

The ISO/MPEG Unified Speech and Audio Coding Standard – Consistent High Quality for all Content Types and at all Bit Rates

MAX NEUENDORF,¹ AES Member, MARKUS MULTRUS,¹ AES Member, NIKOLAUS RETTELBACH,¹ GUILLAUME FUCHS,¹ JULIEN ROBILLIARD,¹ JÉRÉMIE LECOMTE,¹ STEPHAN WILDE,¹ STEFAN BAYER,¹⁰ AES Member, SASCHA DISCH,¹ CHRISTIAN HELMRICH,¹⁰ ROCH LEFEBVRE,² AES Member, PHILIPPE GOURNAY,² BRUNO BESSETTE,² JIMMY LAPIERRE,² AES Student Member, KRISTOFER KJÖRLING,³ HEIKO PURNHAGEN,³ AES Member, LARS VILLEMOES,³ AES Associate Member, WERNER OOMEN,⁴ AES Member, ERIK SCHUIJERS,⁴ KEI KIKURI,⁵ TORU CHINEN,⁶ TAKESHI NORIMATSU,¹ KOK SENG CHONG,⁷ EUNMI OH,⁸ AES Member, MIYOUNG KIM,⁸ SCHUYLER QUACKENBUSH,⁹ AES Fellow, AND BERNHARD GRILL¹

¹*Fraunhofer Institute for Integrated Circuits IIS, Erlangen, 91058 Germany*

²*Université de Sherbrooke/VoiceAge Corp., Sherbrooke, QC, J1K 2R1, Canada*

³*Dolby Sweden, 113 30, Stockholm, Sweden*

⁴*Philips Research Laboratories, Eindhoven, 5656AE, The Netherlands*

⁵*NTT DOCOMO, INC., Yokosuka, Kanagawa, 239-8536, Japan*

⁶*Sony Corporation, Shinagawa, Tokyo, 141-8610, Japan*

⁷*Panasonic Corporation*

⁸*Samsung Electronics, Suwon, Korea*

⁹*Audio Research Labs, Scotch Plains, NJ, 07076, USA*

¹⁰*International Audio Laboratories Erlangen, 91058 Germany*

In early 2012 the ISO/IEC JTC1/SC29/WG11 (MPEG) finalized the new MPEG-D Unified Speech and Audio Coding standard. The new codec brings together the previously separated worlds of general audio coding and speech coding. It does so by integrating elements from audio coding and speech coding into a unified system. The present publication outlines all aspects of this standardization effort, starting with the history and motivation of the MPEG work item, describing all technical features of the final system, and further discussing listening test results and performance numbers which show the advantages of the new system over current state-of-the-art codecs.

0 INTRODUCTION

With the advent of devices that unite a multitude of functionalities, the industry has an increased demand for an audio codec that can deal equally well with all types of audio content including both speech and music at low bit rates. In many use cases, e.g., broadcasting, movies, or audio books, the audio content is not limited to only speech or only music. Instead, a wide variety of content must be processed including mixtures of speech and music. Hence, a unified audio codec that performs equally well on all types of audio content is highly desired. Even though the largest potential for improvements is expected at the lower

end of the bit rate scale, a unified codec requires, of course, to retain or even exceed the quality of presently available codecs at higher bit rates.

Audio coding schemes, such as MPEG-4 High Efficiency Advanced Audio Coding (HE-AAC) [1,2], are advantageous in that they show a high subjective quality at low bit rates for music signals. However, the spectral domain models used in such audio coding schemes do not perform equally well on speech signals at low bit rates.

Speech coding schemes, such as Algebraic Code Excited Linear Prediction (ACELP) [3], are well suited for representing speech at low bit rates. The time domain source-filter model of these coders closely follows the human

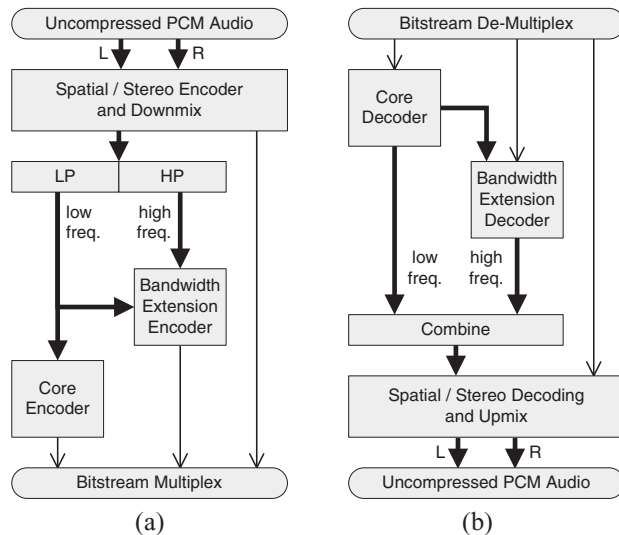


Fig. 1. General structure of a modern audio codec with a core codec accompanied by parametric tools for coding of bandwidth extension and stereo signals. USAC closely follows this coding paradigm. Fig. 1 (a) shows the encoder. Fig. 1(b) shows the corresponding decoder structure. Bold arrows indicate audio signal flow. Thin arrows indicate side information and control data.

speech production process. State-of-the-art speech coders, such as the 3GPP Adaptive Multi-Rate Wideband (AMR-WB) [4,5], perform very well for speech even at low bit rates but show a poor quality for music. Therefore, the source-filter model of AMR-WB was extended by transform coding elements in the 3GPP AMR-WB+ [6,7]. Still, for music signals AMR-WB+ is not able to provide an audio quality similar to that of HE-AAC(v2).

The following sections will first introduce the reader to the above-mentioned state-of-the-art representatives of modern audio and speech coding, HE-AACv2 and AMR-WB+. Re-iterating our contribution to the 132nd AES Convention last year [8], the ISO/IEC MPEG work item is described, followed by a technical description of the standardized coding system at a much higher level of detail than in previous publications [9,10]. Concluding, performance figures from the MPEG Verification Tests are presented and potential applications of the new technology are discussed.

1 STATE OF THE ART

1.1 General Codec Structure

Modern audio codecs and speech codecs typically exhibit a structure as shown in Fig. 1. This scheme consists of three main components: (1) a core-coder (i.e., transform or speech coder) that provides a high quality and largely wave-form preserving representation of low and intermediate frequency signal components; (2) a parametric bandwidth extension, such as Spectral Band Replication (SBR) [2], which reconstructs the high frequency band from replicated low frequency portions through the control of additional parameters; and (3) a parametric stereo coder, such as “Parametric Stereo” [1,11], which represents stereo signals by means of a mono downmix and a corresponding

set of spatial parameters. For low bit rates the parametric tools are able to reach much higher coding efficiency with a good quality / bit rate trade-off. At higher bit rates, where the core coder is able to handle a wider bandwidth and also discrete coding of multiple channels, the parametric tools can be selectively disabled. A general introduction to these concepts can be found in [12].

1.2 HE-AACv2

General transform coding schemes, such as AAC [13,1,2], rely on a sink model motivated by the human auditory system. By means of this psychoacoustic model, temporal and simultaneous masking is exploited for irrelevance removal. The resulting audio coding scheme is based on three main steps: (1) a time/frequency conversion; (2) a subsequent quantization stage, in which the quantization error is controlled using information from a psychoacoustic model; and (3) an encoding stage, in which the quantized spectral coefficients and corresponding side information are entropy-encoded. The result is a highly flexible coding scheme, which adapts well to all types of input signals at various operating points.

To further increase the coding efficiency at low bit rates, HE-AACv2 combines an AAC core in the low frequency band with a parametric bandwidth and stereo extension. Spectral Band Replication (SBR) [2] reconstructs the high frequency content by replicating the low frequency signal portions, controlled by parameter sets containing level, noise, and tonality parameters. “Parametric Stereo” [1,11] is capable of representing stereo signals by a mono downmix and corresponding sets of inter-channel level, phase, and correlation parameters.

1.3 AMR-WB+

Speech coding schemes, such as AMR-WB [4,5], rely on a source model motivated by the mechanism of human speech production. These schemes typically have three major components: (1) a short-term linear predictive coding scheme (LPC), which models the vocal tract; (2) a long-term predictor (LTP) or “adaptive codebook,” which models the periodicity in the excitation signal from the vocal chords; and (3) an “innovation codebook,” which encodes the non-predictable part of the speech signal. AMR-WB follows the ACELP approach that uses an algebraic representation for the innovative codebook: a short block of excitation signal is encoded as a sparse set of pulses and associated gain for the block. The pulse codebook is represented in algebraic form. The encoded parameters in a speech coder are thus: the LPC coefficients, the LTP lag and gain, and the innovative excitation. This coding scheme can provide high quality for speech signals even at low bit rates.

To properly encode music signals, in AMR-WB+ the time domain speech coding modes were extended by a transform coding mode for the excitation signal (TCX). The AMR-WB+ standard also features a low rate parametric high frequency extension as well as parametric stereo capabilities.

1.4 Other Audio Coding Schemes

For matters of completeness the reader should be pointed to other audio coding variants such as MPEG Spatial Audio Object Coding (SAOC), which addresses highly efficient audio coding based on separate object input and rendering [14]. Other examples focus on discriminating between speech and music signals and feeding the signal to specialized codecs depending on the outcome of the classification [15,16,17].

2 THE MPEG UNIFIED SPEECH AND AUDIO CODING WORK ITEM

Addressing the obvious need for an audio codec that can code speech and music equally well, ISO/IEC MPEG issued a Call for Proposal (CfP) on Unified Speech and Audio Coding (USAC) within MPEG-D [18] at the 82nd MPEG Meeting in October 2007. The responses to the Call were evaluated in an extensive listening test, with the result that the joint contribution from Fraunhofer IIS and VoiceAge Corp. was selected as reference model zero (RM0) at the 85th MPEG meeting in summer 2008 [10]. Even at that point the system fulfilled all requirements for the new technology, as listed in the CfP [19].

In the subsequent collaborative phase the RM0 based system was further refined and improved within the MPEG Audio Subgroup until early 2011, when the technical development was essentially finished. The mentioned improvements were introduced by following a well defined core experiment process. In this manner further enhancements from Dolby Labs., Philips, Samsung, Panasonic, Sony, and NTT Docomo were integrated into the system.

After technical completion of the standard, the MPEG Audio Subgroup conducted another comprehensive subjective Verification Test in summer 2011. The results of these tests are summarized in Section 4.

The standard reached International Standard (IS) stage in early 2012 by achieving a positive balloting vote from ISO's National Bodies voting for the standard [20].

3 TECHNICAL DESCRIPTION

3.1 System Overview

USAC preserves the same overall structure of HE-AACv2 as depicted in Fig. 1. An enhanced SBR (eSBR) tool serves as a bandwidth extension module, while MPEG Surround 2-1-2 supplies parametric stereo coding functionality. The core coder consists of an AAC based transform coder enhanced by speech coding technology.

Fig. 2 gives a more detailed insight into the workings of the USAC core decoder. Since in MPEG the encoder is not normatively specified, implementers are free to choose their own encoder architecture as long as it produces valid bitstreams. As a result, USAC provides complete freedom of encoder implementation and—just like any MPEG codec—permits continuous performance improvement even years after finalization of the standardization process.

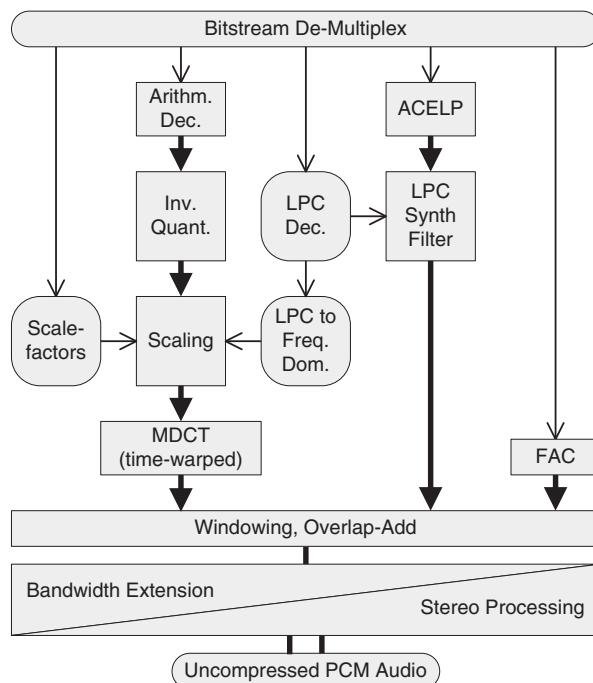


Fig. 2. Overview of USAC core decoder modules. The main decoding path features a Modified Discrete Cosine Transform (MDCT) domain coding part with scalefactor based or LPC based noise shaping. An ACELP path provides speech coder functionality. The Forward Aliasing Cancellation (FAC) enables smooth and flawless transitions between transform coder and ACELP. Following the core decoder, bandwidth extension and stereo processing is provided. Bold black lines indicate audio signal flow. Thin arrows indicate side information and control data.

USAC retains all capabilities of AAC. In Fig. 2 the left signal path resembles the AAC coding scheme. It comprises the function of entropy decoding (arithmetic decoder), inverse quantization, scaling of the spectral coefficients by means of scalefactors, and an inverse MDCT transform. With respect to the MDCT, all flexibility inherited from AAC regarding the choice of the transform window, such as length, shape, and dynamic switching is maintained. All AAC tools for discrete stereo or multichannel operation are included in USAC. As a consequence, USAC can be operated in a mode equivalent to AAC.

