

---

---

**Information technology — Coding of  
audio-visual objects —**

**Part 1:  
Systems**

*Technologies de l'information — Codage des objets audiovisuels —*

*Partie 1: Systèmes*

**PDF disclaimer**

This PDF file may contain embedded typefaces. In accordance with Adobe's licensing policy, this file may be printed or viewed but shall not be edited unless the typefaces which are embedded are licensed to and installed on the computer performing the editing. In downloading this file, parties accept therein the responsibility of not infringing Adobe's licensing policy. The ISO Central Secretariat accepts no liability in this area.

Adobe is a trademark of Adobe Systems Incorporated.

Details of the software products used to create this PDF file can be found in the General Info relative to the file; the PDF-creation parameters were optimized for printing. Every care has been taken to ensure that the file is suitable for use by ISO member bodies. In the unlikely event that a problem relating to it is found, please inform the Central Secretariat at the address given below.

© ISO/IEC 2004

All rights reserved. Unless otherwise specified, no part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office  
Case postale 56 • CH-1211 Geneva 20  
Tel. + 41 22 749 01 11  
Fax + 41 22 749 09 47  
E-mail [copyright@iso.org](mailto:copyright@iso.org)  
Web [www.iso.org](http://www.iso.org)

Published in Switzerland

# Contents

Page

Foreword .....	vi
0 Introduction .....	viii
0.1 Overview .....	viii
0.2 Architecture .....	viii
0.3 Terminal Model: Systems Decoder Model .....	x
0.4 Multiplexing of Streams: The Delivery Layer .....	x
0.5 Synchronization of Streams: The Sync Layer .....	x
0.6 The Compression Layer .....	x
0.7 Application Engine .....	xii
0.8 Extensible MPEG-4 Textual Format (XMT) .....	xii
1 Scope .....	1
2 Normative references .....	1
3 Additional reference .....	2
4 Terms and definitions .....	2
4.1 Access Unit (AU) .....	2
4.2 Alpha Map .....	2
4.3 Atom .....	2
4.4 Audio-visual Object .....	2
4.5 Audio-visual Scene (AV Scene) .....	2
4.6 AVC Parameter Set .....	3
4.7 AVC Access Unit .....	3
4.8 AVC Parameter Set Access Unit .....	3
4.9 AVC Parameter Set Elementary Stream .....	3
4.10 AVC Video Elementary Stream .....	3
4.11 Binary Format for Scene (BIFS) .....	3
4.12 Buffer Model .....	3
4.13 Byte Aligned .....	3
4.14 Chunk .....	3
4.15 Clock Reference .....	3
4.16 Composition .....	3
4.17 Composition Memory (CM) .....	3
4.18 Composition Time Stamp (CTS) .....	3
4.19 Composition Unit (CU) .....	3
4.20 Compression Layer .....	4
4.21 Container Atom .....	4
4.22 Control Point .....	4
4.23 Decoder .....	4
4.24 Decoding buffer (DB) .....	4
4.25 Decoder configuration .....	4
4.26 Decoding Time Stamp (DTS) .....	4
4.27 Delivery Layer .....	4
4.28 Descriptor .....	4
4.29 DMIF Application Interface (DAI) .....	4
4.30 Elementary Stream (ES) .....	4
4.31 Elementary Stream Descriptor .....	4
4.32 Elementary Stream Interface (ESI) .....	4
4.33 M4Mux Channel (FMC) .....	4
4.34 M4Mux Packet .....	5
4.35 M4Mux Stream .....	5
4.36 M4Mux tool .....	5

4.37	Graphics Profile .....	5
4.38	Hint Track .....	5
4.39	Hinter .....	5
4.40	Inter .....	5
4.41	Interaction Stream .....	5
4.42	Intra .....	5
4.43	Initial Object Descriptor .....	5
4.44	Intellectual Property Identification (IPI) .....	5
4.45	Intellectual Property Management and Protection (IPMP) System .....	5
4.46	IPMP Information .....	6
4.47	IPMP System .....	6
4.48	IPMP Tool .....	6
4.49	IPMP Tool Identifier .....	6
4.50	IPMP Tool List .....	6
4.51	Media Node .....	6
4.52	Media stream .....	6
4.53	Media time line .....	6
4.54	Movie Atom .....	6
4.55	Movie Data Atom .....	6
4.56	MP4 File .....	6
4.57	Object Clock Reference (OCR) .....	6
4.58	Object Content Information (OCI) .....	7
4.59	Object Descriptor (OD) .....	7
4.60	Object Descriptor Command .....	7
4.61	Object Descriptor Profile .....	7
4.62	Object Descriptor Stream .....	7
4.63	Object Time Base (OTB) .....	7
4.64	Parametric Audio Decoder .....	7
4.65	Parametric Description .....	7
4.66	Quality of Service (QoS) .....	7
4.67	Random Access .....	7
4.68	Reference Point .....	7
4.69	Rendering .....	7
4.70	Rendering Area .....	7
4.71	Sample .....	8
4.72	Sample Table .....	8
4.73	Scene Description .....	8
4.74	Scene Description Stream .....	8
4.75	Scene Graph Elements .....	8
4.76	Scene Graph Profile .....	8
4.77	Seekable .....	8
4.78	SL-Packetized Stream (SPS) .....	8
4.79	Stream object .....	8
4.80	Structured Audio .....	8
4.81	Sync Layer (SL) .....	8
4.82	Sync Layer Configuration .....	8
4.83	Sync Layer Packet (SL-Packet) .....	8
4.84	Syntactic Description Language (SDL) .....	9
4.85	Systems Decoder Model (SDM) .....	9
4.86	System Time Base (STB) .....	9
4.87	Terminal .....	9
4.88	Time Base .....	9
4.89	Timing Model .....	9
4.90	Time Stamp .....	9
4.91	Track .....	9
4.92	Interaction Stream .....	9
5	Abbreviations and Symbols .....	9
6	Conventions .....	11

