16.2-3)$$

14

15

1 **5.3.22 LSFs Quantization**

2 Routine name: LSF_Q_lsfqnt

3

4 **Input:**

- 5 • Smoothed LSFs, $\underline{lsf}_1(i)$
- 6 • SMV rate, $Rate$

7 **Output:**

- 8 • Quantized LSFs, $\underline{lsf}_1^q(i)$

9 The LSFs are quantized with 25 bits, 21 bits, 20 bits, and 11 bits, for the Rate 1, the Rate
10 1/2, the Rate 1/4, and the Rate 1/8, respectively. The quantization is performed using a
11 weighted-mean-squared error measure, and a multi-stage vector quantization scheme.

12

13 **5.3.22.1 Calculation of LSFs Quantization Weights**

14 The LSFs quantization weights are given by:

$$15 \quad w_{LSF}(i) = \left(P_{LP}^2(\underline{lsf}_1(i)) \right)^{-0.2} \quad (5.3.22.1-1)$$

16 Where $P_{LP}^2(\underline{lsf}_1(i))$ is the LP power spectrum at the $\underline{lsf}_1(i)$, which can be calculated, for
17 example, by:

$$18 \quad P_{LP}^2(\underline{lsf}_1(i)) = 4 \cdot \left(2 - 2 \cos(2\pi \cdot \underline{lsf}_1(i)) \right) \left(\cos(5 \cdot 2\pi \cdot \underline{lsf}_1(i)) + p'_1 \cdot \cos(4 \cdot 2\pi \cdot \underline{lsf}_1(i)) + \right. \\ \left. p'_2 \cdot \cos(3 \cdot 2\pi \cdot \underline{lsf}_1(i)) + p'_3 \cdot \cos(2 \cdot 2\pi \cdot \underline{lsf}_1(i)) + p'_4 \cdot \cos(2\pi \cdot \underline{lsf}_1(i)) + p'_5 / 2 \right) \\ 19 \quad i = 2, 4, 6, 8, 10 \quad (5.3.22.1-2)$$

$$20 \quad P_{LP}^2(\underline{lsf}_1(i)) = 4 \cdot \left(2 + 2 \cos(2\pi \cdot \underline{lsf}_1(i)) \right) \left(\cos(5 \cdot 2\pi \cdot \underline{lsf}_1(i)) + q'_1 \cdot \cos(4 \cdot 2\pi \cdot \underline{lsf}_1(i)) + \right. \\ \left. q'_2 \cdot \cos(3 \cdot 2\pi \cdot \underline{lsf}_1(i)) + q'_3 \cdot \cos(2 \cdot 2\pi \cdot \underline{lsf}_1(i)) + q'_4 \cdot \cos(2\pi \cdot \underline{lsf}_1(i)) + q'_5 / 2 \right) \\ 21 \quad i = 1, 3, 5, 7, 9 \quad (5.3.22.1-3)$$

22 where $[q'_1, q'_2, q'_3, q'_4, q'_5]$ and $[p'_1, p'_2, p'_3, p'_4, p'_5]$ are as calculated in Section 5.3.1.2.4 for the
23 fourth quarter of the encoding frame.

24 The power spectrum can be also calculated from the LSFs.

25 The weighting parameters are further modified according to

$$26 \quad w_{LSF}(i) \Leftarrow (1.64 - 0.16 \cdot i) \cdot w_{LSF}(i) \quad i = 4, \dots, 10 \quad (5.3.22.1-4)$$

27

28

1 **5.3.22.2 LSFs Mean Removal**

2 A vector that represents the mean value of the LSFs is subtracted from the LSFs,
3 generating the mean-removed LSFs vector:

$$4 \quad \underline{lsf}'_1(i) = \underline{lsf}_1(i) - M_{LSF}(i) \quad i = 1, \dots, 10 \quad (5.3.22.2-1)$$

5 **5.3.22.3 LSFs Prediction**

6 The LSFs are predicted using a moving average (MA) prediction scheme. The prediction
7 error, ε_k , is generated by:

$$8 \quad \varepsilon_j(i) = \underline{lsf}'_1(i) - \sum_{p=1}^{O_{Rate}^{MA}} b_{j,p}^{Rate}(i) \cdot \varepsilon_{j^*,p}^q(i) \quad i = 1, \dots, 10 \quad (5.3.22.3-1)$$

9 where O_{Rate}^{MA} , the order of the MA prediction, is 2 for Rate 1 and 4 for all other rates. The
10 prediction coefficients $b_{j,p}^{Rate}(i)$, for each of the possible switch prediction j ($j = 0, 1$ for Rate
11 1 and Rate 1/2, and $j = 0$ Rate 1/4 and Rate 1/8), multiply the previously chosen and
12 quantized prediction error (index by j^*) for O_{Rate}^{MA} frames (indexed by p). Two prediction-
13 error vectors, ε_1 and ε_2 , are generated for Rate 1 and Rate 1/2, and a single prediction
14 vector, ε_1 , is generated for Rate 1/4 and Rate 1/8. Different prediction coefficients are used
15 for each rate, except of the prediction coefficients for Rate 1/4, which use the second set of
16 prediction coefficients of Rate 1/2.

17 **5.3.22.4 Quantization of the LSFs Prediction Error**

18 The LSFs prediction error is quantized using a multi-stage vector quantization scheme. The
19 number of bits used for each rate is different, but the quantization tables are shared
20 between all rates and the quantization schemes are similar. A generic quantization scheme
21 is depicted in the following diagram, while the LSFs quantization for each rate is described
22 in Sections 5.3.22.4.1 to 5.3.22.4.4.
23

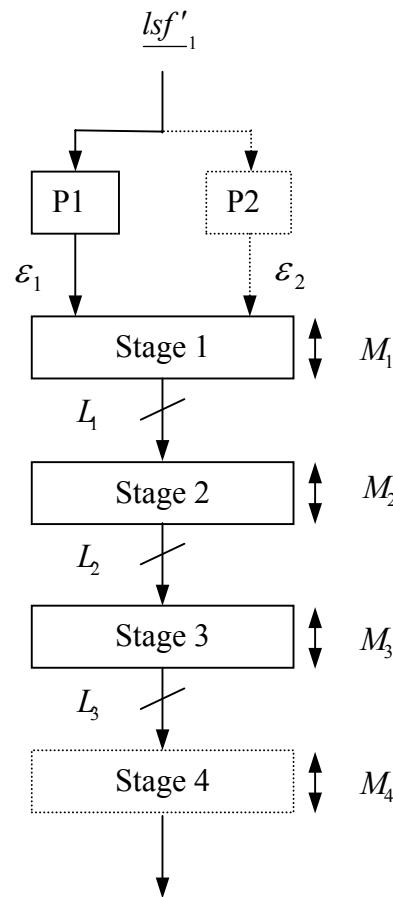


Figure 5.3-3: A Schematic Diagram of (Switched) Predictive Multi-Stage LSFs Quantization

The input to the first stage is the LSFs prediction error, and the input for each consecutive stage is the quantization error of the previous stage. Each stage has several inputs. Since two predictors are used for Rate 1 and Rate 1/2, stage 1 has two inputs for these rates. Since a single predictor is used for Rate 1/4 and Rate 1/8, stage 1 has a single input for these rates. The number of inputs for stage 2 is L_1 , the number of inputs for stage 3 is L_2 , and the number of inputs for stage 4 is L_3 . The number of code vectors in the j^{th} stage is M_j . For the j^{th} stage, denoting the input vector by $I_{LSF}^{j,l}(i)$, where l is the index of the input vector, the codebook vectors $C_{LSF}^{j,r}(i)$, the output vectors are chosen as the L_j best combination of an input vector and a code vector, minimizing the weighted mean squared error measure:

$$WMSE_{LSF}^{j,l,r} = \sum_{i=1}^{10} w_{LSF}(i) \left(I_{LSF}^{j,l}(i) - C_{LSF}^{j,r}(i) \right)^2 \quad (5.3.22.4-1)$$

The output vectors from the j^{th} stage are generated by the L_j best combination of an input vector and a codebook entry, and are given by

$$I_{LSF}^{j+1,l'}(i) = I_{LSF}^{j,l^*}(i) - C_{LSF}^{j,r^*}(i) \quad i = 1, \dots, 10 \quad (5.3.22.4-2)$$

where l' is counting the input vectors to $j^{th} + 1$ stage, and l^* and r^* denote the chosen input vector and the chosen codebook vector, respectively.

4

5.3.22.4.1 LSFs Quantization for Rate 1

Two predictors are used for Rate 1, and both ε_1 and ε_2 are the inputs to stage 1. The following table summarizes the parameters used for Rate 1:

8

M_1	M_2	M_3	M_4	L_1	L_2	L_3
128	128	64	16	8	7	5

9

Table 5.3-10: Parameters for Rate 1 LSFs Quantization

The total number of bits used for LSFs quantization for Rate 1 is 25, comprising of 1 bit for the switched prediction, 7 bits for the first stage, 7 bits for the second stage, 6 bits for the third stage, and 4 bits for the fourth stage.

14

5.3.22.4.2 LSFs Quantization for Rate 1/2

Two predictors are used for Rate 1/2, and both ε_1 and ε_2 are the inputs to stage 1. The following table summarizes the parameters used for Rate 1/2:

18

M_1	M_2	M_3	L_1	L_2
128	128	64	9	7

Table 5.3-11: Parameters for Rate 1/2 LSFs Quantization

The total number of bits used for LSFs quantization for Rate 1/2 is 21, comprising, 7 bits for the first stage, 7 bits for the second stage, and 6 bits for the third stage.

22

5.3.22.4.3 LSFs Quantization for Rate 1/4

A single predictor is used for Rate 1/4, and only ε_1 is the input to stage 1. The predictor is identical to one of the two predictors of Rate 1/2. The following table summarizes the parameters used for Rate 1/4:

27

M_1	M_2	M_3	L_1	L_2
128	128	64	9	7

1

Table 5.3-12: Parameters for Rate 1/4 LSFs Quantization

2

The total number of bits used for LSFs quantization for Rate 1/4 is 20, comprising of 7 bits for the first stage, 7 bits for the second stage, and 6 bits for the third stage.

3

4

5.3.22.4.4 LSFs Quantization for Rate 1/8

5

A single predictor is used for Rate 1/8, and only ε_1 is the input to stage 1. The following table summarizes the parameters used for Rate 1/8:

6

7

8

M_1	M_2	M_3	L_1	L_2
16	16	8	16	16

9

Table 5.3-13: Parameters for Rate 1/8 LSFs Quantization

10

The code vectors for each of the 3 stages are generated as an average of two code vectors in the respective stage used for the LSFs quantization of the other rates. The total number of bits used for LSFs quantization for Rate 1/8 is 11, comprising of 4 bits for the first stage, 4 bits for the second stage, and 3 bits for the third stage.

11

12

13

14

15

5.3.23 Interpolation of Quantized and Unquantized LSFs

The quantized LSFs are interpolated for each subframe. For Rate 1 type-0 frames the interpolation contour is quantized with 2 bits. For all other cases, no bits are transmitted for the interpolation contour, and the LSFs are linearly interpolated between the previous frame quantized LSFs and the current frame quantized LSFs. The interpolated LSFs for each subframe are converted to the quantized prediction coefficients for each subframe, $a_j^q(i)$. The range of j is [1,2,3,4] for Rate 1 frames, [1,2] for Rate 1/2 type-0 frames, and [1,2,3] for Rate 1/2 type-1 frames. The quantized LSFs are converted to a set of quantized prediction coefficients, $\hat{a}_m^q(i)$.

5.3.23.1 Interpolation of Quantized LSFs for Rate 1 Type-0 Frames

For Rate 1 type-0 frames the encoder searches for the best interpolation path and quantizes the best interpolation path with 2 bits. The search is carried by comparing the 4 possible interpolations between the current frame quantized LSFs, $lsf_1^q(i)$, and the previous frame quantized LSFs, $lsf_{1,prv}^q(i)$, with respect the LSFs with were calculated of the second quarter of the coding frame, $lsf_0(i)$.

5.3.23.1.1 Calculating of Distance Weights

A set of weights, w_{INTR} , is calculated from $lsf_0(i)$ according to:

$$\begin{aligned} w_{INTR}(1) &= (1 - lsf_0(1)) \cdot (1 - (lsf_0(2) - lsf_0(1))) \\ w_{INTR}(i) &= (1 - lsf_0(i)) \cdot (1 - \min(lsf_0(i+1) + lsf_0(i), lsf_0(i) + lsf_0(i-1))) \quad i = 2, \dots, 9 \\ w_{INTR}(10) &= (1 - lsf_0(10)) \cdot (1 - (lsf_0(10) - lsf_0(9))) \end{aligned}$$

(5.3.23.1-1)

5.3.23.1.2 Determination of Best Interpolation Path

The best interpolation path is determined by calculating the four possible weighted averages between $lsf_1^q(i)$ and $lsf_{1,prv}^q(i)$, and finding the one that is closest to $lsf_0(i)$. The weighted average, $\overline{lsf}_0^{j,q}$, for $j = 1, \dots, 4$ is calculated by:

$$\overline{lsf}_0^{j,q}(i) = C_j \cdot lsf_1^q(i) + (1 - C_j) \cdot lsf_{1,prv}^q(i) \quad i = 1, \dots, 10 \quad (5.3.23.1-2)$$

where $C_1 = 0.5$, $C_2 = 0.6$, $C_3 = 0.7$, and $C_4 = 0.8$. For each j , the distance between $\overline{lsf}_0^{j,q}(i)$ and $lsf_0(i)$ is calculated as:

$$D_{INTR}^j = \sum_{i=1}^{10} \left| \overline{lsf}_0^{j,q}(i) - lsf_0(i) \right| \cdot w_{INTR}(i) \quad (5.3.23.1-3)$$

The C_{j^*} that minimizes D_{INTR}^j is chosen, and the interpolated LSFs for each of the 4 subframes, $\hat{lsf}_m^q(i)$ are obtained by:

$$1 \quad \hat{l}sf_m^q(i) = \begin{cases} 0.5 \cdot C_{j^*} \cdot lsf_1^q(i) + (1 - 0.5 \cdot C_{j^*}) \cdot lsf_{1,prv}^q(i) & m = 0 \\ C_{j^*} \cdot lsf_1^q(i) + (1 - C_{j^*}) \cdot lsf_{1,prv}^q(i) & m = 1 \\ 0.5 \cdot (1 + C_{j^*}) \cdot lsf_1^q(i) + 0.5 \cdot (1 - C_{j^*}) \cdot lsf_{1,prv}^q(i) & m = 2 \\ lsf_1^q(i) & m = 3 \end{cases} \quad i = 1, \dots, 10$$

2 (5.3.23.1-4)

3 5.3.23.2 Interpolation of Quantized LSFs for Rate 1 Type-1 frames, Rate 1/4 4 frames, and Rate 1/8 Frames

5 For Rate 1 type-1 frames, as well as for Rate 1/4 and Rate 1/8 frames, the LSFs are
6 linearly interpolated for each subframe, and the LSFs for each of the 4 subframes, $\hat{l}sf_m^q(i)$
7 are obtained by:

$$8 \quad \hat{l}sf_m^q(i) = \begin{cases} 0.25 \cdot lsf_1^q(i) + 0.75 \cdot lsf_{1,prv}^q(i) & m = 0 \\ 0.5 \cdot lsf_1^q(i) + 0.5 \cdot lsf_{1,prv}^q(i) & m = 1 \\ 0.75 \cdot lsf_1^q(i) + 0.25 \cdot lsf_{1,prv}^q(i) & m = 2 \\ lsf_1^q(i) & m = 3 \end{cases} \quad i = 1, \dots, 10 \quad (5.3.23.2-1)$$

9 5.3.23.3 Interpolation of Quantized and Unquantized LSFs for Rate 1/2 Type-0 10 Frames

11 For Rate 1/2 type-0 frames the LSFs are linearly interpolated for each subframe, and the
12 LSFs for each of the 2 subframes, $\hat{l}sf_m^q(i)$ are obtained by:

$$13 \quad \hat{l}sf_m^q(i) = \begin{cases} 0.375 \cdot lsf_1^q(i) + 0.625 \cdot lsf_{1,prv}^q(i) & m = 0 \\ 0.875 \cdot lsf_1^q(i) + 0.125 \cdot lsf_{1,prv}^q(i) & m = 1 \end{cases} \quad i = 1, \dots, 10 \quad (5.3.23.3-1)$$

14 The unquantized LSFs are interpolated according to:

$$15 \quad \hat{l}sf_m(i) = \begin{cases} 0.75 \cdot lsf_0(i) + 0.25 \cdot lsf_{1,prv}(i) & m = 0 \\ 0.75 \cdot lsf_1^q(i) + 0.25 \cdot lsf_0(i) & m = 1 \end{cases} \quad i = 1, \dots, 10$$

16 The interpolated unquantized LSFs are converted to interpolated unquantized prediction
17 coefficients $\hat{a}_m(i)$.

18

19 5.3.23.4 Interpolation of Quantized and Unquantized LSFs for Rate 1/2 Type-1 20 Frames

21 For Rate 1/2 type-1 frames the LSFs are linearly interpolated for each subframe, and the
22 LSFs for each of the 3 subframes, $\hat{l}sf_m^q(i)$ are obtained by:

$$23 \quad \hat{l}sf_m^q(i) = \begin{cases} 0.290625 \cdot lsf_1^q(i) + 0.709375 \cdot lsf_{1,prv}^q(i) & m = 0 \\ 0.621875 \cdot lsf_1^q(i) + 0.378125 \cdot lsf_{1,prv}^q(i) & m = 1 \\ 0.953125 \cdot lsf_1^q(i) + 0.046875 \cdot lsf_{1,prv}^q(i) & m = 2 \end{cases} \quad i = 1, \dots, 10$$

24 (5.3.23.4-1)

1 The unquantized LSFs are interpolated according to:

$$2 \quad \hat{l}sf_m(i) = \begin{cases} 0.58125 \cdot lsf_0(i) + 0.41875 \cdot lsf_{1,prv}(i) & m = 0 \\ 0.25 \cdot lsf_1(i) + 0.75 \cdot lsf_0(i) & m = 1 \\ 0.9125 \cdot lsf_1(i) + 0.0875 \cdot lsf_0(i) & m = 2 \end{cases} \quad i = 1, \dots, 10$$

3

(5.3.23.4-2)

4 The interpolated unquantized LSFs are converted to interpolated unquantized prediction
5 coefficients $\hat{a}_m(i)$.

6

7

1 **5.3.24 Recalculating of Weighting Filter Coefficients for Rate 1/2**

2 Since Rate 1/2 frames have 3 subframes for type-1 frame and 2 subframes for type-0
 3 frame, the weighing filter coefficients for each subframe are recalculated. For type-1 frames,
 4 the interpolated first element of the reflection coefficients are recalculated as:

$$\begin{aligned}
 \hat{k}_0(0) &\Leftarrow 0.8375 \cdot \hat{k}_0(0) + 0.1625 \cdot \hat{k}_1(0) \\
 \hat{k}_1(0) &\Leftarrow 0.5 \cdot \hat{k}_1(0) + 0.5 \cdot \hat{k}_2(0) \\
 \hat{k}_2(0) &\Leftarrow 0.8375 \cdot \hat{k}_2(0) + 0.1625 \cdot \hat{k}_3(0)
 \end{aligned}
 \tag{5.3.24-1}$$

6 Similarly, the interpolated LPC prediction gains are recalculated as:

$$\begin{aligned}
 \hat{G}_0^{LPC} &\Leftarrow 0.8375 \cdot \hat{G}_0^{LPC} + 0.1625 \cdot \hat{G}_1^{LPC} \\
 \hat{G}_1^{LPC} &\Leftarrow 0.5 \cdot \hat{G}_1^{LPC} + 0.5 \cdot \hat{G}_2^{LPC} \\
 \hat{G}_2^{LPC} &\Leftarrow 0.8375 \cdot \hat{G}_2^{LPC} + 0.1625 \cdot \hat{G}_3^{LPC}
 \end{aligned}
 \tag{5.3.24-2}$$

8

9 For type-0 frames, the interpolated first element of the reflection coefficients are
 10 recalculated as:

$$\begin{aligned}
 \hat{k}_0(0) &\Leftarrow 0.5 \cdot \hat{k}_0(0) + 0.5 \cdot \hat{k}_1(0) \\
 \hat{k}_1(0) &\Leftarrow 0.5 \cdot \hat{k}_2(0) + 0.5 \cdot \hat{k}_3(0)
 \end{aligned}
 \tag{5.3.24-3}$$

12 Similarly, the interpolated LPC gain is recalculated as:

$$\begin{aligned}
 \hat{G}_0^{LPC} &\Leftarrow 0.5 \cdot \hat{G}_0^{LPC} + 0.5 \cdot \hat{G}_1^{LPC} \\
 \hat{G}_1^{LPC} &\Leftarrow 0.5 \cdot \hat{G}_2^{LPC} + 0.5 \cdot \hat{G}_3^{LPC}
 \end{aligned}
 \tag{5.3.24-4}$$

14 The weighing filter coefficients, $\hat{a}_m^{num}(i)$ and, $\hat{a}_m^{den}(i)$ for each subframe are recalculated for
 15 each subframe, using the recalculated and interpolated prediction coefficients, $\hat{a}_m(i)$, first
 16 element of reflection coefficients, $\hat{k}_m(0)$, and LPC gains, \hat{G}_m^{LPC} , as described in Section
 17 5.3.5.

18

19

1 **5.3.25 Identifying Long-Term Spectral Characteristic**

2 Routine name: CLA_Identify_Input

3

4 **Input:**

- 5 • Frame type, $Type$
- 6 • Interpolated quantized prediction coefficients for the first subframe, $\hat{a}_0^q(i)$
- 7 • First element of quantized LSF vector, $lsf_1^q(1)$
- 8 • Quantized adaptive-codebook gain for the last subframe of the previous frame, g_a^{prv}
- 9 • SMV rate, $Rate$

10 **Output:**

- 11 • Spectral flatness flag, F_{flat}
- 12 • LPC prediction gain from quantized parameters of the first subframe, $G_{0,q}^{LPC}$

13

14 The identification of the long-term spectral characteristics is performed by the analysis of
15 several quantized spectral and pitch characteristics over a long time period.

16 The quantized prediction coefficients of the first subframe are converted to quantized
17 reflection coefficients, $\hat{k}_{0,q}(i)$, and the quantized LPC prediction gain is calculated by:

$$18 \quad G_{0,q}^{LPC} = -10 \cdot \log_{10} \left(\prod_{i=1}^{10} (1 - k_{0,q}^2(i)) \right) \quad (5.3.25-1)$$

19 A counter of the number of consecutive Rate 1/8 frames, $C_{1/8}$, is obtained. A counter of the
20 type-1 frames, C_{type1} , is maintained and used in the calculation of a long term parametric
21 average, \bar{X} as follows:

$$22 \quad X = 2.5 \cdot \left(\min \left(\frac{G_{0,q}^{LPC}}{35}, 1 \right) - 2 \cdot k_{0,q}(1) + \min(g_a^{prv}, 1) + 15 \cdot lsf_1^q(1) \right) \quad (5.3.25-2)$$

$$23 \quad \bar{X} = \begin{cases} 0.75 \cdot X & C_{type1} = 0 \\ 0.9 \cdot \bar{X} + 0.1 \cdot X & 0 < C_{type1} < 10 \\ 0.99 \cdot \bar{X} + 0.01 \cdot X & 10 \leq C_{type1} < 100 \\ 0.998 \cdot \bar{X} + 0.002 \cdot X & 100 \leq C_{type1} \end{cases} \quad (5.3.25-3)$$

24 C_{type1} is updated after the calculation of Equation (5.3.25-2) and Equation (5.3.25-3) above,
25 and is limited to 10000. The final long-term spectral characteristic flag, F_{flat} , is updated
26 only if the frame type is 0, $C_{1/8} > 10$, and $C_{type1} > 200$, as follows:

$$F_{flat} = \begin{cases} 1 & \bar{X} > 5.8 \\ 0 & \bar{X} < 5.3 \end{cases}$$

(5.3.25-4)

2 If none of the conditions is met, F_{flat} is not changed from its previous value.

3

3

1 5.4 Excitation Coding at Rate 1/8

2 5.4.1 Excitation and Unquantized Gain Determination

3 Routine name: GCB_gauss_excit

4 **Input:**

- 5 • Seed, *seed*
- 6 • Energies of weighted residual for two halves of the frame, $E_{rw}(0)$ and $E_{rw}(1)$

7 **Output:**

- 8 • Unquantized fixed-codebook gain, g_c
- 9 • Random excitation, $r(n)$

10

11 A Gaussian random excitation $r(n)$ is generated using the random noise generator
 12 described in Section 5.6.19. For this rate, the seed used is a signed 16-bit static variable
 13 initialized to 0x0E76. The seed is re-set to its initial value whenever the rate is not the Rate
 14 1/8. The gain g_c is determined as follows:

$$15 \quad g_c = 0.7 \cdot \sqrt{\frac{\max(E_{ref} - 960, 0)}{\sum_{n=0}^{L_{frm}} e^2(n)}} \quad (5.4.1-1)$$

16 where E_{ref} is set according to the rates of the previous two frames. If either previous two
 17 frames was not coded using Rate 1/8:

$$18 \quad E_{ref} = E_{wr}(0) + E_{wr}(1) \quad (5.4.1-2)$$

19 otherwise

$$20 \quad E_{ref} = 0.5 \cdot E_{ref} + 0.5 \cdot (E_{wr}(0) + E_{wr}(1)) \quad (5.4.1-3)$$

21

22 5.4.2 Gain quantization

23 Routine name: GEQ_gainNSQMA_1_5

24 **Input:**

- 25 • Random Excitation, $r(n)$
- 26 • Unquantized fixed-codebook gain, g_c

27 **Output:**

- 28 • Quantized fixed-codebook gain, g_c

29

1 A single fixed-codebook gain g_c is quantized, using a prediction scheme, with 5 bits per
 2 frame. The quantized \hat{g}_c is found as the entry that satisfies:

$$3 \quad \arg \min \{ |g_c - \hat{g}_c| \} \quad \text{such that } \hat{g}_c \leq 5g_{mem} \quad (5.4.2-1)$$

4 where g_{mem} is the estimate of the average Rate 1/8 frame gain, and is updated after the
 5 above quantization for each Rate 1/8 frame as

$$6 \quad g_{mem} = \max(\text{round}(\max(g_r, 100)), \text{round}(0.05 \cdot \text{round}(\max(g_r, 100)) + 0.95 \cdot g_{mem})) \quad (5.4.2-2)$$

8 The entries of the codebook contain the correction factor for the predicted fixed-codebook
 9 gain. The prediction of the fixed-codebook gain is based on a 4th order MA prediction of the
 10 fixed-codebook energy.