In addition, USAC introduces new technologies that offer increased flexibility and enhanced efficiency. The AAC Huffman decoder was replaced by a more efficient context-adaptive arithmetic decoder. The scalefactor mechanism as known from AAC can control the quantization noise shaping with a fine spectral granularity. If appropriate, it can be substituted by a Frequency Domain LPC Noise Shaping (FDNS) mechanism that consumes fewer bits. The USAC MDCT features a larger set of window lengths. The 512 and 256 MDCT block sizes complement the AAC 1024 and 128 sizes, providing a more suitable time-frequency decomposition for many signals.

3.2 Core-Coder

In the following subsections each of the technologies employed in the core coder are described in more detail.

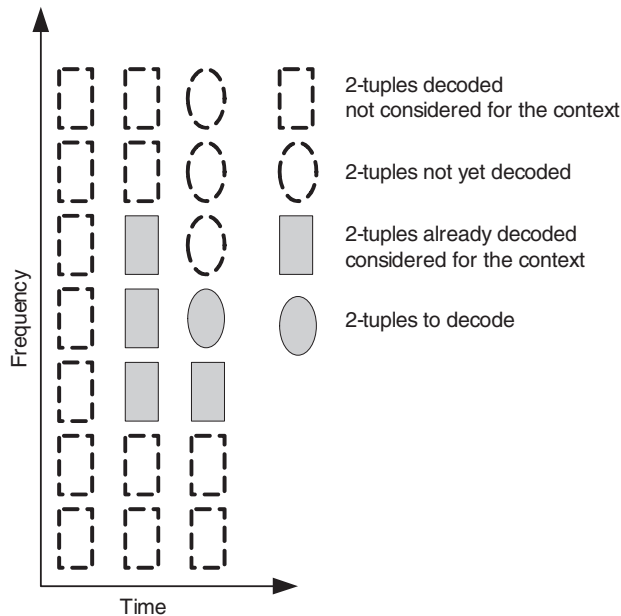


Fig. 3. Context Adaptive Arithmetic Coding.

3.2.1 Arithmetic Coder

In the transform-coder path, a context adaptive arithmetic coder is used for entropy coding of the spectral coefficients. The arithmetic coder works on pairs of two adjacent spectral coefficients (2-tuples). These 2-tuples are split into three parts: (1) the sign bits; (2) the two most significant bit-planes; (3) the remaining least significant bit-planes. For the coding of the two most significant bit-planes, one out of 64 cumulative frequency tables is selected. This selection is derived from a context, which is modeled by previously coded 2-tuples (see Fig. 3).

The remaining least significant bits are coded using one out of three cumulative frequency tables. This cumulative frequency table is chosen depending on the magnitude of the most significant bits in the two uppermost bit-planes.

The signs are transmitted separately at the end of the spectral data. This algorithm allows a saving from 3 to more than 6% of the overall bitrate over AAC Huffman coding while showing comparable complexity requirements [21].

3.2.2 Quantization Module

A scalar quantizer is used for the quantization of spectral coefficients. USAC supports two different quantization schemes, depending on the applied noise shaping: (1) a non-uniform quantizer is used in combination with scalefactor based noise shaping. The scalefactor based noise shaping is performed on the granularity of pre-defined scalefactor bands. To allow for an additional noise shaping within a scalefactor band, a power-law quantization scheme is used [22]. In this non-uniform quantizer the quantization intervals get larger with higher amplitude. Thus, the increase in signal-to-noise ratio with rising signal energy is lower than in a uniform quantizer. (2) A uniform quantizer is used in combination with LPC-based noise shaping. The LPC based noise shaping is able to model the spectral envelope continuously and without subdivision in fixed scalefactor

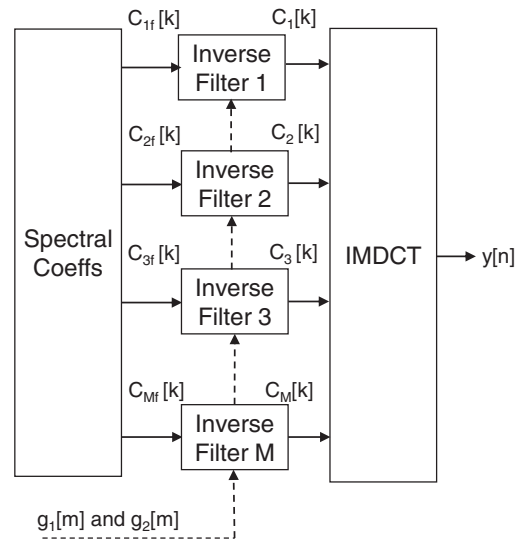


Fig. 4. FDNS processing.

bands. This alleviates the need for an extra intra-band noise shaping.

3.2.3 Noise Shaping Using Scalefactors or LPC

USAC relies on two tools to shape the coding noise when encoding the MDCT coefficients. The first tool is based on a perceptual model and uses a set of scalefactors applied to frequency bands. The second tool is based on linear predictive modeling of the spectral envelope combined with a first-order filtering of the transform coefficients that achieves both frequency-domain noise shaping and sample-by-sample time-domain noise shaping. This second noise shaping tool, called FDNS for Frequency-Domain Noise Shaping, can be seen as a combination of perceptual weighting from speech coders and Temporal Noise Shaping (TNS). Both noise shaping tools are applied in the MDCT domain. The scalefactor approach is more adapted to stationary signals because the noise shaping stays constant over the whole MDCT frame whereas FDNS is more adapted to dynamic signals because the noise shaping evolves smoothly over time. Since the perceptual model using scalefactors is already well documented [22], only FDNS is described below.

When LPC based coding is employed, one LPC filter is decoded for every window within a frame. Depending on the decoded mode, there may be one up to four LPC filters per frame, plus another filter when initiating LPC based coding. Using these LPC coefficients, FDNS operates as follows: for every window, the LPC parameters are converted into a set of $M = 64$ gains $g_k[m]$ in the frequency domain, defining a coarse spectral noise shape at the overlap point between two consecutive MDCT windows. Then, in each of the M bands, a first-order inverse filtering is performed on the spectral coefficients $C_{mf}[k]$, as shown in Fig. 4, to interpolate the noise level within the window boundaries noted as instants A and B in Fig. 5. Therefore, instead of the conventional LPC coefficient interpolation and time-domain filtering as done in speech codecs, the process of noise shaping is applied only in the

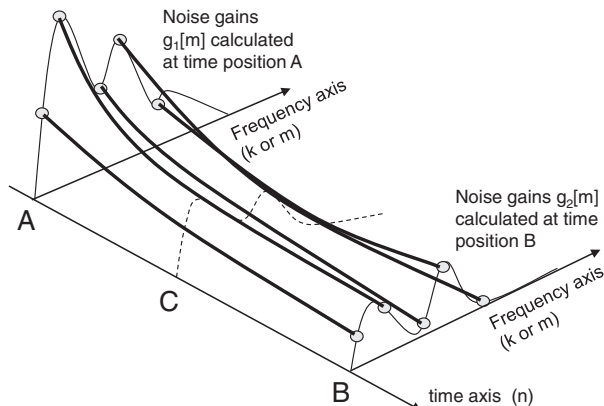


Fig. 5. Effect on noise shape of FDNS processing.

frequency domain. This provides two main advantages: first, the MDCT can be applied to the original signal (rather than the weighted signal as in speech coding), allowing proper time domain aliasing cancellation (TDAC) on the transition between scalefactor and LPC based noise shaping; and second, because of the “TNS-like” feature of FDNS, the noise shape is finely controlled on a sample-by-sample basis rather than on a frame-by-frame basis.

3.2.4 (Time-Warped) MDCT

The Modified Discrete Cosine Transform (MDCT) is well suited for harmonic signals with a constant fundamental frequency F_0 . In this case only a sparse spectrum with a limited number of relevant lines has to be coded. But when F_0 is rapidly varying, typically for voiced speech, the frequency modulation of the individual harmonic lines leads to a smeared spectrum and therefore a loss in coding gain. The Time Warped MDCT (TW-MDCT) [23] overcomes this problem by applying a variable resampling within one block prior to the transform. This resampling reduces or, ideally, completely removes the variation of F_0 . The reduction of this variation causes a better energy compaction of the spectral representation and consequently an increased coding gain compared to the classic MDCT. Furthermore, a careful adaptation of the window functions and of the average sampling frequency retain the perfect reconstruction property and the constant framing of the classic MDCT. The necessary warp information needed for the inverse resampling at the decoder is efficiently coded and part of the side information in the bitstream.

3.2.5 Windowing

In terms of windows and transform block sizes, USAC combines the well-known advantages of the 50% overlap MDCT windows of length 2048 and 256 (transform core of 1024 and 128) from AAC with the higher flexibility of TCX with additional transform sizes of 512 and 256. The long transform windows allow optimal coding of distinctly tonal signals, while the shorter windows with shorter overlaps allow coding of signals with an intermediate and highly varying temporal structure. With this set of windows the codec can adapt its coding mode much more closely to the

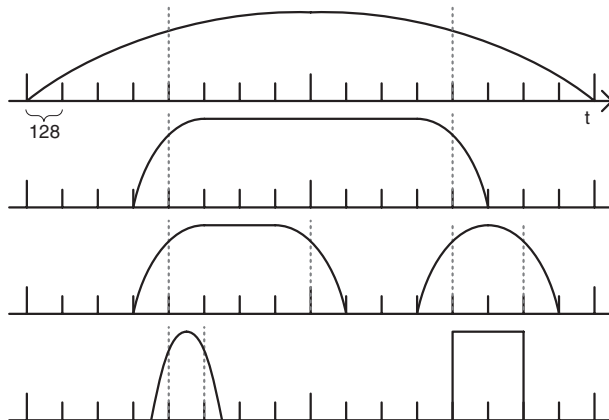


Fig. 6. Schematic overview over the allowed MDCT windows in USAC. Dotted lines indicate transform core boundaries. Bottom right shows ACELP window for reference. Transitional windows are not shown.

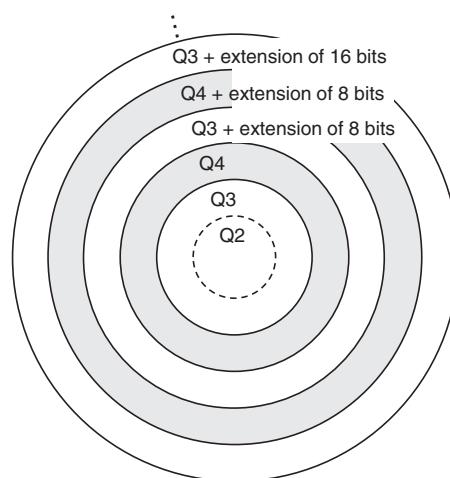


Fig. 7. Embedded structure of the AVQ quantizer.

signal than possible before. Fig. 6 shows the window and transform lengths and general shapes.

Similar to the start and stop windows of AAC, transitional windows accommodate the variation in transform length and coding mode [24]. In the special case of transitions to and from ACELP, the Forward Aliasing Cancellation takes effect (see Section 3.2.9).

Further flexibility is achieved by allowing a 768 sample based windowing scheme. In this mode all transform and window sizes are reduced to $\frac{3}{4}$ th of the above-mentioned numbers. This allows even higher temporal resolution, which is particularly useful in situations where the codec runs on a reduced core sampling rate. This mode is combined with an 8:3 QMF filterbank upsampling in eSBR (see Section 3.3.2) such that a higher audio bandwidth can be achieved at the same time.

3.2.6 Quantization of LPC coefficients

USAC includes a new variable bit rate quantizer structure for the LPC filter coefficients. Rather than using trained codebooks that are memory-consuming, an extremely memory-efficient 2-stage approach based on

algebraic vector quantization (AVQ, see Section 3.2.8) is used. An additional advantage of this approach is that the spectral distortion can be essentially maintained below a pre-set threshold by implicit bit allocation, thus making LPC quantization much less signal dependent. Another aspect of this variable bit rate quantizer is the application of LPC prediction within a frame. Specifically, if more than one LPC filter is transmitted within a frame, a subset of these filters are quantized differentially. This decreases significantly the bit consumption of the LPC quantizer in particular for speech signals. In this 2-stage quantizer, the first stage uses a small trained codebook as a first coarse approximation and the second stage uses a variable-rate AVQ quantizer in a split configuration (16-dimensional LPC coefficients quantized in 2 blocks of 8 dimensions).

3.2.7 ACELP

The time domain encoder in USAC is based on state-of-the-art ACELP speech compression technology. Several speech coding standards, in particular in cellular systems, integrate ACELP. The ACELP module in USAC uses essentially the same components as in AMR-WB+ [6] but with some improvements. LPC quantization was modified such that it is variable in bit rate (as described in Section 3.2.6). And the ACELP technology is more tightly integrated with other components of the codec. In ACELP mode, every quarter frame of 256 samples is split into 4 subframes of 64 samples (or for quarter frames of 192 samples it is split into 3 subframes of 64 samples). Using the LPC filter for that quarter frame (either a decoded filter or an interpolated filter depending on the position in a frame) each subframe is encoded as an excitation signal passed through the LPC filter. The excitation signal is encoded as the sum of two components: a pitch (or LTP) component (delayed, scaled version of the past excitation with properly chosen delay; also called adaptive codebook (ACB)) and an innovative component. The latter is encoded as a sparse vector formed by a series of properly placed non-zero impulses and corresponding signs and global gain. Depending on the available bit rate, the ACELP innovation codebook (ICB) size can be either of 12, 16, 20, 28, 36, 44, 52, or 64 bits. The more bits are spent for the codebook, the more impulses can be described and transmitted. Besides the LPC filter coefficients, the parameters transmitted in an ACELP quarter frame are:

Mean energy	2 bits
LTP pitch	9 or 6 bits
LTP filter	1 bit
ICB	12, 16, 20, 28, 36, 44, 52, or 64 bits
Gains	7 bits

All parameters are transmitted every subframe (every 64 samples), except the Mean energy which is transmitted once every ACELP quarter frame.

