

---

---

**Information technology — High  
efficiency coding and media delivery  
in heterogeneous environments —**

**Part 3:  
3D audio**

*Technologies de l'information — Codage à haute efficacité et livraison  
des médias dans des environnements hétérogènes —*

*Partie 3: Audio 3D*

*ITeH STANDARD PREVIEW  
(standards.iteh.ai)  
Full standard: <https://standards.iteh.ai/catalog/standards/sist/e4315f20-cb86-4e52-81ea-8ad020d7daf4/iso-iec-23008-3-2015>*

**ITeH STANDARD PREVIEW**  
**(standards.iteh.ai)**

Full standard:  
<https://standards.iteh.ai/catalog/standards/sist/e4315f20-cb86-4e52-81ea-8dd020d7daf4/iso-iec-23008-3-2015>



**COPYRIGHT PROTECTED DOCUMENT**

© ISO/IEC 2015, Published in Switzerland

All rights reserved. Unless otherwise specified, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office  
Ch. de Blandonnet 8 • CP 401  
CH-1214 Vernier, Geneva, Switzerland  
Tel. +41 22 749 01 11  
Fax +41 22 749 09 47  
copyright@iso.org  
www.iso.org

# Contents

Page

Foreword .....	viii
Introduction.....	ix
<b>1</b> <b>Scope</b> .....	<b>1</b>
<b>2</b> <b>Normative references</b> .....	<b>1</b>
<b>3</b> <b>Terms, definitions and mnemonics</b> .....	<b>1</b>
3.1 <b>Terms and Definitions</b> .....	<b>1</b>
3.2 <b>Mnemonics</b> .....	<b>1</b>
<b>4</b> <b>Technical Overview</b> .....	<b>2</b>
4.1 <b>Decoder block diagram</b> .....	<b>2</b>
4.2 <b>Overview over the codec building blocks</b> .....	<b>3</b>
4.3 <b>Efficient combination of decoder processing blocks in time domain and QMF domain</b> .....	<b>4</b>
4.4 <b>Rule set for determining processing domains</b> .....	<b>5</b>
4.4.1 <b>Audio Core Codec, Processing Domain</b> .....	<b>5</b>
4.4.2 <b>Mixing</b> .....	<b>6</b>
4.4.3 <b>Audio Core Codec, Interface Domain to Rendering</b> .....	<b>6</b>
4.4.4 <b>Rendering Context</b> .....	<b>6</b>
4.4.5 <b>Post-Processing Context</b> .....	<b>6</b>
4.4.6 <b>End-of-Chain Context</b> .....	<b>7</b>
<b>5</b> <b>MPEG-H 3D Audio Core decoder</b> .....	<b>7</b>
5.1 <b>Terms and Definitions</b> .....	<b>7</b>
5.1.1 <b>Joint Stereo</b> .....	<b>7</b>
5.1.2 <b>MPEG Surround based stereo (MPS 212)</b> .....	<b>7</b>
5.2 <b>Syntax</b> .....	<b>7</b>
5.2.1 <b>General</b> .....	<b>7</b>
5.2.2 <b>Decoder configuration</b> .....	<b>7</b>
5.2.3 <b>MPEG-H 3D Audio Core bitstream payloads</b> .....	<b>22</b>
5.3 <b>Data Structure</b> .....	<b>30</b>
5.3.1 <b>General</b> .....	<b>30</b>
5.3.2 <b>General Configuration Data Elements</b> .....	<b>30</b>
5.3.3 <b>Loudspeaker Configuration Data Elements</b> .....	<b>32</b>
5.3.4 <b>Core Decoder Configuration Data Elements</b> .....	<b>34</b>
5.3.5 <b>Downmix Matrix Data Elements</b> .....	<b>37</b>
5.3.6 <b>HOA Rendering Matrix Data Elements</b> .....	<b>40</b>
5.4 <b>Configuration Element Descriptions</b> .....	<b>42</b>
5.4.1 <b>General</b> .....	<b>42</b>
5.4.2 <b>Downmix configuration</b> .....	<b>43</b>
5.4.3 <b>HOA rendering matrix configuration</b> .....	<b>47</b>
5.5 <b>Tool Descriptions</b> .....	<b>51</b>
5.5.1 <b>General</b> .....	<b>51</b>
5.5.2 <b>Quad Channel Element</b> .....	<b>52</b>
5.5.3 <b>Transform Splitting</b> .....	<b>53</b>
5.5.4 <b>MPEG Surround for Mono to Stereo upmixing</b> .....	<b>60</b>
5.5.5 <b>Enhanced Noise Filling</b> .....	<b>62</b>
5.5.6 <b>Audio Pre-Roll</b> .....	<b>82</b>
5.6 <b>Buffer requirements</b> .....	<b>86</b>
5.6.1 <b>Minimum decoder input buffer</b> .....	<b>86</b>
5.6.2 <b>Bit reservoir</b> .....	<b>86</b>
5.6.3 <b>Maximum bit rate</b> .....	<b>87</b>
5.7 <b>Stream Access Point requirements and inter-frame dependency</b> .....	<b>87</b>
<b>6</b> <b>Dynamic Range Control and Loudness Processing</b> .....	<b>88</b>

