

# **Signals and Communication Technology**

More information about this series at <http://www.springer.com/series/4748>

Tom Bäckström

# Speech Coding

with Code-Excited Linear Prediction

 Springer

Tom Bäckström  
International Audio Laboratories Erlangen  
(AudioLabs)  
Friedrich-Alexander University  
Erlangen-Nürnberg (FAU)  
Erlangen  
Germany

ISSN 1860-4862                      ISSN 1860-4870 (electronic)  
Signals and Communication Technology  
ISBN 978-3-319-50202-1            ISBN 978-3-319-50204-5 (eBook)  
DOI 10.1007/978-3-319-50204-5

Library of Congress Control Number: 2017930277

© Springer International Publishing AG 2017

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by Springer Nature  
The registered company is Springer International Publishing AG  
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Preface

Depending on your source of reference, there are between 7 and 8 billion mobile phones in use at the time of writing this work. It is a breathtaking technological success. In less than 20 years, mobile telephony has reached all parts of the world and there is an almost universal access to mobile phones for all people in the world. That is not to say that all 7 billion people of the world would have their own phone since many have multiple phones, but a large majority of people does have *access* to a mobile phone. Almost all people in the world have the option of making a (mobile) phone call.

The consequences of this technological leap are as astonishing as the speed with which this leap has happened. In the industrial world, practically every adult has a mobile phone. It is a defining characteristic of our age—to be always reachable, to be always online. It has fundamentally changed our attitude towards communication. Being within reach is no more a thing that needs to be planned, gone are the days when we would agree the time and place when we would talk the next time. We are simply always there.

Improved telecommunication solutions have the potential of increasing the amount of communication between people. This is almost only a good thing. Increased communication improves understanding between people, it increases awareness of society and your surroundings, and has thus an impact on the level of democracy.

In the industrial world, democracy can of course be *nice*, but where mobile telephony has had a much larger impact on democracy is the developing world. It is fairly common for small, poor villages and communities in Africa to share one mobile phone for the whole village. It has proved to be more expensive to build and maintain land-lines than mobile networks, and thus hardly any new land-lines are being built. There are plenty of villages that never had a phone, never had an instantaneous communication channel, before the mobile phone. Where before, people had to walk to the next village to get news from their relatives, they can now simply make a phone call. They can even talk to their relatives who have moved abroad, perhaps even overseas, and get the latest news. Imagine the societal and human impact of that: For the first time in human history, any one person in the

world can connect with anybody else in an instant, independent of geographic location!

As with every success, there is a drawback; the huge steps forward have proved to be the biggest enemy of the technology. For example, it would be easy to make the conclusion that since 7 billion people have access to phones, *it works already* and there is thus not much to be improved. Alas, how often have I not heard this claim, even from scientists in the field!

Fortunately, the same argument can be turned upside down to demonstrate the importance of further improvements to mobile telephony. Imagine the impact that the tiniest of improvements would have, when multiplied by 7 billion users! The largest fixed costs of mobile networks surely lie in the network, where efficiency and costs are directly proportional. A 1% improvement in coding efficiency, in the actual amount of data sent, does not have a large impact on the individual phone, but worldwide, already the reduced energy consumption due to a 1% improvement in efficiency would have a huge impact.

At the same time, as the number of devices capable of speech coding is exponentially increasing, the amount of energy required is also increasing rapidly despite improvements in energy efficiency of CPUs. To make it feasible to use speech codecs on such a large number of devices, it will be necessary to reduce the complexity of speech coding algorithms. This is especially evident with the emerging concept of the internet-of-things, where it is envisioned that all electronic devices would be connected and capable of communication. It is obvious that we need to limit the complexity of codecs to reduce power consumption if they are to be implemented on such a large number of devices.

Moreover, even if only a fraction of the improvement yields lower costs for the end-users, it is not difficult to see its impact on that small African village. 50 cents saved per month is good money.