3.2.8 Algebraic Vector Quantization

Algebraic Vector Quantization (AVQ) is a structured quantization technique requiring very little memory and

is intended to quantize signals with uniform distribution. The AVQ tool used in USAC is another component taken from AMR-WB+. It is used to quantize LPC coefficients and FAC parameters (see Section 3.2.9).

The AVQ quantizer is based on the RE8 lattice [25], which has a nice densely packed structure in 8 dimensions. An 8-dimensional vector in RE8 can be represented by a so-called “leader” along with a specific permutation of the leader components. Using an algebraic process, a unique index for each possible permutation can be calculated. Leaders with statistical equivalence can be grouped together to form base codebooks that will define the layers of the indexing. Three base codebooks have been defined: Q2, Q3, and Q4 where indexing all permutations of the selected leaders consumes 8, 12, and 16 bits respectively. To extend the quantizer to even greater size, instead of continuing to add larger base codebooks (Q5 and over), a Voronoi extension has been added to extend algebraically the base codebook. With each additional 8 bits (1 bit per dimension), the Voronoi extension doubles the size of the codebook. Therefore Q3 and Q4 extended by a factor of 2 will use 20 and 24 bits respectively, and for a factor of 4, they will use 28 and 32 bits respectively. Hence, although the first layer (Q2) requires 8 bits, each additional layer in the AVQ tool adds 4 bits to the indexing (1/2 bit resolution). It should be noted that Q2 is a subset of Q3. In the USAC bitstream, the layer number (Qn) is indexed separately using an entropy code since small codebooks are more probable than large codebooks.

3.2.9 Transition Handling

The USAC core combines two domains of quantization, the frequency domain, which uses MDCT with overlapped windows, and the time domain, which uses ACELP with rectangular non-overlapping windows. To compute the synthesis signal, a decoded MDCT frame relies on TDAC of adjacent windows whereas the decoded ACELP excitation uses the LPC filtering. To handle transitions in an effective way between the two modes, a new tool, called “Forward Aliasing Cancellation” (FAC) has been developed. This tool “Forwards” to the decoder the “Aliasing Cancellation” data required to retrieve the signal from the MDCT frame usually accomplished by TDAC. Hence, at transitions between the two domains, additional parameters are transmitted, decoded, and processed to obtain the FAC synthesis as shown in Fig. 8. To recover the complete decoded signal, the FAC synthesis is merely combined with the windowed output of the MDCT. In the specific case of transitions from ACELP to MDCT, the ACELP synthesis and following zero-input response (ZIR) of the LPC filter is windowed, folded, and used as a predictor to reduce the FAC bit consumption.

3.3 Enhanced SBR Bandwidth Extension

3.3.1 Basic Concept of SBR

The Spectral Band Replication (SBR) technology was standardized in MPEG-4 in 2003, as an integral part of High Efficiency AAC (HE-AAC). The tool is a high frequency reconstruction tool that operates on a core coder signal and

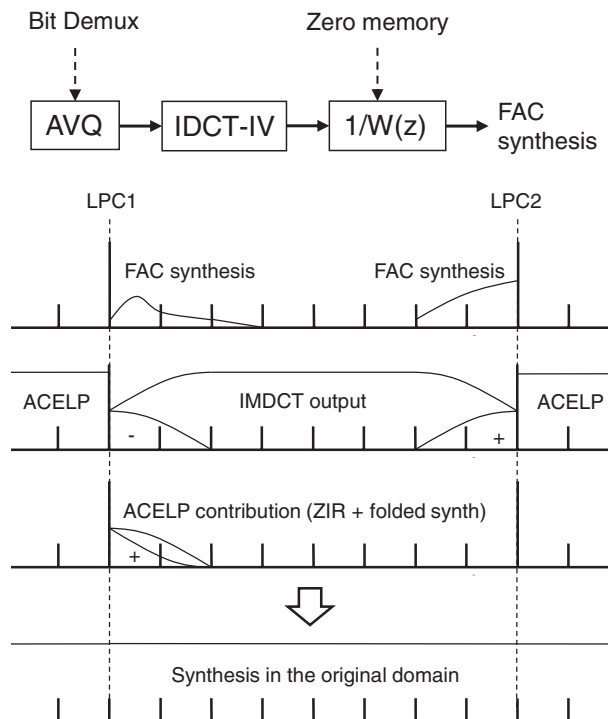


Fig. 8. Forward Aliasing Cancellation applied at transitions between ACELP and MDCT-encoded modes.

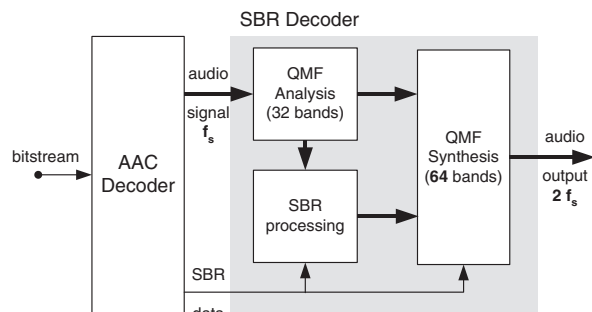


Fig. 9. Basic outline of SBR as used in MPEG-4 in combination with AAC.

extends the bandwidth of the output based on the available lowband signal and control data from the encoder. The principle of SBR and HE-AAC is elaborated on in [2,26, 27].

The SBR decoder operating on the AAC as standardized in MPEG-4 is depicted in Fig. 9. The system shown is a dual rate system where the SBR algorithm operates in a QMF domain and produces an output of wider bandwidth than and twice the sampling rate of the core coded signal going into the SBR module.

The SBR decoder generates a high frequency signal by copy-up methods in the QMF domain as indicated in Fig. 10. An inverse filtering is carried out within each QMF subband in order to adjust the tonality of the subband in accordance with parameters sent from the encoder.

The high frequency regenerated signal is subsequently envelope adjusted based on time/frequency tiles of enve-

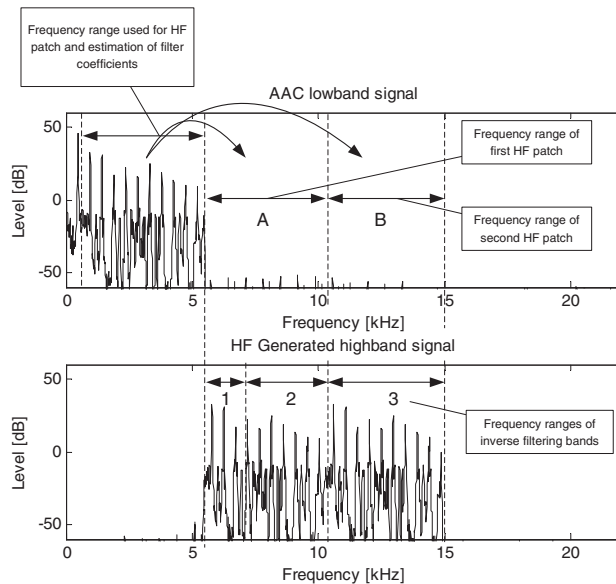


Fig. 10. Basic principle of copy-up based SBR as used in MPEG-4 in combination with AAC.

lope data transmitted from the encoder. During the envelope adjustment, additional noise and sinusoids are optionally added according to parametric data sent from the encoder.

3.3.2 Alternative Sampling Rate Ratios

MPEG-4 SBR was initially designed as a 2:1 system. Here, typically 1024 core coder samples are fed into a 32 band analysis QMF filterbank. The SBR tool performs a 2:1 upsampling in the QMF domain. After reconstructing the high frequency content, the signal is transformed back to time domain by means of a 64 band synthesis QMF filterbank. This results in 2048 time domain samples at twice the core coder sampling rate.

For USAC, the traditional 2:1 system was extended by two additional operating modes. First, to cope with low core coder sampling rates, which are usually used at very low bitrates, a variation of the SBR module similar as standardized in DRM (Digital Radio Mondiale) has been adopted into the USAC standard. In this mode, the 32 band analysis QMF filterbank is replaced by a 16 band QMF analysis filterbank. Hence, the SBR module is also capable of operating as a 4:1 system, where SBR runs at four times the core coder sampling rate. In this case, the maximum output audio bandwidth the system can produce at low sampling rates is increased by a factor of two compared to that of the traditional 2:1 system. This increase in audio bandwidth results in a substantial improvement in subjective quality at very low bit rates.

Second, USAC is also capable of operating in an 8:3 operating mode. In this case, a 24 band analysis QMF filterbank is used. In combination with a 768 core coder frame size, this mode allows for the best trade-off between optimal core-coder sampling rate and high temporal SBR resolution at medium bitrates, e.g., 24 kbit/s.

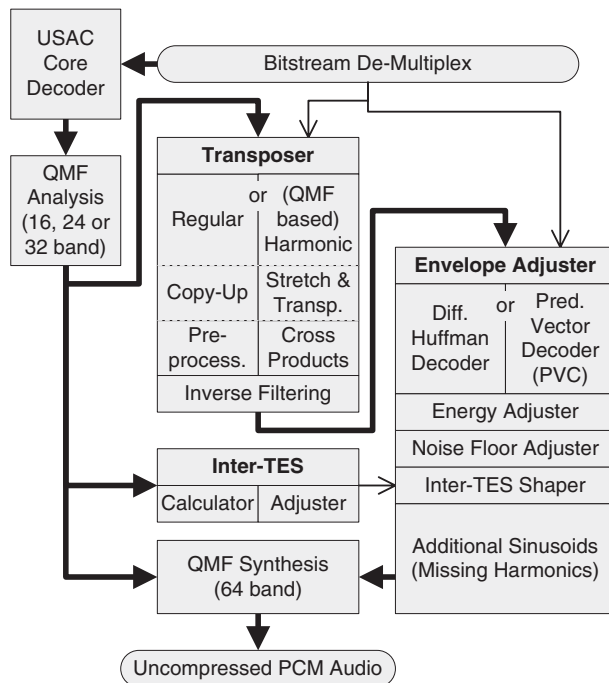


Fig. 11. Complete overview over the enhanced SBR of the USAC system. The figure shows the optional lower complexity QMF domain harmonic transposer, the PVC decoder, and the Inter-TES decoder modules.

3.3.3 Harmonic Transposition

In USAC a harmonic transposer of integer order T maps a sinusoid with frequency ω into a sinusoid with frequency $T\omega$, while preserving signal duration. This concept was originally proposed for SBR in [28], and the quality advantage over the frequency shift method, especially for complex stationary music signals, was verified in [29].

Three orders, $T = 2, 3, 4$, are used in sequence to produce each part of the desired output frequency range using the smallest possible transposition order. If output above the fourth order transposition range is required, it is generated by frequency shifts. When possible, near critically sampled baseband time domains are created for the processing to minimize computational complexity.

The benchmark quality transposer is based on a fine resolution sine windowed DFT. The algorithmic steps for $T = 2$ consist of complex filterbank analysis, subband phase multiplication by two, and filterbank synthesis with time stride twice of that of the analysis. The resulting time stretch is converted into transposition by a sampling rate change. The higher orders $T = 3, 4$ are generated in the same filterbank framework. For a given target subband, inputs from two adjacent source subbands are combined by interpolating phases linearly and magnitudes geometrically. Controlled by one bit per core coder frame, the DFT transposer adaptively invokes a frequency domain oversampling by 50% based on the transient improvement method of [30].

To allow the use of USAC in low-power applications such as portable devices, an alternate, low complexity, transposer that closely follows the bandwidth extension principle of the DFT transposer can be used as shown in Fig. 11. This

low complexity transposer operates in a QMF domain that allows for direct interfacing with the subsequent SBR processing. The coarse resolution QMF transposer suppresses intermodulation distortion by using overlapping block processing [31]. The finer time resolution of the QMF bank itself allows for a better transient response than that of the DFT without oversampling. Moreover, a geometrical magnitude weighting inside the subband blocks reduces potential time smearing.

The inherent spectral stretching of harmonic transposition can lead to a perceptual detachment of single overtones from periodic waveforms having rich overtone spectra. This effect can be attributed to the sparse overtone structure in the stretched spectral portions, since, e.g., a stretching by a factor of two only preserves every other overtone. This is mitigated by the addition of cross products. These consist of contributions from pairs of source subbands separated by a distance corresponding to the fundamental frequency [30]. The control data for cross products is transmitted once per core coder frame and consists of an on/off flag and seven bits indicating the fundamental frequency in the case that the flag is set.

3.3.4 Predictive Vector Coding

Adding the Predictive Vector Coding (PVC) scheme to the eSBR tool introduces a new coding scheme for the SBR spectral envelopes. Whereas in MPEG-4 SBR the spectral envelope is transmitted by means of absolute energies, PVC predicts the spectral envelope in high frequency bands from the spectral envelope in low frequency bands. The coefficient matrices for the prediction are coded using vector quantization. This improves the subjective quality of the eSBR tool, in particular for speech content at low bit rates. Generally, for speech signals, there is a relatively high correlation between the spectral envelopes of low frequency bands and high frequency bands, which can be exploited by PVC. The block diagram of the eSBR decoder including the PVC decoder is shown in Fig. 11.