6.1	Introduction .....	88
6.2	Description .....	88
6.3	Syntax .....	88
6.3.1	Loudness Metadata .....	88
6.3.2	Dynamic Range Control Metadata .....	89
6.3.3	Data Elements .....	90
6.4	Decoding Process .....	91
6.4.1	General.....	91
6.4.2	Dynamic Range Control .....	93
6.4.3	Usage of downmixId in MPEG-H .....	93
6.4.4	DRC Set Selection Process .....	94
6.4.5	DRC-1 for SAOC 3D Content .....	95
6.4.6	DRC-1 for HOA Content .....	96
6.4.7	Loudness Normalization .....	98
6.4.8	Peak Limiter.....	98
6.4.9	Time-Synchronization of DRC gains.....	98
7	Object Metadata Decoding.....	98
7.1	Introduction .....	98
7.2	Description .....	98
7.3	Syntax .....	99
7.3.1	Object Metadata Configuration .....	99
7.3.2	Top level object metadata syntax .....	100
7.3.3	Subsidiary payloads for efficient object metadata decoding .....	100
7.3.4	Subsidiary payloads for object metadata decoding with low delay .....	104
7.4	Data Structure .....	108
7.4.1	Definition of ObjectMetadataConfig() payloads .....	108
7.4.2	Efficient Object Metadata Decoding .....	108
7.4.3	Object Metadata Decoding with Low Delay .....	113
8	Object Rendering .....	117
8.1	Description .....	117
8.2	Terms and Definitions .....	117
8.3	Input data .....	117
8.4	Processing.....	119
8.4.1	Imaginary Loudspeakers .....	119
8.4.2	Dividing the Loudspeaker Setup into a Triangle Mesh.....	120
8.4.3	Rendering Algorithm .....	121
9	SAOC 3D .....	125
9.1	Description .....	125
9.2	Definitions .....	125
9.3	Delay and synchronization .....	127
9.4	Syntax .....	127
9.4.1	Payloads for SAOC 3D .....	127
9.4.2	Definition of SAOC 3D payloads .....	131
9.5	SAOC 3D processing.....	133
9.5.1	Compressed data stream decoding and dequantization of SAOC 3D data.....	133
9.5.2	Time/frequency tranforms .....	133
9.5.3	Signals and parameters .....	133
9.5.4	SAOC 3D decoding.....	135
9.5.5	Dual mode.....	140
10	Generic Loudspeaker Rendering/Format Conversion .....	141
10.1	Description .....	141
10.2	Definitions .....	142
10.2.1	General remarks.....	142
10.2.2	Variable definitions.....	142
10.3	Processing.....	143
10.3.1	Application of transmitted downmix matrices.....	143
10.3.2	Application of transmitted equalizer settings.....	148
10.3.3	Downmix processing involving multiple channel groups.....	148