The other obstacles for mobile telephony imposed by its own success are the economic incentives within standardisation processes. Regrettably, however, the current author is not authorised to pursue an academic discussion of the economic incentives of standardisation any further.

It is then not surprising that during the last decade academic research in speech coding has been scarce. It seems that most of the cutting-edge research happens behind closed, corporate doors and little if any information is leaked out. All that is visible to the outside are the (sometimes) brilliant engineering art-works of speech coding standards, which include hardly any new science. Speech codecs have become very finely tuned machineries which are so complex, that the tiniest modification in any part inevitably breaks something somewhere else.

The purpose of this book is to take a step towards science. The objective is to formalise the most common speech coding tools into a scientific framework, such that their strengths and weaknesses can be assessed in isolation, without the complex interconnections of a full-scale codec. The hope is that this work will give a stable scientific framework which allows for new innovations and paves the way for new break-through technologies.

The emphasis in this book is in understanding *why* and *how* the commonly used methods work. Specific details about standards have for the most parts been intentionally left out to keep the overall structure clear, but ample references are provided for the experts.

The book was originally created as a compendium for the course “Speech Coding” which I taught at the Friedrich-Alexander University Erlangen-Nürnberg (FAU). It is therefore designed to be useful as learning material for students working on a master’s degree in speech and audio processing, as well as information or communications technology. In parallel, the book is also meant to provide in-depth information for engineers and scientists working in the field and can thus be used as a handbook on speech coding.

Erlangen, Germany  
October 2016

Tom Bäckström

# Acknowledgements

What I know about practical methods for speech coding I learned in most parts during my time at the International Audio Laboratories Erlangen<sup>1</sup> from 2008 until 2016. Here, I had the privilege of being closely involved in the standardisation of both the MPEG Unified Speech and Audio Codec (standardised 2012) as well as the 3GPP Enhanced Voice Services speech codec (standardised 2014). This environment provided me demanding development tasks in speech coding but simultaneously sufficient time to learn and develop theory as well. With this book I wish to show my gratefulness for my colleagues at the AudioLabs, FAU as well as Fraunhofer.

I have had the pleasure of working on this book with four co-authors, Sascha Disch, Jérémié Lecomte, Christian Uhle and Guillaume Fuchs, each submitting a guest-authored chapter. All of them are the best experts in their respective fields, complementing those areas where I have less experience. I am grateful for their contributions and I think it is obvious that their participation made the book much better and a more balanced whole.

In writing this book, I have received helpful comments and reviews from Tim Fingscheidt, Christian Fischer Pedersen, Goran Markovic, as well as many students of my courses and other colleagues at the AudioLabs. I remain indebted to you all.

Disclaimer: A chapter of this book was removed from the published version due to a claim that it infringes a non-disclosure agreement.

Erlangen, Germany  
October 2016

Tom Bäckström

---

<sup>1</sup>International Audio Laboratories Erlangen is a joint institute between Fraunhofer Institute of Integrated Circuits (IIS) and Friedrich-Alexander University Erlangen-Nürnberg (FAU).



# Contents

|  |  |    |
|--|--|----|
| <b>1</b>   | <b>Introduction</b> . . . . .  | 1  |
| 1.1  | Objectives of Speech Coding . . . . .  | 1  |
| 1.2  | Perceptual Quality . . . . .   | 3  |
| 1.3  | Speech Signals . . . . .   | 4  |
| 1.4  | Perceptual and Source Models . . . . .   | 5  |
|  | References. . . . .  | 8  |
| <br><b>Part I Basic Properties of Speech Signals</b> |  |    |
| <b>2</b>   | <b>Speech Production and Modelling</b> . . . . .                                     | 11 |
| 2.1  | Introduction . . . . .   | 11 |
| 2.2  | Physiology and Articulation . . . . .  | 11 |
| 2.3  | Phonemes . . . . .   | 16 |
| 2.3.1  | Vowels . . . . .   | 16 |
| 2.3.2  | Consonants . . . . .   | 18 |
| 2.4  | Intonation, Rhythm and Intensity . . . . .   | 19 |
| 2.5  | Vocal Tract Models . . . . .   | 20 |
| 2.6  | Models of the Glottal Excitation. . . . .  | 22 |
| 2.7  | Obstruent Modelling. . . . .   | 24 |
| 2.8  | Nasal Cavities. . . . .  | 25 |
| 2.9  | Lips . . . . .   | 25 |
| 2.10   | System Modelling. . . . .  | 26 |
|  | References. . . . .  | 29 |
| <b>3</b>   | <b>Principles of Entropy Coding with Perceptual<br/>Quality Evaluation</b> . . . . . | 31 |
| 3.1  | Source and Entropy Coding . . . . .  | 31 |
| 3.2  | Quantisation . . . . .   | 34 |