The analysis and synthesis QMF banks and HF generator remain unchanged, but the HF envelope adjuster is modified to process the high frequency envelopes generated by the PVC decoder. In the PVC decoder, the high frequency envelopes are generated by multiplying a prediction coefficient matrix with the low frequency envelopes. A prediction codebook in the PVC decoder holds 128 coefficient matrices. An index of the prediction coefficient matrix that provides the lowest difference between predicted and actual envelopes is transmitted as a 7 bit value in the bitstream.

3.3.5 Inter-Subband-Sample Temporal Envelope Shaping (Inter-TES)

For transient input signals audible distortion (pre/post-echoes) can occur in the high frequency components generated by eSBR due to its limited temporal resolution. Although splitting a frame into several shorter time segments can avoid the distortion, this requires more bits for eSBR information. In contrast, Inter-TES can reduce the distortion with smaller number of bits by taking advantage

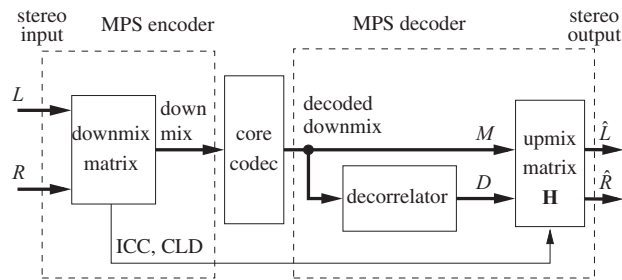


Fig. 12. Basic structure of the MPS 2-1-2 parametric stereo encoder and decoder used in USAC. Bold lines denote audio signal paths, whereas the thin arrow denotes the flow of parametric side information.

of the correlation between temporal envelopes in the low and high frequency bands. Inter-TES requires 1 bit for its activation and 2 bits for an additional parameter described below.

Fig. 11 shows the Inter-TES module as part of the eSBR block diagram. When Inter-TES is activated, the temporal envelope of the low frequency signal is first calculated, and the gain values are then computed by adjusting the temporal envelope of the low frequency signal according to a transmitted parameter. Finally, the gain values are applied to the transposed high frequency signal including noise components. As shown in Fig. 11, the shaped high frequency signal and the independent sinusoids are added if necessary, and then fed to the synthesis QMF bank in conjunction with the low frequency signal.

3.4 Stereo Coding

3.4.1 Discrete vs. Parametric Stereo Coding

There are two established approaches for coding of stereophonic audio signals. “Discrete stereo coding” schemes strive to represent the individual waveforms of each of the two channels of a stereo signal. They utilize joint stereo coding techniques such as mid/side (M/S) coding [32] to take inter-channel redundancy and binaural masking effects into account. “Parametric stereo coding” schemes [33,34, 35], on the other hand, are designed to represent the perceived spatial sound image of the stereo signal. They utilize a compact parametric representation of the spatial sound image that is conveyed as side information in addition to a mono downmix signal and used in the decoder to recreate a stereo output signal. Parametric stereo coding is typically used at low target bit rates, where it achieves a higher coding efficiency than discrete stereo coding. USAC extends, combines, and integrates these two stereo coding schemes, thus bridging the gap between them.

3.4.2 Parametric Stereo Coding with MPEG Surround 2-1-2

Parametric stereo coding in USAC is provided by an MPEG Surround 2-1-2 (MPS 2-1-2) downmix/upmix module that was derived from MPEG Surround (MPS) [11,36, 37]. The signal flow of the MPS 2-1-2 processing is depicted in Fig. 12.

At the encoder, MPS calculates a downmix signal and parameters that capture the essential spatial properties of the input channels. These spatial parameters, namely the inter-channel level differences (CLDs) and inter-channel cross-correlations (ICCs), are only updated at a relatively low time-frequency resolution based on the limits of the human auditory system to perceive spatial phenomena, thus requiring a bit rate of only a few kbit/s.

In the decoder, a decorrelated signal D , generated from the downmixed input signal M , is fed along with M into the upmixing matrix \mathbf{H} , as depicted in the right dashed box in Fig. 12. The coefficients of \mathbf{H} are determined by the parametric spatial side information generated in the encoder.

Stereo sound quality is enhanced by utilizing phase parameters in addition to CLDs and ICCs. It is well-known that inter-channel phase differences (IPDs) can play an important role in stereo image quality, especially at low frequencies [33]. In contrast to parametric stereo coding in MPEG-4 HE-AAC v2 [1], phase coding in USAC only requires the transmission of IPD parameters, since it has been shown that the overall phase differences parameters (OPDs) can be analytically derived from the other spatial parameters on the decoder side [38,39]. The USAC parametric stereo phase coding can handle anti-phase signals by applying an unbalanced weighting of the left and right channels during downmixing and upmixing processes. This improves stability for stereo signals where out-of-phase signal components would otherwise cancel each other in a simple mono downmix.

3.4.3 Unified Stereo Coding

In a parametric stereo decoder, the stereo signal is reconstructed by an upmix matrix from the mono downmix signal and a decorrelated version of the downmix, as shown in the right part of Fig. 12. MPS enhances this concept by optionally replacing parts of the decorrelated signal with a residual waveform signal. This ensures scalability up to the same transparent audio quality achievable by discrete stereo coding, whereas the quality of a parametric stereo coder without residual coding might be limited by the parametric nature of the spatial sound image description. Unlike MPS, where the residual signals are coded independently from the downmix signals, USAC tightly couples the coding of the downmix and residual signals.

As described above, in addition to CLD and ICC parameters, USAC also employs IPD parameters for coding the stereo image. The combination of parametric stereo coding involving IPD parameters and integrated residual coding is referred to as “unified stereo coding” in USAC. In order to minimize the residual signal, an encoder as shown in the left half of Fig. 13 is used. In each frequency band, the left and right signals L and R are fed into a traditional mid/side (i.e., sum/difference) transform. The resulting signals are gain normalized by a factor c . A prediction of the scaled difference signal is made by multiplication of the mid (i.e., sum) signal M with a complex-valued parameter α . Both c and α are a function of the CLD, ICC, and IPD

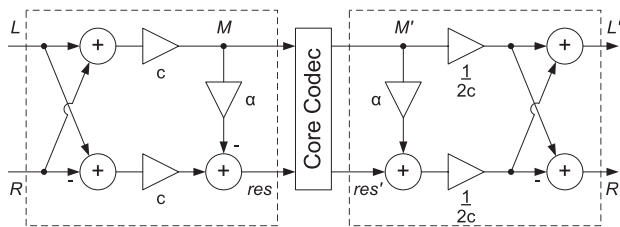


Fig. 13. Block diagram of unified stereo encoder (left) and decoder (right).

parameters. The resulting M and res signals are then fed into the 2-channel USAC core encoder that includes a correspondingly modified psychoacoustic model and can either encode the downmix and residual signal directly or can encode a mid/side transformed version known as “pseudo L/R” signal.

The decoder follows the inverse path, as depicted in the right half of Fig. 13. The optimal order of MPS and SBR processing in the USAC decoder depends on the bandwidth of the residual signal. If no or a bandlimited residual signal is used, it is advantageous to apply mono SBR decoding followed by MPS 2-1-2 decoding. At higher bit rates, where the residual signal can be coded with the same bandwidth as the downmix signal, it is beneficial to apply MPS 2-1-2 decoding prior to stereo SBR decoding.

3.4.4 Transient Steering Decorrelator

Applause signals are known to be a challenge for parametric stereo coding. In a simple model, applause signals can be thought of as being composed of a quasi-stationary noise-like background sound originating from the dense, far-off claps, and a collection of single, prominently exposed claps. Both components have very different properties that need to be addressed in the parametric upmix [40].

Upmixed applause signals usually lack spatial envelopment due to the insufficiently restored transient distribution and are impaired by temporally smeared transients. To preserve a natural and convincing spatio-temporal structure, a decorrelating technique is needed that can handle both of the extreme signal characteristics as described by the applause model. The Transient Steering Decorrelator (TSD) is an implementation of such a decorrelator [41]. TSD basically denotes a modification of the MPS 2-1-2 processing within USAC.

The block diagram of the TSD embedded in the upmix box of the MPS 2-1-2 decoder module is shown in Fig. 14. The mono downmix is split by a transient separation unit with fine temporal granularity into a transient signal path and a non-transient signal path. Decorrelation is achieved separately within each signal path through specially adapted decorrelators. The outputs of these are added to obtain the final decorrelated signal. The non-transient signal path M_1 utilizes the MPS 2-1-2 late-reverb-type decorrelator. The transient signal path M_2 comprises a parameter-controlled transient decorrelator. Two frequency independent parameters that entirely guide the TSD process are transmitted in

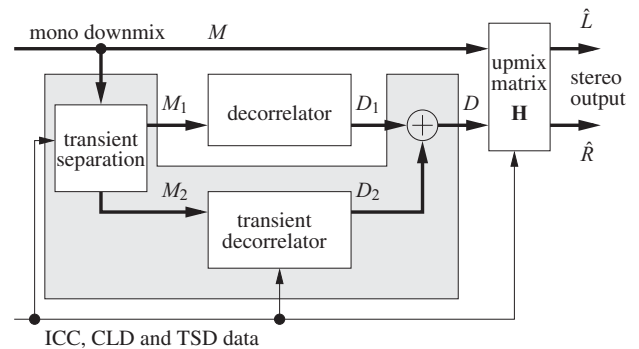


Fig. 14. TSD (highlighted by gray shading) within the MPS 2-1-2 module of the USAC decoder.

the TSD side information: a binary decision that controls the transient separation in the decoder and phase values that spatially steer the transients in the transient decorrelator. Spatial reproduction of transient events does not require fine spectral granularity. Hence, if TSD is active, MPS 2-1-2 may use broadband spatial cues to reduce side information.

3.4.5 Complex Prediction Stereo Coding

MPS 2-1-2 and unified stereo coding employ complex QMF banks, which are shared with the SBR bandwidth extension tool. At high bit rates, however, the SBR tool is typically not operated, while unified stereo coding would still provide an improved coding efficiency compared to traditional joint stereo coding techniques such as mid/side coding. In order to achieve this improved coding efficiency without the computational complexity caused by the QMF banks, USAC provides a complex prediction stereo coding tool [42] that operates directly in the MDCT domain of the underlying transform coder.

Complex prediction stereo coding applies linear predictive coding principles to minimize inter-channel redundancy in mid signal M and side signal S . The prediction technique is able to compensate for inter-channel phase differences as it employs a complex-valued representation of either M or S in combination with a complex-valued prediction coefficient α . The redundant coherent portions between M and S signal are minimized—and the signal compaction maximized—by subtracting from the smaller of the two a weighted and phase-adjusted version of the larger one—the downmix spectrum D —leading to a residual spectrum E . Downmix and residual are then perceptually coded and transmitted along with prediction coefficients. Fig. 15 shows the block diagram of a complex prediction stereo encoder and decoder, where L and R represent the MDCT spectra of the left and right channel, respectively.

The key here is to utilize a complex-valued downmix spectrum D obtained from a modulated complex lapped transform (MCLT) [43] representation for which the MDCT is the real part and whose imaginary part is the modified discrete sine transform (MDST). Given that in USAC, M and S are obtained via real-valued MDCTs, an additional real-to-imaginary (R2I) transform is required so that D can be constructed in both encoder and decoder [42]. In USAC, an efficient approximation of the R2I transform is utilized

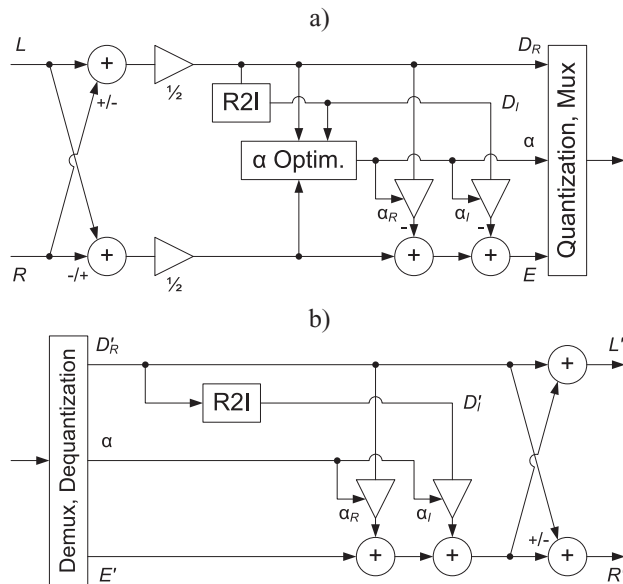


Fig. 15. Block diagram of complex prediction stereo encoder a) and decoder b).

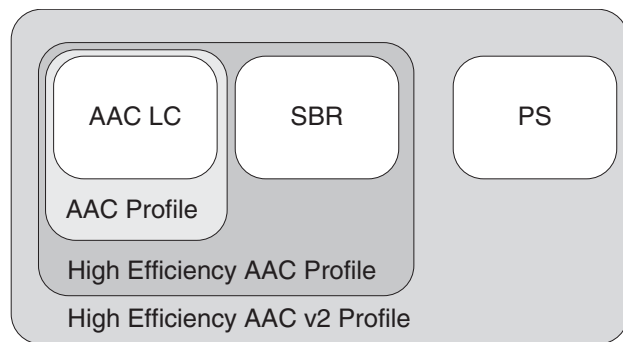


Fig. 16. The AAC family of profiles.

that operates directly in the frequency domain and does not increase the algorithmic delay of the coder.

3.5 System Aspects

3.5.1 Profiles

MPEG defines profiles as a combination of standardized tools. While all tools are always retained in the standard, the profile provides a subset or combination of tools that serve specific industry needs. Although there are many profiles defined in MPEG-4 Audio, the most successful and widely adopted ones are the “AAC family” of profiles, i.e., the “AAC Profile,” the “HE-AAC v1 Profile,” and the “HE-AAC v2 Profile.”