11 The relation between the correction factor γ and the quantized fixed-codebook gain is given
 12 by

$$13 \quad \hat{g}_c = \gamma \cdot \tilde{g}_c \quad (5.4.2-3)$$

14 where \hat{g}_c is the quantized fixed-codebook gain, and \tilde{g}_c is the predicted fixed-codebook
 15 gain. The predicted fixed-codebook gain is based on a 4th order MA prediction of the fixed-
 16 codebook energy, and is given by

$$17 \quad \tilde{g}_c = 10^{\frac{1}{20}(\tilde{E} - E_c + \bar{E})} \quad (5.4.2-4)$$

18 where the $\bar{E} = 30\text{dB}$ is the mean energy,

$$19 \quad E_c = 10 \log_{10} \left(\frac{1}{L_{frm}} \sum_{n=0}^{L_{frm}-1} c_{unf}^{fix}(n)^2 \right) \quad (5.4.2-5)$$

20 and

$$21 \quad \tilde{E} = \sum_{i=1}^4 b_i \cdot (20 \log_{10} \gamma_{frm-i}^q) \quad (5.4.2-6)$$

22 where the prediction coefficients of the MA prediction are $\{b_1, b_2, b_3, b_4\} = \{0.7, 0.6, 0.4, 0.2\}$,
 23 and γ_{frm-i}^q are the quantized correction factors for the previous 4 frames.

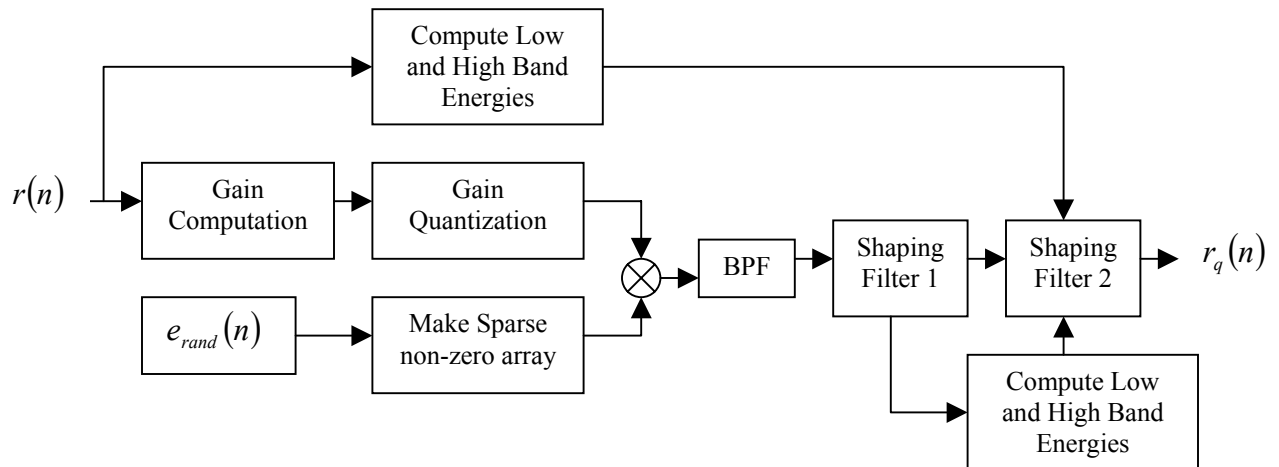
24 To avoid production of a NULL (all 1s) Rate 1/8 packet at the encoder, the Rate 1/8 gain
 25 search shall be restricted to the first 31 entries of the table when the quantized LSF indices
 26 are all 1s.

27
 27

1 5.5 Excitation Coding at Rate 1/4

2 Noise Excited Linear Prediction (NELP) is used for Rate 1/4 frames. A pseudo-random noise
 3 is colored using a set of filters, and scaled appropriately by gains in 10 subframes of 16
 4 samples each to model the excitation.

5



6

7

Figure 5.5-1: Excitation Coding for Rate 1/4 Frames

8

9 5.5.1 Computation and Quantization of Gains

10 Routine name: nelp_encoder

11 **Input:**

- 12 • Pre-processed speech, $\tilde{s}(n)$
- 13 • Quantized LSFs, $lsf_1^q(i)$

14 **Output:**

- 15 • Quantized excitation, $r_q(n)$

16

17

18 The pre-processed speech, $\tilde{s}(n)$, is converted to the residual signal $r(n)$ using the
 19 quantized LSFs, $lsf_1^q(i)$ (by converting the LSFs to LPCs, and zero filtering the pre-
 20 processed speech with the LPCs). The energy E_i of the residual signal is computed for each
 21 of the 10 subframes ($0 \leq i < 10$) as

$$22 \quad E_i = \frac{1}{16} \sum_{n=16i}^{16i+15} r^2(n) \quad (5.5.1-1)$$

23 The 10 subframe gains G_i are computed from the above energies as

$$1 \quad G_i = \sqrt{E_i} \quad (5.5.1-2)$$

2 Two normalization factors for two splits ($0 \leq j < 2$) of the above 10 gains are determined by

$$3 \quad G'[j] = \sqrt{\frac{1}{5} \sum_{i=0}^4 G_{j*5+i}^2} \quad (5.5.1-3)$$

4 The above two normalization factors are vector quantized using a two dimensional
5 codebook $UVG1CB[32][2]$ to provide an index $i\hat{G}'$ as

$$6 \quad i\hat{G}' = \arg \min_k \left(\sum_{j=0}^1 (\log_{10}(G'[j]) - UVG1CB[k][j])^2 \right), \quad 0 \leq k < 32 \quad (5.5.1-4)$$

7 The quantized normalization factors are then computed for ($0 \leq j < 2$) using

$$8 \quad \hat{G}'[j] = 10^{UVG1CB[i\hat{G}'][j]} \quad (5.5.1-5)$$

9 The 10 subframe gains are then normalized using the quantized normalization factors
10 according to

$$11 \quad G''[2j+i] = \frac{G[2j+i]}{\hat{G}'[j]}, \quad 0 \leq i < 5; 0 \leq j < 2 \quad (5.5.1-6)$$

12 The 10 normalized gains are then vector quantized using a 3-dimensional codebook
13 $UVG2CB[2][64][5]$ to provide two indices $iG''[0]$, and $iG''[1]$ such that

$$14 \quad i\hat{G}''[j] = \arg \min_k \left(\sum_{i=0}^4 (G''[2j+i] - UVG2CB[j][k][i])^2 \right), \quad 0 \leq k < 64; 0 \leq j < 2$$

15 (5.5.1-7)

16 The quantized, normalized gains are then given by

$$17 \quad \hat{G}''[2j+i] = UVG2CB[j][i\hat{G}''[j]] \quad , \quad 0 \leq i < 5; 0 \leq j < 2 \quad (5.5.1-8)$$

18 The quantized gains used to scale the NELP excitation then are

$$19 \quad \hat{G}[2j+i] = \hat{G}'[j] \hat{G}''[2j+i], \quad 0 \leq i < 5; 0 \leq j < 2 \quad (5.5.1-9)$$

20

21 **5.5.2 Random number generation**

22 A pseudo-random vector e_{rand} of 160 samples is generated as follows,

$$23 \quad NELPseed = (521 \times NELPseed + 259) / 32768 \quad (5.5.2-1)$$

$$24 \quad e_{rand}(n) = NELPseed, \quad 0 \leq n < 160 \quad (5.5.2-2)$$

25 For Rate 1/4, the seed is initialized at the beginning of each frame to a signed 16-bit
26 number that is derived from the bits representing the LSF.
27

1

2 **5.5.3 Creation of sparse non-zero excitation**

3 The above pseudo-random vector $e_{rand}(n)$ is divided into 10 subframes, and each subframe
 4 is made sparse by retaining the 4 largest absolute magnitude samples, and zeroing out the
 5 remaining 12. Each of these 10 subframes of the sparse vectors is denoted by $e_{sparse}(i, n)$,
 6 $0 \leq i < 10$; $0 \leq n < 16$. These sparse vectors are then scaled by their appropriate gains to
 7 create 160 samples of a sparse, excitation signal

$$8 \quad r_1(16i + n) = \sqrt{3}\hat{G}(i)e_{sparse}(i, n), \quad 0 \leq i < 10; 0 \leq n < 16 \quad (5.5.3-1)$$

9 **5.5.4 Shaping the excitation**

10 The above excitation $r_1(n)$ is filtered by a 12th order bandpass filter, in order to remove the
 11 low band and high components to yield the signal $r_2(n)$. This bandpass filter is given by the
 12 following transfer function

$$13 \quad BP(z) = \frac{BP_{n0} + BP_{n1}z^{-1} + \dots + BP_{n12}z^{-12}}{1 + BP_{d1}z^{-1} + \dots + BP_{d12}z^{-12}} \quad (5.5.4-1)$$

14 where the coefficients BP_{ni} , and BP_{di} are given in the tables `bp1_num_coef`, and
 15 `bp1_den_coef` respectively in the C-code.

16 $r_2(n)$ is then filtered by the 10th order shaping filter $SF1(z)$ to give the signal $r_3(n)$. This
 17 shaping filter is given by the following transfer function

$$18 \quad SF1(z) = \frac{SF1_{n0} + SF1_{n1}z^{-1} + \dots + SF1_{n12}z^{-10}}{1 + SF1_{d1}z^{-1} + \dots + SF1_{d12}z^{-10}} \quad (5.5.4-2)$$

19 where the coefficients $SF1_{ni}$, and $SF1_{di}$ are given in the tables `shapel_num_coef`, and
 20 `shapel_den_coef` respectively in the C-code.

21 $r_3(n)$ is then scaled according to the following equation to yield

$$22 \quad r_4(n) = \left(\frac{\sqrt{\sum_{n=0}^{159} [r(n)]^2}}{\sqrt{\sum_{n=0}^{159} [r_3(n)]^2}} \right) r_3(n) \quad (5.5.4-3)$$

23 $r(n)$ and $r_4(n)$ are filtered by the 10th order lowpass filter $LP(z)$ to give $r_L(n)$ and $r_{4L}(n)$.
 24 Similarly, $r(n)$ and $r_4(n)$ are filtered by the 10th order highpass filter $HP(z)$ to give $r_H(n)$
 25 and $r_{4H}(n)$.

$$26 \quad LP(z) = \frac{LP_{n0} + LP_{n1}z^{-1} + \dots + LP_{n12}z^{-10}}{1 + LP_{d1}z^{-1} + \dots + LP_{d12}z^{-10}} \quad (5.5.4-4)$$

$$27 \quad HP(z) = \frac{HP_{n0} + HP_{n1}z^{-1} + \dots + HP_{n12}z^{-10}}{1 + HP_{d1}z^{-1} + \dots + HP_{d12}z^{-10}} \quad (5.5.4-5)$$

1 where the coefficients LP_{ni} , and LP_{di} are given in the tables **txlpf1_num_coef**, and
 2 **txlpf1_den_coef** respectively in the C-code, and the coefficients HP_{ni} , and HP_{di} are given
 3 in the tables **txhpf1_num_coef**, and **txhpf1_den_coef** respectively in the C-code.

5 Two factors R_L , and R_H are computed as follows

$$6 \quad R_L = 10 \log_{10} \left(\frac{\sum_{n=0}^{n=159} [r_L(n)]^2}{\sum_{n=0}^{n=159} [r_{4L}(n)]^2} \right) \quad (5.5.4-6)$$

7 and

$$8 \quad R_H = 10 \log_{10} \left(\frac{\sum_{n=0}^{n=159} [r_H(n)]^2}{\sum_{n=0}^{n=159} [r_{4H}(n)]^2} \right) \quad (5.5.4-7)$$

9 A filter index F_{idx} is computed as follows from R_L , and R_H , and sent to the decoder.

$$10 \quad F_{idx} = \begin{cases} 1 & ; \quad R_L < -3 \\ 2 & ; \quad R_L \geq -3 \text{ and } R_H < -3 \\ 0 & ; \quad \textit{otherwise} \end{cases} \quad (5.5.4-8)$$

11 The shaping filter $SF(z)$ is determined from F_{idx} as follows

$$12 \quad SF(z) = \begin{cases} SF1(z) & ; \quad F_{idx} = 0 \\ SF1(z)SF2(z) & ; \quad F_{idx} = 1 \\ SF1(z)SF3(z) & ; \quad F_{idx} = 2 \end{cases} \quad (5.5.4-9)$$

13 where, $SF2(z)$, $SF3(z)$ are

$$14 \quad SF2(z) = \frac{SF2_{n0} + SF2_{n1}z^{-1} + \dots + SF2_{n12}z^{-10}}{1 + SF2_{d1}z^{-1} + \dots + SF2_{d12}z^{-10}} \quad (5.5.4-10)$$

$$15 \quad SF3(z) = \frac{SF3_{n0} + SF3_{n1}z^{-1} + \dots + SF3_{n12}z^{-10}}{1 + SF3_{d1}z^{-1} + \dots + SF3_{d12}z^{-10}} \quad (5.5.4-11)$$

16 where the coefficients $SF2_{ni}$, $SF2_{di}$, $SF3_{ni}$, and $SF3_{di}$ are given in the tables
 17 **shape2_num_coef**, **shape2_den_coef**, **shape3_num_coef**, and **shape3_den_coef**
 18 respectively in the C-code.

19 Finally, $r_3(n)$ is filtered by $SF(z)$ to provide the quantized formant residual $r_q(n)$ for the
 20 Rate 1/4 frame.

21
 21

5.6 Excitation Coding at Rate 1 and Rate 1/2

5.6.1 Introduction to Excitation Coding for Rate 1 and Rate 1/2

The coding of the excitation for Rate 1 and Rate 1/2 is performed for each subframe. Since the coding approach for both rates is the same, the subframe processing loop in the SMV encoder is shared between Rate 1 and Rate 1/2, while some internal sections of the subframe processing loop are distinct for each rate and for each frame type.

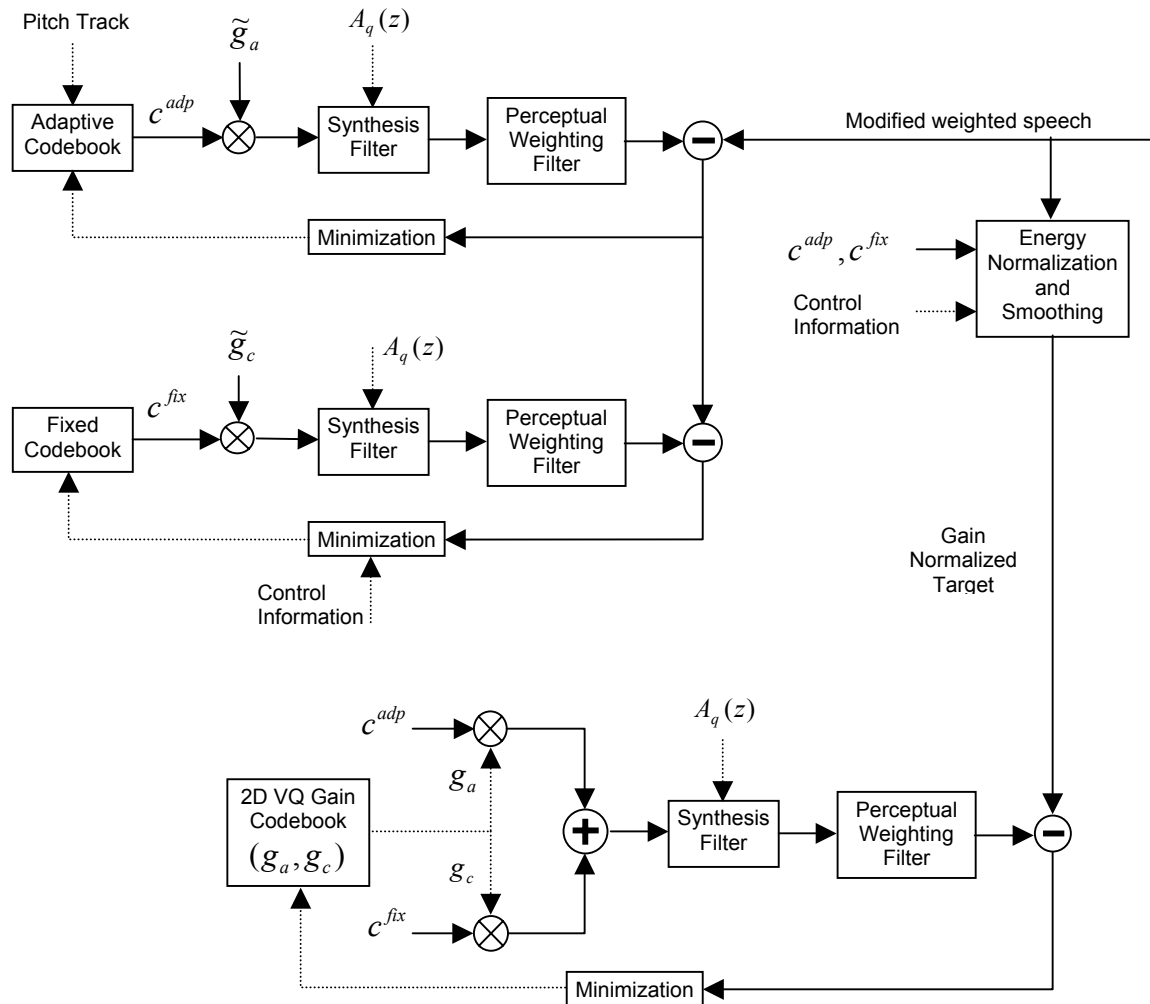


Figure 5.6-1: Subframe Excitation Coding for Rate 1 and Rate 1/2 Type-0 Frames

The subframe excitation coding for type-0 frames is depicted in Figure 5.6-1. The subframe processing includes a search for an adaptive-codebook contribution, a fixed-codebook contribution, and the joint quantization of the adaptive-codebook gain and the fixed-codebook gain, using gain-normalized target signal. The choice of the codebook vector and

the quantization of the gains are completed for each subframe before the next subframe is processed.

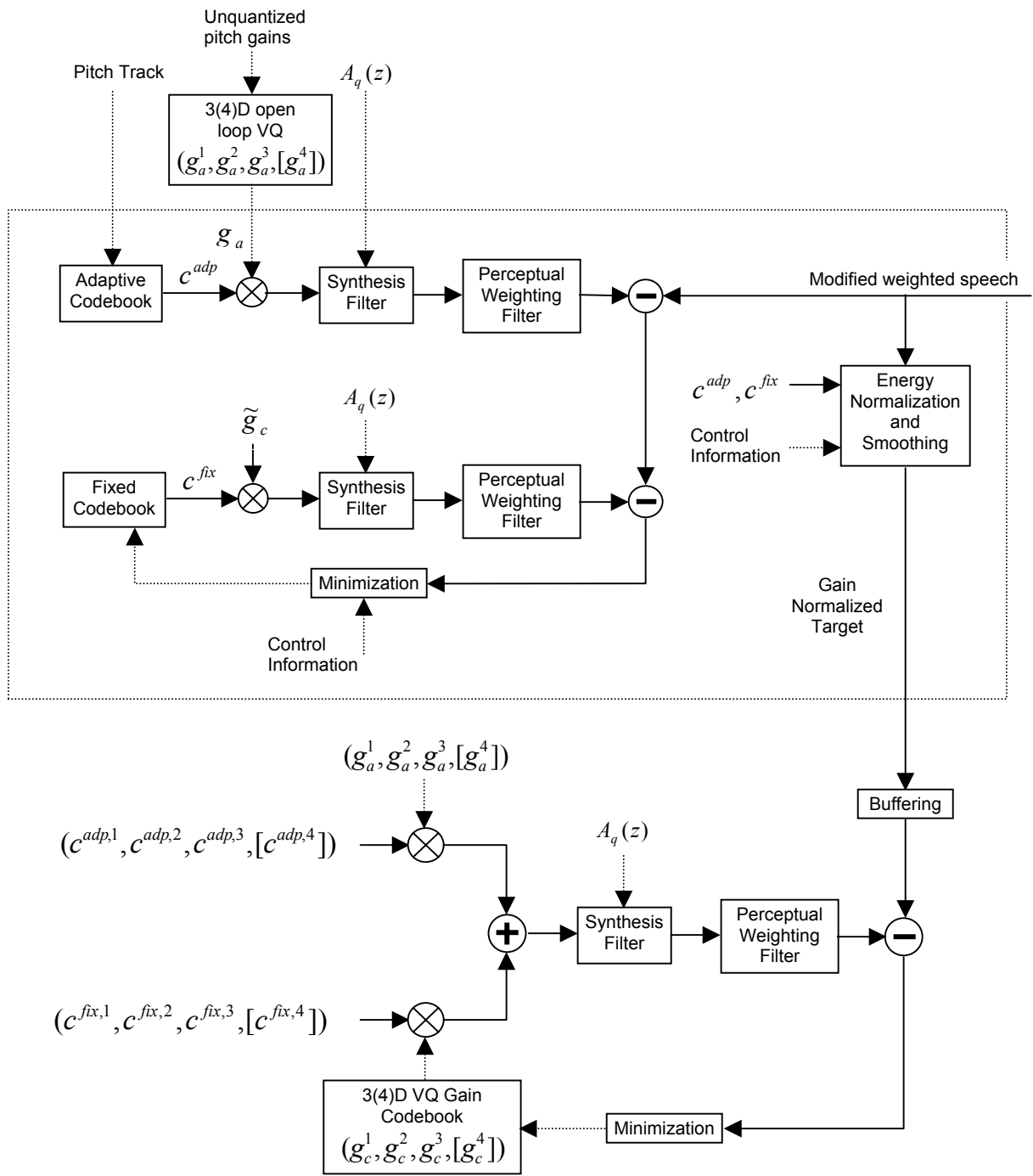


Figure 5.6-2: Subframe excitation coding for Rate 1 and Rate 1/2 type-1 frames

1 The subframe excitation coding for type-1 frames is depicted in Figure 5.6-2. Before the
 2 subframe is processed, the adaptive-codebook gains for all of the subframes are jointly
 3 quantized. The subframe processing includes only the search for the fixed-codebook
 4 contribution, since the adaptive-codebook contribution is obtained from the interpolated
 5 pitch contour. After the fixed-codebook contribution is determined for all subframes, the
 6 fixed-codebook gains are jointly quantized, using gain-normalized target signal.

7 Before the subframe processing loop is executed, the number of subframes is set, according
 8 to the rate and the frame type. For Rate 1 the number of subframes is 4, for Rate 1/2 type-
 9 0 frames the number of subframes is 2, and for Rate 1/2 type-1 frames the number of
 10 subframes is 3. At the beginning of the subframe loop the number of samples in the
 11 subframe, L_{SF} , is set, according to the following table:

	Rate 1		Rate 1/2	
	Type-0 frames	Type-1 frames	Type-0 frames	Type-1 frames
1 st subframe	40	40	80	53
2 nd subframe	40	40	80	53
3 rd subframe	40	40	—	54
4 th subframe	40	40	—	—

13 **Table 5.6-1: L_{SF} - Number of Samples in a Subframe for Rate 1 and Rate 1/2**

14 Sections 5.6.2 to 5.6.17 describe the subframe processing for Rate 1 and Rate 1/2, and are
 15 repeated for all 2,3, or 4 subframes. Note that the subframe processing loop is performed
 16 even when the selected rate is Rate 1/4 or Rate 1/8; however, the loop is performed
 17 without the excitation search and is needed only to keep the proper memory updates. For
 18 some of the symbols in the description of the excitation in the following description, the
 19 index of the subframe is omitted for simplicity.

20
 21

1 **5.6.2 Calculation of Impulse Response of the Combined Synthesis and**
 2 **Weighting Filter**

3 Routine name: LPC_ImpulseResonse

4 **Input:**

- 5 • Weighing filter numerator coefficients for the subframe, $\hat{a}_m^{num}(i)$
- 6 • Weighing filter denominator coefficients for the subframe, $\hat{a}_m^{den}(i)$
- 7 • Interpolated quantized prediction coefficients for the subframe, $\hat{a}_m^q(i)$

8 **Output:**

- 9 • Impulse response of the combined filter for the subframe, $h_{imp}(n)$

10

11 The synthesis and weighing filter are combined in the form:

$$12 \quad F(z) = \frac{A^{num}(z)}{A(z)A^{den}(z)} \quad (5.6.2-1)$$

13 The filter $A^{num}(z)$ uses the coefficients $\hat{a}_m^{num}(i)$ and the filter $A^{den}(z)$ uses the coefficients
 14 $\hat{a}_m^{den}(i)$, calculated according to Section 5.3.5 for Rate 1 and according to Section 5.3.24 for
 15 Rate 1/2. The filter $A(z)$ uses the interpolated quantized prediction coefficients, $\hat{a}_m^q(i)$,
 16 obtained according to Section 5.3.23. The impulse response of the combined filter for the
 17 length of the subframe is calculated and stored in $h_{imp}(n)$.