|                               |  |           |
|-------------------------------|--|-----------|
| 3.3                           | Rate-Distortion Optimisation . . . . .                           | 36        |
| 3.4                           | Evaluation Models and Bitrate . . . . .                          | 37        |
|                               | Appendix . . . . .   | 39        |
|                               | References. . . . .  | 44        |
| <br><b>Part II Core Tools</b> |  |           |
| <b>4</b>                      | <b>Spectral Envelope and Perceptual Masking Models . . . . .</b> | <b>47</b> |
| 4.1                           | Introduction . . . . .   | 47        |
| 4.2                           | Linear Predictive Coding . . . . .                               | 49        |
| 4.2.1                         | Definition of the Linear Predictive Model . . . . .              | 51        |
| 4.2.2                         | Estimation of the Autocorrelation and Windowing . . . . .        | 53        |
| 4.2.3                         | Pre-processing Tools . . . . .                                   | 56        |
| 4.3                           | Perceptual Modelling . . . . .                                   | 59        |
| 4.4                           | Quantisation and Coding of Linear Predictive Models . . . . .    | 61        |
| 4.4.1                         | Representations . . . . .  | 61        |
| 4.4.2                         | Evaluation of Quantisation Quality . . . . .                     | 64        |
| 4.4.3                         | Vector Quantisation and Coding . . . . .                         | 65        |
| 4.4.4                         | Interpolation of Envelope Models . . . . .                       | 68        |
| 4.5                           | Summary and Discussion . . . . .                                 | 69        |
|                               | Appendix . . . . .   | 70        |
|                               | References. . . . .  | 75        |
| <b>5</b>                      | <b>Windowing and the Zero Input Response . . . . .</b>           | <b>77</b> |
| 5.1                           | Introduction and Motivation . . . . .                            | 77        |
| 5.2                           | Filter Windowing . . . . .                                       | 79        |
| 5.2.1                         | Modelling a Stationary Signal . . . . .                          | 80        |
| 5.2.2                         | Adaptive Modelling. . . . .                                      | 81        |
| 5.3                           | Source Model Windowing . . . . .                                 | 82        |
| 5.4                           | Perceptual Windowing . . . . .                                   | 83        |
| 5.5                           | Evaluation and Minimisation of Perceptual Error . . . . .        | 84        |
| 5.6                           | Frames and Subframes . . . . .                                   | 87        |
| 5.7                           | Summary . . . . .  | 88        |
|                               | Appendix . . . . .   | 88        |
|                               | References. . . . .  | 90        |
| <b>6</b>                      | <b>Fundamental Frequency . . . . .</b>                           | <b>91</b> |
| 6.1                           | Source Modelling . . . . .                                       | 91        |
| 6.2                           | Long-Term Prediction. . . . .                                    | 92        |
| 6.3                           | Codebook Formulation . . . . .                                   | 92        |
| 6.4                           | Stability . . . . .  | 93        |
| 6.5                           | Lag Estimation . . . . .   | 94        |
| 6.6                           | Lag Quantisation. . . . .  | 95        |
|                               | References. . . . .  | 96        |