10.3.4	Initialization of the format converter .....	149
10.3.5	Audio signal processing.....	165
11	Immersive Loudspeaker Rendering / Format Conversion .....	171
11.1	Description .....	171
11.2	Syntax.....	172
11.3	Definitions .....	173
11.3.1	General remarks .....	173
11.3.2	Variable definitions .....	173
12	Higher Order Ambisonics (HOA) .....	221
12.1	Technical Overview .....	221
12.1.1	Block Diagram .....	221
12.1.2	Overview of the decoder tools .....	222
12.2	Syntax.....	223
12.2.1	Configuration of HOA elements.....	223
12.2.2	Payloads of HOA elements.....	224
12.3	Data Structure.....	229
12.3.1	Definitions of HOA Config .....	229
12.3.2	Definitions of HOA payload.....	231
12.4	HOA Tool Description .....	234
12.4.1	HOA Frame Converter.....	234
12.4.2	Spatial HOA decoding.....	243
12.4.3	HOA Renderer.....	255
13	Binaural Renderer .....	263
13.1	Introduction.....	263
13.2	Frequency-Domain Binaural Renderer.....	264
13.2.1	Introduction.....	264
13.2.2	Definitions .....	266
13.2.3	Parameterization of Binaural Room Impulse Responses .....	270
13.2.4	Frequency-Domain Binaural Processing .....	282
13.3	Time-Domain Binaural Renderer .....	289
13.3.1	Introduction.....	289
13.3.2	Definitions .....	290
13.3.3	Parameterization of Binaural Room Impulse Responses .....	291
13.3.4	Time-Domain Binaural Processing.....	296
14	MPEG-H 3D audio stream (MHAS).....	297
14.1	Overview.....	297
14.2	Syntax.....	297
14.2.1	Main MHAS syntax elements.....	297
14.2.2	Subsidiary MHAS syntax elements .....	299
14.3	Semantics.....	299
14.3.1	mpeghAudioStreamPacket() .....	299
14.3.2	MHASPacketPayload() .....	300
14.4	Description of MHASPacketTypes.....	300
14.4.1	PACTYP_FILLDATA .....	300
14.4.2	PACTYP_MPEGH3DACFG .....	300
14.4.3	PACTYP_MPEGH3DAFRAME.....	301
14.4.4	PACTYP_SYNC .....	301
14.4.5	PACTYP_SYNGAP .....	301
14.4.6	PACTYP_MARKER .....	301
14.4.7	PACTYP_CRC16 and PACTYP_CRC32 .....	302
14.4.8	PACTYP_DESCRIPTOR .....	302
14.4.9	PACTYP_USERINTERACTION .....	302
14.4.10	PACTYP_LOUDNESS_DRC .....	302
14.4.11	PACTYP_BUFFERINFO.....	303
14.5	Application Examples .....	303
14.5.1	Light-weighted broadcast.....	303
14.5.2	MPEG-2 Transport Stream.....	303
14.6	Multi-Stream Delivery and Interface .....	304

15	Metadata Audio Elements (MAE).....	306
15.1	Introduction .....	306
15.2	Syntax .....	307
15.3	Semantics .....	311
15.4	Definition of mae_metaDataElementIDs.....	319
16	Loudspeaker Distance Compensation .....	319
17	Interfaces to the MPEG-H 3D audio decoder .....	320
17.1	General.....	320
17.2	Interface for local setup information .....	321
17.2.1	General.....	321
17.2.2	WIRE output .....	321
17.2.3	Syntax for local setup information.....	321
17.2.4	Semantics for local setup information .....	321
17.3	Interface for local loudspeaker setup and rendering.....	322
17.3.1	General.....	322
17.3.2	Syntax for local loudspeaker signaling.....	322
17.3.3	Semantics for local loudspeaker signaling.....	323
17.4	Interface for binaural room impulse responses (BRIRs).....	324
17.4.1	Introduction .....	324
17.4.2	Syntax of Binaural Renderer Interface .....	324
17.4.3	Semantics .....	327
17.5	Interface for local screen size information .....	332
17.5.1	General.....	332
17.5.2	Syntax .....	332
17.5.3	Semantics .....	332
17.6	Interface for signaling of local zoom area.....	333
17.6.1	General.....	333
17.6.2	Syntax .....	333
17.6.3	Semantics .....	334
17.7	Interface for user interaction .....	334
17.7.1	Introduction .....	334
17.7.2	Definition of User Interaction Categories.....	334
17.7.3	Definition of an Interface for User Interaction.....	335
17.7.4	Syntax of interaction interface .....	335
17.7.5	Semantics of interaction interface .....	336
17.8	Interface for loudness normalization and dynamic range control (DRC).....	338
18	Application and processing of local setup information and interaction data.....	338
18.1	<i>Element Metadata</i> Preprocessing .....	338
18.2	Interactivity Limitations and Restrictions .....	341
18.2.1	General Information.....	341
18.2.2	WIRE Interactivity .....	341
18.2.3	Position Interactivity .....	342
18.2.4	Screen-Related Element Remapping and Object Remapping for Zooming .....	342
18.2.5	Closest Speaker Playback .....	342
18.3	Screen-Related Element Remapping.....	342
18.4	Object Remapping for Zooming.....	344
18.5	Determination of the Closest Speaker.....	345
Annex A	(normative) Tables for arithmetic decoding of IGF information .....	347
A.1	cf_se01[27] .....	347
A.2	cf_se10[27] .....	347
A.3	cf_se02[ 7][27] .....	347
A.4	short cf_se20[ 7][27] .....	347
A.5	short cf_se11[ 7][ 7][27] .....	348
A.6	cf_off_se01 .....	349
A.7	cf_off_se10 .....	349
A.8	cf_off_se02[ 7].....	349
A.9	short cf_off_se20[ 7].....	350
A.10	cf_off_se11[ 7][ 7] .....	350