18

19

1 5.6.3 Generating the Target Signal

2 Routine name: PRC_TargetSignal

3 **Input:**

- 4 • Weighing filter numerator coefficients for the subframe, $\hat{a}_m^{num}(i)$
- 5 • Weighing filter denominator coefficients for the subframe, $\hat{a}_m^{den}(i)$
- 6 • Interpolated quantized prediction coefficients for the subframe, $\hat{a}_m^q(i)$
- 7 • Modified speech, $s_{mod}(n)$

8 **Output:**

- 9 • Target signal for closed-loop pitch search or pitch interpolation, $T_{gs}(n)$
- 10 • Residual modified signal, $r_{mod}(n)$, concatenated to the adaptive codebook, $C^{adp}(n)$

11

12 For each subframe (with a size defined by the rate and by the frame type), the modified
 13 speech signal, $s_{mod}(n)$ (Section 5.3.19), is filtered by the weighting filter to generate the
 14 target signal for the closed-loop search (type-0 frames) and the pitch interpolation (type-1
 15 frames). The zero-input response of the previous frame is removed from the synthesis filter.
 16 The filtering is performed in 2 steps, and result in the modified residual signal, as well as
 17 the target signal.

18 At the first step, the residual signal $r_{mod}(n)$ is generated from the modified speech $s_{mod}(n)$
 19 by the all-zero filter $A(z)$, where the filter coefficients are the interpolated prediction
 20 coefficients for each subframe. The residual for the current subframe is concatenated to the
 21 adaptive codebook from the previous subframes, to generate the adaptive codebook,
 22 $C^{adp}(n)$, for the current subframe.

23 At the second step, the target signal $T_{gs}(n)$ is generated from the residual signal $r_{mod}(n)$ by
 24 the combined filter:

$$25 \quad F(z) = \frac{A^{num}(z)}{A(z)A^{den}(z)} \quad (5.6.3-1)$$

26 The zero-input response removal is performed by the manipulation of the states of the $A(z)$
 27 and $A^{den}(z)$ filters.

28

29

1 **5.6.4 Generating the Ideal Excitation**

2 Routine name: PRC_Ideal_Excit

3 **Input:**

- 4 • Weighing filter numerator coefficients for the subframe, $\hat{a}_m^{num}(i)$
- 5 • Weighing filter denominator coefficients for the subframe, $\hat{a}_m^{den}(i)$
- 6 • Interpolated quantized prediction coefficients for the subframe, $\hat{a}_m^q(i)$
- 7 • Target signal for closed-loop pitch search or pitch interpolation, $T_{gs}(n)$

8 **Output:**

- 9 • Ideal excitation, $I_{ext}(n)$

10

11 The ideal excitation $I_{ext}(n)$ is generated from the target signal $T_{gs}(n)$ as the zero-state
12 response of the filter:

13
$$\hat{F}(z) = \frac{A(z)A_{den}(z)}{A_{num}(z)} \quad (5.6.4-1)$$

14

15

1 **5.6.5 Adaptive-codebook Contribution for Type-1 Frames**

2 Routine name: LTP_PP_pitch_ext

3

4 **Input:**

- 5 • Pitch contour, $Pit(n)$
- 6 • Target signal for closed-loop pitch search or pitch interpolation, $T_{gs}(n)$
- 7 • Adaptive codebook, $C^{adp}(n)$
- 8 • Impulse response of the combined filter for the subframe, $h_{imp}(n)$

9 **Output:**

- 10 • Adaptive-codebook contribution, $c_{unf}^{adp}(n)$
- 11 • Filtered adaptive-codebook contribution, $c_f^{adp}(n)$
- 12 • Fractional pitch for the subframe, l_{FRC}^p
- 13 • Integer pitch for the subframe, l_{INT}^p
- 14 • Normalized closed-loop pitch correlation, R_p^{wgt}

15

16 For type-1 frames, the adaptive-codebook vector contribution for the subframe is generated
 17 using the interpolated pitch contour, $Pit(n)$ (Section 5.3.14.2). The filtered adaptive-
 18 codebook contribution is obtained by the convolution of the adaptive-codebook vector, and
 19 the normalized adaptive-codebook correlation is derived.

20

21 **5.6.5.1 The Adaptive-codebook Vector Contribution**

22 The adaptive codebook vector contribution for the subframe, $c_{unf}^{adp}(n)$, is obtained by pitch-
 23 filtering the adaptive codebook $C^{adp}(n)$ according to:

$$24 \quad c_{unf}^{adp}(n) = C^{adp}(n) = C^{adp}(n - Pit(n)) \quad n = 0, \dots, L_{SF} - 1 \quad (5.6.5.1-1)$$

25 Since $Pit(\cdot)$ can have fractional values, a set of Sinc windows (sampled at 21 different
 26 phases) is used to obtain $c_{unf}^{adp}(n)$ values at fractional time locations. The fractional pitch for
 27 the subframe, l_{FRC}^p , is set to $Pit(L_{SF} / 2)$ and l_{INT}^p is set to the integer part of the fractional
 28 pitch, i.e. $I(l_{FRC}^p)$.

29

30

1 **5.6.5.2 The Filtered Adaptive Vector Contribution and Normalized Pitch Correlation**

2 The filtered adaptive-codebook vector, $c_f^{adp}(n)$, is obtained by a zero-state convolution of
 3 $c_{unf}^{adp}(n)$ with $h_{imp}(n)$. The normalized pitch correlation is obtained by:

4
$$R_p^{wgt} = \frac{\max\left(\sum_{n=0}^{L_{SF}-1} T_{gs}(n) \cdot c_f^{adp}(n), 0\right)}{\sqrt{\sum_{n=0}^{L_{SF}-1} (T_{gs}(n))^2 \cdot \sum_{n=0}^{L_{SF}-1} (c_f^{adp}(n))^2}}$$
 (5.6.5.2-1)

5

6

5.6.6 Adaptive-codebook Contribution For Rate 1 Type-0 Frames

Routine name: LTP_close_8_5k_pitch

Input:

- Target signal for closed-loop pitch search or pitch interpolation, $T_{gs}(n)$
- Adaptive codebook, $C^{adp}(n)$
- Impulse response of the combined filter for the subframe, $h_{imp}(n)$
- Integer open-loop pitch for the subframe, l_p^{SF}
- Signal modification classifier, C_{SM}
- Fractional pitch lag of previous subframe, $l_{FRC}^{p,prv}$

Output:

- Adaptive-codebook contribution, $c_{unf}^{adp}(n)$
- Filtered adaptive-codebook contribution, $c_f^{adp}(n)$
- Fractional pitch for the subframe, l_{FRC}^p
- Integer pitch for the subframe, l_{INT}^p
- Normalized closed-loop pitch correlation, R_p^{wgt}
- Adaptive-codebook gain (closed-loop pitch gain), g_a

For Rate 1 type-0 frames the pitch is represented by 8 bits for the first and third subframes, and by 5-bit differential coding for the second and the fourth subframes. The closed-loop pitch search interval is $l_p^{SF} \pm 3$ for the first subframe and $l_p^{SF} \pm 3$ for the third subframe. For the second and the fourth subframes the closed-loop pitch search interval is $l_{FRC}^{p,prv} \pm 4.328125$ around the closed-loop pitch of the previous subframe, as defined in the differential-pitch codebook. For the fourth subframe, if $C_{SM} > 0$, the closed-loop pitch search interval is $[l_{FRC}^{p,prv}(2) - 4.328125, \dots, l_{FRC}^{p,prv}(2) + 4.328125] \cap [l_p^{SF} - 1, \dots, l_p^{SF} + 1]$ where $l_{FRC}^p(2)$ denotes the closed-loop pitch of the third subframe. (If the intersection is found to be an empty set, a single point, $l_{FRC}^{p,prv}(2) - 4.328125$ or $l_{FRC}^{p,prv}(2) + 4.328125$ is used.) Note that the resulting search pitch interval is further constrained so that it is always a subset of the allowable pitch interval $[17, \dots, 148]$.

5.6.6.1 Normalized Cross Correlation Between the Target and the Filtered Excitation

The normalized correlations between the filtered adaptive codebook and target signal at the integer values of the search interval are computed. The number of normalized correlation values is the number of integer points in the interval, with extra 5 integer points in each direction for further interpolation. For each lag l in the search interval, the filtered adaptive codebook, $C_{flt,l}^{adp}(n)$ is obtained by convolving the adaptive codebook $C^{adp}(n)$ with $h_{imp}(n)$, using a recursive filtering procedure. The normalized correlation for each lag, $R_{nrm}^{adp}(l)$, is given by:

$$R_{nrm}^{adp}(l) = \frac{\sum_{n=0}^{L_{SF}-1} T_{gs}(n) \cdot C_{flt,l}^{adp}(n)}{\sqrt{\sum_{n=0}^{L_{SF}-1} (C_{flt,l}^{adp}(n))^2}} \quad (5.6.6.1-1)$$

where $L_{SF} = 40$.

5.6.6.2 Fractional Pitch Search and Quantization

The integer closed-loop pitch is obtained as the lag l that maximizes the normalized correlation, $R_{nrm}^{adp}(l)$. For the last subframe, if the maximum of the normalized correlation is negative, and $C_{SM} > 0$, the integer closed-loop pitch is set to the open-loop pitch lag.

The fractional pitch is searched on the search interval of $[-0.9, \dots, 0.9]$ around the integer closed-loop pitch lag. The allowable fractional pitch lag values on that interval are defined by the direct pitch quantization table for the first and the third subframes, and by the differential quantization table for the second and the fourth subframes. The fractional pitch is obtained by a band-limited interpolation of the normalized correlation $R_{nrm}^{adp}(l)$, calculated in Section 5.6.6.1, at the allowable fractional pitch lags on the search interval, and choosing the pitch by the lag that maximizes the interpolated normalized correlation. The adaptive-codebook contribution is constructed in the vector $c_{unf}^{adp}(n)$. l_{INT}^p is simply the rounded version of l_{FRC}^p .

5.6.6.3 Adaptive-codebook Gain (Closed-Loop Pitch Gain) and Normalized Correlation

The adaptive-codebook vector $c_{unf}^{adp}(n)$ is convolved with $h_{imp}(n)$, to generate the filtered adaptive-codebook vector $c_f^{adp}(n)$. The closed-loop pitch gain is given by:

$$g_a = \frac{\sum_{n=0}^{L_{SF}-1} T_{gs}(n) \cdot c_f^{adp}(n)}{\sum_{n=0}^{L_{SF}-1} c_f^{adp}(n) \cdot c_f^{adp}(n)} \quad (5.6.6.3-1)$$

1 where $L_{SF} = 40$. The closed-loop normalized correlation is given by:

$$2 \quad R_p^{wgt} = \frac{\sum_{n=0}^{L_{SF}-1} T_{gs}(n) \cdot c_f^{adp}(n)}{\sqrt{\sum_{n=0}^{L_{SF}-1} c_f^{adp}(n) \cdot c_f^{adp}(n) \cdot \sum_{n=0}^{L_{SF}-1} T_{gs}(n) \cdot T_{gs}(n)}} \quad (5.6.6.3-2)$$

3

4

1 **5.6.7 Adaptive-codebook Contribution For Rate 1/2 Type-0 Frames**

2 Routine name: LTP_close_7b_pitch

3

4 **Input:**

- 5 • Target signal for closed-loop pitch search or pitch interpolation, $T_{gs}(n)$
- 6 • Adaptive codebook, $C^{adp}(n)$
- 7 • Impulse response of the combined filter for the subframe, $h_{imp}(n)$
- 8 • Integer open-loop pitch for the subframe, l_p^{SF}
- 9 • Signal modification classifier, C_{SM}

10 **Output:**

- 11 • Adaptive-codebook contribution, $c_{unf}^{adp}(n)$
- 12 • Filtered adaptive-codebook contribution, $c_f^{adp}(n)$
- 13 • Fractional pitch for the subframe, l_{FRC}^p
- 14 • Integer pitch for the subframe, l_{INT}^p
- 15 • Normalized closed-loop pitch correlation, R_p^{wgt}
- 16 • Adaptive-codebook gain, g_a

17

18 For Rate 1/2 type-0 frames the pitch lag is represented by 7 bits for the first and the
 19 second subframes, which represent the integer pitch lags on the range $[21, \dots, 148]$. The
 20 closed-loop pitch search for Rate 1/2 type-0 is similar to the closed-loop pitch search for
 21 Rate 1 type-0 frames, except only integer pitch values are searched. The closed-loop search
 22 interval is $l_p^{SF} \pm 3$ around the open-loop pitch lags. If the open-loop pitch is smaller than
 23 21, it is multiplied by 2 so that it fits into the allowable pitch range. For the second
 24 subframe, if $C_{SM} > 0$ and if $l_p^{SF}(1)$ is larger than 60, the closed-loop pitch search interval is
 25 reduced to $l_p^{SF}(1) \pm 1$.

26

27 **5.6.7.1 Normalized Cross Correlation Between the Target and the Filtered Excitation**

28 The normalized cross correlation for each lag, $R_{nm}^{adp}(l)$, between the filtered adaptive
 29 codebook and target signal at the integer values of the search interval are computed, as
 30 described for Rate 1 type-0 frames in Section 5.6.6.

1 **5.6.7.2 Integer Pitch Search and Quantization**

2 The integer closed-loop pitch lag l_{INT}^p is obtained as the lag that maximizes the normalized
 3 correlation. For the last subframe, if the maximum normalized correlation is negative and
 4 $C_{SM} > 0$, the integer closed-loop pitch lag is set to the open-loop pitch. The adaptive-
 5 codebook contribution, $c_{unf}^{adp}(n)$, is constructed according to the integer closed-loop pitch
 6 l_{INT}^p . l_{FRC}^p in this case is set to l_{INT}^p .

7

8 **5.6.7.3 Adaptive-codebook Gain (Closed-Loop Pitch Gain) and Normalized Correlation**

9 Similar to Rate 1 type-0 frames, the adaptive-codebook vector $c_{unf}^{adp}(n)$ is convolved with
 10 $h_{imp}(n)$, to generate the filtered adaptive-codebook vector $c_f^{adp}(n)$. The closed-loop pitch
 11 gain, g_a , and the normalized pitch correlation is calculated according to Equations
 12 (5.6.6.3-1) and (5.6.6.3-2) in Section 5.6.6.3, with $L_{SF} = 80$.

13

14

1 **5.6.8 Scaling of Adaptive-Codebook Gain**

2 In order to enhance the perceptual performance of the fixed-codebook search, the adaptive-
 3 codebook gain is scaled down before it is used to generate the target signal for the fixed-
 4 codebook search, $T_g(n)$. In this way, $T_g(n)$ would have a small amount of periodicity,
 5 which encourages placing of pulses around high-energy segments during the pulse-
 6 codebook search (i.e., around the pulse spike in the case of voiced speech).

7 For type-1 frames, if $F_{noisyV} < 2$, the open-loop pitch gain, g_a , is scaled down according to:

$$8 \quad g_a \leftarrow g_a \cdot (0.3 \cdot R_p^{wgt} + 0.7) \quad (5.6.8-1)$$

9 For type-0 frames, the closed-loop pitch gain, g_a , is scaled down according to:

$$10 \quad g_a \leftarrow g_a \cdot (0.5 \cdot R_p^{wgt} + 0.5) \quad (5.6.8-2)$$

11

12

13 **5.6.9 Update of Target and Ideal Excitation**

14 The target for the fixed-codebook search, $T_g(n)$, is generated from the target for the
 15 adaptive-codebook search, $T_{gs}(n)$, the filtered adaptive-codebook contribution, $c_f^{adp}(n)$,
 16 and the closed-loop pitch gain, g_a , according to:

$$17 \quad T_g(n) = T_{gs}(n) - g_a \cdot c_f^{adp}(n) \quad (5.6.9-1)$$

18 The ideal excitation is updated, using the unfiltered adaptive-codebook contribution,
 19 $c_{unf}^{adp}(n)$, and the closed-loop pitch gain, g_a , according to:

$$20 \quad I_{ext}(n) \leftarrow I_{ext}(n) - g_a \cdot c_{unf}^{adp}(n) \quad (5.6.9-2)$$

21

1 5.6.10 Calculation of Fixed-codebook Smoothing Parameters

2 Routine name: SMO_excitation_analysis

3

4 **Input:**

- 5 • Modified speech, $s_{\text{mod}}(n)$
- 6 • Ideal excitation, $I_{\text{ext}}(n)$
- 7 • Fractional pitch for the subframe, l_{FRC}^p
- 8 • Voice activity decision, VAD
- 9 • First element of interpolated reflection coefficient vector for the sub frame, $\hat{k}_m(0)$

10 **Output:**

- 11 • Subframe mode factor, $C_{\text{sub_mod}}$
- 12 • Gain normalization factor, β_E

13

14 The routine generates two factors for the subframe, one used for sub-codebook selection,
 15 and the other for gain normalization. A counter of the consecutive inactive voiced
 16 subframes (according to the frame VAD decision), C_{VAD} , is kept and is limited to 8. An
 17 internal counter of the number of noise updates, $C_{\text{update_noise}}$, is obtained and limited by 20.
 18 An internal counter of the number of speech updates, $C_{\text{update_speech}}$, are obtained and limited
 19 by 60. Using the noise update counter and the speech update counter, the maximum of the
 20 SNR, SNR_{max} , is updated, using an average of the maximum of the absolute values during
 21 noise segments, \bar{N}_{max} and the average of the maximum of the absolute values during
 22 speech segments, \bar{S}_{max} . The maximum of the SNR is used to set a level of adaptation, C_{adpt} .

23 The normalized correlation is calculated on the modified speech for the integer portion of
 24 the pitch lag, $\lfloor l_{\text{FRC}}^p \rfloor$, and for $\lfloor l_{\text{FRC}}^p \rfloor + 1$ are calculated according to:

$$25 \quad R_1 = \frac{\sum_{n=0}^{79} s_{\text{mod}}(n) \cdot s_{\text{mod}}(n - \lfloor l_{\text{FRC}}^p \rfloor)}{\sqrt{\sum_{n=0}^{79} s_{\text{mod}}^2(n) \cdot s_{\text{mod}}(n) \sum_{n=0}^{79} s_{\text{mod}}^2(n - \lfloor l_{\text{FRC}}^p \rfloor)}} \quad (5.6.10-1)$$

$$26 \quad R_2 = \frac{\sum_{n=0}^{79} s_{\text{mod}}(n) \cdot s_{\text{mod}}(n - \lfloor l_{\text{FRC}}^p \rfloor + 1)}{\sqrt{\sum_{n=0}^{79} s_{\text{mod}}^2(n) \cdot s_{\text{mod}}(n) \sum_{n=0}^{79} s_{\text{mod}}^2(n - \lfloor l_{\text{FRC}}^p \rfloor + 1)}} \quad (5.6.10-2)$$

27 (If $\lfloor l_{\text{FRC}}^p \rfloor = 148$, the second correlation is calculated for the lag 147.)

1 The maximum between both correlations is kept, and is denoted R_{L^p} . A moving average of
 2 R_{L^p} is obtained by:

$$3 \quad \bar{R}_{L^p} = 0.9 \cdot \bar{R}_{L^p} + 0.1 \cdot R_{L^p} \quad (5.6.10-3)$$

4 A buffer of maximum of absolute values in the subframe, and a buffer of sum of absolute
 5 values in the frame, for the last 15 subframes, are maintained and undated. The buffer of
 6 maximum absolute values is divided into 5 groups, each of 3 subframes, and the derivative
 7 of the maximum of absolute sum for each group is obtained.

8 Using the maximum of absolute values for the current subframe, the level of adaptation
 9 C_{adpt} , a tracking level of the noise, the average pitch correlation R_{L^p} , and the VAD decision,
 10 a binary speech mode, F_{speech} , is obtained.

11 Update of the noise or the speech parameters is performed, using the parameters described
 12 above for a decision between the update of noise parameters or the update of speech
 13 parameters, and the rate of the updates.

14 The maximum of the absolute value of the ideal excitation subframe is obtained, and
 15 denoted by I_{max} . The ratio between I_{max} and a long term average of it, \bar{I}_{max} , is denoted
 16 R_{max} . The speech mode, F_{speech} , and R_{max} are used to generate the binary excitation mode
 17 factor, C_{ext_mod} , and the gain smoothing factor, β_E .

18 The parameters described above are used for the proper update of the long-term average
 19 \bar{I}_{max} .

20 For each subframe, the frame class C_{FR} is multiplied by C_{ext_mod} to generate the subframe-
 21 modified class, C_{sub_mod} . Further, if $F_{noisyV} = 2$, C_{sub_mod} is also set to 2.

22

23

5.6.11 Principles and Algorithms for Pulse Codebook Search

The fixed-codebook contribution is obtained by a search through several sub-codebooks (with the exception of Rate 1 type-1 frames), and by choosing the best contribution (pulse location and signs, or random vector) and the sub-codebook that contains that best contribution. Three sub-codebooks of pulse structure are used for Rate 1 type-0 frames. A single codebook of pulse structure is used for Rate 1 type-1 frame. Two sub-codebooks of pulse structure and a Gaussian sub-codebook are used for Rate 1/2 type-0 frames. Three sub-codebooks of pulse structure are used for Rate 1/2 type-1 frames.

5.6.11.1 Analysis-by-Synthesis Approach for Fixed-Codebook Contribution

The sub-codebooks are searched for the best fixed-codebook contribution. The best sub-codebook, and the best combination of pulses (or random excitation) in the chosen sub-codebook, are selected for the fixed-codebook contribution. The selection between the sub-codebooks is based on a distance measure that combines the minimization of the mean squared error between the filtered excitation and the target signal (known as closed-loop analysis-by-synthesis approach) and a set of perceptual criteria.

The analysis-by-synthesis approach require finding the fixed-codebook excitation vector, c , which minimizes the mean-square error measure (the time domain index, n , is omitted in the following equations):

$$\varepsilon = \|T_g - gHc\|^2 \quad (5.6.11.1-1)$$

where H is the lower-triangular matrix, where each column includes the impulse response, with the first element on the diagonal. It is well known that by solving Equation (5.6.11.1-1) for the optimal gain and substituting the optimal gain into Equation (5.6.11.1-1), it can be shown that the minimizing ε_i is equivalent to maximizing:

$$\tilde{\varepsilon} = \frac{\left(\left(T_g\right)^T Hc\right)^2}{\|Hc\|^2} = \frac{\left(\left(H^T T_g\right)^T c\right)^2}{c^T H^T Hc} \quad (5.6.11.1-2)$$

If we define $b_{T_g} = \left(H^T T_g\right)^T$ and $y = Hc$:

$$\tilde{\varepsilon} = \frac{\left(b_{T_g} c\right)^2}{y^T y} \quad (5.6.11.1-3)$$

5.6.11.2 Principles of Iterative Search Procedure for Pulse Sub-Codebooks

The fixed-codebook search for the pulse sub-codebooks (or a codebook) is based on an iterative algorithm. The algorithm searches for the best pulse position of each pulse or a pair of pulses, while keeping all the other pulses at their previously determined positions. After the locations for all the pulses in a particular sub-codebook are searched, the procedure can be repeated for further refinement of the pulse locations. The repeated refinement is called “a turn”. During turn 0, the initial locations for the pulses in the sub-codebook are determined.

The fast-search iterative algorithm is based on the linearity of the elements in the numerator and the denominator. Let the vector c^- to include all the other pulses in the sub-codebook, except of the i^{th} pulse, i.e., $c_i = c^- + p_i$, where p_i is a vector that includes

1 a single pulse at the location j_i and zeros elsewhere. Define $y^- = Hc^-$, then $y = y^- + Hp_i$.
 2 Equation (5.6.11.1-3), for the additional i^{th} pulse, can be written as:

$$\begin{aligned} \tilde{\varepsilon}_i &= \frac{(b_{T_g}c)^2}{\|y\|^2} = \frac{[b_{T_g}(c^- + p_i)]^2}{(y^- + Hp_i)^T (y^- + Hp_i)} = \\ &= \frac{[(b_{T_g}c)^- + b_{T_g}p_i]^2}{\|y^-\|^2 + 2(y^-)^T Hp_i + \|Hp_i\|^2} \end{aligned} \quad (5.6.11.2-1)$$

4 The fast search algorithm is based on the pre-calculation of the terms $b_{T_g}p_i$, Hp_i , and
 5 $\|Hp_i\|^2$ for any possible pulse location in the subframe. The evaluation of numerator and
 6 the denominator of $\tilde{\varepsilon}_i$ requires only the calculation of the correlation between y^- and Hp_i ,
 7 and 3 additions.