The AAC family of profiles, as outlined in Fig. 16, is hierarchical. The structure of the profiles ensures that (a) an AAC decoder plays AAC LC (Low Complexity), (b) an HE-AAC decoder plays AAC LC and SBR, and (c) an HE-AAC v2 decoder plays AAC LC, SBR, and PS.

In the MPEG USAC standard, two profiles are defined:

- 1) Extended HE-AAC Profile,
- 2) Baseline USAC Profile.

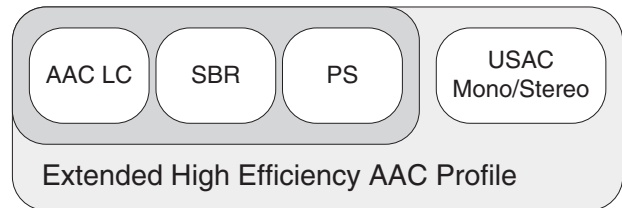


Fig. 17. The Extended HE-AAC Profile.

The Baseline USAC Profile contains the complete USAC codec except for the DFT transposer, the time-warped filterbank, and the MPS fractional delay decorrelator.

The Extended High Efficiency AAC Profile contains all of the tools of the High Efficiency AAC v2 Profile and is as such capable of decoding all AAC family profile streams. In order to exploit the consistent performance across content types at low bit rates, the profile additionally incorporates mono/stereo capability of the Baseline USAC Profile as outlined in Fig. 17.

As a result the Extended High Efficiency AAC Profile represents a natural evolution of one of the most successful families of profiles in MPEG Audio.

On the other hand, the Baseline USAC Profile provides a clear stand-alone profile for applications where a universally applicable codec is required but the capability of supporting the existing MPEG-4 AAC profiles is not relevant.

The worst case decoding complexity of both profiles is listed in Tables 1 and 2. The complexity numbers are indicated in terms of Processor Complexity Units (PCU) and RAM Complexity Units (RCU). PCUs are specified in MOPS and RCUs are expressed in kWords (1000 words).

Each profile typically consists of several levels. The levels are defined hierarchically and denote the worst case complexity for a given decoder configuration. A higher level indicates an increased decoder complexity, which goes along with support for a higher number of channels or a higher output sampling rate.

First implementations of an Extended High Efficiency AAC Profile decoder indicate comparable complexity and memory requirements as for High Efficiency AAC v2 Profile when operated at the same level.

3.5.2 Transport

The way of signaling and transport of the USAC payload is very similar to MPEG-4 HE-AACv2. As for HE-AACv2, the concept of a signaling within MPEG-4 Audio is supported. For this purpose, a new Audio Object Type (AOT) for USAC is defined within the MPEG-4 *AudioSpecificConfig*. The *AudioSpecificConfig* can also carry the *UsacConfig* data, which is needed to properly set up the decoder.

The mandatory explicit signaling of all USAC decoder tools, such as SBR and PS, avoids several problems of HE-AACv2. For the reason of backward compatibility to decoders not supporting SBR or Parametric Stereo, an implicit signaling was introduced in HE-AACv2. As a consequence, a decoder at start-up was not able to clearly determine output sampling rate, channel configuration, or number of

Table 1. Baseline USAC Profile processor and RAM complexity depending on decoder level.

Level	Max. channels	Max. sampling rate [kHz]	Max. PCU	Max. RCU
1	1	48	7	6
2	2	48	12	11
3	5.1	48	31	28
4	5.1	96	62	28

samples per frame. In contrast to HE-AACv2, a USAC decoder unambiguously determines its configuration by reading the *UsacConfig* data at start-up.

A set of *audioProfileLevelIndication* values allows for the signaling of the required decoder profile and level.

As for HE-AACv2, the frame-wise payload (*UsacFrame*) directly corresponds to MPEG-4 access units. In combination with the MPEG-4 signaling, a multitude of transport formats natively supports the carriage of USAC. For streaming applications, the use of, e.g., LATM/LOAS [1], IETF RFC 3016 [44], and RFC 3640 [45] is possible. For broadcasting applications, MPEG-2 transport stream [46] may be used. Finally, the MP4 and 3GPP file formats [47,48] provide support for file-based applications.

MP4 and 3GPP file format-based applications can now benefit from mandatory edit list support in USAC decoders to provide an exact time alignment: a decoder can reconstruct the signal with the exact starting and ending times, as compared to the original signal. Thus, additional samples at the beginning or end, introduced by frame-based processing and other buffering within the codec, are removed on the decoder side, ensuring, e.g., gapless playback.

3.5.3 Random Access

Various tools in USAC may exploit inter-frame correlation to reduce the bit demand. In SBR, MPS 2-1-2 and complex prediction stereo coding, time differential coding relative to the previous frame may be used. The arithmetic coder may refer to a context based on the previous frame. Though these techniques improve coding efficiency for the individual frames, they come at the cost of introducing a source of inter-frame dependencies. This means that a given frame may not be decoded without the knowledge of the previous frame.

In case of transmission over an error prone channel or in case of broadcasting where a continuously transmitted

stream is received and shall be decoded starting with a randomly received first frame, these inter-frame dependencies can make the tune-in phase challenging.

For the reasons listed above, USAC audio streams contain random access frames that can be decoded entirely independent from any previous frame (“independent frames”). The information whether a frame acts as an “independent frame” is conveyed in the first bit of the USAC frame and can be easily retrieved.

The frame independence is achieved by resetting the arithmetic coder context and forcing SBR, MPS 2-1-2, and complex prediction stereo coding to frequency-differential coding only. The independent frame serves as safe starting points for random access decoding, also after a frame loss.

In addition to the indication of independent frames, great importance was attached to an explicit signaling of potential core-coder frame dependent information. Wherever window size, window shape, or the need for FAC data is usually derived from the previous frame, this information can be unambiguously determined from the payload of any given independent frame.

4 PERFORMANCE

4.1 Listening Test Description

Three listening tests were performed to verify the quality of USAC. The objective of these verification tests was to confirm that the goals set out in the original Call for Proposals are met by the final standard [18,49]. ISO/IEC National Bodies could then take the test results as documented in the Verification Test report [50] into account when casting their final vote for USAC. Since the goal of USAC was the development of an audio codec that performs at least as good as the better of the best speech codec (AMR-WB+) and the best audio codec (HE-AACv2) around, the verification tests compared USAC to these codecs for mono

Table 2. Extended High Efficiency AAC Profile processor and RAM complexity depending on decoder level.

Level	Max. channels	Max. sampling rate [kHz]	Max. PCU	Max. RCU
1	n/a	n/a	n/a	n/a
2	2	48 (Note 1)	12	11
3	2	48 (Note 1)	15	11
4	5.1 (Note 2)	48	25	28
5	5.1 (Note 2)	96	49	28
6	7.1 (Note 2)	48	34	37
7	7.1 (Note 2)	96	67	53

Note 1: Level 2 and Level 3 differ for the decoding of HE-AACv2 bitstreams with respect to the max. AAC sampling rate in case Parametric Stereo data is present [1]. Note 2: USAC is limited to mono or stereo.

Table 3. Conditions for Test 1 (mono at low bitrates).

Condition	Label
Hidden reference	HR
Low pass anchor at 3.5 kHz ¹	LP3500
Low pass anchor at 7 kHz	LP7000
USAC at 8 kbit/s	USAC-8
USAC at 12 kbit/s	USAC-12
USAC at 16 kbit/s	USAC-16
USAC at 24 kbit/s	USAC-24
HE-AAC v2 at 12 kbit/s	HE-AAC-12
HE-AAC v2 at 24 kbit/s	HE-AAC-24
AMR-WB+ at 8 kbit/s	AMR-8
AMR-WB+ at 12 kbit/s	AMR-12
AMR-WB+ at 24 kbit/s	AMR-24

Table 4. Conditions for Test 2 (stereo at low bitrates).

Condition	Label
Hidden reference	HR
Low pass anchor at 3.5 kHz	LP3500
Low pass anchor at 7 kHz	LP7000
USAC at 16 kbit/s	USAC-16
USAC at 20 kbit/s	USAC-20
USAC at 24 kbit/s	USAC-24
HE-AAC v2 at 16 kbit/s	HE-AAC-16
HE-AAC v2 at 24 kbit/s	HE-AAC-24
AMR-WB+ at 16 kbit/s	AMR-16
AMR-WB+ at 24 kbit/s	AMR-24

and stereo at several bit rates. The results also provide the possibility to create a quality versus bit rate curve (a.k.a. rate-distortion curve) showing how the perceived quality of USAC progresses at different bit rates. The conditions included in each test are given in Tables 3 to 5.

Two further audio codecs were included as references: HE-AACv2 and AMR-WB+. In order to assess whether the new technology exceeds even the combined references codec performance, the concept of the Virtual Codec (VC) was introduced. The VC score is derived from the two reference codecs by choosing the better of the HE-AACv2 or AMR-WB+ score for each item at each operating point. Consequently the VC as a whole is always at least as good as the reference codecs and often better. It thus constitutes an even higher target to beat.

4.2 Test Items

Twenty-four test items were used in the test, consisting of eight items from each of three content categories: Speech, Speech mixed with Music, and Music. Test items were stereo signals sampled at 48 kHz and were approximately 8 seconds in duration. A large number of relatively short test items were used so that the items could encompass a greater diversity of content. Mono items were derived from the stereo items by either averaging left and right channel signals or taking only the left channel if averaging would result in significant comb filtering or phase cancellation.

¹ Bandlimited but keeping the same stereo width as the original (hidden reference)

Table 5. Conditions for Test 3 (stereo at high bitrates).

Condition	Label
Hidden reference	HR
Low pass anchor at 3.5 kHz	LP3500
Low pass anchor at 7 kHz	LP7000
USAC at 32 kbit/s	USAC-32
USAC at 48 kbit/s	USAC-48
USAC at 64 kbit/s	USAC-64
USAC at 96 kbit/s	USAC-96
HE-AAC v2 at 32 kbit/s	HE-AAC-32
HE-AAC v2 at 64 kbit/s	HE-AAC-64
HE-AAC v2 at 96 kbit/s	HE-AAC-96
AMR-WB+ at 32 kbit/s	AMR-32

Table 6. MUSHRA Subjective Scale.

Descriptor	Range
EXCELLENT	80 to 100
GOOD	60 to 80
FAIR	40 to 60
POOR	20 to 40
BAD	0 to 20

All items were selected to be challenging for all codecs under test.

4.3 Test Methodology

All tests followed the MUSHRA methodology [51] and were conducted in an acoustically controlled environment (such as a commercial sound booth) using reference quality headphones.

All items were concatenated to form a single file for processing by the systems under test. USAC processing was done using the Baseline USAC Profile encoder and decoder.

Fifteen test sites participated in the three tests. Of these, 13 test sites participated in test 1, 8 test sites in test 2, and 6 test sites in test 3. Listeners were post-screened and only those that showed consistent assessments were used in the statistical analysis. This post-screening consisted of checking whether, for a given listener in a given test, the Hidden Reference (HR) was always given a score larger than or equal to 90 and whether the anchors are scored monotonic ($LP3500 \leq LP7000 \leq HR$). Only the scores of listeners having met these two post-screening conditions were retained for statistical analysis. After post-screening, tests 1, 2 and 3 had 60, 40 and 25 listeners, respectively.

4.4 Test Results

Figs. 18 to 20 show the average absolute scores for each codec, including the VC, at the different operating points tested. Note that the scores between the tested points for a given codec are linearly interpolated in these Figures to show the trend of the quality/bit rate curve. The scores from all test sites, after listener post-screening, are pooled for this analysis. Vertical bars around each average score indicate the 95% confidence intervals using a t-distribution. The

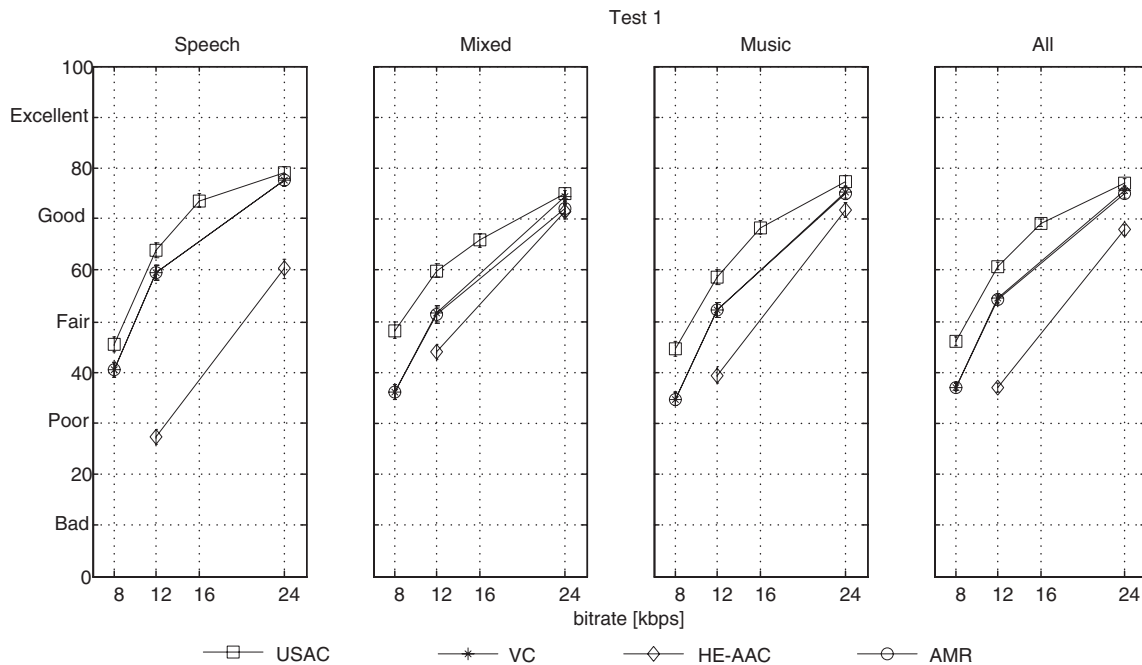


Fig. 18. Average absolute scores in Test 1 (mono at low bitrates) for USAC, HE-AACv2 (HE-AAC in the legend), AMR-WB+ (AMR in the legend), and the Virtual Codec (VC).

vertical axis in Figs. 18 to 20 uses the MUSHRA subjective scale, shown in Table 6.