- 7 Residual Coding** . . . . . 97
  - 7.1 Background . . . . . 97
  - 7.2 Codebook Design . . . . . 98
  - 7.3 Algebraic Codebooks . . . . . 101
    - 7.3.1 Quantisation and Codebook Design. . . . . 102
    - 7.3.2 Codebook Search. . . . . 104
    - 7.3.3 Encoding . . . . . 107
  - 7.4 Other Codebooks . . . . . 113
  - References. . . . . 115
- 8 Signal Gain and Harmonics to Noise Ratio** . . . . . 117
  - 8.1 Basic Gain Quantisation . . . . . 117
  - 8.2 Harmonics to Noise Ratio. . . . . 119
  - Reference . . . . . 120

**Part III Advanced Tools and Extensions**

- 9 Pre- and Postfiltering** . . . . . 123
  - 9.1 Introduction . . . . . 123
  - 9.2 Formant Enhancement and Excitation Shaping . . . . . 124
  - 9.3 Bass Postfilter. . . . . 126
  - 9.4 Pitch Sharpening. . . . . 127
  - 9.5 Pulse Dispersion and Phase Scrambling . . . . . 128
  - References. . . . . 129
- 10 Frequency Domain Coding**. . . . . 131
  - 10.1 Introduction . . . . . 131
  - 10.2 Overlapping Windowing. . . . . 133
  - 10.3 Windowing with Critical Sampling. . . . . 137
  - 10.4 Time-Frequency Transforms . . . . . 140
  - 10.5 Perceptual Modelling and Quantisation. . . . . 143
  - 10.6 Entropy Coding and the Rate Loop . . . . . 144
  - 10.7 Overall Codec Structure and Integration. . . . . 147
  - 10.8 Summary . . . . . 149
  - References. . . . . 150
- 11 Bandwidth Extension** . . . . . 151
  - 11.1 Introduction . . . . . 151
  - 11.2 Bandwidth Extension Basics. . . . . 152
  - 11.3 Bandwidth Extension Implementation. . . . . 153
    - 11.3.1 Bandwidth Extension Encoder. . . . . 153
    - 11.3.2 Bandwidth Extension Decoder. . . . . 154

|           |  |            |
|-----------|--|------------|
| 11.4      | Common Bandwidth Extension Methods . . . . .                               | 156        |
| 11.4.1    | Spectral Band Replication . . . . .  | 157        |
| 11.4.2    | Speech Bandwidth Extension . . . . .                                       | 157        |
| 11.4.3    | Intelligent Gap Filling . . . . .  | 158        |
| 11.5      | Conclusion . . . . .   | 158        |
|           | References. . . . .  | 159        |
| <b>12</b> | <b>Packet Loss and Concealment . . . . .</b>                               | <b>161</b> |
| 12.1      | Introduction . . . . .   | 161        |
| 12.1.1    | Error Types . . . . .  | 162        |
| 12.1.2    | Methods for Concealment of Errors in Multimedia<br>Transmissions . . . . . | 163        |
| 12.2      | Sender Based Concealment. . . . .  | 163        |
| 12.2.1    | Interleaving . . . . .   | 164        |
| 12.2.2    | Forward Error Correction . . . . .   | 165        |
| 12.2.3    | Multiple Description Coding and Randomised<br>Quantisation . . . . .       | 168        |
| 12.3      | Receiver Based Methods. . . . .  | 168        |
| 12.3.1    | Concealment . . . . .  | 169        |
| 12.3.2    | Recovery . . . . .   | 176        |
| 12.3.3    | Jitter Buffer . . . . .  | 178        |
| 12.4      | Interactive Methods . . . . .  | 180        |
| 12.4.1    | Retransmission . . . . .   | 181        |
| 12.4.2    | Adaptive or Channel-Aware Modes. . . . .                                   | 182        |
| 12.5      | Conclusion . . . . .   | 182        |
|           | References. . . . .  | 183        |
| <b>13</b> | <b>Voice Activity Detection . . . . .</b>                                  | <b>185</b> |
| 13.1      | Introduction . . . . .   | 185        |
| 13.1.1    | Applications . . . . .   | 185        |
| 13.1.2    | Requirements. . . . .  | 187        |
| 13.2      | Methods . . . . .  | 188        |
| 13.2.1    | Feature Extraction . . . . .   | 188        |
| 13.2.2    | Detection, Classification and Supervised Learning. . . . .                 | 192        |
| 13.2.3    | Example Implementations . . . . .  | 195        |
| 13.3      | Evaluation. . . . .  | 197        |
| 13.3.1    | Performance Measures. . . . .  | 198        |
| 13.3.2    | Challenges of VAD. . . . .   | 198        |
|           | References. . . . .  | 201        |
| <b>14</b> | <b>Relaxed Code-Excited Linear Prediction (RCELP). . . . .</b>             | <b>205</b> |
| 14.1      | Generalised Analysis-by-Synthesis . . . . .                                | 205        |
| 14.2      | RCELP for Voiced Coding . . . . .  | 207        |
| 14.2.1    | Pitch Estimation and Pulse Localisation . . . . .                          | 208        |
| 14.2.2    | Continuous Linear Delay Contour. . . . .                                   | 209        |