<b>7</b>	<b>Streaming Framework .....</b>	<b>11</b>
<b>7.1</b>	<b>Systems Decoder Model.....</b>	<b>11</b>
<b>7.2</b>	<b>Object Description Framework.....</b>	<b>17</b>
<b>7.3</b>	<b>Synchronization of Elementary Streams .....</b>	<b>72</b>
<b>7.4</b>	<b>Multiplexing of Elementary Streams .....</b>	<b>83</b>
<b>8</b>	<b>Syntactic Description Language .....</b>	<b>92</b>
<b>8.1</b>	<b>Introduction .....</b>	<b>92</b>
<b>8.2</b>	<b>Elementary Data Types.....</b>	<b>92</b>
<b>8.3</b>	<b>Composite Data Types .....</b>	<b>95</b>
<b>8.4</b>	<b>Arithmetic and Logical Expressions.....</b>	<b>99</b>
<b>8.5</b>	<b>Non-Parsable Variables .....</b>	<b>99</b>
<b>8.6</b>	<b>Syntactic Flow Control .....</b>	<b>99</b>
<b>8.7</b>	<b>Built-In Operators.....</b>	<b>101</b>
<b>8.8</b>	<b>Scoping Rules .....</b>	<b>101</b>
<b>9</b>	<b>Profiles .....</b>	<b>101</b>
<b>Annex A</b>	<b>(informative) Time Base Reconstruction .....</b>	<b>103</b>
<b>A.1</b>	<b>Time Base Reconstruction.....</b>	<b>103</b>
<b>A.2</b>	<b>Temporal aliasing and audio resampling .....</b>	<b>104</b>
<b>A.3</b>	<b>Reconstruction of a Synchronised Audio-visual Scene: A Walkthrough .....</b>	<b>105</b>
<b>Annex B</b>	<b>(informative) Registration procedure.....</b>	<b>106</b>
<b>B.1</b>	<b>Procedure for the request of a Registration ID (RID) .....</b>	<b>106</b>
<b>B.2</b>	<b>Responsibilities of the Registration Authority.....</b>	<b>106</b>
<b>B.3</b>	<b>Contact information for the Registration Authority .....</b>	<b>106</b>
<b>B.4</b>	<b>Responsibilities of Parties Requesting a RID .....</b>	<b>107</b>
<b>B.5</b>	<b>Appeal Procedure for Denied Applications.....</b>	<b>107</b>
<b>B.6</b>	<b>Registration Application Form.....</b>	<b>107</b>
<b>Annex C</b>	<b>(informative) The QoS Management Model for ISO/IEC 14496 Content .....</b>	<b>110</b>
<b>Annex D</b>	<b>(informative) Conversion Between Time and Date Conventions .....</b>	<b>111</b>
<b>D.1</b>	<b>Conversion Between Time and Date Conventions .....</b>	<b>111</b>
<b>Annex E</b>	<b>(informative) Graphical Representation of Object Descriptor and Sync Layer Syntax .....</b>	<b>113</b>
<b>E.1</b>	<b>Length encoding of descriptors and commands .....</b>	<b>113</b>
<b>E.2</b>	<b>Object Descriptor Stream and OD commands.....</b>	<b>114</b>
<b>E.3</b>	<b>OCI stream .....</b>	<b>114</b>
<b>E.4</b>	<b>Object descriptor and its components .....</b>	<b>115</b>
<b>E.5</b>	<b>OCI Descriptors .....</b>	<b>117</b>
<b>E.6</b>	<b>Sync layer configuration and syntax .....</b>	<b>120</b>
<b>Annex F</b>	<b>(informative) Elementary Stream Interface .....</b>	<b>121</b>
<b>Annex G</b>	<b>(informative) Upstream Walkthrough.....</b>	<b>123</b>
<b>G.1</b>	<b>Introduction .....</b>	<b>123</b>
<b>G.2</b>	<b>Configuration.....</b>	<b>123</b>
<b>G.3</b>	<b>Content access procedure with DAI .....</b>	<b>124</b>
<b>G.4</b>	<b>Example.....</b>	<b>124</b>
<b>Annex H</b>	<b>(informative) Scene and Object Description Carousel .....</b>	<b>128</b>
<b>Annex I</b>	<b>(normative) Usage of ITU-T Recommendation H.264   ISO/IEC 14496-10 AVC .....</b>	<b>129</b>
<b>I.1</b>	<b>SL packet encapsulation of AVC Access Unit .....</b>	<b>129</b>
<b>I.2</b>	<b>Handling of Parameter Sets .....</b>	<b>129</b>
<b>I.3</b>	<b>Usage of ISO/IEC 14496-14 AVC File Format in MPEG-4 Systems .....</b>	<b>130</b>
<b>Annex J</b>	<b>(informative) Patent statements .....</b>	<b>131</b>
<b>J.1</b>	<b>General .....</b>	<b>131</b>
<b>J.2</b>	<b>Patent Statements for Version 1.....</b>	<b>131</b>
<b>J.3</b>	<b>Patent Statements for Version 2.....</b>	<b>132</b>
	<b>Bibliography .....</b>	<b>134</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.