Figs. 18 to 20 show that, when averaging over all content types, the average score of USAC is significantly above that of the VC, with 95% confidence intervals not overlapping by a wide margin. Two exceptions are at 24 kbit/s mono and 96 kbit/s stereo where USAC and the VC have overlapping confidence intervals but with the average score of USAC above that of the VC. Furthermore, Figs. 18 to 20 show that when considering each signal content type individually

(speech, music or speech mixed with music), the absolute score for USAC is always greater than the absolute score of the VC and often by a large margin. This is most apparent in Test 2 (stereo operation between 16 and 24 kbit/s), with a 6 to 18 point advantage for USAC on the 100-point scale. A third observation from Figs. 18 to 20 is that the quality for USAC is much more consistent across signal content types than the two state-of-the-art codecs considered (HE-AACv2 and AMR-WB+). This is especially apparent at medium and low rate operation (Figs. 18 and 19).

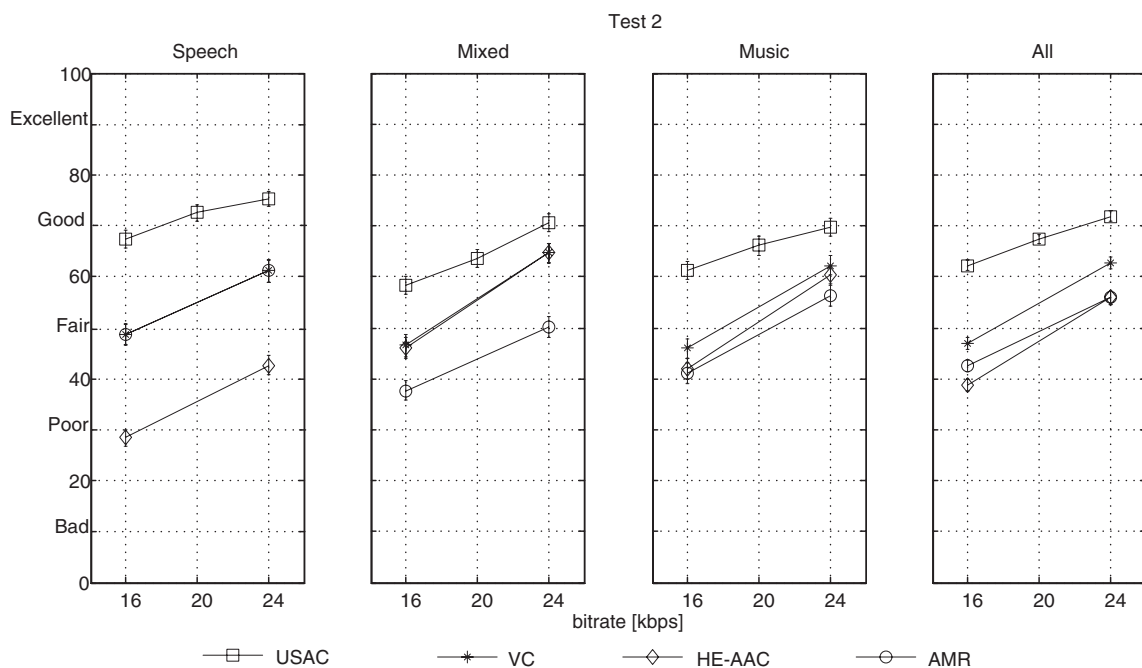


Fig. 19. Average absolute scores in Test 2 (stereo at low bit rates) for USAC, HE-AACv2 (HE-AAC in the legend), AMR-WB+ (AMR in the legend), and the Virtual Codec (VC).

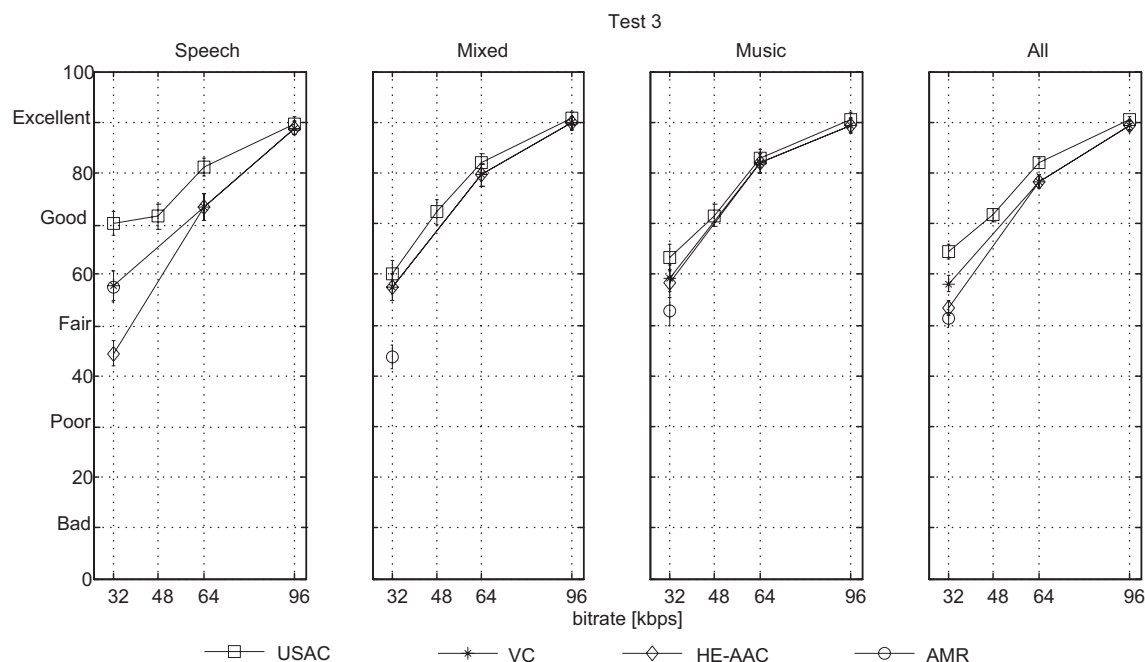


Fig. 20. Average absolute scores in Test 3 (stereo at high bit rates) for USAC, HE-AACv2 (HE-AAC in the legend), AMR-WB+ (AMR in the legend), and the Virtual Codec (VC).

The USAC verification test results show that USAC not only matches the quality of the better of HE-AACv2 and AMR-WB+ on all signal content types and at all bit rates tested (from 8 mono to 96 kbit/s stereo), but USAC actually exceeds that sound quality, and often by a large margin, in the bit rate range from 8 kbit/s mono to 64 kbit/s stereo. At higher bit rates, the quality of USAC converges to that of HE-AACv2 with or without SBR.

5 APPLICATIONS

USAC extends the HE-AACv2 range of use towards lower bit rates. As it additionally delivers at least the same quality as HE-AACv2 at higher rates, it also allows for applications requiring scalability over a large bit rate range. This makes USAC especially interesting for applications where bandwidth is limited or varying. Although mobile bandwidth is increasing with the upcoming 4G mobile standards, at the same time mobile data bandwidth usage increases dramatically. Moreover, multimedia streaming is accounting for a major part of today's growth in mobile bandwidth traffic.

In applications such as streaming multimedia to mobile devices, bandwidth scalability is a key requirement to ensure a pleasant user experience also under non-optimal conditions. Users want to receive the content without dropouts not only when being the only user in a cell and not moving. They also want to listen to their favorite Internet radio station when sitting in a fast traveling car or train, or while waiting for the very same train in a crowded station.

In digital radio, saving on transmission bandwidth reduces distribution costs and allows for a greater diversity of programs. Coding efficiency is most relevant for mobile reception, where robust channel coding schemes add

to the needed transmission bandwidth. Even in mobile TV, where video occupies the largest share of the transmission bandwidth, adding additional audio tracks like simulcasting stereo and multichannel audio or adding additional services like audio descriptive channels will significantly increase bandwidth demand. This raises the need for a highly efficient compression scheme, delivering good audio quality for both music and spoken material at low bit rates. The situation is similar for audio books. Even though these contain mostly speech content, which may justify using dedicated speech codecs, background music and effects should be reproduced in high quality as well.

For all of the above-mentioned applications, the new USAC standard seems perfectly suited because of its extended bit rate range, quality consistency, and unprecedented efficiency.

6 CONCLUSION

The ISO/IEC 23003-3:2012 MPEG-D Unified speech and audio coding standard is the first codec that reliably merges the world of general audio coding and the world of speech coding into one solid design. At the same time, USAC can be seen as the true successor of a long line of successful MPEG general audio codecs that started with MPEG-1 Audio and its most famous member, mp3. This was followed by AAC and HE-AAC(v2) that commercially share the success of mp3, as both codecs are present in virtually every mobile phone and many TV sets and currently available digital audio players.

USAC now further builds on the technologies in mp3 and AAC and takes these one step further: It includes all the essential components of its predecessors in a further evolved form. It can, therefore, do everything mp3, AAC,

and HE-AAC can do but is more efficient than its predecessors. Through the integration of the ACELP and TCX elements of AMR-WB+, USAC also represents a new state-of-the-art in low rate speech and mixed content coding. This makes USAC today the most efficient codec for all signal categories, including speech signals. Starting at bit rates of around 8 kbit/s and up, it will deliver the best speech, music, and mixed signal quality possible today for a given bit rate. Similar to AAC, USAC will scale toward perceptual transparency for higher bit rates.

During standardization, care was taken to keep the codec as lean as possible. As a result, the increase in implementation complexity over its predecessor is moderate and implementations for typical AAC and HE-AAC processing platforms are already available. All in all, USAC can be considered the true 4th generation MPEG Audio codec, again setting a new state-of-the-art like its predecessors.

7 ACKNOWLEDGMENT

A standardization project of this dimension can only ever be realized by a joint effort of a considerable number of highly skilled experts. Therefore, the authors would like to extend their appreciation to the following list of contributors for their important and substantial work over the long duration of creation of the USAC standard:

T. Bäckström, P. Carlsson, Z. Dan, F. de Bont, B. den Brinker, S. Döhla, B. Edler, P. Ekstrand, D. Fischer, R. Geiger, P. Hedelin, J. Herre, M. Hildenbrand, J. Hilpert, J. Hirschfeld, J. Kim, J. Koppens, U. Krämer, A. Kuntz, F. Nagel, M. Neusinger, A. Niedermeier, M. Nishiguchi, M. Ostrovskyy, B. Resch, R. Salami, F. Lavoie, J. Samuelsson, G. Schuller, L. Sehlström, V. Subbaraman, M. Luis Valero, S. Wabnik, P. Warmbold, Y. Yokotani, H. Zhong, and H. Zhou.

Furthermore, the authors would like to express their gratitude to all members of the MPEG Audio Subgroup for their support and collaborative work during the standardization process.