- 14.2.3 Signal Modification . . . . . 210
- 14.2.4 Pitch Coding and Subsequent Coding . . . . . 211
- 14.3 Parametric Unvoiced Coding . . . . . 212
- 14.4 RCELP in Standards. . . . . 214
- References. . . . . 215

**Part IV Standards and Specifications**

- 15 Quality Evaluation . . . . . 219**
  - 15.1 Introduction . . . . . 219
    - 15.1.1 Applications . . . . . 219
    - 15.1.2 Aspects of Quality. . . . . 220
    - 15.1.3 Test Material. . . . . 222
  - 15.2 Subjective Evaluation. . . . . 222
    - 15.2.1 ITU-T Recommendation P.800 . . . . . 224
    - 15.2.2 Recommendation ITU-T Recommendation P.835 . . . . . 228
    - 15.2.3 Recommendation ITU-R BS.1534 (MUSHRA). . . . . 228
  - 15.3 Objective Evaluation. . . . . 231
    - 15.3.1 ITU-T Recommendation P.862 (PESQ). . . . . 232
    - 15.3.2 ITU-T Recommendation P.863 (POLQA) . . . . . 233
    - 15.3.3 ITU-R Recommendation BS.1387 (PEAQ) . . . . . 234
  - References. . . . . 235
- Index . . . . . 237**

# About the Authors

**Tom Bäckström, D.Sc. (tech)**, is now professor at Aalto University, in Helsinki, Finland, but during the preparation of this book he was professor at International Audio Laboratories Erlangen, which is a co-operation unit between the Friedrich-Alexander University Erlangen-Nürnberg (FAU) and Fraunhofer Institute of Integrated Circuits (IIS), both in Erlangen, Germany.

**Sascha Disch, Dr.-Ing.**, is a senior scientist at International Audio Laboratories Erlangen, Fraunhofer IIS.

**Jérémie Lecomte** was a researcher at Fraunhofer IIS during the preparation of this book, but is now at Amazon.

**Christian Uhle, Dr.-Ing.**, is a senior scientist at International Audio Laboratories Erlangen, Fraunhofer IIS.

**Guillaume Fuchs, Ph.D.**, is a senior scientist at International Audio Laboratories Erlangen, Fraunhofer IIS.

# Acronyms

Throughout this work, acronyms and abbreviations have been avoided whenever possible. The motivation is that since many abbreviations are known only within a very specific professional field, using them makes it difficult for the uninitiated to join or even follow a discussion. Since acronyms are, however, still used frequently among engineers and scientists, the following list of acronyms is provided as a reference.