8

9 **5.6.11.3 Updates of Pulse Contributions**

10 When the best pulse is chosen, its contribution should be added to the overall pulse
 11 contribution, prior to the search for the next pulse to create the terms y_i^{update} and
 12 $(b_{T_g}c)_i^{\text{update}}$. (If all the pulses in the sub-codebook are included, $y_i^{\text{update}} = y$ and
 13 $(b_{T_g}c)_i^{\text{update}} = (b_{T_g}c)$). On the other hand, at each search for a new possible location for a
 14 specific pulse, the contribution of that pulse at its old location should be removed. The
 15 addition or the subtraction of the pulse contribution is done, simply, by:

$$16 \quad y_i^{\text{update}} = y^- + Hp_i \quad (5.6.11.3-1)$$

$$17 \quad y^- = y_i^{\text{update}} - Hp_i \quad (5.6.11.3-2)$$

$$18 \quad (b_{T_g}c)_i^{\text{update}} = (b_{T_g}c)^- + b_{T_g}p_i \quad (5.6.11.3-3)$$

$$19 \quad (b_{T_g}c)^- = (b_{T_g}c)_i^{\text{update}} - b_{T_g}p_i \quad (5.6.11.3-4)$$

20

21 **5.6.11.4 Concept of Pitch Enhancement for Pulse Codebooks**

22 Each pulse in the pulse codebook is called a main pulse. Each main pulse can be repeated
 23 within the subframe, with an appropriate gain factor, on intervals corresponding to the
 24 pitch lag of the subframe. Pulse insertion after the main pulse is called forward pitch
 25 enhancement, and pulse insertion before the main pulse is called backward pitch
 26 enhancement. The pre-calculation of $b_{T_g}p_i$, Hp_i , and $\|Hp_i\|^2$ include both the forward
 27 pitch enhancement and the backward pitch enhancement. Note, in particular, that if only
 28 forward pitch enhancement was used, a single vector, which includes the forward pitch-
 29 enhanced impulse response of a pulse at the beginning of the frame, would be sufficient to
 30 represent Hp_i with a shift of that vector, and the term $(y^-)^T Hp_i$ could be easily calculated.
 31 However, since the backward pitch enhancement includes non-causal contributions, up to

1 4 different vectors are needed to be pre-calculated to represent H_{p_i} , depending on the pitch
 2 lag of the subframe. (The maximal subframe length is 80, but in that case the minimum pitch
 3 lag is 21, hence the maximal number of different vectors is 4.) The appropriate
 4 representative vector is selected according to the index of the main pulse and the pitch lag
 5 of the subframe. The representative vectors are denoted by H_p^m $m=0,\dots,3$ where m
 6 indicates the number of backward pitch enhancement pulses. Note that not all H_p^m need to
 7 be calculated, and the number depends on the subframe size and the pitch lag for the
 8 subframe
 9

10 5.6.11.5 Pre-Calculation Procedure

11 Routine name: FCS_PreCalc_FBCB

13 Input:

- 14 • Target signal, $T_g(n)$
- 15 • Impulse response vector (possibly modified), $\hat{h}_{imp}(n)$
- 16 • Integer pitch lag for the subframe, l_{INT}^p
- 17 • Modified pitch gain, \hat{g}_a (see Section 5.6.13.2 and 5.6.14.2)

18
 19 The pitch enhancement coefficients are calculated by:

$$20 \quad P_{mag}(m) = (\hat{g}_a)^m \quad m = 0, \dots, 4 \quad (5.6.11.5-1)$$

21 The representative vectors for H_{p_i} are calculated as follows. The vector H_p^0 is obtained by:

$$22 \quad H_p^0(j) = \sum_{m=0}^4 P_{mag}(m) \cdot \hat{h}_{imp}(j - m \cdot l_{INT}^p) \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.5-2)$$

23 where $\hat{h}_{imp}(j)$ is a causal impulse response. The vectors H_p^m $m=1,\dots,4$ are obtained
 24 recursively using:

$$25 \quad H_p^m(j) = H_p^{m-1}(j - l_{INT}^p) + P_{mag}(m) \cdot \hat{h}_{imp}(j) \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.5-3)$$

26 Note that H_p^m $m=1,\dots,4$ are not evaluated if $m \cdot l_{INT}^p \geq L_{SF}$.

27 An auxiliary function, targeted for computation reduction during the search, is calculated
 28 as:

$$29 \quad P_m(j) = m \quad \text{for } m \cdot l_{INT}^p \leq j < (m+1) \cdot l_{INT}^p \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.5-4)$$

30 The weighted filtered target signal is obtained by:

$$31 \quad b_{Tg}(j) = \sum_{n=j}^{L_{SF}-1} T_g(n) \cdot \hat{h}_{imp}(n-j) \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.5-5)$$

1 The pitch enhanced $b_{Tg} p_i$ is calculated for each index in the subframe (p_i such that
2 $j_i = i$) by:

$$3 \quad b_{Tg} p_i = \sum_{m=-5}^5 P_{mag}(m) \cdot b_{Tg}(i + m \cdot l_{INT}^p), \quad i = 0, \dots, L_{SF} - 1 \quad ; 0 \leq i + m \cdot l_{INT}^p < L_{SF} .$$

4 (5.6.11.5-6)

5 The signs of $b_{Tg} p_i$ are saved in a sign array for each i in the subframe according to:

$$6 \quad \begin{aligned} S_i &= 1 & b_{Tg} p_i &> 0 \\ S_i &= -1 & b_{Tg} p_i &\leq 0 \end{aligned}$$

(5.6.11.5-7)

7 The pitch enhanced $b_{Tg} p_i$ is then modified to its absolute value for each i in the subframe.

8 The energy of the pitch-enhanced and filtered pulse contribution, $\|Hp_i\|^2$, is obtained for
9 each index in the subframe (p_i such that $j_i = i$) by:

$$10 \quad \|Hp_i\|^2 = \sum_{l=0}^{L_{SF}-i-1+P_m(i) \cdot l_{INT}^p} \left(H_p^{P_m(i)}(l) \right)^2 \quad i = 0, \dots, L_{SF} - 1$$

(5.6.11.5-8)

11 using a recursive calculation procedure.

12 An alternate method of computing the filtered codebook energy factor, $\|y\|^2$ in equation
13 5.6.11.2-1, uses the correlation matrix of the impulse response, Φ , whose elements are
14 defined by

$$15 \quad \phi(i, j) = \sum_{l=MAX((i-P_m(i) \cdot l_{INT}^p), (j-P_m(j) \cdot l_{INT}^p))}^{L_{SF}-1} H_p^{P_m(i)}(l-i) H_p^{P_m(j)}(l-j) \quad ; \quad 0 \leq i < L_{SF} \text{ and } 0 \leq j < L_{SF}$$

16 (5.6.11.5-9)

17

18

19 **5.6.11.6 An Iterative Algorithm Based on Position Search for a Single Pulse**

20 Routine name: FCS_Fast_Search_FBCB

21

22 The iterative algorithm based on position search for a single pulse achieves an initial
23 positioning or re-evaluation of all the pulse locations in the sub-codebook by searching for
24 the best location of a single pulse at a time. The algorithm can perform any number of
25 turns for further refinements of all pulse locations. At each turn, the best location for each
26 pulse in the sub-codebook is searched while keeping the previous locations for all the other
27 pulses in the sub-codebook. At the first turn the pulses are added sequentially, where the
28 first pulse is positioned independently of the other pulses, the second pulse is added to the
29 sub-codebook vector while retaining the first pulse, and so on, until the last pulse in the
30 sub-codebook is positioned.

31

32 5.6.11.6.1 Initialization of the Search Parameters

1 If the iterative algorithm is called for an initial positioning of the first pulse in the sub-
 2 codebook (first turn), the search parameters, $(b_{Tg}c)^{update}$ and y^{update} , are initialized to:

$$3 \quad (b_{Tg}c)^{update} = 0 \quad (5.6.11.6-1)$$

$$4 \quad y^{update}(j) = 0 \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.6-2)$$

5 If the iterative algorithm is called for a refinement of previously set of positions, the search
 6 parameters are initialized to:

$$7 \quad (b_{Tg}c)^{update} = \sum_{P_i} b_{Tg} p_i \quad (5.6.11.6-3)$$

$$8 \quad y^{update}(j) = \sum_{P_i} S_i \cdot H_p^{P_m(i)}(j - i + P_m(i) \cdot l_{INT}^p) \quad j = 0, \dots, L_{SF} - 1 \quad (5.6.11.6-4)$$

9 where i runs over all the previously chosen pulses in the sub-codebook.

10

11 5.6.11.6.2 The Search Procedure

12 The search procedure is structured of three nested loops. The external loop progresses the
 13 turns for further refinement of an already established set of pulse locations. The
 14 intermediate loop progresses sequentially from the first pulse to the last pulse in the sub-
 15 codebook, and refines the location of each individual pulse (tested pulse). Before the
 16 internal loop is entered, the contribution of the tested pulse in its previous location is
 17 removed from the search parameters as described in Section 5.6.11.3 by Equations
 18 (5.6.11.3-2) and (5.6.11.3-4).

19 In the internal search loop, $\tilde{\varepsilon}_i$ is evaluated for all the possible track locations for the tested
 20 pulse. The error measure $\tilde{\varepsilon}_i$ is evaluated using the procedure described in Section 5.6.11.2
 21 by Equation (5.6.11.2-1). The location i^* that maximizes $\tilde{\varepsilon}_i$ is selected.

22 At the end of inner loop, the contribution of the tested pulse at the selected position is
 23 added to the search parameters as described in Section 5.6.11.3 by Equations (5.6.11.3-1)
 24 and (5.6.11.3-3).

25 Alternatively, the filtered codebook energy factor, $\|y\|^2$ in Equation 5.6.11.2-1 can be
 26 computed as

$$27 \quad \|y\|^2 = \sum_{k=0}^{N_p-1} \phi(p_k, p_k) + 2 \sum_{k=0}^{N_p-2} \sum_{l=k+1}^{N_p-1} S_k S_l \phi(p_k, p_l) \quad (5.6.11.6-5)$$

28 where N_p is the number of pulses in the code-vector. It can be noted that the above
 29 equation is a general equation and the actual computation can be efficiently performed
 30 using a set of partial summations as the pulses are searched iteratively.

31

32 5.6.11.6.3 Building the Pulse Excitation

33 After the final turn of sub-codebook refinement has been completed, the pulse excitation,
 34 including the forward and the backward pitch enhancement pulse, is constructed,

1 according to the selected pulse locations, the subframe pitch lag l_{INT}^p , and the pitch
 2 enhancement coefficients P_{mag} .

3

4 **5.6.11.7 An Iterative Algorithm Based on Joint Position Search for Two Pulses**

5 Routine name: FCS_Fast_Search_FBCB2

6 The iterative algorithm that is based on a joint position search for two pulses achieves an
 7 initial positioning or re-evaluation of all the pulse locations in the sub-codebook by jointly
 8 searching for the best locations of two pulses at a time. The algorithm can perform any
 9 number of turns for further refinements of all pulse locations. At each turn, the best
 10 locations for two pulses in the sub-codebook are searched while keeping the previous
 11 locations for all the other pulses in the sub-codebook. At the first turn pairs of pulses are
 12 added to the fixed-codebook contribution sequentially, where the first pair of pulses is
 13 positioned independently of the other pulses, the second pair is added, and so on, until the
 14 last pair in the sub-codebook is positioned.

15 Note that each turn in the joint search for two pulses is of higher computational complexity
 16 than is the search for a single pulse, but pre-search procedure can be used to limit the
 17 search complexity while keeping the advantages of the joint search.

18

19 5.6.11.7.1 Initialization of the Search Parameters

20 The initialization of the search parameters, for the case of the first turn or any sequential
 21 turn is performed exactly as described in Section 5.6.11.6.1 above.

22

23 5.6.11.7.2 Pre-Search for 2-Pulse Sub-Codebook

24 To reduce the search complexity in the case of 2-pulse sub-codebook (used for Rate 1/2
 25 frames, see Section 5.6.14.1), a low-complexity pre-search procedure is performed to
 26 identify a set of candidates for further search. The number of pre-search candidates
 27 depends on the frame type and the subframe pitch lag, l_{INT}^p , according to:

$$28 \quad N_{pre} = \begin{cases} 16 & \text{type-1 frame, } l_{INT}^p < 30 \\ 19 & \text{type-1 frame, } l_{INT}^p \geq 30 \\ 18 & \text{type-0 frame, } l_{INT}^p < 45 \\ 22 & \text{type-0 frame, } l_{INT}^p \geq 45 \end{cases} \quad (5.6.11.7-1)$$

29 The pre-search identifies the N_{pre} pulse locations in each of the two tracks that maximize
 30 the error criterion:

$$31 \quad \hat{\varepsilon}_i = \frac{(b_{Tg} p_i)^2}{\|Hp_i\|^2} \quad (5.6.11.7-2)$$

32

33 5.6.11.7.3 Pre-Search for 3-Pulse Sub-Codebook

1 The pre-search for the 3-pulse sub-codebook positions only the main pulse by maximizing
 2 the same error criterion of Equation (5.6.11.7-2) above. After the main pulse is positioned,
 3 the tracks for the second and the third pulse are dynamically allocated. See Section
 4 5.6.14.1 for the details of the 3-pulse dynamic tracks. Note that the main pulse location is
 5 not modified any further.

6

7 5.6.11.7.4 The Search Procedure

8 The search procedure is structured of four nested loops. The external loop progresses the
 9 turns for further refinement of an already established set of pulse locations. The second
 10 loop progresses on an arranged set of two pulses in the sub-codebook at a time, and refines
 11 the locations of the two pulses in the set (tested pulse pair). Before the internal loop is
 12 entered, for turn 1 or higher, the contributions of the tested pulse pair in their previous
 13 locations are removed from the search parameters in Section 5.6.11.3 by Equations
 14 (5.6.11.3-2) and (5.6.11.3-4).

15 The two internal search loops are performed on either all the possible track locations for
 16 the tested pulse pair, or on the pre-selected possible track locations for the tested pulse
 17 pair. The order of the two loops is structured such that the internal loop is performed for
 18 the pulse with the larger number of possible track locations. At the outer loop for the first
 19 pulse, which runs on all possible locations for this pulse, indexed by i_1 , the search
 20 parameters are temporarily updated with the update procedure described in Section
 21 5.6.11.2 by Equation (5.6.11.2-1), using $b_{Tg}p_{i_1}$, Hp_{i_1} , and $\|Hp_{i_1}\|^2$. In the inner loop for the
 22 second pulse, indexed i_2 , the search parameters are updated using $b_{Tg}p_{i_2}$, Hp_{i_2} , and
 23 $\|Hp_{i_2}\|^2$. The error measure $\tilde{\varepsilon}_{i_1, i_2}$ is evaluated using the temporarily updated search
 24 parameters, and the locations i_1^* and i_2^* that maximize $\tilde{\varepsilon}_{i_1, i_2}$ are selected.

25 After the inner two loops are completed, in order to prepare the search parameters for the
 26 next pair search, the contribution of the tested pulse pair at their newly selected positions
 27 is added to the previous search parameters as described in Section 5.6.11.3 by Equations
 28 (5.6.11.3-1) and (5.6.11.3-3).

29 Alternatively, the filtered codebook energy factor, $\|y\|^2$ in Equation 5.6.11.2-1 can be
 30 computed as

$$31 \quad \|y\|^2 = \sum_{k=0}^{N_p-1} \phi(p_k, p_k) + 2 \sum_{k=0}^{N_p-2} \sum_{l=k+1}^{N_p-1} S_k S_l \phi(p_k, p_l) \quad (5.6.11.7-3)$$

32 where N_p is the number of pulses in the code-vector. It can be noted that the above
 33 equation is a general equation and the actual computation can be efficiently performed
 34 using a set of partial summations as the pulses are searched iteratively.

35

36 5.6.11.7.5 Building the Pulse Excitation

37 After the final turn of sub-codebook refinement has been completed the pulse excitation,
 38 including the forward and the backward pitch enhancement pulse, is constructed,
 39 according to the selected pulse locations, the subframe integer pitch lag l_{INT}^p , and the pitch
 40 enhancement coefficients P_{mag} .

1

2 **5.6.11.8 Joint Position Search for Five Pulses**3 Routine name: FCS_Full_Search_FBCB_5p

4

5 For Rate 1/2 type-1 frames, a special 5-pulse sub-codebook, described in Table 5.6-8, is
6 used. The locations for the 5 pulses are searched with a joint positioning search. The joint
7 positioning search includes 5 nested loops, which test all possible locations for the 5
8 pulses. For each pulse, indexed by i , the contribution to search parameters is added and
9 temporarily with the update procedure described in Section 5.6.11.2 by Equation
10 (5.6.11.2-1), using $b_{Tg}p_i$, Hp_i , and $\|Hp_i\|^2$. The error measure $\tilde{\mathcal{E}}_{i_1, i_2, i_3, i_4, i_5}$ is evaluated using
11 the updated search parameters, and the locations i_1^* , i_2^* , i_3^* , i_4^* and i_5^* that maximize
12 $\tilde{\mathcal{E}}_{i_1, i_2, i_3, i_4, i_5}$ are selected.

13

14

5.6.12 Gaussian Codebook Search

Routine name: GCB_gauss_cb

The Gaussian codebook is searched in a closed-loop manner in order to minimize the perceptual mean-squared error between the original and the reconstructed speech. The Gaussian excitation vector is constructed from two orthogonal basis vectors, where the first basis vector contains all of the even sample points in the sub-frame, and the second basis vector contains all of the odd sample points; hence, each basis vector is of dimension 40, and the same pseudo Gaussian codebook is used for the two basis vectors. For each basis vector, 54 candidates are considered, and 1 bit is used to specify the sign. This results in 11664 entries, which are specified by 13 bits plus the unused index space of the 2-pulse codebook, i.e. 3472 entries.

In order to reduce the complexity two candidates are pre-selected in an open-loop approach for each basis vector, by maximizing the cross correlation functions $|R_1^G(l)|$ and $|R_2^G(l)|$, respectively:

$$R_1^G(l) = \sum_{n=0}^{39} I_{exc}(2 \cdot n) \cdot c_{Ga,unf}^l(n) \quad l = 0, \dots, 53 \quad (5.6.12-1)$$

$$R_2^G(l) = \sum_{n=0}^{39} I_{exc}(2 \cdot n + 1) \cdot c_{Ga,unf}^l(n) \quad l = 0, \dots, 53 \quad (5.6.12-2)$$

where $I_{exc}(n)$ is the ideal excitation after removal of the contribution from the adaptive codebook, and $c_{Ga,unf}^l(n)$ $l = 0, \dots, 53$ is the candidate excitation vectors from the Gaussian codebook. The same codebook is used for both basis vectors. The signs of the candidate vectors $sign(R_1^G(l))$ and $sign(R_2^G(l))$ are determined as the sign of the correlation functions $R_1^G(l)$ and $R_2^G(l)$ for the respective candidate vectors. The final candidate is determined among the four possible combinations of the pre-selected candidates for the two basis vectors by maximizing

$$R^G(l_1, l_2) = \frac{\left(\sum_{j=0}^{L_{SF}-1} T_g(n) \cdot c_{Ga,f}^{l_1, l_2}(n) \right)^2}{\sum_{j=0}^{L_{SF}-1} \left(c_{Ga,f}^{l_1, l_2}(n) \right)^2} \quad l_1 = L_1(0), L_1(1) \quad l_2 = L_2(0), L_2(1) \quad (5.6.12-3)$$

where $L_1(0), L_1(1)$ and $L_2(0), L_2(1)$ specifies the candidate vectors for the two basis vectors and $c_{Ga,f}^{l_1, l_2}(n)$ is the result of the convolution between $h_{imp}(n)$ and the Gaussian code vector $c_{Ga,unf}^{l_1, l_2}(n)$ reconstructed according to:

$$c_{Ga,unf}^{l_1, l_2}(n) = \begin{cases} sign(R_1^G(l_1)) \cdot c_{Ga,unf}^{l_1}\left(\frac{n}{2}\right) & n \text{ even} \\ sign(R_2^G(l_2)) \cdot c_{Ga,unf}^{l_2}\left(\frac{n-1}{2}\right) & n \text{ odd} \end{cases} \quad n = 0, \dots, L_{SF} - 1 \quad (5.6.12-4).$$

5.6.13 Fixed-codebook Contribution for Rate 1

5.6.13.1 Fixed-codebook Structure for Rate 1

The fixed-codebook contribution for Rate 1 type-0 frames comprises 3 sub-codebooks, each of 5 pulses. The 3 sub-codebooks differ in the tracks for the 5 pulses, and only one of the 3 sub-codebooks is selected and used. The fixed-codebook contribution for Rate 1 type-1 frames is a single codebook of 8 pulses.

The following tables give the tracks for the three 5-pulse sub-codebooks for type-0 frames, and the tracks for the single 8-pulse codebook for type-1 frames.

	Track locations in a 40-sample subframe (type-0 frames)	# of bits
Pulse 1	1, 3, 6, 8, 11, 13, 16, 18, 21, 23, 26, 28, 31, 33, 36, 38	4
Pulse 2	4, 9, 14, 19, 24, 29, 34, 39	3
Pulse 3	1, 3, 6, 8, 11, 13, 16, 18, 21, 23, 26, 28, 31, 33, 36, 38	4
Pulse 4	4, 9, 14, 19, 24, 29, 34, 39	3
Pulse 5	0, 2, 5, 7, 10, 12, 15, 17, 20, 22, 25, 27, 30, 32, 35, 37	4

Table 5.6-2: Tracks for the 1st 5-Pulse Sub-Codebook for Rate 1 Type-0 Frames

Since Pulse 1 and Pulse 3 share the same track, the sign information for both tracks can be represented with a single bit (using the order of the pulse locations). Since Pulse 2 and Pulse 4 also share the same track, their sign information can also be represented with a single bit. The sign of Pulse 5 is represented by a single bit. The number of bits for the representation of the 1st 5-pulse sub-codebook for Rate 1 type-0 frames is 21.

	Track locations in a 40-sample subframe (type-0 frames)	# of bits
Pulse 1	0, 1, 2, 3, 4, 6, 8, 10	3
Pulse 2	5, 9, 13, 16, 19, 22, 25, 27	3
Pulse 3	7, 11, 15, 18, 21, 24, 28, 32	3
Pulse 4	12, 14, 17, 20, 23, 26, 30, 34	3
Pulse 5	29, 31, 33, 35, 36, 37, 38, 39	3

Table 5.6-3: Tracks for the 2nd 5-Pulse Sub-Codebook for Rate 1 Type-0 Frames

Since each pulse has an individual sign, the number of bits used to represent the 2nd 5-pulse sub-codebook for Rate 1 type-0 frames is 20.

	Track locations in a 40-sample subframe (type-0 frames)	# of bits
Pulse 1	0, 1, 2, 3, 4, 5, 6, 7	3
Pulse 2	8, 9, 10, 11, 12, 13, 14, 15	3
Pulse 3	16, 17, 18, 19, 20, 21, 22, 23	3
Pulse 4	24, 25, 26, 27, 28, 29, 30, 31	3

Pulse 5	32, 33, 34, 35, 36, 37, 38, 39	3
---------	--------------------------------	---

1

2

Table 5.6-4: Tracks for the 3rd 5-Pulse Sub-Codebook for Rate 1 Type-0 Frames

3

Since each pulse has an individual sign, the number of bits used to represent the 2nd 5-pulse sub-codebook for Rate 1 type-0 frames is 20.

4

5

	Track locations in a 40-sample subframe (type-1 frames)	# of bits
Pulse 1	0, 2, 5, 7, 10, 12, 15, 17, 20, 22, 25, 27, 30,32, 35, 37	4
Pulse 2	1, 6, 11, 16, 21, 26, 31, 36	3
Pulse 3	3, 8, 13, 18, 23, 28, 33, 38	3
Pulse 4	4, 9, 14, 19, 24, 29, 34, 39	3
Pulse 5	0, 2, 5, 7, 10, 12, 15, 17, 20, 22, 25, 27, 30, 32, 35, 37	4
Pulse 6	1, 6, 11, 16, 21, 26, 31, 36	3
Pulse 7	3, 8, 13, 18, 23, 28, 33, 38	3
Pulse 8	4, 9, 14, 19, 24, 29, 34, 39	3

6

7

Table 5.6-5: Tracks for the 8-Pulse Codebook for Rate 1 Type-1 Frames

8

Since the each pulse shares a track with another pulse, the sign information for both pulses can be represented with a single bit, and the number of bits used to represent the 8-pulse codebook for Rate 1 type-1 frames is 30 bits.

9

10

11

5.6.13.2 Fixed-codebook Search for Rate 1

12

Routine name: FCS_cdbk_search_8500BSP

13

14

Input:

15

16

- Target signal, $T_g(n)$

17

- Impulse response, $h_{imp}(n)$

18

- Ideal excitation, $I_{ext}(n)$

19

- Integer pitch lag, l_{INT}^p

20

- Pitch gain, g_a (for type-1 frames the pitch gain is quantized pitch gain for the subframe, for type-0 frames the pitch gain is the quantized pitch gain of the previous subframe)

21

22

23

- Subframe-modified class, C_{sub_mod}

24

- Noise-to-signal ratio, NSR

25

- Normalized pitch correlation, R_p^{wgt}

1 **Output:**

- 2 • Unfiltered fixed-codebook excitation, $c_{unf}^{fix}(n)$

3

4 For Rate 1, the impulse response is not modified, and $\hat{h}_{imp}(n)$ used in Section 5.6.11.5 is
 5 identical to $h_{imp}(n)$, as defined in Section 5.6.2.

6

7 5.6.13.2.1 Pitch Enhancement Gain and Pre-Calculation of Search Parameters

8 The pitch gain used for pitch enhancement is obtained as:

$$9 \quad \hat{g}_a = \begin{cases} 0.75 \cdot \min\{g_a, 1\} & l_{INT}^p < 20 \text{ and } g_a > 0.5 \\ g_a & \text{otherwise} \end{cases} \quad (5.6.13.2-1)$$

10 \hat{g}_a is further modified according to the frame type according to:

$$11 \quad \hat{g}_a \leftarrow \begin{cases} \min\{\hat{g}_a, 0.8\} & \text{type-0 frame} \\ \min\{\hat{g}_a, 1.0\} & \text{type-1 frame} \end{cases} \quad (5.6.13.2-2)$$

12 .The search parameters are pre-calculated, as described in Section 5.6.11.5.