<b>Annex B</b> (normative) <b>SAOC 3D Decorrelator pre-mixing matrices</b> .....	<b>351</b>
<b>B.1</b> Premixing matrix for output configurations with small number of output channels.....	<b>351</b>
<b>B.2</b> Premixing matrix for 22.2 output configuration.....	<b>351</b>
<b>B.3</b> Algorithm for generating pre-mixing matrices.....	<b>352</b>
<b>B.3.1</b> Input to the algorithm and representations.....	<b>352</b>
<b>B.3.2</b> Algorithm steps.....	<b>352</b>
<b>Annex C</b> (informative) <b>Encoder Tools</b> .....	<b>356</b>
<b>C.1</b> General Overview.....	<b>356</b>
<b>C.1.1</b> Encoder block diagram.....	<b>356</b>
<b>C.1.2</b> Overview of the encoder and decoder building blocks.....	<b>356</b>
<b>C.2</b> Core Encoder Tools.....	<b>357</b>
<b>C.2.1</b> Quad Channel Element.....	<b>357</b>
<b>C.2.2</b> Transform Splitting.....	<b>358</b>
<b>C.2.3</b> Calculation of Residual Signal for MPEG Surround with Hybrid Residual Coding.....	<b>359</b>
<b>C.2.4</b> Enhanced Noise Filling.....	<b>359</b>
<b>C.3</b> Object Metadata Encoding.....	<b>360</b>
<b>C.3.1</b> Pre-Processing of the Object Metadata.....	<b>360</b>
<b>C.3.2</b> Efficient Object Metadata Encoding.....	<b>361</b>
<b>C.3.3</b> Object Metadata Encoding with Low Delay.....	<b>361</b>
<b>C.3.4</b> Spatially skipping objects.....	<b>361</b>
<b>C.4</b> SAOC 3D Encoder.....	<b>361</b>
<b>C.4.1</b> Overview.....	<b>361</b>
<b>C.4.2</b> Calculation of the SAOC 3D parameters.....	<b>361</b>
<b>C.4.3</b> Time/frequency transform.....	<b>362</b>
<b>C.4.4</b> Framing.....	<b>362</b>
<b>C.4.5</b> Parameter quantization and coding.....	<b>362</b>
<b>Annex D</b> (informative) <b>Peak limiter for unguided clipping prevention</b> .....	<b>384</b>
<b>Annex E</b> (normative) <b>Compact Template Downmix Matrices</b> .....	<b>385</b>
<b>Annex F</b> (normative) <b>HOA Tables</b> .....	<b>386</b>
<b>F.6</b> 32 Uniformly Distributed Positions in Spherical Coordinates.....	<b>390</b>
<b>F.12</b> Table of 256x8 weighting values, WeightValCdbk.....	<b>404</b>
<b>Annex G</b> (informative) <b>Low Complexity HOA Rendering</b> .....	<b>421</b>
<b>G.1</b> Tool Description.....	<b>421</b>
<b>G.2</b> Predominant Sound Rendering.....	<b>421</b>
<b>G.3</b> Ambient Sound Rendering.....	<b>422</b>
<b>G.4</b> Output Signal Composition.....	<b>422</b>
<b>Annex H</b> (informative) <b>Information on delay and complexity of Time-Domain binauralization</b> .....	<b>423</b>
<b>H.1</b> Complexity and latency.....	<b>423</b>
<b>H.1.1</b> Algorithm description.....	<b>423</b>
<b>H.1.2</b> Complexity.....	<b>423</b>
<b>H.1.3</b> Latency.....	<b>424</b>
<b>H.2</b> Experimental results.....	<b>425</b>
<b>H.3</b> Alternative low-delay implementations.....	<b>426</b>
<b>Bibliography</b> .....	<b>428</b>

## Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work. In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see [www.iso.org/directives](http://www.iso.org/directives)).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see [www.iso.org/patents](http://www.iso.org/patents)).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation on the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the WTO principles in the Technical Barriers to Trade (TBT) see the following URL: [Foreword - Supplementary information](#)

The committee responsible for this document is ISO/IEC JTC 1, *Information technology*, Subcommittee SC 29, *Coding of audio, picture, multimedia and hypermedia information*.