|         |   |
|---------|---|
| 3GPP    | 3rd generation partnership project (standardisation organisation)   |
| AbS     | Analysis by synthesis (optimisation method)   |
| ACELP   | Algebraic code-excited linear prediction (speech codec)   |
| ACR     | Absolute category rating (subjective evaluation method)   |
| ADPCM   | Adaptive pulse code modulation (quantisation method with source model)                                    |
| AMR     | Adaptive multirate (speech coding standard)   |
| AMR-NB  | Adaptive multirate narrowband (speech coding standard)  |
| AMR-WB  | Adaptive multirate wideband (speech coding standard)  |
| AMR-WB+ | Adaptive multirate wideband Plus (speech coding standard)   |
| ARQ     | Automatic retransmission request (method for concealment of lost packets)                                 |
| BBWE    | Blind bandwidth extension (advanced speech and audio coding tool)   |
| BPF     | Bass postfilter (advanced speech and audio coding tool)   |
| BWE     | Bandwidth extension (advanced speech and audio coding tool)   |
| CBR     | Constant bitrate (transmission rate)  |
| CCR     | Comparison category rating (subjective evaluation method)   |
| CELP    | Code-excited linear prediction (speech codec)   |
| CNG     | Comfort noise generator (advanced speech coding tool)   |
| DCR     | Degradation category rating (subjective evaluation method)  |
| DPCM    | Differential pulse code modulation (quantisation method which assumes that signal has low-pass character) |

|        |   |
|--------|---|
| DTX    | Discontinuous transmission (coding mode)  |
| ERB    | Equivalent rectangular bandwidth (frequency scale corresponding to perceptual accuracy)   |
| EVS    | 3GPP Enhanced voice services (speech coding standard)   |
| FEC    | Forward error correction (method for concealment of lost packets)   |
| FIR    | Finite impulse response (filter)  |
| G.718  | Frame error robust narrowband and wideband embedded variable bitrate coding of speech and audio from 8-32 kbit/s (speech coding standard) |
| HNR    | Harmonics to noise ratio (measure of signal characteristic)   |
| IID    | Independent and identically distributed (random variable)   |
| IIR    | Infinite impulse response (filter)  |
| IPA    | International phonetic alphabet (phonetic notation system)  |
| ISF    | Immittance spectrum frequency (representation)  |
| ISP    | Immittance spectrum pair (polynomials)  |
| ITU    | International telecommunication union (standardisation organisation)  |
| ITU-D  | ITU's Development sector (standardisation organisation)   |
| ITU-R  | ITU's Radio communication sector (standardisation organisation)   |
| ITU-T  | ITU's Telecommunication standardisation sector (standardisation organisation)   |
| LPC    | Linear predictive coding (method)   |
| LSD    | Log-spectral distortion (distance measure)  |
| LSF    | Line spectrum frequency (representation)  |
| LSP    | Line spectrum pair (polynomials)  |
| LTP    | Long-term prediction (method)   |
| MDCT   | Modified discrete cosine transform (lapped time-frequency transform)  |
| MPEG   | Moving picture experts group (standardisation organisation)   |
| MUSHRA | Multiple stimuli with hidden reference and anchor (subjective quality evaluation method)  |
| PCM    | Pulse code modulation (quantisation method)   |
| PEAQ   | Method for objective measurements of perceived audio quality (objective quality evaluation method)  |
| PESQ   | Perceptual evaluation of speech quality (objective quality evaluation method)   |
| PLC    | Packet loss concealment (advanced speech coding tool)   |
| POLQA  | Perceptual objective listening quality assessment (objective quality evaluation method)   |
| RCELP  | Relaxed code-excited linear prediction (low-bitrate speech codec)   |
| SAMPA  | Speech assessment methods phonetic alphabet (phonetic notation system)  |
| SBR    | Spectral band replication (method for bandwidth extension)  |



|      |  |
|------|--|
| USAC | MPEG Unified speech and audio codec (speech coding standard) |
| VAD  | Voice activity detection (advanced speech coding tool)       |
| VBR  | Variable bitrate (transmission rate)                         |
| ZIR  | Zero input response (part of the filter windowing method)    |