13

14 5.6.13.2.2 Fixed-codebook Search for Type-0 Frames

15 The two 5-pulse sub-codebooks, described in Table 5.6-2 and Table 5.6-3, are searched,
 16 each with a single search turn, using the iterative search algorithm based on a position
 17 search for a single pulse at a time, described in Section 5.6.11.6. The sub-codebook that
 18 results in the maximum error criterion (minimum error) $\tilde{\epsilon}$ is selected, with its combination
 19 of 5 pulses.

20 The criterion is then weighted, based on the sharpness of the ideal excitation, the signal-to-
 21 noise ratio, and the subframe-modified class.

22 The third 5-pulse sub-codebook described in Table 5.6-4 is searched, using the iterative
 23 search algorithm based on a position search for a single pulse at a time, described in
 24 Section 5.6.11.6, and its error criterion is compared to the weighted criterion from the
 25 previous step. If the criterion from the third sub-codebook is larger than the weighted
 26 criterion, the third sub-codebook is chosen; otherwise, the best of the two first sub-code-
 27 books is chosen.

28 The pulse locations in the selected sub-codebook are further refined with 3 more turns,
 29 using the iterative search algorithm based a on joint position search for two pulses,
 30 described in Section 5.6.11.7. If the 1st 5-pulse sub-codebook, described in Table 5.6-2, is
 31 chosen, the redundancy in the representation of two pulses with two signs on identical
 32 tracks is removed between track 1 and track 3 and between track 2 and track 4, by
 33 exchanging the pulse locations according to the individual pulse signs, after the final
 34 refinement turn.

35 The fixed-codebook contribution, $c_{unf}^{fix}(n)$, is reconstructed by placing the pulses at the
 36 selected positions with the selected signs, and the addition of pitch enhancements pulses.

37

1 5.6.13.2.3 Fixed-codebook Search for Type-1 Frames

2 The 8-pulse codebook is first searched with 2 turns of the iterative search algorithm that is
3 based on a joint position search for two pulses described in Section 5.6.11.7. The pulse
4 locations are further refined with 3 more turns of the iterative search algorithm based on
5 position search for a single pulse, described in Section 5.6.11.6. The redundancy in the
6 representation of two pulses with two signs on identical tracks is removed between track 1
7 and track 5, track 2 and track 5, track 3 and track 6, and track 4 and track 8, by
8 exchanging the pulse locations according to the individual pulse signs, after the final
9 refinement turn.

10

11

12

13

14

5.6.14 Fixed-codebook Contribution for Rate 1/2

5.6.14.1 Fixed-codebook Structure for Rate 1/2

The fixed-codebook contribution for Rate 1/2 type-0 frames comprises 3 sub-codebooks, one of 2 pulses, the second of 3 pulses, and the third of random excitation. The fixed-codebook contribution for Rate 1/2 type-1 frames also comprises 3 sub-codebooks; the first of 2 pulses, the second is of 3 pulses, and the third of 5 pulses.

The following tables give the tracks for the 2-pulse sub-codebook for type-0 frames, and the tracks for the 2-pulse and the 5-pulse sub-codebook for type-1 frames.

	Track locations in an 80-sample subframe (type-0 frames)	# of bits
Pulse 1	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79	6.5
Pulse 2	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79	6.5

Table 5.6-6: Tracks for the 2-Pulse Sub-Codebook for Rate 1/2 Type-0 Frames

Since $\log_2(80) = 6.322... < 6.5$, the locations of both pulses are combined and coded with $2 \cdot 6.5 = 13$ bits. Moreover, since Pulse 1 and Pulse 2 share the same track, the sign information for both tracks can be represented with a single bit. The number of bits for the representation of the 2-pulse sub-codebook for Rate 1/2 type-0 frames is 14.

	Track locations in a 53/54-sample subframe (type-1 frames)	# of bits
Pulse 1	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30, 32, 34, 36, 38, 40, 42, 44, 46, 48, 50, 52	5
Pulse 2	1, 3, 5, 7, 9, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 25, 27, 29, 31, 33, 35, 37, 39, 41, 43, 45, 47, 49, 51	5

Table 5.6-7: Tracks for the 2-Pulse Sub-Codebook for Rate 1/2 type-1 frames

The signs for each pulse are represented individually, and the total number of bits for the representation of the 2-pulse sub-codebook for Rate 1/2 type-1 frames is 12.

	Track locations in a 53/54-sample subframe (type-1 frames)	# of bits
Pulse 1	0, 15, 30, 45	2

Pulse 2	0, 5	1
Pulse 3	10, 20	1
Pulse 4	25, 35	1
Pulse 5	40, 50	1

1

2 **Table 5.6-8: Tracks for the 5-pulse Sub-Codebook for Rate 1/2 Type-1 Frames**

3 The signs for each pulse are represented individually, and the total number of bits for the
4 representation of the 5-pulse sub-codebook for Rate 1/2 type-1 frames is 11.

5 The 13-bit 3-pulse sub-codebook for type-0 frames and the 11-bit 3-pulse sub-codebook for
6 type-1 frames are constructed dynamically, where the tracks for the second and the third
7 pulse depend on the location of the first pulse. The following tables give the fixed tracks for
8 the first pulse for both sub-codebooks.

9

	Track locations in a 80-sample subframe (type-0 frames)	# of bits
Pulse 1	0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75	4

10

11 **Table 5.6-9: Track for the First Pulse in the 3-Pulse Sub-Codebook for Rate 1/2 type-0**
12 **Frames**

13

	Track locations in a 53/54-sample subframe (type-1 frames)	# of bits
Pulse 1	3, 6, 9, 12, 15, 18, 21, 24, 27, 30, 33, 36, 39, 42, 45, 48	4

14

15 **Table 5.6-10: Track for the First Pulse in the 3-Pulse Sub-Codebook for Rate 1/2**
16 **Type-1 Frames**

17 The tracks for the 2nd pulse and the 3rd pulse are constructed in a small interval around the
18 location of the first pulse. The length of the interval is 16 samples for type-0 frames and 8
19 samples for type-1 frames. The interval is positioned with the first pulse in the middle, but
20 if the first pulse is closer to the subframe boundaries, the interval is adjusted to be fully
21 included in the subframe. The 2nd pulse and the 3rd pulse are located on alternating even
22 and odd locations to the left of the first pulse with reverse odd and even locations to the
23 right of the first pulse.

24 Since 3 bits are used for the 2nd pulse and the 3rd pulse for type-0 frames, and since the
25 signs of each pulse are represented individually, the total number of bits for the
26 representation of the 3-pulse sub-codebook for Rate 1/2 type-0 frames is 13.

27 Since 2 bits are used for the 2nd pulse and the 3rd pulse for type-1 frames, and since the
28 signs of each pulse are represented individually, the total number of bits for the
29 representation of the 3-pulse sub-codebook for Rate 1/2 type-1 frames is 11.

30

31 **5.6.14.2 Fixed-codebook Search for Rate 1/2**

32

33 Routine name: FCS_cdbk_search_4000BPS

34

1 **Input:**

- 2 • Adaptive codebook, $C^{adp}(n)$
- 3 • Quantized and interpolated prediction coefficients, $\hat{a}_m^q(n)$
- 4 • Target signal, $T_g(n)$
- 5 • Impulse response, $h_{imp}(n)$
- 6 • Ideal excitation, $I_{ext}(n)$
- 7 • Integer pitch lag, l_{INT}^p
- 8 • Pitch gain, g_a (for type-1 frames the pitch gain is quantized pitch gain for the
- 9 subframe, for type-0 frames the pitch gain is the quantized pitch gain of the
- 10 previous subframe)
- 11 • Subframe-modified class, C_{sub_mod}
- 12 • Noise-to-signal ratio, NSR
- 13 • Normalized pitch correlation, R_p^{wgt}
- 14 • Previous frame rate, $Rate^{prv}$
- 15 • LPC prediction gain from quantized parameters of the first subframe, $G_{0,q}^{LPC}$

16 **Output:**

- 17 • Unfiltered fixed-codebook excitation, $c_{unf}^{fix}(n)$

18

19 The fixed-codebook search for Rate 1/2 incorporates three codebook enhancements, in
 20 addition to the pitch enhancement described in Section 5.6.11.4. The codebook
 21 enhancement is incorporated into the impulse response, by modifying $h_{imp}(n)$ (Section
 22 5.6.2) to form $\hat{h}_{imp}(n)$, used in Section 5.6.11.

23

24 5.6.14.2.1 Parameters of Formant Correction Filter and Stability Index

25 A weighing parameter, α^{corr} , is set according to the LPC prediction gain from quantized
 26 parameters of the first subframe, $G_{0,q}^{LPC}$ (Section 5.3.25)

$$27 \quad \alpha^{corr} = \begin{cases} 0.0125 & G_{0,q}^{LPC} < 10 \\ 0.025 & 10 \leq G_{0,q}^{LPC} < 20 \\ 0.05 & 20 \leq G_{0,q}^{LPC} \end{cases} \quad (5.6.14.2-1)$$

28 The parameters for an all-pole filter are obtained (for the sub-frame indexed by m) from
 29 the quantized and interpolated prediction coefficients $\hat{a}_m^q(j)$ by:

$$1 \quad a^{corr}(j) = \hat{a}_m^q(j) \cdot (\alpha^{corr})^j \quad j = 1, \dots, 10 \quad (5.6.14.2-2)$$

2 At the first subframe, a stability index for the entire frame is calculated by:

$$3 \quad S_t = 1 - \frac{|G_{0,q}^{LPC} - G_{0,q}^{LPC,prv}|}{|G_{0,q}^{LPC} + G_{0,q}^{LPC,prv}|} \quad (5.6.14.2-3)$$

4 where $G_{0,q}^{LPC,prv}$ is $G_{0,q}^{LPC}$ of the previous frame. If the coding rate is not Rate 1/2, S_t is set to
5 0.

6

7 5.6.14.2.2 Pulse Addition Based on Adaptive-codebook Correlations

8 The weighted quantized prediction coefficients $\tilde{a}_m^q(j)$ (for the sub-frame indexed by m) are
9 obtained from the quantized and interpolated prediction coefficients $\hat{a}_m^q(j)$ according to:

$$10 \quad \tilde{a}_m^q(j) = (0.85)^j \cdot \hat{a}_m^q(j) \quad j = 1, \dots, 10 \quad (5.6.14.2-4)$$

11 The 170 samples of the previously quantized and reconstructed excitation are filtered by
12 the filter:

$$13 \quad \tilde{W}(z) = \frac{1}{\tilde{A}^q(z)} \quad (5.6.14.2-5)$$

14 where $\tilde{a}_m^q(j)$ are the coefficients of $\tilde{A}^q(z)$, to generate the weighted quantized excitation
15 $c_w^q(n)$. Note the previously quantized and reconstruction excitation for type-1 frames is
16 comprised only of samples from the previous frames, while the previously quantized and
17 reconstructed excitation for type-0 is comprised of samples from previous frames, as well as
18 samples from previous subframes of the current frame.

19 A set of correlation lags on the previously quantized and reconstructed excitation is
20 constructed. If $17 < \frac{l_{INT}^p}{3} < L_{SF}$, the lags:

$$21 \quad L^C = \left\lfloor \frac{j \cdot l_{INT}^p}{3} + 0.5 \right\rfloor \quad j = 1, \dots, 5 \quad (5.6.14.2-6)$$

22 are included in the set of correlation lags. Similarly, if $17 < \frac{l_{INT}^p}{4} < L_{SF}$, the lags:

$$23 \quad L^C = \left\lfloor \frac{j \cdot l_{INT}^p}{4} + 0.5 \right\rfloor \quad j = 1, \dots, 5 \quad (5.6.14.2-7)$$

24 are included in the set of correlation lags. The set is then limited to include lags that are
25 smaller than $l_{INT}^p - 6$ and smaller than L_{SF} . The number of correlation lags is denoted N_w^C ,
26 and the correlation lags by L_m^C , $m = 1, \dots, N_w^C$. The normalized correlation on the weighted
27 previously reconstructed excitation is calculated by:

$$R_w^C(m) = \frac{\sum_{n=0}^{79} c_w^q(90+n) \cdot c_w^q(90+n-L_m^C)}{\frac{1}{2} \left(\sum_{n=0}^{79} c_w^q(90+n) \cdot c_w^q(90+n) + \sum_{n=0}^{79} c_w^q(90+n-L_m^C) \cdot c_w^q(90+n-L_m^C) \right)}$$

1
2
3

(5.6.14.2-8)

4 The normalized correlations are further limited to the range $0 \leq R_w^C(m) \leq 1$.

5 If the current frame is type-1 and the previous frame was also Rate 1/2 type-1 frame, the correlations are normalized by:

$$R_w^C(m) \Leftarrow 0.5 \cdot R_w^C(m) \cdot S_t \cdot \min(\hat{g}_p, 1) \quad m = 1, \dots, N_w^C$$

6
7

(5.6.14.2-9)

8 Otherwise, the correlation is normalized by

$$R_w^C(m) \Leftarrow 0.25 \cdot R_w^C(m) \cdot S_t \cdot \min(\hat{g}_p, 0.75) \quad m = 1, \dots, N_w^C$$

9

(5.6.14.2-10)

10 5.6.14.2.3 Generating the Modified Impulse Response

11 The modified impulse response $\hat{h}_{imp}(n)$ is generated from the impulse response $h_{imp}(n)$.
12 The correlation enhancement is performed by

$$h'_{imp}(n) = h_{imp}(n) + \sum_{m=1}^{N_w^C} R_w^C(m) \cdot h_{imp}(n - L_m^C) \quad n = 0, \dots, L_{SF} - 1$$

13

(5.6.14.2-11)

14 where $h_{imp}(n)$ is zero for $i < 0$.

15 The formant correction is performed by the filtering of $h'_{imp}(n)$ by the zero-state all-pole
16 filter:

$$H^{corr}(z) = \frac{1}{A^{corr}(z)}$$

17

(5.6.14.2-12)

18 where the coefficients of $A^{corr}(z)$ were calculated in Section 5.6.14.2.1.

19 The last modification of the impulse response is obtained by the convolution of output of
20 the formant correction filter with a constant 20-sample dispersion vector, to generate
21 $\hat{h}_{imp}(n)$.

22

23 5.6.14.2.4 Pitch Enhancement Gain and Pre-Calculation of Search Parameters

24 The pitch gain used for pitch enhancement is obtained as:

$$\hat{g}_a = \begin{cases} \min\{g_a, 1\} & \text{type -1 frame} \\ 0.75 & \text{type -0 frame} \end{cases}$$

25

(5.6.14.2-13)

26 The search parameters are pre-calculated, as described in Section 5.6.11.5.

1

2 5.6.14.2.5 Fixed-codebook Search for Type-0 Frames

3 The 2-pulse sub-codebook, described in Table 5.6-6, is searched with a first search turn
 4 with the iterative search algorithm based on joint position search for two pulses (Section
 5 5.6.11.7), and a second turn with iterative search algorithm based on position search for a
 6 single pulse (Section 5.6.11.6). The redundancy in the representation of two pulses with
 7 two signs on identical tracks is removed between track 1 and track 2, by exchanging the
 8 pulse locations according to the individual pulse signs.

9 The criterion $\tilde{\varepsilon}$ from the 2-pulse sub-codebook is weighted, as a function of the pitch lag
 10 and the normalized pitch correlation.

11 Before the 3-pulse sub-codebook is searched, the modified pitch gain \hat{g}_a is set to 0.25, and
 12 search parameters are pre-calculated again, as described in Section 5.6.11.5. The 3-pulse
 13 sub-codebook is then searched, using a single turn of the iterative search algorithm based
 14 on joint position search for two pulses.

15 The maximum between the weighted criterion from the 2-pulse sub-codebook and the
 16 criterion from the 3-pulse sub-codebook selects which sub-codebook is chosen at this
 17 stage, with its best pulse locations and signs.

18 The chosen criterion is weighted again, using the subframe-modified class, the sharpness
 19 of the ideal excitation, the normalized pitch correlation, the information of the pre-selected
 20 sub-codebook (3-pulse or 2-pulse), the noise-to-signal ratio, and the pitch lag.

21 The Gaussian sub-codebook is searched last (Section 5.6.12), using the original impulse
 22 response, $h_{imp}(n)$ (Section 5.6.2).

23 The weighted criterion from the previous search is compared to the criterion from the
 24 Gaussian sub-codebook, and the best sub-codebook is selected.

25 The fixed-codebook contribution, $c_{unf}^{fix}(n)$, is reconstructed. If the 2-pulse sub-codebook or
 26 the 3-pulse sub-codebook were selected, the codebook enhancement contributions are
 27 added to the reconstructed fixed-codebook contribution.

28

29 5.6.14.2.6 Fixed-codebook Search for Type-1 Frames

30 The 2-pulse sub-codebooks, described in Table 5.6-7, is searched with a first search turn
 31 with the iterative search algorithm based on a joint position search for two pulses (Section
 32 5.6.11.7), and a second turn with iterative search algorithm based on a position search for
 33 a single pulse (Section 5.6.11.6).

34 The 3-pulse sub-codebook is searched next, using a single turn of the iterative search
 35 algorithm based on a joint position search for two pulses (Section 5.6.11.7).

36 The maximum between the criterion from the 2-pulse sub-codebook and the criterion from
 37 the 3-pulse sub-codebook selects which sub-codebook is chosen at this stage, with its best
 38 pulse locations and signs.

39 The chosen criterion is then weighted, using the subframe-modified class, the sharpness of
 40 the ideal excitation, and the noise-to-signal ratio.

41 The 5-pulse sub-codebook is searched last using the joint position search described in
 42 Section 5.6.11.8.

43 The weighted criterion from the previous search is compared to the criterion from the 5-
 44 pulse sub-codebook, and the best sub-codebook is selected.

1 The fixed-codebook contribution, $c_{unf}^{fix}(n)$, is reconstructed, and the codebook enhancement
2 contributions are added to the reconstructed fixed-codebook contribution.
3

4 **5.6.15 Filtering of Fixed-Codebook Contribution**

5 The fixed-codebook vector, $c_{unf}^{fix}(n)$, is convolved with the impulse response $h_{imp}(n)$ to
6 generate the filtered fixed-codebook vector, $c_f^{fix}(n)$.

7
8

1 **5.6.16 Gain Re-optimization and Normalization for Rate 1/2 and Rate 1**

2 Routine name: GEQ_gain_reopt_2

3 **Input:**

- 4 • Filtered adaptive-codebook vector, $c_f^{adp}(n)$
- 5 • Filtered fixed-codebook vector, $c_f^{fix}(n)$
- 6 • Target vector, $T_{gs}(n)$

7 **Output:**

- 8 • Quantized fixed-codebook gain, g_c

9

10 The purpose of this module is to modify the target signal prior to the gain quantization.
 11 First a re-optimization process takes places as follows. If the normalized pitch correlation
 12 is below 0.75 and the frame is of type 0, the adaptive and fixed-codebook gains are jointly
 13 optimized by minimizing the WMSE between the original and the reconstructed speech
 14 according to

$$15 \quad \{g_a, g_c\} = \arg \min \left\{ \sum_{n=0}^{L_{SF}-1} \left(T_{gs}(n) - \left(g_a c_f^{adp}(n) + g_c c_f^{fix}(n) \right) \right)^2 \right\} \quad (5.6.16-1)$$

16 This results in the following estimates of the two gains:

$$17 \quad g_a = \frac{R_{a,t} \cdot R_{c,c} - R_{a,c} \cdot R_{c,t}}{R_{a,a} \cdot R_{c,c} - R_{a,c} \cdot R_{a,c}} \quad (5.6.16-2)$$

18 and

$$19 \quad g_c = \frac{R_{c,t} - g_a R_{a,c}}{R_{c,c}}, \quad (5.6.16-3)$$

20 where

$$21 \quad R_{a,t} = \sum_{n=0}^{L_{SF}-1} c_f^{fix}(n) \cdot T_{gs}(n), \quad (5.6.16-4)$$

$$22 \quad R_{c,t} = \sum_{n=0}^{L_{SF}-1} c_f^{fix}(n) \cdot T_{gs}(n), \quad (5.6.16-5)$$

$$23 \quad R_{a,a} = \sum_{n=0}^{L_{SF}-1} \left(c_f^{adp}(n) \right)^2, \quad (5.6.16-6)$$

$$24 \quad R_{c,c} = \sum_{n=0}^{L_{SF}-1} \left(c_f^{fix}(n) \right)^2, \quad (5.6.16-7)$$

$$25 \quad R_{a,c} = \sum_{n=0}^{L_{SF}-1} c_f^{adp}(n) \cdot c_f^{fix}(n). \quad (5.6.16-8)$$

1 where L_{SF} is the subframe length.

2 If the normalized pitch correlation is above 0.75 or the frame is of type 1, the adaptive and
3 fixed-codebook gains are disjointedly optimized. The adaptive-codebook gain was already
4 calculated according to Equation (5.6.6.3-1). Subsequently, the fixed-codebook gain is
5 estimated according to

$$6 \quad g_c = \arg \min \left\{ \sum_{n=0}^{L_{SF}-1} \left(T_g(n) - (g_c \cdot c_f^{fix}(n)) \right)^2 \right\}, \quad (5.6.16-9)$$

7 where

$$8 \quad T_g(n) = T_{gs}(n) - g_a \cdot c_f^{adp}(n). \quad (5.6.16-10)$$

9 The gain is given by

$$10 \quad g_c = \frac{\sum_{n=0}^{L_{SF}-1} T_g(n) \cdot (c_f^{fix}(n))}{\sum_{n=0}^{L_{SF}-1} (c_f^{fix}(n))^2} \quad (5.6.16-11)$$

11 Second, a gain normalization process takes place as follows. First, we describe the gain
12 normalization process for type-0 frames. The target energy of both the quantized excitation
13 and reconstructed speech is estimated according to the smoothing factor β_E (Section
14 5.6.10).

15

16 Routine name: PRC_GainsNorm_Gc_Gp

17 **Input:**

- 18 • Filtered adaptive-codebook vector, $c_f^{adp}(n)$
- 19 • Filtered fixed-codebook vector, $c_f^{fix}(n)$
- 20 • Unfiltered adaptive-codebook vector, $c_{unf}^{adp}(n)$
- 21 • Unfiltered fixed-codebook vector, $c_{unf}^{fix}(n)$
- 22 • Target vector, $T_{gs}(n)$
- 23 • Noise-to-signal ratio, NSR
- 24 • Gain smoothing factor, β_E
- 25 • Excitation mode factor, C_{ext_mod}
- 26 • Weighted residual signal, $r_w(n)$
- 27 • Residual of modified signal, $r_{mod}(n)$

28 **Output:**

- 29 • Normalized gains, g_a and g_c

1

2 **5.6.16.1 Normalization for $NSR > 0.125$** 3 A scaling factor g_s for the codebook gains is estimated in order to match the energy
4 targets. It is based on the smoothed open and closed loop energy targets.5 It should be noted that the smoothing is variable and may be zero. The open loop scaling
6 factor is given by

7
$$g_s = 0.75 \cdot \sqrt{\frac{\max(\tilde{E}_e - 5, 0)}{\sum_{n=0}^{L_{SF}-1} (g_a \cdot c_{unf}^{adp}(n) + g_c \cdot c_{unf}^{fix}(n))^2}}, \quad (5.6.16.1-1)$$

8 where it is bounded by $g_s < 1.2/g_a$ and the target energy of the quantized excitation (also
9 referred as the open loop target energy) is given by

10
$$\tilde{E}_e = \beta_E \cdot \tilde{E}_e + (1 - \beta_E) \cdot \sum_{n=0}^{L_{SF}-1} r_w(n)^2. \quad (5.6.16.1-2)$$

11 The target energy of the reconstructed speech (also referred as the closed loop target
12 energy) is given by

13
$$\tilde{E}_s = \beta_E \cdot \tilde{E}_s + (1 - \beta_E) \cdot \sum_{n=0}^{L_{SF}-1} T_{gs}(n)^2. \quad (5.6.16.1-3)$$

14

15

16

17 **5.6.16.2 Normalization for $NSR \leq 0.125$** 18 The determination of g_s is derived as follows. First, an open-loop scaling factor and a
19 closed-loop scaling factor are computed as:

20
$$g_{ol} = \sqrt{\frac{\max(\tilde{E}_e - 5, 0)}{\sum_{n=0}^{L_{SF}-1} (g_a \cdot c_{unf}^{adp}(n) + g_c \cdot c_{unf}^{fix}(n))^2}}, \quad (5.6.16.2-1)$$

21 and

22
$$g_{cl} = \sqrt{\frac{\max(\tilde{E}_s - 5, 0)}{\sum_{n=0}^{L_{SF}-1} (g_a \cdot c_f^{adp}(n) + g_c \cdot c_f^{fix}(n))^2}}, \quad (5.6.16.2-2)$$

23 Here,

24
$$\tilde{E}_e = \beta_E \cdot \tilde{E}_e + (1 - \beta_E) \cdot \sum_{n=0}^{L_{SF}-1} r_{mod}(n)^2 \quad (5.6.16.2-3)$$

25 and

$$\tilde{E}_s = \beta_E \cdot \tilde{E}_s + (1 - \beta_E) \cdot \sum_{n=0}^{L_{SF}-1} T_{gs}(n)^2 \quad (5.6.16.2-4)$$

2 An attenuation of the LPC filter is also calculated as:

$$g^{LPC} = \sqrt{\frac{\sum_{n=0}^{L_{SF}-1} r_{\text{mod}}(n)^2}{\sum_{n=0}^{L_{SF}-1} T_{gs}(n)^2}} \quad (5.6.16.2-5)$$

4 where g^{LPC} is bounded by $g^{LPC} < 0.8$.