8 REFERENCES

- [1] ISO/IEC 14496-3:2009, "Coding of Audio-Visual Objects, Part 3: Audio," Aug. 2009.
- [2] M. Wolters, K. Kjörning, D. Himm, and H. Purnhagen, "A Closer Look into MPEG-4 High Efficiency AAC," presented at the *115th Convention* of the Audio Engineering Society (2003 Oct.), convention paper 5871.
- [3] C. Laflamme, J.-P. Adoul, R. Salami, S. Morissette, and P. Mabilieu, "16 kbps Wideband Speech Coding Technique Based on Algebraic Celp," in *Proc. IEEE ICASSP 1991*, Toronto (1991 Apr.), IEEE, vol. 1, pp. 13–16.
- [4] 3GPP, "Adaptive Multi-Rate - Wideband (AMR-WB) Speech Codec; General Description," 2002, 3GPP TS 26.171.
- [5] B. Bessette, R. Salami, R. Lefebvre, M. Jelinek, J. Rotola-Pukkila, J. Vainio, H. Mikkola, and K. Jarvinen, "The Adaptive Multirate Wideband Speech Codec (AMR-WB)," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 8, pp. 620–636 (2002 Nov.).
- [6] 3GPP, "Audio Codec Processing Functions; Extended Adaptive Multi-Rate - Wideband (AMR-WB+) Codec; Transcoding Functions," 2004, 3GPP TS 26.290.
- [7] J. Makinen, B. Bessette, S. Bruhn, P. Ojala, R. Salami, and A. Taleb, "AMR-WB+: A New Audio Coding Standard for 3rd Generation Mobile Audio Services," in *Proc. IEEE ICASSP 2005*, Philadelphia (2005 Mar.), vol. 2, pp. 1109–1112.
- [8] M. Neuendorf, M. Multus, N. Rettelbach, G. Fuchs, J. Robilliard, J. Lecomte, S. Wilde, S. Bayer, S. Disch, C. R. Helmrich, R. Lefebvre, P. Gournay, B. Bessette, J. Lapierre, K. Kjörning, H. Purnhagen, L. Villemoes, W. Oomen, E. Schuijers, K. Kikuri, T. Chinen, T. Norimatsu, K. S. Chong, E. Oh, M. Kim, S. Quackenbush, and B. Grill, "MPEG Unified Speech and Audio Coding - The ISO/MPEG Standard for High-Efficiency Audio Coding of All Content Types," presented at the *132nd Convention* of the Audio Engineering Society (2012 Apr.), convention paper 8654.
- [9] M. Neuendorf, P. Gournay, M. Multus, J. Lecomte, B. Bessette, R. Geiger, S. Bayer, G. Fuchs, J. Hilpert, N. Rettelbach, R. Salami, G. Schuller, R. Lefebvre, and B. Grill, "Unified Speech and Audio Coding Scheme for High Quality at Low Bitrates," in *Proc. IEEE ICASSP 2009*, Taipei (2009 Apr.), IEEE, pp. 1–4.
- [10] M. Neuendorf, P. Gournay, M. Multus, J. Lecomte, B. Bessette, R. Geiger, S. Bayer, G. Fuchs, J. Hilpert, N. Rettelbach, F. Nagel, J. Robilliard, R. Salami, G. Schuller, R. Lefebvre, and B. Grill, "A Novel Scheme for Low Bitrate Unified Speech and Audio Coding - MPEG RM0," presented at the *126th Convention* of the Audio Engineering Society (2009 May), convention paper 7713.
- [11] J. Breebaart and C. Faller, *Spatial Audio Processing: MPEG Surround and Other Applications* (John Wiley & Sons Ltd, West Sussex, England, 2007).
- [12] W.-H. Chiang, C. Hwang, and Y. Hsu, "Advances in Low Bit-Rate Audio Coding: A Digest of Selected Papers from Recent AES Conventions," *J. Audio Eng. Soc.*, vol. 51, pp. 956–964 (2003 Oct.).
- [13] K. Brandenburg and M. Bosi, "Overview of MPEG Audio: Current and Future Standards for Low Bit-Rate Audio Coding," *J. Audio Eng. Soc.*, vol. 45, pp. 4–21 (1997 Jan./Feb.).
- [14] J. Herre, H. Purnhagen, J. Koppens, O. Hellmuth, J. Engdegård, J. Hilpert, L. Villemoes, L. Terentiv, C. Falch, A. Hölzer, M. L. Valero, B. Resch, H. Mundt, and H. Oh, "MPEG Spatial Audio Object Coding - The ISO/MPEG Standard for Efficient Coding of Interactive Audio Scenes," *J. Audio Eng. Soc.*, vol. 60, pp. 655–673 (2012 Sep.).
- [15] R. M. Aarts and R. T. Dekkers, "A Real-Time Speech-Music Discriminator," *J. Audio Eng. Soc.*, vol. 47, pp. 720–725 (1999 Sep.).
- [16] J. G. A. Barbedo and A. Lopes, "A Robust and Computationally Efficient Speech/Music Discriminator," *J. Audio Eng. Soc.*, vol. 54, pp. 571–588 (2006 Jul./Aug.).
- [17] S. Garcia Galan, J. E. Muñoz Exposito, N. Ruiz Reyes, and P. Vera Candeas, "Design and Implementation of a Web-Based Software Framework for Real Time

Intelligent Audio Coding Based on Speech/Music Discrimination,” presented at the 122nd *Convention* of the Audio Engineering Society (2007 May), convention paper 7005.

[18] ISO/IEC JTC1/SC29/WG11, “Call for Proposals on Unified Speech and Audio Coding,” Shenzhen, China, Oct. 2007, MPEG2007/N9519.

[19] ISO/IEC JTC1/SC29/WG11, “MPEG Press Release,” Hannover, Germany, July 2008, MPEG2008/N9965.

[20] ISO/IEC 23003-3:2012, “MPEG-D (MPEG audio technologies), Part 3: Unified Speech and Audio Coding,” 2012.

[21] G. Fuchs, V. Subbaraman, and M. Multrus, “Efficient Context Adaptive Entropy Coding for Real-Time Applications,” in *Proc. IEEE ICASSP 2011*, Prague (2011 May), pp. 493–496.

[22] M. Bosi, K. Brandenburg, S. Quackenbush, L. Fielder, K. Akagiri, H. Fuchs, M. Dietz, J. Herre, G. Davidson, and Y. Oikawa, “ISO/IEC MPEG-2 Advanced Audio Coding,” *J. Audio Eng. Soc.*, vol. 45, pp. 789–814 (1997 Oct.).

[23] B. Edler, S. Disch, S. Bayer, G. Fuchs, and R. Geiger, “A Time-Warped MDCT Approach to Speech Transform Coding,” presented at the 126th *Convention* of the Audio Engineering Society (2009 May), convention paper 7710.

[24] J. Lecomte, P. Gournay, R. Geiger, B. Bessette, and M. Neuendorf, “Efficient Cross-Fade Windows for Transitions between LPC-Based and Non-LPC Based Audio Coding,” presented at the 126th *Convention* of the Audio Engineering Society (2009 May), convention paper 7712.

[25] S. Ragot, M. Xie, and R. Lefebvre, “Near-Ellipsoidal Voronoi Coding,” *IEEE Transactions on Information Theory*, vol. 49, no. 7, pp. 1815–1820 (2003 July).

[26] M. Dietz, L. Liljeryd, K. Kjörling, and O. Kunz, “Spectral Band Replication, a Novel Approach in Audio Coding,” presented at the 112th *Convention* of the Audio Engineering Society (2002 May), convention paper 5553.

[27] A. C. den Brinker, J. Breebaart, P. Ekstrand, J. Engdegård, F. Henn, K. Kjörling, W. Oomen, and H. Purnhagen, “An Overview of the Coding Standard MPEG-4 Audio Amendments 1 and 2: HE-AAC, SSC, and HE-AAC v2,” *EURASIP J. Audio, Speech, and Music Processing*, vol. 2009, Article ID 468971 (2009).

[28] L. Liljeryd, P. Ekstrand, F. Henn, and K. Kjörling, “Source Coding Enhancement Using Spectral Band Replication,” 2004, EP0940015B1/ WO98/57436.

[29] F. Nagel and S. Disch, “A Harmonic Bandwidth Extension Method for Audio Codecs,” in *Proc. IEEE ICASSP 2009*, Taipei (2009 Apr.), pp. 145–148.

[30] L. Villemoes, P. Ekstrand, and P. Hedelin, “Methods for Enhanced Harmonic Transposition,” in *IEEE Workshop on Appl. of Signal Proc. to Audio and Acoustics*, New Paltz (2011 Oct.), pp. 161–164.

[31] H. Zhong, L. Villemoes, P. Ekstrand, S. Disch, F. Nagel, S. Wilde, C. K. Seng, and T. Norimatsu, “QMF

Based Harmonic Spectral Band Replication,” presented at the 131st *Convention* of the Audio Engineering Society (2011 Oct.), convention paper 8517.

[32] J. D. Johnston and A. J. Ferreira, “Sum-Difference Stereo Transform Coding,” in *Proc. IEEE ICASSP 1992*, San Francisco (1992 Mar.), vol. 2, pp. 569–572.

[33] F. Baumgarte and C. Faller, “Binaural Cue Coding - Part I: Psychoacoustic Fundamentals and Design Principles,” *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 6, pp. 509–519 (2003 Nov.).

[34] C. Faller and F. Baumgarte, “Binaural Cue Coding - Part II: Schemes and Applications,” *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 6, pp. 520–531 (2003 Nov.).

[35] E. Schuijers, J. Breebaart, H. Purnhagen, and J. Engdegård, “Low Complexity Parametric Stereo Coding,” presented at the 116th *Convention* of the Audio Engineering Society (2004 May), convention paper 6073.

[36] J. Herre, K. Kjörling, J. Breebaart, C. Faller, S. Disch, H. Purnhagen, J. Koppens, J. Hilpert, J. Rödén, W. Oomen, K. Linzmeier, and K. S. Chong, “MPEG Surround – The ISO/MPEG Standard for Efficient and Compatible Multichannel Audio Coding,” *J. Audio Eng. Soc.*, vol. 56, pp. 932–955 (2008 Nov.).

[37] ISO/IEC 23003-1:2007, “MPEG-D (MPEG audio technologies), Part 1: MPEG Surround,” 2007.

[38] J. Lapiere and R. Lefebvre, “On Improving Parametric Stereo Audio Coding,” presented at the 120th *Convention* of the Audio Engineering Society (2006 May), convention paper 6804.

[39] J. Kim, E. Oh, and J. Robilliard, “Enhanced Stereo Coding with Phase Parameters for MPEG Unified Speech and Audio Coding,” presented at the 127th *Convention* of the Audio Engineering Society (2009 Oct.), convention paper 7875.

[40] S. Disch and A. Kuntz, “A Dedicated Decorrelator for Parametric Spatial Coding of Applause-like Audio Signals,” in *Microelectronic Systems: Circuits, Systems and Applications*, A. Heuberger, G. Elst, and R. Hanke, Eds. (Springer Berlin Heidelberg, 2011), pp. 355–363.

[41] A. Kuntz, S. Disch, T. Bäckström, J. Robilliard, and C. Uhle, “The Transient Steering Decorrelator Tool in the Upcoming MPEG Unified Speech and Audio Coding Standard,” presented at the 131st *Convention* of the Audio Engineering Society (2011 Oct.), convention paper 8533.

[42] C. R. Helmrich, P. Carlsson, S. Disch, B. Edler, J. Hilpert, M. Neusinger, H. Purnhagen, N. Rettelbach, J. Robilliard, and L. Villemoes, “Efficient Transform Coding of Two-Channel Audio Signals by Means of Complex-Valued Stereo Prediction,” in *Proc. IEEE ICASSP 2011*, Prague (2011 May), pp. 497–500.

[43] F. Küch and B. Edler, “Aliasing Reduction for Modified Discrete Cosine Transform Domain Filtering and its Application to Speech Enhancement,” in *IEEE Workshop on Appl. of Signal Proc. to Audio and Acoustics*, New Paltz (2007 Oct.), pp. 131–134.

[44] Y. Kikuchi, T. Nomura, S. Fukunaga, Y. Matsui, and H. Kimata, "RTP Payload Format for MPEG-4 Audio/Visual Streams," Nov. 2000, IETF RFC 3016.

[45] J. van der Meer, D. Mackie, V. Swaminathan, D. Singer, and P. Gentric, "RTP Payload Format for Transport of MPEG-4 Elementary Streams," Nov. 2003, IETF RFC 3640.

[46] ISO/IEC 13818-1:2007, "Generic Coding of Moving Pictures and Associated Audio Information, Part 1: Systems," 2007.

[47] ISO/IEC 14496-14:2003, "Coding of Audio-Visual Objects, Part 14: MP4 File Format," 2003.

[48] 3GPP, "Transparent End-to-End Packet Switched Streaming Service (PSS); 3GPP File Format (3GP)," 2011, 3GPP TS 26.244.

[49] ISO/IEC JTC1/SC29/WG11, "Evaluation Guidelines for Unified Speech and Audio Proposals," Antalya, Turkey, Jan. 2008, MPEG2008/N9638.

[50] ISO/IEC JTC1/SC29/WG11, "Unified Speech and Audio Coding Verification Test Report," Torino, Italy, July 2011, MPEG2011/N12232.

[51] International Telecommunication Union, "Method for the Subjective Assessment of Intermediate Sound Quality (MUSHRA)," 2003, ITU-R, Recommendation BS. 1534-1, Geneva.

THE AUTHORS



Max Neuendorf



Markus Multrus



Nikolaus Rettelbach



Guillaume Fuchs



Julien Robilliard



Jérémie Lecomte



Stephan Wilde



Stefan Bayer



Sascha Disch



Christian R. Helmrich



Roch Lefebvre



Philippe Gournay



Bruno Bessette



Jimmy Lapierre



Kristofer Kjörling



Heiko Purnhagen



Lars Villemoes



Werner Oomen



Erik Schuijers



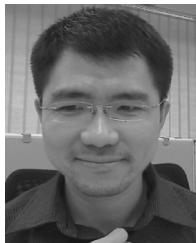
Kei Kikuri



Toru Chinen



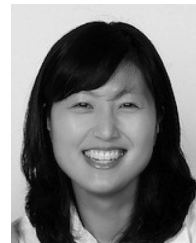
Takeshi Norimatsu



Chong Kok Seng



Eunmi Oh



Miyoung Kim



Schuyler Quackenbush



Bernhard Grill

Max Neuendorf is group manager of the audio and speech coding group at Fraunhofer IIS. He joined Fraunhofer IIS after acquiring his Diploma in electrical engineering and information technology at the Technische Universität München (TUM) in 2002. Mr. Neuendorf has been working in the field of modern audio codec development for 11 years and is the main editor of the ISO/IEC MPEG-D USAC standard document.

Markus Multus received his Diplom-Ingenieur degree in electrical engineering from the Friedrich-Alexander University Erlangen-Nuremberg in 2004. Since then he has been working on audio codec development at Fraunhofer IIS. His main research areas include waveform and parametric audio/speech coding and bandwidth extension. Mr. Multus was an active member of the ISO MPEG committee during the USAC standardization and is co-editor of the ISO/IEC MPEG-D USAC standard document.

Nikolaus Rettelbach joined Fraunhofer IIS in 2000 after receiving his Diplom-Ingenieur degree in electrical engineering from the Friedrich-Alexander University of Erlangen-Nuremberg, Germany. Since joining IIS he has been working in the field audio coding and in 2010 became group manager of the high quality audio coding group at Fraunhofer IIS.