ISO/IEC 23008 consists of the following parts, under the general title *Information technology — High efficiency coding and media delivery in heterogeneous environments*:

- *Part 1: MPEG media transport (MMT)*
- *Part 2: High efficiency video coding*
- *Part 3: 3D audio*
- *Part 4: MMT Reference and Conformance Software*
- *Part 5: Reference software for high efficiency video coding*
- *Part 8: HEVC conformance testing*
- *Part 10: MPEG media transport forward error correction (FEC) codes*
- *Part 11: MPEG media transport composition information*
- *Part 12: Image file format*
- *Part 13: MMT Implementation Guidelines*



## Introduction

3D sound systems are able to realize a significantly enhanced sound experience relative to current widespread 5.1 channel audio programs and playback systems. These systems demand high quality audio coding and error-free transmission in order to keep the timbre, sound localization and sound envelopment of the original audio program. Presentation over headphones with suitable spatialization are also considered.

This part of ISO/IEC 23008-3 “High Efficiency Coding and Media Delivery in Heterogeneous Environments — Part 3: 3D Audio” provides means for all scenarios where there is a need to compress a multi-channel audio program (e.g. 22.2 channel program) and to render it to the native target number of loudspeakers. In order to reach a wide market, a 3D Audio program is able to be downmixed to a lower hierarchy of loudspeakers, for example 10.1 or 8.1 channels. In addition, all scenarios support a level of random access to facilitate broadcast break-in, and “trick modes” such as fast forward when playing from stored media.

The main focus of this specification are applications such as audio for Home Theatres where the audio presentation is immersive, involving many loudspeakers (e.g. from 10 to more than 20) surrounding the listener and placed below, at and above ear-level. Moreover applications as Personal TV, TV for SmartPhones and Multi-channel Audio-only Programs are envisioned. These require that 3D Audio encoding/decoding systems are able to operate at low bitrates appropriate for efficient transmission over a cellular channel. At the same time the sense of envelopment and accurate sonic localization even for systems having a tablet-sized visual displays with speakers built into the device and headphone listening are maintained.

**iTeh STANDARD PREVIEW**  
(standards.iteh.ai)  
Full standard:  
<https://standards.iteh.ai/catalog/standards/iso-iec-23008-3-2015/cb86-4e52-81ea-8dd020d7daf4/iso-iec-23008-3-2015>

**iTeh STANDARD PREVIEW**  
**(standards.iteh.ai)**

Full standard:  
<https://standards.iteh.ai/catalog/standards/sist/e4315f20-cb86-4e52-81ea-8dd020d7daf4/iso-iec-23008-3-2015>

# Information technology — High efficiency coding and media delivery in heterogeneous environments — Part 3: 3D audio

## 1 Scope

This part of ISO/IEC 23008-3 specifies technology which supports the efficient transmission of 3D audio signals and flexible rendering for the playback of 3D audio in a wide variety of listening scenarios. These include 3D home theater setups, 22.2 loudspeaker systems, automotive entertainment systems and playback over headphones connected to a tablet or smartphone.

## 2 Normative references

The following documents, in whole or in part, are normatively referenced in this document and are indispensable for its application. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 13818-1:2013, *Information technology — Generic Coding of moving pictures and associated audio information: Systems*

ISO/IEC 14496-3:2009, *Information technology — Coding of audio-visual objects — Part 3: Audio*

ISO/IEC 14496-11, *Information technology — Coding of audio-visual objects — Part 11: Scene description and application engine*

ISO/IEC 23001-8:2013, *Information technology — MPEG systems technologies — Part 8: Coding-independent code-points*

ISO/IEC 23001-8:2013/Amd.1, *Information technology — MPEG systems technologies — Part 8: Coding-independent code-points, AMENDMENT 1: New audio code points*

ISO/IEC 23003-1:2007, *Information technology — MPEG audio technologies — Part 1: MPEG Surround*