5

6 Then g_s is calculated in the following steps if the flag $C_{\text{ext_mod}}$ is not 0 as:

$$g_s = \min\left\{\frac{g^{LPC}}{0.8} \cdot g_{ol} + \left(1 - \frac{g^{LPC}}{0.8}\right) \cdot g_{cl}, g_{cl}\right\} \quad (5.6.16.2-6)$$

8 Then it is refined in two steps:

$$g_s = \min\left\{0.7 \cdot g_s, \frac{1.2}{g_a}\right\} \quad (5.6.16.2-7)$$

10 and $g_s = \min\{1 + g^{LPC}, g_s\}$.

11 Otherwise, $g_s = 0.75 \cdot \min\{g_{ol}, g_{cl}\}$. It is further refined in this case as $g_s = \min\{g_s, \frac{1.2}{g_a}\}$

12 In any case g_s is always lower-bounded by 1.0.

13

14 Based on the final scaling factor the unquantized gains are modified according to

$$g'_a = g_s \cdot g_a \quad (5.6.16.2-8)$$

$$g'_c = g_s \cdot g_c \quad (5.6.16.2-9)$$

17 and the target is corrected accordingly to

$$T_{gs}(n) = g'_a \cdot c_f^{adp}(n) + g'_c \cdot c_f^{fix}(n) \quad (5.6.16.2-10)$$

19

20 **5.6.16.3 Fixed-Codebook Gain Normalization for Type-1 Frames and $NSR > 0.25$**

21 Routine name: PRC_GainNorm_Gc

22 **Input:**

23 • Filtered adaptive-codebook vector, $c_f^{adp}(n)$

24 • Filtered fixed-codebook vector, $c_f^{fix}(n)$

- 1 • Unfiltered adaptive-codebook vector, $c_{unf}^{adp}(n)$
- 2 • Unfiltered fixed-codebook vector, $c_{unf}^{fix}(n)$
- 3 • Target vector, $T_g(n)$
- 4 • Residual of modified signal, $r_{mod}(n)$

5 **Output:**

- 6 • Normalized gain, g_c

7

8 The different scale factors are computed as:

9

$$10 \quad g_{ol} = \sqrt{\frac{\sum_{n=0}^{L_{SF}-1} r_{mod}^2(n)}{\sum_{n=0}^{L_{SF}-1} (g_c \cdot c_{unf}^{fix}(n))^2}} \quad (5.6.16.3-1)$$

11 and

$$12 \quad g_{cl} = \sqrt{\frac{\sum_{n=0}^{L_{SF}-1} T_g^2(n)}{\sum_{n=0}^{L_{SF}-1} (g_c \cdot c_f^{fix}(n))^2}} \quad (5.6.16.3-2)$$

13 g^{LPC} is computed:

$$14 \quad g^{LPC} = \sqrt{\frac{\sum_{n=0}^{L_{SF}-1} r_{mod}(n)^2}{\sum_{n=0}^{L_{SF}-1} T_g(n)^2}} \quad (5.6.16.3-3)$$

15 The gain-scaling factor g_s is computed as:

$$16 \quad g_s = 0.75 \cdot g^{LPC} \cdot g_{ol} + (1 - 0.75 \cdot g^{LPC}) \cdot g_{cl} \quad (5.6.16.3-4)$$

17 It is refined as:

$$18 \quad g_s = \min(1 + 0.75 \cdot g^{LPC}, \max(1.0, 0.5 \cdot g_s)) \quad (5.6.16.3-5)$$

19 Based on the final scaling factor, the unquantized gains are modified according to:

$$20 \quad g'_a = g_a, \quad (5.6.16.3-6)$$

$$21 \quad g'_c = g_s \cdot g_c, \quad (5.6.16.3-7)$$

22 and the target is corrected accordingly to

$$T_g(n) = g'_a \cdot c_f^{adp}(n) + g'_c \cdot c_f^{fix}(n)$$

(5.6.16.3-8)

1
2
3

1
2

3 **5.6.17 Adaptive-Codebook and Fixed-Codebook Gain Quantization for** 4 **Type-0 Frames**

5 Routine name: GEQ_gainVQMA_2

6

7 **Input:**

- 8 • Filtered adaptive-codebook vector, $c_f^{adp}(n)$
- 9 • Filtered fixed-codebook vector, $c_f^{fix}(n)$
- 10 • Unfiltered adaptive-codebook vector, $c_{unf}^{adp}(n)$
- 11 • Unfiltered fixed-codebook vector, $c_{unf}^{fix}(n)$
- 12 • Target vector, $T_g(n)$

13 **Output:**

- 14 • Quantized gains, g_a and g_c

15

16 For type-0, the adaptive and fixed-codebook gains are jointly vector quantized with 7 bits
17 per subframe. The 2-dimensional codebook is searched exhaustively for the entry that
18 minimizes the MSE between the target and the reconstructed speech signal, i.e. minimizing

$$19 \quad E = \sum_{n=0}^{L_{SF}-1} \left(T_g(n) - (g_a c_f^{adp}(n) + g_c c_f^{fix}(n)) \right)^2 \quad (5.6.17-1)$$

20 where L_{SF} is 40 and 80 for Rate 1 and Rate 1/2 respectively. The quantized adaptive and
21 fixed-codebook gains are derived from the 7 bits codebook. The entries of the codebook
22 contain the adaptive-codebook gain and the correction factor for the predicted fixed-
23 codebook gain.

24 The prediction of the fixed-codebook gain is based on a 4th and 2nd order MA prediction of
25 the fixed-codebook energy for Rate 1 and Rate 1/2 respectively.

26 The relation between the correction factor γ^m and the quantized fixed-codebook gain is
27 given by

$$28 \quad g_c = \gamma^m \cdot \tilde{g}_c \quad (5.6.17-2)$$

29 where g_c is the quantized fixed-codebook gain and \tilde{g}_c is the predicted fixed-codebook gain.
30 The predicted fixed-codebook gain is given by

$$31 \quad \tilde{g}_c = 10^{\frac{1}{20}(\tilde{E}_k - E_c + \bar{E})} \quad (5.6.17-3)$$

32 where the $\bar{E} = 30\text{dB}$ is the mean energy,

$$E_c = 10 \log_{10} \left(\frac{1}{L_{SF}} \sum_{n=0}^{L_{SF}-1} c_{unf}^{fix}(n)^2 \right), \quad (5.6.17-4)$$

2 and

$$\tilde{E} = \sum_{i=1}^p b_i \cdot (20 \log_{10}(\gamma^{m-i})), \quad (5.6.17-5)$$

4 where p is 4 and 2 for Rate 1 and Rate 1/2 respectively. The prediction coefficients of the
5 MA prediction are $\{b_1, b_2, b_3, b_4\} = \{0.7, 0.6, 0.4, 0.2\}$ and $\{b_1, b_2\} = \{0.6, 0.3\}$ for Rate 1 and
6 Rate 1/2 respectively.

7

8 **5.6.18 Fixed-Codebook Gain Quantization for Type-1 Frames**

9 Routine name: GEQ_gainVQMA_3, GEQ_gainVQMA_4

10 **Input:**

- 11 • Filtered adaptive-codebook vector, $c_f^{adp}(n)$
- 12 • Filtered fixed-codebook vector, $c_f^{fix}(n)$
- 13 • Unfiltered adaptive-codebook vector, $c_{unf}^{adp}(n)$
- 14 • Unfiltered fixed-codebook vector, $c_{unf}^{fix}(n)$
- 15 • Target vector, T_g
- 16 • Quantized gains, $g_a(m)$, $m = 1, 2, 3, [4]$

17 **Output:**

- 18 • Quantized gains, $g_c(m)$, $m = 1, 2, 3, [4]$

19

20 This section describes the fixed-codebook gain quantization for Rate 1 and Rate 1/2. For
21 type-1 frames the adaptive-codebook gain quantization is described in Section 5.3.18. The
22 fixed-codebook gain vector $g_c(m)$ for type-1 frame is quantized to 10 and 8 bits for Rate 1
23 and Rate 1/2, respectively. The quantization is predictive delayed quantization; i.e. it takes
24 place after all of subframes have been processed with unquantized fixed-codebook gains.
25 The excitation signals, the target signals, and the quantized adaptive- codebook gains, are
26 buffered during the subframe processing and used to perform delayed joint quantization of
27 the fixed-codebook gain vector.

28 For Rate 1, the 4-dimensional codebook is searched in order to minimize

$$E = \sum_{m=1}^4 \sum_{n=0}^{40} \left((T_g^m(n) - g_a(m) c_f^{adp,m}(n) - g_c(m) c_f^{fix,m}(n))^2 \right), \quad (5.6.18-1)$$

30 where the quantized adaptive-codebook gains $\{g_a(m)\}$ originate from the original frame
31 based processing, and $\{T_g^m(n)\}$, $\{c_f^{adp,m}(n)\}$, and $\{c_f^{fix,m}(n)\}$ are buffered during the
32 subframe processing. Note that the superscript m is the m^{th} subframe.

1 Similarly for Rate 1/2, we have

$$\begin{aligned}
 E &= \alpha_1 \sum_{n=0}^{52} \left((T_g^1(n) - g_a(1) \cdot c_f^{adp,1}(n)) - g_c(1) \cdot c_f^{fix,1}(n) \right)^2 \\
 &+ \alpha_2 \sum_{n=0}^{52} \left((T_g^2(n) - g_a(2) \cdot c_f^{adp,2}(n)) - g_c(2) \cdot c_f^{fix,2}(n) \right)^2, \\
 &+ \alpha_3 \sum_{n=0}^{53} \left((T_g^3(n) - g_a(3) \cdot c_f^{adp,3}(n)) - g_c(3) \cdot c_f^{fix,3}(n) \right)^2
 \end{aligned} \tag{5.6.18-2}$$

3 where the quantized adaptive-codebook gains $\{g_a(m)\}$ originate from the original frame
 4 based processing, and $\{T_g^m(n)\}$, $\{c_f^{adp,m}(n)\}$, and $\{c_f^{fix,m}(n)\}$ are buffered during the
 5 subframe processing and

$$\alpha_m = \frac{\sum_{n=0}^{L_{sf}(m)} (T_g^m(n))^2}{\sum_{n=0}^{L_{fjm}} T_g(n)^2} \tag{5.6.18-3}$$

7 For Rate 1, the fixed-codebook gains $\{g_c(1), g_c(2), g_c(3), g_c(4)\}$ are derived from a 2-stage
 8 10 bit codebook (7 and 3 bits respectively), where the entries of the codebook contain a 4-
 9 dimensional correction factor for the predicted fixed-codebook gains. For Rate 1/2, the
 10 fixed-codebook gains $\{g_c(1), g_c(2), g_c(3)\}$ are derived from an 8-bit codebook where the
 11 entries of the codebook contain a 3-dimensional correction factor for the predicted fixed-
 12 codebook gains. The prediction of the fixed-codebook gains is based on an MA prediction of
 13 the fixed-codebook energy.

14 The relation between the correction factors γ^m and the quantized fixed-codebook gains is
 15 given by

$$g_c^m = \gamma^m \cdot \tilde{g}_c^m \tag{5.6.18-4}$$

17 where g_c^m is the quantized fixed-codebook gain and \tilde{g}_c^m is the predicted fixed-codebook gain
 18 of the m^{th} subframe of current frame.

19 The predicted fixed-codebook gains are based on an MA prediction of the fixed-codebook
 20 energy, and it is given by

$$\tilde{g}_c^m = 10^{\frac{1}{20}(\tilde{E}^m - E_c^m + \bar{E})} \tag{5.6.18-5}$$

22 where the $\bar{E} = 34\text{dB}$ is the mean energy,

$$E_c^m = 10 \log_{10} \left(\frac{1}{L_{SF}} \sum_{n=0}^{L_{SF}-1} (c_{unf}^{fix,j}(n))^2 \right), \tag{5.6.18-6}$$

24 and

$$\tilde{E}^m = \sum_{i=1}^N b_{m,i} \cdot (20 \log_{10}(\gamma_{prev}^{N+1-i})) \tag{5.6.18-7}$$

1 where N is 4 and 3 for Rate 1 and Rate 1/2 respectively. The predictor coefficients for Rate
2 1 are

$$3 \quad \{b_{m,1}, b_{m,2}, b_{m,3}, b_{m,4}\} = \begin{cases} \{0.7, 0.6, 0.4, 0.2\} & m = 1 \text{ (1st subframe)} \\ \{0.4, 0.2, 0.1, 0.05\} & m = 2 \text{ (2nd subframe)} \\ \{0.3, 0.2, 0.075, 0.025\} & m = 3 \text{ (3rd subframe)} \\ \{0.2, 0.075, 0.025, 0.0\} & m = 4 \text{ (4th subframe)} \end{cases}, \quad (5.6.18-8)$$

4 and the predictor coefficients for Rate 1/2 are

$$5 \quad \{b_{m,1}, b_{m,2}, b_{m,3}\} = \begin{cases} \{0.6, 0.30, 0.100\} & m = 1 \text{ (1st subframe)} \\ \{0.4, 0.25, 0.100\} & m = 2 \text{ (2nd subframe)} \\ \{0.3, 0.15, 0.075\} & m = 3 \text{ (3rd subframe)} \end{cases}. \quad (5.6.18-9)$$

6 Consequently, the prediction of the energies of the 4/3 subframes is based on the same
7 past memory.

8

9 **5.6.19 Random Number Generator**

10 Routine name: GCB_quickran

11 **Input:**

- 12 • Seed, *seed*

13

14 **Output:**

- 15 • Uniformly distributed variable, *rand(seed)*
- 16 • Updated Seed, *seed*

$$17 \quad r(n) = \sum_{i=1}^{12} rand(seed)$$

18 where *rand()* generates a zero mean, uniformly distributed variable. For each sample,
19 *seed* is updated as follows:

$$20 \quad seed = extract_l(seed \cdot 25173 + 13849), \quad (5.6.19-1)$$

21 where *extract_l(.)* extracts the 16 least significant bits of a 32 bit argument.

22

23

23

23

1 **6 SMV DECODER**

2 **6.1 Rate and Frame Type Retrieval and Frame Error Detection**

3

4 **6.1.1 Handling NULL Rate 1/8 Frame**

5 If the decoder receives an all 1's (null) Rate 1/8 packet, frame erasure processing is
6 performed for the current frame. (The encoder bit-packing routine ensures that, for 1/8 rate,
7 an all 1's bit stream is replaced by a valid bit stream that includes a single 0.)

8 **6.1.2 Handling All Zero Frames**

9 If the decoder receives an all 0's packet (for all rates), the current frame is processed as
10 frame erasure. (The encoder bit-packing routine ensures that, for all rates, an all 0's bit
11 stream is replaced by a valid bit stream that includes a single 1.)

12 **6.1.3 Handling 'Rate 1 with Bit Errors' Frame**

13 If the decoder receives 'Rate 1 with Bit Errors' frame, it will process the frame as an erasure
14 (Insufficient Frame Quality) as in Section 6.8.

15 **6.1.4 Channel Rate-Determination Algorithm Error Detection**

16 The multiplex sub-layer occasionally provides bad rates to the SMV decoder due to errors in
17 the Channel rate-determination algorithm (CRDA).

18 Due to constraints imposed on the SMV encoder the following frame rate, and type
19 transitions between successive frames are illegal:

- 20 a) Type-1 Rate 1 frame following a Rate 1/8 frame
- 21 b) Type-1 Rate 1/2 frame following a Rate 1/8 frame
- 22 c) Type-1 Rate 1 frame following a Rate 1/4 frame
- 23 d) Type-1 Rate 1/2 frame following a Rate 1/4 frame
- 24 e) Rate 1/8 frame following a Type-1 Rate 1 frame
- 25 f) Rate 1/8 frame following a Type-1 Rate 1/2 frame
- 26 g) Rate 1/4 frame following a Type-1 Rate 1 frame
- 27 h) Rate 1/4 frame following a Type-1 Rate 1/2 frame

28 Once the frame rate and the type of the current frame have been determined from the
29 packet, a check for the above transitions shall be made. If any of the above transitions is
30 detected, the current frame shall be declared a bad frame. Frame erasure processing
31 should be performed as described in Section 6.8. In addition, if the transitions of a, b, c, or
32 d are detected, the following decoder memories are reset to their initialization values:

- 33 i. Excitation memory
- 34 ii. Frame Erasure Concealment excitation memory
- 35 iii. LPC synthesis memory
- 36 iv. Postfilter synthesis memory
- 37 v. Past fixed codebook energy
- 38 vi. Previous frame LSFs

39 Further, if the transitions of a, b, c, or d are detected, the fixed codebook gain moving
40 average vector quantization memories are extended by their old values.

41 If the Rate Sanity Flag in the Rate 1 SMV packet is 1, the frame shall be declared a bad
42 frame. Frame erasure processing should be performed as described in Section 6.8.

1 If the Rate Sanity Flag in the Rate 1/4 SMV packet is 1, the frame shall be declared a bad
2 frame. Frame erasure processing should be performed as described in Section 6.8.

3 If the two shaping filter bits (SF_IDX[0] and SF_IDX[1]) in the Rate 1/4 SMV packet are 1,
4 the frame shall be declared a bad frame. Frame erasure processing should be performed as
5 described in Section 6.8.

6 If the current frame is a type 1, Rate 1 frame following a type 1 (Rate 1 or Rate 1/2) frame,
7 and the pitch lag of the current frame is different from the pitch lag of the previous frame
8 by more than 15, the frame shall be declared a bad frame. Frame erasure processing
9 should be performed as described in Section 6.8.

10 If the current frame is a Type-1, Rate 1/2 frame following a Type-1 (Rate 1 or Rate 1/2)
11 frame, and the pitch lag of the current frame is greater than 40, and is different from the
12 pitch lag of the previous frame by more than 15, the frame shall be declared a bad frame.
13 Frame erasure processing should be performed as described in Section 6.8.

14 If the current frame is a Rate 1/8 frame whose gain factor g_r is such that

$$15 \quad g_r > 5g_{mem} \quad (6.1.4-1)$$

16 the frame shall be declared a bad frame. Frame erasure processing should be performed as
17 described in Section 6.8, where g_{mem} is the estimate of the average frame gain for Rate 1/8,
18 and is updated after the above test in each Rate 1/8 frame as

$$19 \quad g_{mem} = \max(\text{round}(\max(g_r, 100)), \text{round}(0.05 \cdot \text{round}(\max(g_r, 100)) + 0.95 \cdot g_{mem}))$$

20 (6.1.4-2)

21 After the fixed codebook excitation, $c_{unf}^{fix}(n)$, and the fixed-codebook gain, g_c , have been
22 determined as described in Section 6.7.4 and 6.7.6 for each subframe of fifty type-0 frames
23 (Rate 1 or Rate 1/2) that immediately follow a type-1 frame, the following test for unusually
24 high fixed-codebook gain shall be performed:

$$25 \quad E_{nrm}^{fix} = \frac{g_c E_{unf}^{fix}}{\text{MAX}\left(0.01, \sqrt{10^{(\hat{E}_{MAX}/10-8)}}\right)}, \quad (6.1.4-3)$$

26 where E_{unf}^{fix} is the energy of the unfiltered fixed codebook excitation, and \hat{E}_{MAX} is the
27 maximum LPC synthesized speech frame energy in decibels in the last fifty type-1 frames.

28 For Rate 1, if $\left\{-2.34(G^{LPC})^3 + 170.16(G^{LPC})^2 - 41431(G^{LPC}) + 36754\right\} < E_{nrm}^{fix}$ and $G^{LPC} < 35$, or
29 $E_{nrm}^{fix} > 21000$, the subframe gain g_c should be reset to zero, and the FCB gain moving
30 average VQ memory should be reset to the initial values to avoid high energy speech
31 artifacts.

32 For Rate 1/2, if $\left\{-1.84(G^{LPC})^3 + 162.8(G^{LPC})^2 - 4796(G^{LPC}) + 47520\right\} < E_{nrm}^{fix}$ and $G^{LPC} < 35$,
33 or $E_{nrm}^{fix} > 18000$ the subframe gain g_c should be reset to zero, and the FCB gain moving
34 average VQ memory should be reset to the initial values to avoid high energy speech
35 artifacts.

36 For Rate 1/2 type-0 frames, if an index in the unused index space for the 2-pulse codebook
37 (with the additional Gaussian codebook and the additional TTY/TTD/DATA/DTMF flag) is
38 detected, the frame is erased.

39

6.2 General Structure of the Decoder

The decoder comprises 3 elements:

1. Excitation generation
2. LPC synthesis filter
3. Post filtering

The operation of each element is different for each rate. In addition, a frame erasure concealment algorithm is used to generate the decoder output during single or multiple frame erasure.

Figure 6.2-1 describes the general structure of the decoder, without the frame concealment module. The specific details of decoder operations are described in the next sections.

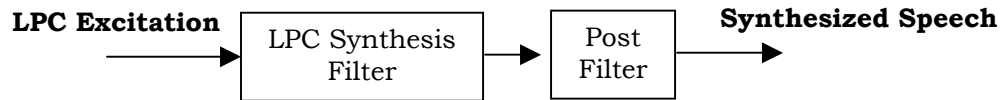


Figure 6.2-1: General Decoder Structure

6.3 Decoding of the LSFs

The indices into the LSFs multi-stage tables shall be retrieved from the bitstream, according to the transmitted rate. The quantized prediction error, $\varepsilon^q(i)$, shall be constructed by the sum of the entries in the quantization tables, as specified in Section 5.3.22.4. The reconstructed LSFs shall be obtained from the combination of the MA prediction, the decoded prediction error obtained from the decoded indices, and the mean vector of the LSFs by:

$$\underline{lsf}^q(i) = \varepsilon^q(i) + \sum_{p=1}^{O_{Rate}^{MA}} b_{j,p}^{Rate}(i) \cdot \varepsilon_p^q(i) + M_{LSF}(i) \quad i = 1, \dots, 10 \quad (6.1.4-1)$$

where j is the decoded or default predictor according to the rate, and O_{Rate}^{MA} , the order of the MA prediction, is 2 for Rate 1 and 4 for all other rates

6.4 Interpolation of Decoded LSFs

The LSFs shall be interpolated the same way specified for the encoder in Section 5.3.23. For Rate 1 type-0 frames, the interpolation path shall be decoded from the bitstream, and the decoded interpolation path shall be used.

1 **6.5 Excitation Decoding for Rate 1/8**

2 The decoding of Rate 1/8 is depicted in Figure 6.5-1. The gain index shall be decoded from
3 the bitstream input to the decoder, and the gain correction factor shall be obtained from
4 the quantization table. The gain value for the frame shall be calculated from the correction
5 factor and the predicted energy, as described in Section 5.4.2. The excitation shall be
6 generated by a random vector generator, described in Section 5.6.19. The excitation shall
7 be multiplied by the gain, to generate the excitation to the LPC synthesis filter.

8

9

10

11

12

13

14

15

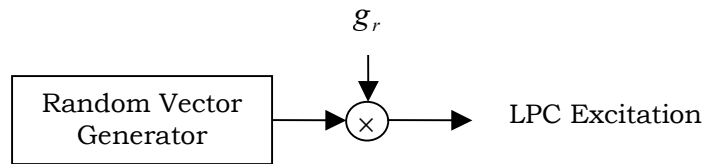


Figure 6.5-1: Excitation Decoding for Rate 1/8

6.6 Decoding of Rate 1/4

The decoding of Rate 1/4 is depicted in.

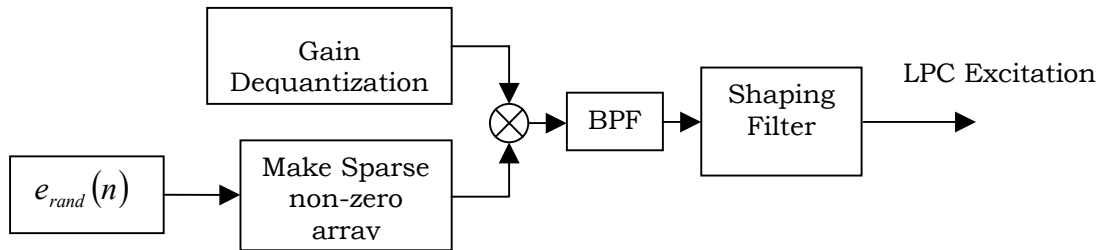


Figure 6.6-1: Excitation Decoding for Rate 1/4

The gain indices $i\hat{G}'$, and $i\hat{G}''[j]$, $0 \leq j < 2$, and the shaping filter index F_{idx} are decoded from the bitstream input to the decoder. The gains are dequantized using Equations (5.5.1-5), (5.5.1-8) and (5.5.1-9). A pseudo-random vector, $e_{rand}(n)$, is generated as described in Section 5.5.2. As described in Section 5.5.3, this pseudo-random vector is used to generate a gain scaled sparse signal, $r_1(n)$, which is then shaped by the bandpass filter $BP(z)$, and the shaping filter $SF(z)$ to yield the Rate 1/4 frame LPC excitation. As shown in Equation (5.5.4-9), the index of the filter, F_{idx} is used to establish $SF(z)$. The Rate 1/8 frame excitation is then filtered by the LP synthesis filter, and the post-processing filter to generate speech.