Guillaume Fuchs received his Engineering Degree from the engineering school INSA of Rennes, France, in 2001 and his Ph.D from the University of Sherbrooke, Canada, in 2006, both in electrical engineering. From 2001 to 2003, he worked on digital image compression as research engineer at Canon Research, France. From 2003 to 2006, he worked with VoiceAge, Canada, as scientific collaborator on speech coding. In 2006, he joined Fraunhofer IIS and since then has been working on developing speech and audio coders for different standards. His main research interests include speech and audio source coding and speech analysis.

Julien Robilliard received an M.Sc. degree in Acoustics from the Aalborg University, Denmark, in 2005. He joined Fraunhofer IIS, in the Audio & Multimedia department, in 2008. His main fields of interest and expertise are spatial hearing, parametric stereo, and multichannel audio coding.

Jérémie Lecomte is a Senior Research Engineer at Fraunhofer IIS. He received his Engineering Degree in electrical engineering from the ESEO engineering school (France) in 2007 and his M.S. degree in signal processing and telecommunications from the Université de Sherbrooke (Canada) in 2008. His main fields of interests are in speech and audio coding and robustness. He has contributed to the MPEG standards on Unified Speech and Audio Coding.

Stephan Wilde received a Dipl.-Ing. degree in information and communication technology from the Friedrich-Alexander University of Erlangen-Nuremberg, Germany, in 2009. After graduation he joined Fraunhofer IIS in Erlangen, Germany, where he works as a research assistant in the field of audio and speech coding. His main research interests include audio/speech coding and bandwidth extension.

Stefan Bayer received a Diplomingenieur degree in electrical engineering and audio engineering from the Graz University of Technology, Austria, in 2005. From 2005 until 2011 he joined the Fraunhofer IIS in Erlangen, Germany, where he worked as a research assistant in the field of perceptual audio coding and speech coding and contributed to the development of MPEG Unified Speech and Audio Coding. Since 2011 he is a researcher at the International Audio Laboratories Erlangen, where his research interests include audio coding and audio signal analysis.

Sascha Disch received his Diplom-Ingenieur degree in electrical engineering from Hamburg University of Technology (TUHH), Germany, in 1999. From 1999 to 2007 he joined the Fraunhofer IIS, Erlangen, Germany. At Fraunhofer he worked in research and development in the field of perceptual audio coding and audio processing. In MPEG standardization of parametric spatial audio coding he contributed as a developer and served as a co-editor of the MPEG Surround standard. From 2007 to 2010 he was a researcher at the Laboratory of Information Technology, Leibniz University Hanover (LUH), Germany, from which he received his Dr.-Ingenieur degree in 2011. During that time, Dr. Disch also participated in the development of the Unified Speech and Audio Coding (USAC) standard at MPEG. Since 2010 Dr. Disch is affiliated with the IIS again. His research interests as a Senior Scientist at Fraunhofer include waveform and parametric audio signal coding, audio bandwidth extension, and digital audio effects.

Christian R. Helmrich received a B.Sc. degree in computer science from Capitol College, Laurel, MD, USA, in 2005 and a M.Sc. degree in information and media technologies from Hamburg University of Technology (TUHH), Germany, in 2008. Since then he has been working on numerous speech and audio coders for Fraunhofer IIS in Erlangen, partly as a Senior Engineer. Mr. Helmrich is currently pursuing a Ph.D. degree in audio analysis and coding at the International Audio Laboratories Erlangen, a joint institution of Fraunhofer IIS and the University of Erlangen-Nuremberg. His main research interests include audio and video coding as well as restoration from analog sources.

Roch Lefebvre is professor at the Faculty of Engineering in Université de Sherbrooke, Quebec, Canada, where he also leads the Research Group on Speech and Audio Coding since 1998. He received the B.Sc. degree in physics from McGill University in 1989, and the M.Sc. and Ph.D. degrees in electrical engineering from the Université de Sherbrooke in 1992 and 1995, respectively. Professor Lefebvre is cofounder of VoiceAge Corporation, a Montreal-based company that commercializes digital speech and audio solutions for communications and consumer electronics. His research areas include low bit rate and robust speech and audio coding, noise reduction, and music analysis for user interaction. Professor Lefebvre has published numerous journal and conference papers and participated in two books on speech coding and signal processing. He also contributed to several standardization activities in speech and audio coding within ITU-T, 3GPP and ISO MPEG. He is a member of the IEEE and the AES.

Philippe Gournay is an Adjunct Professor at the Université de Sherbrooke, Quebec, Canada, and a Senior Researcher for VoiceAge, Montreal (QC), Canada. He received his Engineering Degree in electrical engineering from the ENSSAT engineering school, France, in 1991 and his Ph.D. in signal processing and telecommunications from the Université de Rennes I, France, in 1994. His current research interests are in speech and audio coding and speech enhancement.

Bruno Bessette was born in Canada in 1969. He received the B.S. degree in electrical engineering from Université de Sherbrooke, Quebec, Canada, in 1992. During his studies and as a graduated engineer, he worked in industry in the field of R&D. He joined the speech coding group of Université de Sherbrooke in 1994. Initially, he has been taking part on the design and implementation of speech coding algorithms where he developed new ideas, several of which are integral parts of standards widely deployed in wireless systems around the world. In 1999, he became co-founder of VoiceAge Corporation, a company based in Montreal, in which he is involved full time. In the past several years he has been actively involved in wideband speech and audio codecs standardization activities within ITU-T and 3GPP. His work in audio coding led him to collaborate with Fraunhofer IIS from 2007 to 2012 in the ISO/MPEG standardization leading to USAC standard.

Jimmy Lapierre received a Master's Degree (M.Sc.A.) in electrical engineering from the Université de Sherbrooke in Québec, Canada, in 2006. He currently works for VoiceAge, while pursuing a Ph.D. in Sherbrooke, focusing his research on improving low bit rate perceptual audio coding.

Kristofer Kjörling received his M.Sc. degree in electrical engineering from the Royal Institute of Technology in Stockholm, Sweden, concluded by a Master's Thesis on early concepts for what would evolve into MPEG-4 SBR as part of High Efficiency AAC. From 1997 to 2007 he worked at Coding Technologies AB in Stockholm, among other things leading the company's standardization effort in MPEG, actively participating in and contributing to the standardization of MPEG-4 HE AAC, MPEG-D MPEG Surround, MPEG-D USAC. The main research work has been on Spectral Band Replication (HE-AAC), MPEG Surround, and MPEG-D USAC. Mr. Kjörling is a co-recipient of the 2013 IEEE Masaru Ibuka Consumer Electronics Award for his work on HE AAC. Since late 2007 he has been employed by Dolby Laboratories, where he holds the position of Principal Member Technical Staff.

Heiko Purnhagen was born in Bremen, Germany, in 1969. He received an M.S. degree (Diplom) in electrical engineering from the University of Hannover, Germany, in 1994, concluded by a thesis project on automatic speech recognition carried out at the Norwegian Institute of Technology, Trondheim, Norway.

From 1996 to 2002 he was with the Information Technology Laboratory at the University of Hannover, where he pursued a Ph.D. degree related to his research on very low-bit-rate audio coding using parametric signal representations. In 2002 he joined Coding Technologies (now Dolby Sweden) in Stockholm, Sweden, where he is working on research, development, and standardization of low bit-rate audio and speech coding systems. He contributed to the development of parametric techniques for efficient coding of stereo and multichannel signals, and his recent research activities include the unified coding of speech and audio signals.

Since 1996 Mr. Purnhagen has been an active member of the ISO MPEG standardization committee and editor or co-editor of several MPEG standards. He is the principal author of the MPEG-4 parametric audio coding specification known as HILN (harmonic and individual lines plus noise) and contributed to the standardization of MPEG-4 High Efficiency AAC, the MPEG-4 Parametric Stereo coding tool, and MPEG Surround. He is a member of the AES and IEEE and has published various papers on low-bit-rate audio coding and related subjects. He enjoys listening to live performances of jazz and free improvised music, and tries to capture them in his concert photography.

Lars Villemoes received the M.Sc. degree in engineering and the Ph.D. degree in mathematics from the Technical University of Denmark in 1989 and 1992, respectively. From the Royal Institute of Technology (KTH) in Sweden he received the TeknD. and Docent degrees in mathematics in 1995 and 2001. As a postdoctoral researcher 1995–1997, he visited the Department of Mathematics at Yale University and the Signal Processing group of the Department of Signals, Systems, and Sensors, KTH. From 1997 to 2001, he was a Research Associate in wavelet theory at the Department of Mathematics, KTH. He then joined Coding Technologies and since 2008 is a Senior Member of Technical Staff at Dolby. His main interest is the development of new tools for the processing and representation of signals in time and frequency. He has contributed to the MPEG standards on Unified Speech and Audio Coding, Spatial Audio Object Coding, MPEG Surround and High Efficiency AAC.

Werner Oomen received the Ingenieur degree in Electronics from the University of Eindhoven, The Netherlands, in 1992. In that same year he joined Philips Research Laboratories in Eindhoven in the Digital Signal Processing group, leading and contributing to a diversity of audio signal processing projects, primarily in the field of audio source coding algorithms. As of 1999, the development aspects in the various projects have become more prominent as well as his supporting role towards Philips' Intellectual Properties and Standardization department. As of 1995, Werner has been involved in standardization bodies, primarily 3GPP and MPEG, where for the latter he has actively contributed to the standardization of o.a. MPEG2-AAC, MPEG4-WB CELP, Parametric (Stereo) Coding, Lossless Coding of 1-bit Over-Samples Audio, MPEG Surround, Spatial Audio Object Coding, and Unified Speech and Audio Coding.

Erik Schuijers was born in the Netherlands in 1976. He received the M.Sc. degree in electrical engineering from the Eindhoven University of Technology, The Netherlands, in 1999. In 2000 he joined Philips, primarily working on research and development of audio coding technology, resulting in contributions to various MPEG audio standards. In 2013 he moved to the Brain, Body & Behavior department of Philips Research, contributing to research on interpretation of physiological sensor data for infants and frail elderly.

Kei Kikuri received his B.E. and M.E. degrees from the Yokohama National University, Kanagawa, Japan, in 1998

and 2000, respectively. In 2000 he joined NTT DOCOMO, Inc. Since joining NTT DOCOMO he has been engaged in research and development of speech and audio coding technology.

Toru Chinen received the M.S. degree in electrical engineering from Tokyo Metropolitan University in 1991. Since then, he has been working for algorithm development and implementation of audio codecs, e.g., MPEG AAC and USAC. Currently he is leading a group of engineers to develop an audio codec handling both channel- and object-based audio signals.

Takeshi Norimatsu, after receiving a master degree of computer science from Kyushu University, Japan, in 1981, worked at Panasonic Corporation as a researcher of speech and audio technologies and an engineer of various audio-visual products. Since 1996, he has contributed to the standardization of MPEG audio coding technologies such as TwinVQ, HE-AAC, and MPEG Surround. From 2006 to 2012, he was a chair of SC29/WG11 audio committee in Information Technology Standards Commission of Japan (ITSCJ). This study was conducted while he worked in Panasonic.

Chong Kok Seng graduated with a Bachelor of Electrical and Computer Systems Engineering degree in 1994 and a Ph.D. in 1997 from Monash University, Melbourne, Victoria, Australia. He was a Principal Engineer at Panasonic R&D Center, Singapore. His area of research is audio coding, acoustics, and signal processing. He has published 50+ patents in his career so far.

Eunmi L. Oh received her Ph.D. in psychology (with an emphasis on psychoacoustics) from the University of Wisconsin-Madison in 1997. She has been with Samsung Electronics since 2000. Her recent researches include perceptual audio coding, super-wideband speech coding, and 3-D audio processing.

Miyoung Kim received M.S degree from Kyungbook National University, Korea in 2002. Since 2002 she has been with Samsung Electronics. Her main research interests are perceptual audio coding and speech coding.

Schuyler Quackenbush received a B.S. degree from Princeton University in 1975 an M.S. degree in electrical engineering in 1980, and a Ph.D. degree in electrical engineering in 1985 from Georgia Institute of Technology. He was Member of Technical Staff at AT&T Bell Laboratories from 1986 until 1996, when he joined AT&T Laboratories as Principal Technical Staff Member. In 2002 he founded Audio Research Labs, an audio technology consulting company. He has been the Chair of the International Standards Organization Moving Picture Experts Group (ISO/MPEG) Audio subgroup since 1998, and was one of the authors of the ISO/IEC MPEG Advanced Audio Coding standard. Dr. Quackenbush is a fellow of the Audio Engineering Society (AES) and a senior member of the Institute of Electrical and Electronics Engineers (IEEE).

●
Dr.-Ing. Bernhard Grill studied electrical engineering at the Friedrich-Alexander-University Erlangen-Nuremberg. From 1988 until 1995 he engaged in the development and implementation of audio coding algorithms at the Fraunhofer IIS in Erlangen. In 1992 he received, together with Jürgen Herre and Ernst Eberlein, the “Joseph-von-Fraunhofer Award” for his contributions to the mp3 technology. Since 2000 Bernhard Grill has been heading the

audio department at Fraunhofer IIS. Together with Harald Popp, he is responsible for the successful development and commercial application of new audio coding algorithms and associated technologies. In 2000 Bernhard Grill, Karlheinz Brandenburg, and Harald Popp, as representatives of the larger team at Fraunhofer IIS, were honored with the “German Future Award” as the inventors of mp3. Since July 2011 Bernhard Grill is Deputy Director of the Fraunhofer IIS.