ISO/IEC 23003-2:2010, *Information technology — MPEG audio technologies — Part 2: Spatial Audio Object Coding (SAOC)*

ISO/IEC 23003-3:2012, *Information technology — MPEG audio technologies — Part 3: Unified speech and audio coding*

ISO/IEC 23003-4:2015, *Information technology — MPEG audio technologies — Part 4: Dynamic range control*

## 3 Terms, definitions and mnemonics

### 3.1 Terms and Definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 14496-3:2009, 1.3 (Terms and definitions), in ISO/IEC 14496-3:2009, 1.4 (Symbols and abbreviations) and in ISO/IEC 23003-3:2012, 3.1 (Terms and definitions) apply.

### 3.2 Mnemonics

The following mnemonics are defined to describe the different data types used in the coded bitstream payload.

- bslbf Bit string, left bit first, where “left” is the order in which bit strings are written in ISO/IEC 14496. Bit strings are written as a string of 1s and 0s within single quote marks, for example '1000 0001'. Blanks within a bit string are for ease of reading and have no significance.
- uimsbf Unsigned integer, most significant bit first.
- vclbfc Variable length code, left bit first, where “left” refers to the order in which the variable length codes are written.
- tcimsbf Two's complement integer, most significant (sign) bit first.

## 4 Technical Overview

### 4.1 Decoder block diagram

The 3D Audio Codec System consists of an MPEG-H 3D Audio Core Codec for coding of channel, object and Higher Order Ambisonics (HOA) signals. The core codec is based on the MPEG-D USAC codec. To increase the efficiency for coding a large amount of objects, MPEG SAOC technology has been adopted. Several types of renderers perform the tasks of rendering objects to channels, rendering channels to a different loudspeaker setup, rendering HOA signals to the loudspeaker setup or rendering virtual loudspeaker channels or HOA components to headphones.

When object signals are explicitly transmitted or parametrically encoded using SAOC, the corresponding Object Metadata information is compressed and multiplexed into the 3D-Audio bitstream.

Figure 1 shows the different algorithmic blocks of the 3D-Audio system.