6.7 Decoding of Rate 1 and Rate 1/2

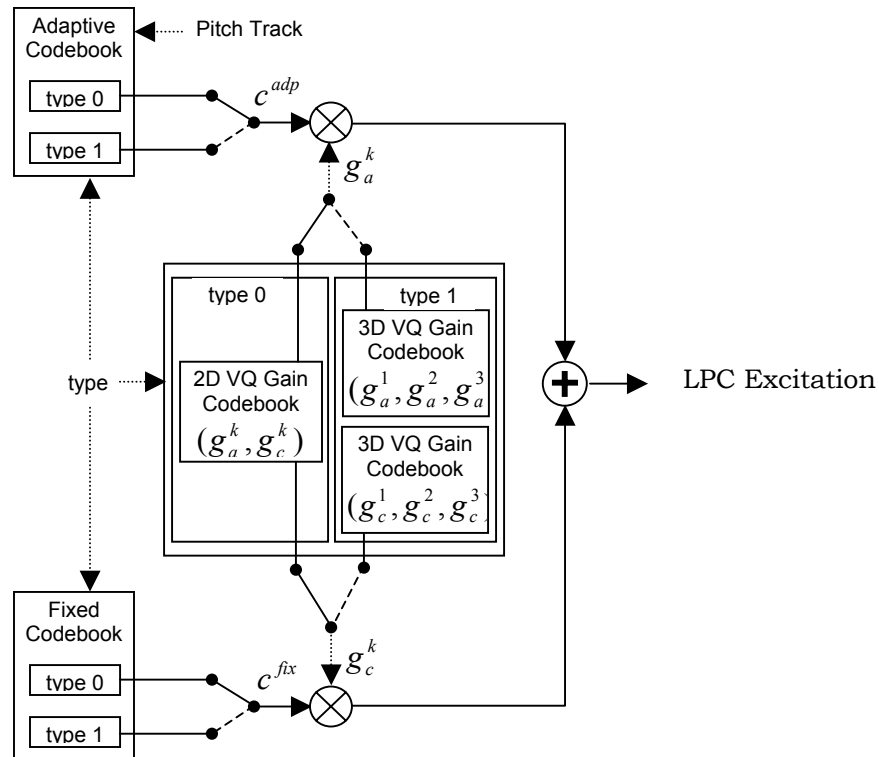


Figure 6.7-1: The SMV Decoder for Rate 1 and Rate 1/2

Figure 6.7-1 describes the excitation generation for Rate 1 and Rate 1/2. The excitation is constructed from the fixed-codebook contribution that is multiplied by the fixed-codebook gain, which is added to the adaptive-codebook contribution that multiplied by the adaptive codebook gain.

For each subframe, the fixed-codebook contribution is constructed from the pulse locations and signs, or from the Gaussian codebook indices. The adaptive-codebook contribution for type-0 frames is constructed from the pitch lag values for each subframe. For type-1 frames the adaptive-codebook contribution is obtained from the interpolated pitch contour. The gains for type-0 frames are obtained from the joint table for the fixed-codebook gain and the adaptive-codebook gain. For type-1 frames, the adaptive-codebook gains for the whole frame are obtained from the table for the adaptive-codebook gains, and the fixed-codebook gains for the whole frame are obtained from the table for fixed-codebook gains.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

6.7.1 Identifying Long-Term Spectral Characteristics

The identification of the long-term spectral characteristics shall be performed as done at the encoder (Section 5.3.25).

6.7.2 Pitch Decoding (Adaptive-codebook Contribution)

For the case of type-1 frames, the pitch track and pitch interpolation shall be performed as done in the encoder (Section 5.3.14.2) using the decoded frame pitch lag and the previous frame pitch lags.

The adaptive-codebook vector contribution for the subframe, $c_{unf}^{adp}(n)$, shall be obtained from the adaptive codebook $C^{adp}(n)$ with fractional interpolation of $1/20$ of a sample, using a 21-phase sampled Sinc window. The fractional pitch for the subframe, l_{FRC}^p , and the integer pitch for the subframe l_{INT}^p , shall be obtained from the interpolated pitch contour in the middle of the subframe. The fractional pitch is the value of the pitch contour at the middle of the subframe, and the integer pitch is the integer part of the fractional pitch. For the case of type-0 frames (Rate 1 and Rate 1/2), the adaptive-codebook contribution shall be derived from the decoded subframe pitch lags in the same way as in the encoder (Sections 5.6.6 and 5.6.7).

6.7.3 Adaptive-codebook Gain Decoding for Rate 1 and Rate 1/2 (Type-1)

Using the decoded index, the adaptive-codebook gains $g_a(m)$ shall be obtained directly from the quantization tables.

6.7.4 Fixed-codebook Excitation Decoding for Rate 1 and Rate 1/2

Using the received indices comprising the information about the codebook type and pulse positions and signs, along with the relevant decoded information such as reconstructed LP coefficients, adaptive-codebook gains (current for type-1 frames, previous for type-0 frames), pitch lags, and fixed-codebook contribution $c_{unf}^{fix}(n)$ shall be reconstructed as it was done in the encoder (Sections 5.6.13 and 5.6.14).

6.7.5 Fixed-codebook Excitation Contribution for Rate 1/8 and Rate 1/4

The fixed-codebook contribution shall be obtained using a Gaussian random excitation. The seed for the random number generator shall be obtained and updated as was done at the encoder (Section 5.5 and 5.6.19)

6.7.6 Gain Decoding

The adaptive-codebook gain $g_a(m)$ for Rate 1 and Rate 1/2 (type 0) shall be obtained from the relevant tables using the decoded index. The fixed-codebook gains $g_c(m)$ for Rate 1, Rate 1/2 and Rate 1/8 shall be obtained using the proper prediction. The prediction shall be complemented with the decoded prediction errors as it was done at the encoder (Sections 5.4.2, 5.6.17 and 5.6.18) in order to obtain the final $g_c(m)$.

1 **6.7.7 Signal-to-Noise Ratio (SNR) at the Decoder**

2 An estimate of the SNR at the decoder is calculated as:

3
$$SNR_{dec} = E_{dec}^S - E_{dec}^N$$
 (6.7.7-1)

4 where

5
$$E_{dec}^S = \max\{0.95 \cdot E_{dec}^S + 0.05 \cdot E_{dec}, E_{dec}\},$$
 (6.7.7-2)

6
$$E_{dec}^N = \min\{0.95 \cdot E_{dec}^N + 0.05 \cdot E_{dec}, E_{dec}\}$$
 (6.7.7-3)

7 and

8
$$E_{dec} = 10 \cdot \log_{10} \left(\frac{1}{L_{fsm}} \sum_{i=0}^{L_{fsm}-1} syn(i)^2 \right).$$
 (6.7.7-4)

9 $syn(\cdot)$ is the synthesized speech at the decoder. Note that E_{dec}^S is updated during Rate 1 or
10 Rate 1/2 and E_{dec}^N is updated during Rate 1/8.

11

12

6.8 Frame Erasure Concealment

A bad frame indicator for the frame $bfi(frm)$ is set if an erased frame indicator was received. A counter N_{bfi} counts the number of consecutive bad frames. It is reset to 0 when a good frame is received. A counter N_{gfi} counts the number of consecutive good frames after a bad frame has been received. N_{gfi} is saturated to 4 frames.

The decoder rate of operation and frame type for the current frame are set as follows if $bfi(frm)$ is 1:

$$Rate(frm) = Rate(frm - 1) \quad (6.7.7-1)$$

$$Type(frm) = Type(frm - 1) \quad (6.7.7-2)$$

6.8.1 LSF Extrapolation

If $bfi(frm) = 1$, use LSF parameters of previous frame i.e. $\underline{lsf}^{frm}(i) = \underline{lsf}^{frm-1}(i)$, $i = 1, \dots, 10$. The state variables of the LSF prediction are recomputed with an attenuation as:

$$\varepsilon_j(i) = 0.6 \cdot (\underline{lsf}'(i) - \sum_{p=1}^{O_{Rate}} b_{j,p}^{Rate}(i) \cdot \varepsilon_{j,p}^q(i)) \quad i = 1, \dots, 10 \quad (6.8.1-1)$$

where \underline{lsf}' is the mean-removed \underline{lsf} .

Special care is catered for good frames that are subsequent to a bad frame if F_c is set and the following is true:

if $1 \leq N_{gfi} \leq 2$ and $lsf(i) - lsf(i-1) < 75$ Hz, set the minimum spacing between consecutive LSFs to 125 Hz, and set an attenuation flag F_a to 1. If $lsf(i) - lsf(i-1) \geq 75$ Hz, set the minimum spacing between consecutive LSFs to 100 Hz. Otherwise, the minimum spacing is set to 50 Hz.

A caution flag F_c is set for the subsequent frame if the following is true:

$$bfi(frm) = 1 \text{ and } (k(1) > 0.4 \text{ or } (Rate = 1/8 \text{ and } SNR_{dec} > 50 \text{ dB})) \quad (6.8.1-2)$$

where $k(1)$ is the first reflection coefficient obtained from the decoded linear prediction coefficients. In this case the state variables of LSF prediction are prepared for the next frame, and are reset to:

$$\varepsilon_j(i) = \underline{lsf}'(i) - \sum_{p=1}^{O_{Rate}} b_{j,p}^{Rate}(i) \cdot \varepsilon_p^q(i) \quad i = 1, \dots, 10 \quad (6.8.1-3)$$

where $\underline{lsf}'(i)$ is the mean-removed vector of a preset $\underline{lsf}^{FLAT}(i)$ derived from a flat spectrum.

The caution flag F_c is reset to 0 if $N_{gfi} > 3$.

6.8.2 Pitch Extrapolation

Different pitch extrapolation algorithms are devised for type-0 and type-1 frames.

For Type-0: if $bfi(frm)$ is 1,

$$\begin{aligned}
& l_p(m) = l_p(m-1) \quad \text{if } k(1) < 0 \\
& l_p(m) = l_p(m-1) + 1 \quad \text{otherwise}
\end{aligned} \tag{6.8.2-1}$$

2 where $k(1)$ is the first reflection coefficient from the decoded prediction coefficients.

3 For type 1, if $bfi(frm)$ is 1,

4 If $\left| l_{pp}^q(frm-1) - l_{pp}^q(frm-2) \right| < 4$ and $\left| l_{pp}^q(frm-2) - l_{pp}^q(frm-3) \right| < 4$:

$$l_{pp}^q(0) = \min(147, \max(17, 3 \cdot A + B)) \tag{6.8.2-2}$$

6 otherwise

$$l_{pp}^q(frm) = l_{pp}^q(frm-1) \tag{6.8.2-3}$$

8 Where A and B are obtained through linear extrapolation as:

$$A = 0.5 \cdot \left(l_{pp}^q(frm-1) - l_{pp}^q(frm-2) \right) \tag{6.8.2-4}$$

10 and

$$B = \frac{l_{pp}^q(frm-3) + l_{pp}^q(frm-2) + l_{pp}^q(frm-1) - 3 \cdot A}{3} \tag{6.8.2-5}$$

12 A special feature is introduced to correct the pitch track and the ACB buffer for isolated
13 frame errors. This takes place if $bfi(frm)$ is 0 and $bfi(frm-1)$ is 1 and $bfi(frm-2)$ is 0.
14 The pitch track correction happens only the current frame and the previous frame are of
15 type 1. In this case, we use the correctly received $l_{pp}^q(frm)$ to correct $l_{pp}^q(frm-1)$ via
16 linear interpolation as:

$$l_{pp}^q(frm-1) = \left(\frac{l_{pp}^q(frm) - l_{pp}^q(frm-2) \cdot 159}{319} \right) + l_{pp}^q(frm-2) \tag{6.8.2-6}$$

18 The correction of the adaptive-codebook buffer happens for either frame type if both the
19 current and previous frame is either Rate 1 or Rate 1/2 and $g_a > 0.5$. The correction of the
20 adaptive-codebook buffer is basically the re-computing of the excitation of the preceding
21 bad frame using the corrected pitch lag and tracks. The gains and fixed-codebook
22 excitations remain as they were extrapolated during the preceding bad frame.
23

24 6.8.3 Gain Extrapolation

25 For type-1 frames:

$$\begin{aligned}
& g_a(m) = 0.95 \quad \text{if } N_{bfi} = 1 \\
& g_a(m) = \alpha \cdot g_a(m-1) \quad \text{if } N_{bfi} > 1
\end{aligned} \tag{6.8.3-1}$$

27 where α is an attenuation factor that is a function of N_{bfi} . An average adaptive-codebook
28 gain \bar{g}_a over the preceding frame is computed.

$$\begin{aligned}
& g_c(m) = 0.0 \quad \text{if } \bar{g}_a > 0.6 \\
& g_c(m) = \beta \cdot \sqrt{\frac{E_{sc}(m-1)}{E_c(m)}} \quad \text{otherwise}
\end{aligned} \tag{6.8.3-2}$$

1 where β is an attenuation factor that is a function of N_{bfi} , $E_{sc}(m-1)$ is the gain scaled
 2 fixed-codebook excitation energy from the preceding subframe, and $E_c(m)$ is fixed-
 3 codebook excitation energy of the current subframe.

4 For Rate 1 and Rate 1/2 type-0 frames, the gains are extrapolated as follows:

5 If $N_{bfi} = 1$, an average adaptive-codebook gain \bar{g}_a over the preceding subframes is
 6 computed in the first frame:

$$7 \quad \bar{g}_a = \frac{1}{L} \sum_{i=1}^L G_a(8-i), \quad (6.8.3-3)$$

8 where G_a is buffer containing the last 8 g_a and L is a function of the pitch lag. A factor
 9 $F_a(m,0)$ measuring the proportion of the adaptive-codebook excitation energy to the total
 10 excitation energy is computed on a subframe basis.

$$11 \quad \begin{aligned} g_a &= 0.95 && k(1) < 0.4 \text{ and } E_f(frm-1) \text{ and } \max(F_a(m,-1)) > 0.7 \\ g_a &= \min(0.95, \bar{g}_a) && (k(1) < 0.4 \text{ and } E_f(frm-1) \text{ and } \max(F_a(m,-1)) > 0.7) \text{ or } F_{ons} = 1 \\ g_a &= 0.2 && \text{otherwise} \end{aligned}$$

$$12 \quad (6.8.3-4)$$

14 where $E_f(frm-1)$ is the energy of the previous frame and F_{ons} is a flag indicating an onset
 15 at the previous frame. For other subframes, $g_a(m) = g_a(m-1)$.

16 If $N_{bfi} > 1$

$$17 \quad g_a(m) = \alpha \cdot g_a(m-1), \quad (6.8.3-5)$$

18 where α is a function of N_{bfi} .

19 The fixed-codebook gain is given as

$$20 \quad g_c(m) = \beta \cdot \sqrt{\frac{E_{sc}(m-1)}{E_c(m)}}, \quad (6.8.3-6)$$

21 where β is function of both N_{bfi} and $g_a(m)$. F_{ons} is set if there is an increase of at least
 22 13 dB between the current and the previous 2 frames or an increase of at least 7 dB in LPC
 23 gain between the current and previous frame.

24
 25 For Rate 1/4,

$$26 \quad \hat{G}[j] = \beta \cdot \hat{G}_{old}[j] \quad 0 \leq j < 10 \quad (6.8.3-7)$$

27 where β is a function of N_{bfi} , and $\hat{G}_{old}[j]$ is the set gains for the previous Rate 1/4 frame.

28
 29 For Rate 1/8,

$$30 \quad g_c(m) = \beta \cdot \sqrt{\frac{E_{sc}(m-1)}{E_c(m)}} \quad (6.8.3-8)$$

31 where β is a function of N_{bfi} .

32

6.8.4 Excitation Generation

A Gaussian random excitation $r(n)$ is generated using the random noise generator described in Section 5.6.19. $r(n)$ is generated regardless of rate in the case of bad frame, The seed used, is a signed 16-bit static variable initialized to 21845. The update of the seed is described in Section 5.6.19. No reset of this seed is performed.

6.9 Post-Processing

Routine name: PPR_post_process

Input:

- Interpolated prediction coefficients, $\hat{a}_m^q(i)$
- Decoded pitch lag value, $l_p(m)$
- Synthesized speech, $\hat{s}(n)$

Output:

- Postfiltered synthesized speech, $s'_p(n)$

Post-processing is made up of three functions: adaptive postfiltering, high-pass filtering and signal up-scaling. The adaptive postfilter is a cascade of three filters: a pitch postfilter $H_p(z)$, a tilt compensation filter $H_t(z)$, and a short-term postfilter $H_s(z)$, followed by and adaptive gain control (AGC). The postfilter is updated every subframe. The postfiltering process is organized as follows. First the synthesized speech $\hat{s}(n)$ is filtered through a filter $A(z/\gamma_{num})$ to produce a weighted residual speech signal $\hat{r}_w(n)$. The signal $\hat{r}_w(n)$ is used to compute the pitch delay l_{pf} and a gain g_{pf} . The signal $\hat{r}_w(n)$ is filtered through the pitch postfilter $H_p(z)$ to produce the signal $r'_w(n)$. $r'_w(n)$ is passed through a tilt compensation filter $H_t(z)$, leading to a signal $r_t(n)$. $r_t(n)$ is filtered through a synthesis filter $1/A(z/\gamma_{den})$ to give a signal $\hat{s}_p(n)$.

Finally, an AGC is then applied between the output signal $\hat{s}_p(n)$ and $\hat{s}(n)$ resulting in the signal $s'_p(n)$. The high-pass filtering and scaling are performed on $s'_p(n)$.

6.9.1 Weighted residual signal

The synthesized signal $\hat{s}(n)$ is filtered through a filter

$$A\left(\frac{z}{\gamma_n}\right) = 1 + \sum_{i=1}^{10} a_i \gamma_{num}^i z^{-i} \quad (6.9.1-1)$$

to produce a weighted residual speech signal $\hat{r}_w(n)$, where $\{a_i\}$ are the interpolated LP coefficients and γ_{num} is a adaptive weighting factor, updated each subframe as:

$$\gamma_{num} = 0.75 \cdot \gamma_{num} + 0.25 \cdot \gamma_0$$

where γ_0 is set to 0.63 for Rate 1, and 0.6 otherwise.

6.9.2 Pitch postfilter

The pitch postfilter is given by:

$$1 \quad H_p(z) = \frac{g_N}{1 + g_{pf} z^{-l_{pf}}} \quad (6.9.2-1)$$

2 where l_{pf} is pitch postfilter delay and g_{pf} is pitch postfilter gain and g_N is an AGC gain.

3 The pitch postfilter delay l_{pf} is obtained as the delay that maximizes the normalized
4 correlation $R_{pf}(D)$. D are the delays in the range of $[l_p - 1.5, l_p + 1.5]$ with 1/8 resolution,
5 l_p is the decoded subframe pitch lag. $R_{pf}(D)$ is given by

$$6 \quad R_{pf}(D) = \frac{\sum_{n=0}^{L_{SF}} \hat{r}_w(n) \hat{r}_w^D(n)}{\sum_{n=0}^{L_{SF}} \hat{r}_w^D(n) \hat{r}_w^D(n)} \quad (6.9.2-2)$$

7 where $\hat{r}_w^D(n)$ is the weighted residual at delay D . The pitch postfilter gain g_{pf} is obtained
8 in the following by first computing g_{pf} as:

$$9 \quad g_{pf} = \frac{\sum_{n=0}^{L_{SF}} \hat{r}_w(n) v(n)}{\sum_{n=0}^{L_{SF}} v(n) v(n)} \quad (6.9.2-3)$$

10 where $v(n)$ is the optimum contribution from $\hat{r}_w(n)$ at the optimum delay D . Then, g_{pf}
11 is limited to be greater than 1 followed by a scaling as:

$$12 \quad g_{pf} = \alpha \cdot g_{pf} \quad (6.9.2-4)$$

13 where α is 0.4 if Rate is 1, or 0.5 if Rate is 1/2, or 0 otherwise.

14 After the extraction of $v(n)$, $v(n)$ is updated as:

$$15 \quad v(n) = r_w(n) + g_{pf} \cdot v(n), n = 0, \dots, L_{SF} - 1 \quad (6.9.2-5)$$

16 The AGC gain g_N is obtained as

$$17 \quad g_N = \sqrt{\frac{\sum_{n=0}^{L_{SF}} \hat{r}_w(n) \hat{r}_w(n)}{\sum_{n=0}^{L_{SF}} v(n) v(n)}} \quad (6.9.2-6)$$

18 The final output $r'_w(n)$ is calculated as:

$$19 \quad r'_w(n) = g_N \cdot v(n), n = 0, \dots, L_{SF} - 1 \quad (6.9.2-7)$$

20 The pitch postfilter is only applied if Rate is 1 or 1/2. In the case of $bfi(0)$ or N_{bfi} is 1, the
21 optimum search for the delay D is not done but rather D is set to the decoded subframe
22 pitch lag l_p .

23

1 **6.9.3 Tilt Compensation**

2 The tilt compensation operation consists of the following steps. First, the impulse response
3 $h_{pf}(n)$ of the pole-zero filter

$$4 \frac{A\left(\frac{z}{\gamma_{num}}\right)}{A\left(\frac{z}{\gamma_{den}}\right)}, \quad (6.9.3-1)$$

5 truncated to 22 samples, is calculated where $\gamma_{den} = 0.75$. Second, the first reflection
6 coefficient $k_{pf}(1)$ of $h_{pf}(n)$ is computed as:

$$7 \quad k_{pf}(1) = -\frac{\sum_{j=0}^{18} h_{pf}(j)h_{pf}(j+1)}{\sum_{j=0}^{20} h_{pf}(j)h_{pf}(j)} \quad (6.9.3-2)$$

8 A tilt factor t_f is obtained from $k_{pf}(1)$ as:

$$9 \quad t_f = \max(\min(0.5 \cdot k_{pf}(1), 0), -0.1) \quad (6.9.3-3)$$

10 Then t_f is modified if $k_{pf}(1)$ is negative as:

$$11 \quad \begin{aligned} t_f &= t_f - 0.1, & \text{if } F_{flat} = 0 \\ t_f &= t_f - 0.05, & \text{otherwise} \end{aligned} \quad (6.9.3-4)$$

12 Finally, the tilt compensation leads to the output $r_t(n)$ as:

$$13 \quad r_t(n) = r'_w(n) + t_f \cdot r'_w(n-1), n = 0, \dots, L_{SF} - 1 \quad (6.9.3-5)$$

15 **6.9.4 Short-Term Postfilter**

16 The output of the tilt compensation filter $r_t(n)$ is passed through a synthesis filter

$$17 \quad \frac{1}{A\left(\frac{z}{\gamma_d}\right)} \quad (6.9.4-1)$$

18 to produce the synthesized post-filtered signal $\hat{s}_p(n)$, where $\gamma_d = 0.75$.

20 **6.9.5 Adaptive Gain Control**

21 Adaptive gain control is used to compensate for gain differences between the reconstructed
22 speech $\hat{s}(n)$ and the post-filtered signal $\hat{s}_p(n)$. The gain-scaling factor G_f for the current
23 subframe is calculated as:

$$G_f = \sqrt{\frac{\sum_{j=0}^{L_{SF}} \hat{s}(j)\hat{s}(j)}{\sum_{j=0}^{L_{SF}} \hat{s}_p(j)\hat{s}_p(j)}}. \quad (6.9.5-1)$$

2 Then G_f is modified as:

$$\begin{aligned} G_f &= 1.1 \cdot G_f, & \text{if } F_{flat} = 0 \\ G_f &= 1.15 \cdot G_f, & \text{otherwise} \end{aligned} \quad (6.9.5-2)$$

4 The gain scaled post-filtered signal $s'_p(n)$ is obtained as

$$s'_p(n) = g_f(n) \cdot \hat{s}_p(n), n = 0, \dots, L_{SF} - 1 \quad (6.9.5-3)$$

6 where $g_f(n)$ is updated on a sample-by-sample basis as

$$g_f(n) = 0.9 \cdot g_f(n) + 0.1 \cdot G_f \quad (6.9.5-4)$$

9 **6.9.6 High-pass filtering and up-scaling**

10 A high-pass filter at cutoff frequency of 80 Hz is applied to the gain scaled post-filtered
11 signal $s'_p(n)$. The filter is given by:

$$H_{pf}^h(z) = \frac{2.0 \cdot (0.92727435 - 1.8544941z^{-1} + 0.92727435z^{-2})}{1 - 1.9059465z^{-1} + 0.9114024z^{-2}} \quad (6.9.6-1)$$

13 The factor 2 in the numerator is for the embedded up-scaling by 2.

14
15
16
17
17

1 7 SMV DATA CAPABILITIES

2 The SMV provides the capability for passing data by replacing some or all of the speech
3 allocated bits by general user data bits. This capability is called in-band or in-bitstream
4 data. In this document, TTY/TDD Baudot Code (45.45 baud and 50 baud) and DTMF are
5 the only data payloads that are defined and supported by the SMV standard.

6 In-band TTY/TDD Baudot Code is mandatory. All SMV implementations shall include in-
7 band TTY/TDD capability. However, in-band DTMF is not always mandatory. Depending
8 upon applications such as wireless local loop (fixed wireless) or mobility, in-band DTMF can
9 be optional. In any case, in-band and out-of-band DTMF shall not be transmitted
10 simultaneously.

11 For mobility applications, in-band DTMF signaling shall not be used.

12 For fixed wireless (also known as wireless local loop) applications, sending DTMF in the
13 forward link is common (e.g., answering machine access). In order to guarantee that in-band
14 DTMF can be decoded correctly in both forward and reverse link, implementing in-band
15 DTMF decoder is mandatory in both directions (forward and reverse), whereas implementing
16 the in-band DTMF encoder is optional.

17 In-band data may be transported in one of three different ways:

- 18 1. Rate 1/2 packet with data (61bits/packet = 3050 bps)
- 19 2. Rate 1 packet with data (144 bits/packet = 7200 bps)
- 20 3. Rate 1 packet with Rate 1/2 speech and data in remaining bits (64 bits/packet = 3200
21 bps)

22 The data rate noted above is the data net rate, due to overhead needed for the detection of
23 the special data packets.