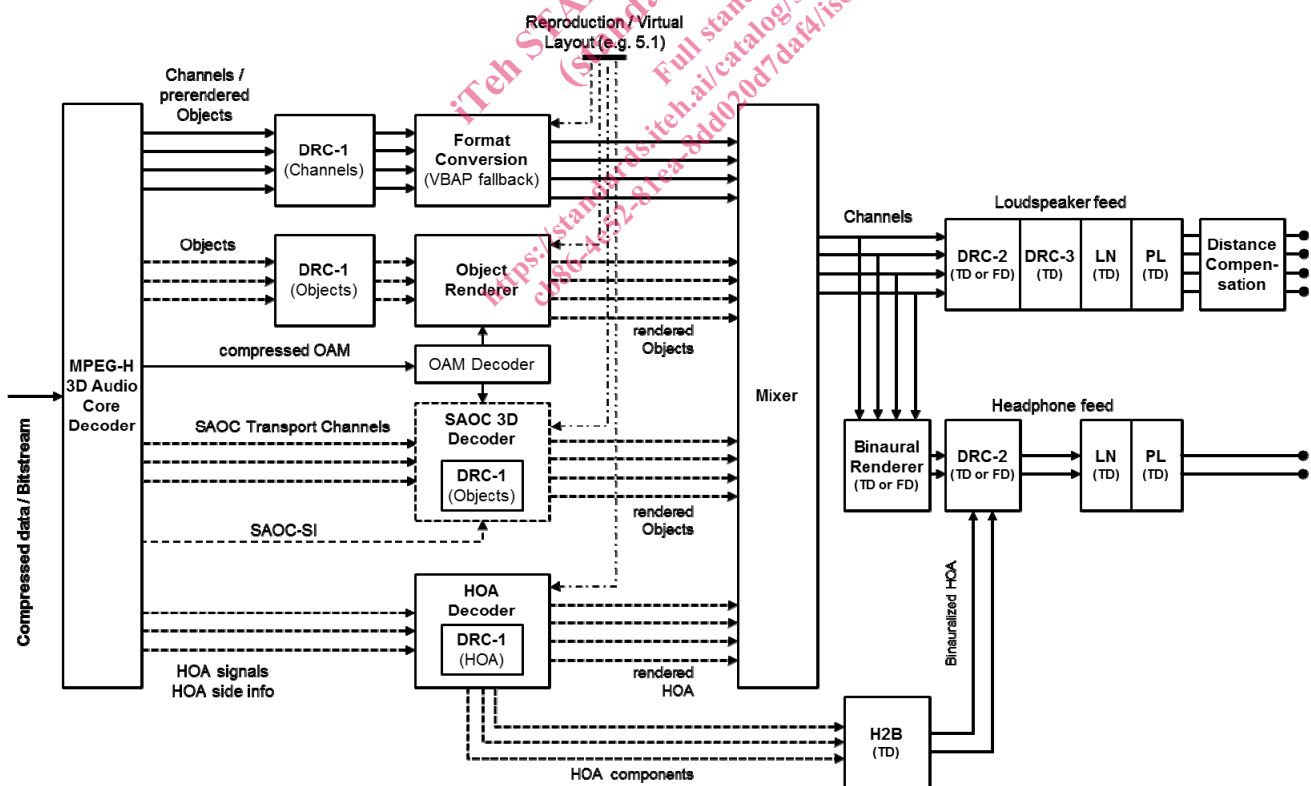


Figure 1 — Block diagram of the 3D-Audio decoder.

(DRC: Dynamic Range Control, SAOC: Spatial Audio Object Coding, HOA: Higher Order Ambisonics, LN: Loudness Normalization, PL: Peak Limiter)

## 4.2 Overview over the codec building blocks

The MPEG-H 3DA Core Codec for loudspeaker-channel signals, discrete object signals, object downmix signals and pre-rendered signals is based on MPEG-D USAC technology. It handles the coding of the multitude of signals by creating channel- and object-mapping information based on the geometric and semantic information of the input's channel and object assignment. This mapping information describes how input channels and objects are mapped to channel elements (CPEs, SCEs, LFEs) and the corresponding information is transmitted to the decoder.

The coding of objects is possible in different ways, depending on the rate/distortion requirements and the interactivity requirements for the renderer. The following object coding variants are possible:

- Prerendered objects: Object signals are pre-rendered and mixed to multi-channel or HOA signals before encoding, as appropriate. The subsequent coding chain then operates on multi-channel or HOA signals.
- Discrete object waveforms: Objects are supplied as monophonic waveforms to the encoder. The encoder uses single channel elements SCEs to transmit the objects in addition to the channel signals. The decoded objects are rendered and mixed at the receiver side. Compressed object metadata information is transmitted to the receiver/renderer alongside.
- Parametric object waveforms: Object properties and their relation to each other are described by means of SAOC parameters. The downmix of the object signals is coded with the MPEG-H 3D Audio Core codec. The parametric information is transmitted alongside. The number of downmix channels is chosen depending on the number of objects and the overall data rate. Compressed object metadata information is transmitted to the SAOC renderer.

The SAOC Encoder and Decoder for object signals are based on MPEG SAOC technology. The system is capable of recreating, modifying and rendering a number of audio objects based on a smaller number of transmitted channels and additional parametric data (OLDs, IOCs, DMGs).

The SAOC decoder reconstructs the object/channel signals from the decoded SAOC transport channels and parametric information, and generates the output audio scene based on the reproduction layout, the decompressed object metadata information and optionally on the user interaction information.

The Object Metadata Codec efficiently codes the associated metadata that specifies the geometrical position and volume of each object in 3D space by quantization of the object properties in time and space. The compressed object metadata is transmitted to the receiver as side information.

The Object Renderer utilizes the compressed object metadata to generate object waveforms according to the given reproduction format. Each object is rendered to certain output channels according to its metadata. The output of this block results from the sum of the partial results.

The Loudspeaker Renderer converts between the transmitted channel configuration and the desired reproduction format. It is thus called 'format converter'. In case of conversions to lower numbers of output channels it creates downmixes. The system automatically generates optimized downmix matrices for the given combination of input and output formats and applies these matrices in a downmix process. The format converter allows for standard loudspeaker configurations as well as for random configurations with non-standard loudspeaker positions.

The Higher Order Ambisonics (HOA) Decoder/Renderer reconstructs the HOA coefficient signals based on the HOA transport channels decoded by the 3D Audio Core Decoder and the HOA specific side information. The coding principle is based on a separate transmission of so-called predominant sounds and ambient sound scene components. Subsequently the HOA renderer generates the loudspeaker channel feeds based on the reproduction layout.