24 The TTY/TDD and DTMF features define their own error recovery schemes. For general
25 data applications, however, the in-band data channel is a raw bit stream and does not
26 currently define protocols for handshaking or schemes for error detection/correction. The
27 decoder, however, is capable of automatically detecting the in-band data by decoding the
28 packet. No out-of-band signaling is required to send in-band data from the encoder to the
29 decoder.

30 When the encoder sends in-band data to the decoder, the data is packed in either Rate 1/2
31 frames or Rate 1 frames and they contain a unique bit pattern that can be identified by the
32 decoder to indicate that the packet contains data. The decoder then reads the header bits
33 to determine which type of data is in the frame and calls the appropriate algorithm to
34 regenerate the signal at the decoder's output, or to send the data to an appropriate
35 application.

36 The bit allocation for Rate 1/2 data frames is summarized in Table 7-1.

37

Bit Positions	# of Bits	Rate 1/2 Data Frame
1	1	SVS (frame type)
2 - 4	3	Data Header
5 - 65	61	Data Block
66 - 80	15	Illegal Codeword

38 **Table 7-1: Bit Allocation for Rate 1/2 Data Frames**

39

40 Rate 1 packets can have one of two different data formats. The first is for transmitting data
41 along with speech parameters (Data + Speech). The second contains only data (Data +
42 Auxiliary Data). The bit allocations for both Rate 1 data formats are shown in Table 7-2.
43

1

Bit Positions	# of Bits	Rate 1 Data + Speech Frames	Rate 1 Data + Auxiliary Data Frames
1	1	SVS (frame type)	SVS (frame type)
2 – 81	80	Rate 1/2 Speech Frame	Auxiliary Data
82 – 84	3	Data Header	Data Header
85 – 148	64	Data Block	Data Block
149 – 170	22	Illegal Codeword	Illegal Codeword

2

Table 7-2: Bit Allocation for Rate 1 Data Frames

3

4 For all 3 in-band data formats, the SVS (frame type) field is set to zero and a codeword is
5 put in the packet that is considered illegal for speech frames. A packet is identified as
6 containing data if these to fields are set to the values below. The values of the SVS (frame
7 type) bit and illegal codewords for each of the in-band data formats are in Table 7-3.

	Rate 1/2 Data	Data + Speech	Data + Auxiliary Data
SVS (frame type)	0	0	0
Illegal Codeword	0x7FC8 (15 bits)	0x2002EA (22 bits)	0x2002EA (22 bits)

8

Table 7-3: SVS (frame type) and Illegal Codeword Values

9

10 Whereas the SVS and illegal codeword identify the packet as containing data, the Data
11 Header is a 3-bit field, which identifies the data format and payload. Table 7-4 contains a
12 description of the Data Header Field for Rate 1 and Rate 1/2 packets.

Data Header Value	Rate 1 (Bits 82-84)	Rate 1/2 (Bits 2-4)
0	Reserved	Reserved
1	Data + Speech	Data
2	Reserved	TTY/TDD
3	Reserved	DTMF
4	Data + Auxiliary Data	Reserved
5 – 7	Reserved	Reserved

13

Table 7-4: Data Header Description

14

7.1 TTY/TDD Payload Format

The TTY/TDD payload shall be transported using the Rate 1/2 packet data format only. Rate 1 data formats shall not be used for this application. For TTY/TDD payloads, the Data Header field shall be set to 2 (see Table 7-4) and the first 8 bits of the data field shall correspond to the TTY Type field. The Baudot code is currently the only TTY Type supported by this standard. Table 7.1-1 summarizes the valid values for the TTY Type field.

TTY Type	Description
0	Reserved
1	Baudot Code
2 – 255	Reserved

Table 7.1-1: TTY Type Field

7.1.1 Baudot Code Payload

Payloads for TTY/TDD Baudot codes at 45.45 baud and 50 baud are defined for SMV. For transmitting the Baudot code, the Data Field shall be set to 2 and TTY Type shall be set to 1. For Baudot, the bits in the Rate 1/2 packet are defined as in Table 7.1-2.

Rate 1/2 Packet Bit Positions	# of Bits	Field	Value
1	1	SVS (frame type)	0
2 – 4	3	Data Header	2
5 – 12	8	TTY Type	1
13 – 20	8	TTY Header	See Table 7.3-1
21 – 36	16	TTY Character	See Table 7.3-2
37 – 52	16	TTY Rate	See Table 7.3-3
53 – 65	13	Reserved	N/A
66 – 80	15	Illegal Codeword	0x7FC8

Table 7.1-2: Rate 1/2 Packet for TTY/TDD Baudot Code Payload

7.2 DTMF Payload Format

The DTMF payload shall be transported using the Rate 1/2 packet data format only. Rate 1 data formats shall not be used for this application. Table 7.2-1 shows a Rate 1/2 packet with the format for a DTMF payload. The Data Header field shall be set to 3 (see Table 7-4) and the first 16 bits of the data field shall correspond to a 16-bit signed value DTMF digit, as specified in Table 7.2-2.

Rate 1/2 Packets Bit Positions	# of Bits	Field	Value
1	1	SVS (frame type)	0
2 – 4	3	Data Header	3
5 – 20	16	DTMF Digit	See Table 7.2-2
21 – 65	45	Reserved	N/A
66 – 80	15	Illegal Codeword	0x7FC8

1
2**Table 7.2-1: Rate 1/2 Packet with DTMF Payload**

Data Field	DTMF Signal
1	1
2	2
3	4
5	5
6	6
7	7
8	8
9	9
10	0
11	*
12	#
13	A
14	B
15	C
16	D
-1	no signal detected

3

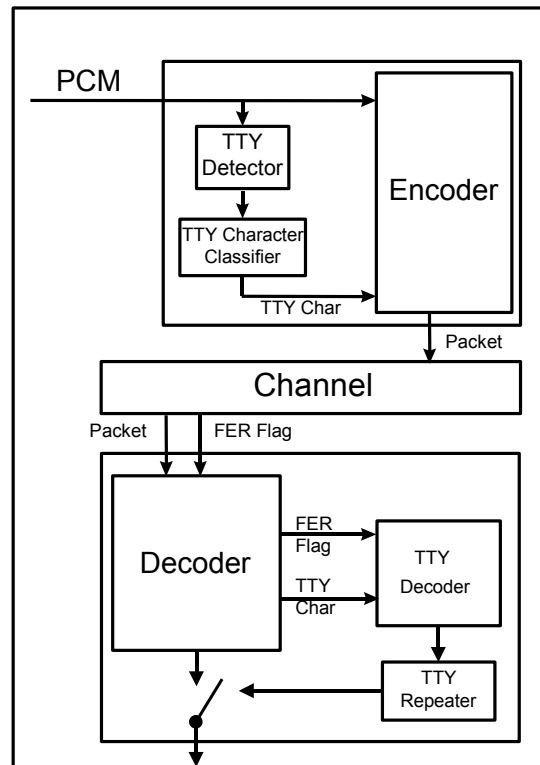
Table 7.2-2: DTMF Data Values

4 **7.3 TTY/TDD Operation**

5 The following feature provides a method for reliably transporting the 45.45 bps and 50 bps
6 Baudot code in the SMV, making digital wireless telephony accessible to TTY/TDD users.
7 The feature supports voice carryover/hearing carryover (VCO/HCO). VCO allows a
8 TTY/TDD user to switch between receiving TTY and talking into the phone. Similarly, HCO
9 allows a user to switch between transmitting TTY characters and picking up the phone to
10 listen. When Baudot tones are not present, the vocoder operates as usual; there is no
11 modification or added delay to the voice path when speech is present.

12 The Baudot code is specified in ITU-T Recommendation V.18. The Baudot code is a
13 carrierless, binary FSK signaling scheme. A 1400 Hz. tone is used to signal a logical "1"
14 and a 1800 Hz. tone is used to signal a logical "0". A 45.45 bps TTY bit has a duration of
15 22 ms. \pm 0.4 and a 50 bps bit has a duration of 20 ms. \pm 0.4, and a character consists of 1
16 start bit, 5 data bits, and 1.5 – 2 stop bits. When a character is not being transmitted,
17 silence, or a noisy equivalent, is transmitted. Hence, a TTY character spans a minimum of
18 6 speech processing frames. When the TTY encoder processing detects a character, it
19 sends the character and its header (see Section 7.3.1 for a description of the header) to the
20 decoder over a minimum of 6 consecutive frames and a maximum of 16 frames. Because of
21 channel impairments, the decoder may not receive all of the packets sent by the encoder
22 because of frame erasures and random bit errors. The decoder uses the redundancy to
23 correct any corrupted TTY information. Once the decoder recognizes the TTY character
24 being sent, the decoder's TTY repeater regenerates the Baudot tones corresponding to that

1 character. The TTY processing in the encoder processes the received PCM input one frame
 2 at a time and labels each frame as either NON_TTY, TTY_ONSET, or as a TTY character.
 3 The vocoder is in one of two states, TTY_MODE or NON_TTY_MODE. In the absence of
 4 Baudot tones, the encoder
 5



6
 7 **Figure 7.3-1: TTY/TDD Processing Block Diagram**
 8

9 and decoder are in the NON_TTY_MODE and the encoder and decoder process the frame as
 10 speech. When Baudot tones are present, the encoder and decoder enter TTY_MODE and
 11 process the TTY information as described below. When the vocoder stops detecting TTY
 12 information, it resumes processing speech.

13 It is recommended that a mechanism exist to disable the TTY/TDD feature in the vocoder.
 14

15 **7.3.1 TTY Header and Character Format**

16 The TTY information passed from the encoder to the decoder contains a header, character
 17 information, and the baud rate. When the encoder is transmitting a TTY character, the
 18 header contains a sequence number to distinguish that character from the characters
 19 immediately before and after. The same TTY information is transmitted for each instance of
 20 a character for a minimum of 6 frames and a maximum of 16 frames. The header cycles
 21 through its range of valid values, one value for each instance of a character. During the
 22 TTY_ONSET state, the encoder sends the decoder a TTY_SILENCE message in the header
 23 and character field. The TTY header uses 2 bits in an 8-bit field (see Table 7.3-1). The TTY
 24 character uses 5 bits in a 16-bit field (see Table 7.3-2). The TTY baud rate uses 1 bit in a
 25 16-bit field (see Table 7.3-3).
 26

Range	Description
-------	-------------

0	TTY_SILENCE
1 – 3	TTY Counter
4 – 255	Reserved

Table 7.3-1: Range of Values for TTY/TDD Baudot Header Field

Range	Description
4	TTY_SILENCE (when TTY Header = 0)
0 – 31	TTY Character
32 – 65535	Reserved

Table 7.3-2: Range of Values for TTY/TDD Baudot Character Field

Range	Description
0	45.45 Baud
1	50 Baud
2 – 65535	Reserved

Table 7.3-3: Range of Values for TTY/TDD Baud Rate Field

7.3.2 TTY Encoder Processing

The TTY encoder processing takes the larger task of detecting TTY characters and divides it into a series of smaller tasks, creating different levels of detection. It is through this divide and conquer approach that the `tty_enc()` routine has low complexity in the absence of Baudot tones.

The first level of detection is to divide the 160 samples in the speech frame into 10 blocks of 16 samples. These blocks are called detection intervals, or dits. Each dit is classified as `NON_TTY`, `LOGIC_0`, or `LOGIC_1`.

The next level of detection is to determine if there are enough `LOGIC_0` or `LOGIC_1` dits in a row to form a TTY bit. The transition from a “0” bit to a “1” bit is used to determine the onset of a TTY character and the `TTY_SILENCE` message is sent to the decoder. The `TTY_SILENCE` message shall continue to be sent until a TTY character is detected, or until a `NON_TTY` frame is detected. When enough bits are detected to form a character, the rate of the character is determined and the TTY character information is sent to the decoder.

7.3.2.1 TTY Encoder Inputs

- TTY character information from the previous frame (header, TTY character, and baud rate).
- 160 PCM samples from the output of the high pass filter.

7.3.2.2 Dit Classification

The 160 samples from the high pass filter’s output are converted into 10 dits. For every block of 16 samples, it computes the spectral energy at the mark and space frequencies using a 16-point DFT at 1400 Hz and 1800 Hz with a rectangular window. The ratio of the maximum energy between the mark and space energy and the total energy is compared to a threshold; i.e.,

$$\frac{\max(\text{mark_energy}, \text{space_energy})}{\text{total_energy}} > \text{THRESH}$$

If that threshold is exceeded, the dit is labeled as either a mark or a space, whichever one has the greater energy. If the threshold is not met, the dit is labeled as `NON_TTY`.

1 A second threshold is applied to the total energy of the dit. If the total energy is less than
 2 the threshold, the dit is classified as NON_TTY. This minimum energy threshold prevents
 3 the TTY detector from being overly sensitive to low level tones that may have been
 4 unintentionally reflected back into the TTY detector. The threshold is set at approximately
 5 -50 dBm.

7 7.3.2.3 Dits to Bits

8 The dits are used to form a TTY bit. Dits are classified as LOGIC_0, LOGIC_1, or
 9 UNKNOWN. A nominal bit consists of 11 dits for 45.45 baud and 10 dits for 50 baud. A bit
 10 is not required to have a continuous run of LOGIC_0's or LOGIC_1's in order to be detected.
 11 Spurious UNKNOWN detections are permitted, up to a threshold.

12 A TTY bit is searched by looking at the dits within a variable length sliding window. The
 13 window varies in length from 8 dits to 13 dits. The number of LOGIC_0's, LOGIC_1's, and
 14 UNKNOWN dit detections are counted in the window. The dits in the search window must
 15 meet the following criteria in order to detect a TTY bit:

- 16 • Minimum of 6 LOGIC_0 (LOGIC_1) dits
- 17 • Maximum of 2 LOGIC_1 (LOGIC_0) dits
- 18 • Maximum of 5 UNKNOWN dits

19 Heuristics are also applied so that the sliding window is centered over the TTY bit being
 20 detected. If all of the thresholds are met, a "0" ("1") TTY bit is declared, and the overall
 21 length of the bit and the number of bad dits within the window are recorded. If the dits in
 22 the search window do not meet the criteria for a TTY bit, a gap between TTY bits is declared
 23 and the gap is recorded. The window is slid by one dit and the search is repeated.

24 The length of the search window is allowed to vary. The length is adjusted depending on
 25 the previous character's baud rate, and where the mark/space transitions occur.

26 Two history buffers are maintained to record the detected TTY bits and gaps. The array
 27 `tty_bit_hist[]` maintains a history of the bits that are detected and `tty_bit_len_hist[]` stores
 28 the lengths, in dits, of the gaps and TTY bits that are detected. Both arrays are 9 elements
 29 long, there is one element for the start bit, five for the data bits, one for the stop bit, and
 30 one for the memory bit, the bit before the start bit, as depicted in [Figure 7.3-2](#).
 31 The last element in the array is used to record partially detected TTY bits.

0	1	2	3	4	5	6	7	8
Memory Bit	Start Bit	LSB Data 0	Data 1	Data 2	Data 3	MSB Data 4	Stop Bit	Next Bit

35 **Figure 7.3-2: TTY Bit History Buffer**

36 When a TTY bit or a gap is detected, its value is recorded in `tty_bit_hist[]` and its length is
 37 stored in `tty_bit_len_hist[]`. The length is used by `get_tty_char()` to threshold the length of a
 38 candidate character, and by `tty_rate()` to determine if a detected character is 45.45 baud or
 39 50 baud. Since the length of the memory bit is irrelevant, `tty_bit_len_hist[0]` is used to
 40 count the number of NON_TTY and TTY_SILENCE dits that were detected within the TTY
 41 bits.

42 Because the speech frames may not coincide with the boundaries of the TTY bits, it is
 43 possible that a bit may straddle two speech frames. It is possible, therefore, that the
 44 sliding window may contain only a partial bit within a frame. These partial detections are
 45 recorded in element 8 of the TTY bit history buffers, labeled "Next Bit" in Figure 7.3-2.

1 **7.3.2.4 TTY Character Classification**

2 The routine `get_tty_char()` checks if the detected bits form a TTY character. The following
3 conditions must be met in order for a character to be declared.

- 4 • The bit preceding the start bit must not be a “0”.
- 5 • The start bit must be a “0”
- 6 • The stop bit must be a “1”

- 7 • The length of the candidate character, in dits, must be within a maximum and
8 minimum threshold.
- 9 • The number of bad dit detections within the character must be within a threshold.

10 If all of the conditions are met, a character is declared.

11 Once a character is found, the character information and its header are sent to the decoder
12 a minimum of 6 frames and a maximum of 16 frames. The constant
13 `FRAMING_HANGOVER` dictates the maximum number of times the information for the
14 same character is sent. If a new character is framed before `FRAMING_HANGOVER` is
15 reached, the information for the old character is terminated and the new information is
16 sent to the decoder.

18 **7.3.2.5 TTY Baud Rate Determination**

19 After a character has been detected, `tty_rate()` determines the baud rate of the character by
20 counting the length, in dits, of the start bit and the five data bits. The length of the stop bit
21 is not used because its length is too variable.

22 A character with a length of 63 dits or greater is declared 45.45 baud, otherwise it is
23 declared 50 baud. A hangover of three characters is maintained before `tty_rate()` will switch
24 from one baud rate to the other. That is to say, if the baud rate changes in the middle of a
25 call, three consecutive characters must be detected at the new baud rate in order for
26 `tty_rate()` to declare the new baud rate.

28 **7.3.2.6 TTY State Machine**

29 The routine `get_tty_state()` is responsible for changing the state of the TTY encoder
30 processing. There are 3 states, `NON_TTY_MODE`, `TTY_ONSET`, and `TTY_MODE`.
31 `Get_tty_state()` is responsible for determining `NON_TTY_MODE` and `TTY_ONSET`.

32 Changing TTY state from `NON_TTY_MODE` to `TTY_ONSET` requires that a “0” bit is followed
33 by a “1” bit. This rule requires the presence of both the space tone and the mark tone and
34 for the tones to be the correct duration. This test must be met in order to declare
35 `TTY_ONSET` for the first time.

36 `NON_TTY_MODE` is declared by `get_tty_state()` whenever a non-TTY bit is detected or when
37 a “0” bit occupies the bit preceding the start bit.

39 **7.3.3 TTY/TDD Decoder Processing**

40 The TTY decoder processing must recognize when TTY/TDD information is in the packet,
41 recover from channel impairments to decode the TTY character being sent, and regenerate
42 the Baudot tones corresponding to that character.

43 When the decoder is in `NON_TTY_MODE`, the packet is decoded in the usual manner for
44 speech. When the decoder makes the transition to `TTY_MODE`, it mutes its output until it
45 receives TTY character information or until it transitions back to `NON_TTY_MODE`.

47 **7.3.3.1 TTY Decoder Inputs**

- 48 • TTY Information (header, TTY character, baud rate).

- Bad frame indicator

7.3.3.2 Decoding the TTY/TDD Information

The task of detecting the presence of TTY information, recovering from frame and bit errors, and decoding the TTY character is performed in `tty_dec()`. If the routine detects that TTY information is being sent, the `tty_dec()` flag is set to non-zero and the PCM buffer is filled with the appropriate Baudot tones. If TTY is not detected, the flag is set to zero and the PCM buffer is returned unmodified.

The routine labels each frame as `NON_TTY`, `TTY_SILENCE`, `FER`, or a TTY character, and maintains a history buffer of these classifications for 11 frames: 9 frames of lookahead, 1 current frame, and 1 frame of lookback (see Table 7.3-4). The most recent packet enters the buffer at location 0, but the decision for the current frame is based on the contents of element 9. The buffer is updated at the end of each frame, shifting its contents to the right by one. The buffer is initialized to `NON_TTY`.

Frame:	0	1	2	3	4	5	6	7	8	9	10
Description:	Lookahead									Current Frame	Look-back

Table 7.3-4: `tty_dec()` History Buffer

At the start of each frame, the most recent information is sanity checked to see if it is consistent with TTY information. If the frame erasure flag is set, the frame is labeled `FER`, otherwise the TTY/TDD information is checked to see if the header and TTY character fields fall within the allowed range of values. If all of the tests pass, Frame 0 is labeled with the TTY character in the history buffer.

The TTY decoder processing can reliably regenerate the TTY characters despite channel impairments because the character information is transmitted a minimum of 6 times from the encoder. Errors are corrected by a voting process. The current frame and nine frames of lookahead are used to determine the correct TTY header, character, and baud rate. Errors are replaced with the winner of the voting process.

Voting is conducted under the following conditions:

- Any time the current frame is labeled as `FER`.
- Every time the current frame contains information for the start of a new character. Because a new character must contain a minimum of 6 frames of the same information, a vote is taken to verify that the information is present before it will generate the tones for that character. Once a character wins the vote, any frame errors, bit errors, or other inconsistencies are corrected in the frame window where the character information is expected, i.e. the current frame and the adjacent frames of lookahead will contain the same header and character information.
- Any time the current frame contains `TTY_SILENCE` or a TTY character, and the frame of lookback contains `NON_TTY`. This makes it harder for the decoder to erroneously go into `TTY_MODE`, thus preventing false alarms when speech is present.

Once `tty_dec()` makes its decision on the current frame, `tty_dec()` calls `tty_gen()` to generate the appropriate PCM samples.

7.3.3.3 Baudot Generator

Once the current frame is labeled by `tty_dec()`, `tty_gen()` is called to fill the PCM buffer with the appropriate Baudot tones. In the case of `NON_TTY`, the PCM buffer is returned unmodified. In the case of `TTY_SILENCE`, the PCM is muted.

1 Generating TTY characters is more involved because one character spans many frames, so
 2 `tty_gen()` must generate the Baudot tones one subframe at a time. When a TTY character
 3 needs to be regenerated, `tty_gen()` puts a subframe's worth of samples in the PCM buffer. It
 4 keeps track of which bit it is in the middle of generating and the number of samples left to
 5 generate for that bit, so that the next time it is called, it can pick up where it left off. Once
 6 `tty_gen()` begins to generate a character, it will generate the entire character before it will
 7 generate the next character. This is done so that the repeater will only generate valid TTY
 8 characters.

9 There exists logic in `tty_gen()` to detect when the next character arrives before the current
 10 one is finished. If the next character arrives before the current one can be regenerated, the
 11 bit duration is reduced by 4 samples and a minimum of 1 stop bit is generated. Otherwise,
 12 nominal bit lengths are generated, i.e. 22 ms bits for 45.45 baud and 20 ms bits for 50
 13 baud, and a minimum of 1.5 stop bits is generated for 45.45 baud, and a minimum of 1.75
 14 stop bits is generated for 50 baud.

15 There exists a provision in the ITU-T Recommendation V.18 for the TTY/TDD device to
 16 extend its stop bit in order to prevent a TTY/TDD device from detecting its own echo. This
 17 routine will extend the stop bit a maximum of 300 ms if a TTY character is followed by
 18 silence. If a new character arrives before 300 ms has elapsed, the extended stop bit is
 19 terminated and the new character is generated immediately.

20 The tones themselves are generated by `tone_gen()`. Before `tty_gen()` returns, it updates the
 21 decoder's lookback field in the TTY history buffer with the information corresponding to the
 22 last samples generated. For example, if `tty_gen()` finished generating a character in the
 23 middle of the subframe and started generating silence, the lookback field is updated with
 24 `TTY_SILENCE`.

26 **7.3.3.4 Tone Generator**

27 The routine `tone_gen()` is a sine wave generator. Given a frequency and the number of
 28 samples, it will generate the PCM samples by using a 2 tap marginally stable IIR filter. The
 29 filter implements the trigonometric identity:

$$30 \quad \cos(k\omega) = 2 \cdot \cos(\omega) \cdot \cos((k-1)\omega) - \cos((k-2)\omega)$$

31 It is a zero excitation filter, using only its past 2 samples and the cosine of the frequency to
 32 be generated, to produce the next sample.

33 **7.4 DTMF**

34 **7.4.1 DTMF Detector**

35 The DTMF (Dual Tone Multiple Frequency) detector detects the presence and the
 36 information content of incoming DTMF tone signals. Its row and its column frequency as
 37 shown in Table 7.4-1 uniquely identify each key on the keypad. It outputs a 16-bit value
 38 associated with each of the 16 possible DTMF tone combinations.
 39

1

		Column Frequency Group			
		1209 Hz	1336 Hz	1477 Hz	1633 Hz
Row Frequency Group	697 Hz	1	2	3	A
	770 Hz	4	5	6	B
	852 Hz	7	8	9	C
	941 Hz	*	0	#	D

2

Table 7.4-1: Touch-tone telephone keypad:

3

A row and a column tone is associated with each digit.

4

7.4.1.1 Description of the DTMF Detector

5

The DTMF detector consists of three main functions:

6

1. AGC – normalizes the input frame.

7

2. Goertzel – performs spectral analysis at the 8 frequencies and at the second harmonic of the column frequencies using an 80 point Goertzel algorithm. The Goertzel algorithm provides an efficient means of calculating the DFT energy at a specific frequency (see Oppenheim and Schaffer, Digital Signal Processing, Prentice-Hall) In addition a 160 point Goertzel algorithm is used at the first and fourth row frequency in order to satisfy the requirements found in Bellcore document GR-506-CORE LSSGR: Signaling and Analog Interfaces Issue 1, June 1996.

8

9

10

11

12

13

14

3. Check– checks the spectral energies to determine if a valid DTMF tone is detected.

15

7.4.1.2 INPUT

16

- Pointer to 16 bit linear time series for current 20 msec. frame and 10 msec. look-ahead.

17

7.4.1.3 OUTPUT

18

- 16 bit valid DTMF tone digit if a tone is detected or 0xffff to indicate no DTMF present.

19

20