International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 2.

The main task of the joint technical committee is to prepare International Standards. Draft International Standards adopted by the joint technical committee are circulated to national bodies for voting. Publication as an International Standard requires approval by at least 75 % of the national bodies casting a vote.

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.

The International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC) draw attention to the fact that it is claimed that compliance with this document may involve the use of a patent.

The ISO and IEC take no position concerning the evidence, validity and scope of this patent right.

The holder of this patent right has assured the ISO and IEC that he is willing to negotiate licences under reasonable and non-discriminatory terms and conditions with applicants throughout the world. In this respect, the statement of the holder of this patent right is registered with the ISO and IEC. Information may be obtained from the companies listed in Annex J.

ISO/IEC 14496-1 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 29, *Coding of audio, picture, multimedia and hypermedia information*.

This third edition cancels and replaces the second edition (ISO/IEC 14496-1:2001). It also incorporates the Amendments ISO/IEC 14496-1:2001/Amd.1:2001, ISO/IEC 14496-1:2001/Amd.3:2004, ISO/IEC 14496-1:2001/Amd.4:2003 and ISO/IEC 14496-1:2001/Amd.7:2004, which have been technically revised.

ISO/IEC 14496 consists of the following parts, under the general title *Information technology — Coding of audio-visual objects*:

- *Part 1: Systems*
- *Part 2: Visual*
- *Part 3: Audio*
- *Part 4: Conformance testing*
- *Part 5: Reference software*
- *Part 6: Delivery Multimedia Integration Framework (DMIF)*
- *Part 7: Optimized reference software for coding of audio-visual objects*

- *Part 8: Carriage of ISO/IEC 14496 contents over IP networks*
- *Part 9: Reference hardware description*
- *Part 10: Advanced Video Coding*
- *Part 11: Scene description and application engine*
- *Part 12: ISO base media file format*
- *Part 13: Intellectual Property Management and Protection (IPMP) extensions*
- *Part 14: MP4 file format*
- *Part 15: Advanced Video Coding (AVC) file format*
- *Part 16: Animation Framework eXtension (AFX)*
- *Part 17: Streaming text format*
- *Part 18: Font compression and streaming*
- *Part 19: Synthesized texture stream*

## 0 Introduction

### 0.1 Overview

ISO/IEC 14496 specifies a system for the communication of interactive audio-visual scenes. This specification includes the following elements:

1. the coded representation of natural or synthetic, two-dimensional (2D) or three-dimensional (3D) objects that can be manifested audibly and/or visually (audio-visual objects) (specified in part 2, 3, 10, 11 and 16 of ISO/IEC 14496);
2. the coded representation of the spatio-temporal positioning of audio-visual objects as well as their behavior in response to interaction (scene description, specified in part 11 of ISO/IEC 14496);
3. the coded representation of information related to the management of data streams (synchronization, identification, description and association of stream content, specified in this part of ISO/IEC 14496);
4. a generic interface to the data stream delivery layer functionality (specified in part 6 of ISO/IEC 14496);
5. an application engine for programmatic control of the player: format, delivery of downloadable Java byte code as well as its execution lifecycle and behavior through APIs (specified in part 11 of ISO/IEC 14496); and
6. a file format to contain the media information of an ISO/IEC 14496 presentation in a flexible, extensible format to facilitate interchange, management, editing, and presentation of the media specified in part 12 (ISO File Format), part 14 (MP4 File Format) and part 15 (AVC File Format) of ISO/IEC 14496.

The overall operation of a system communicating audio-visual scenes can be paraphrased as follows:

At the sending terminal, the audio-visual scene information is compressed, supplemented with synchronization information and passed to a delivery layer that multiplexes it into one or more coded binary streams that are transmitted or stored. At the receiving terminal, these streams are demultiplexed and decompressed. The audio-visual objects are composed according to the scene description and synchronization information and presented to the end user. The end user may have the option to interact with this presentation. Interaction information can be processed locally or transmitted back to the sending terminal. ISO/IEC 14496 defines the syntax and semantics of the bitstreams that convey such scene information, as well as the details of their decoding processes.

This part of ISO/IEC 14496 specifies the following tools:

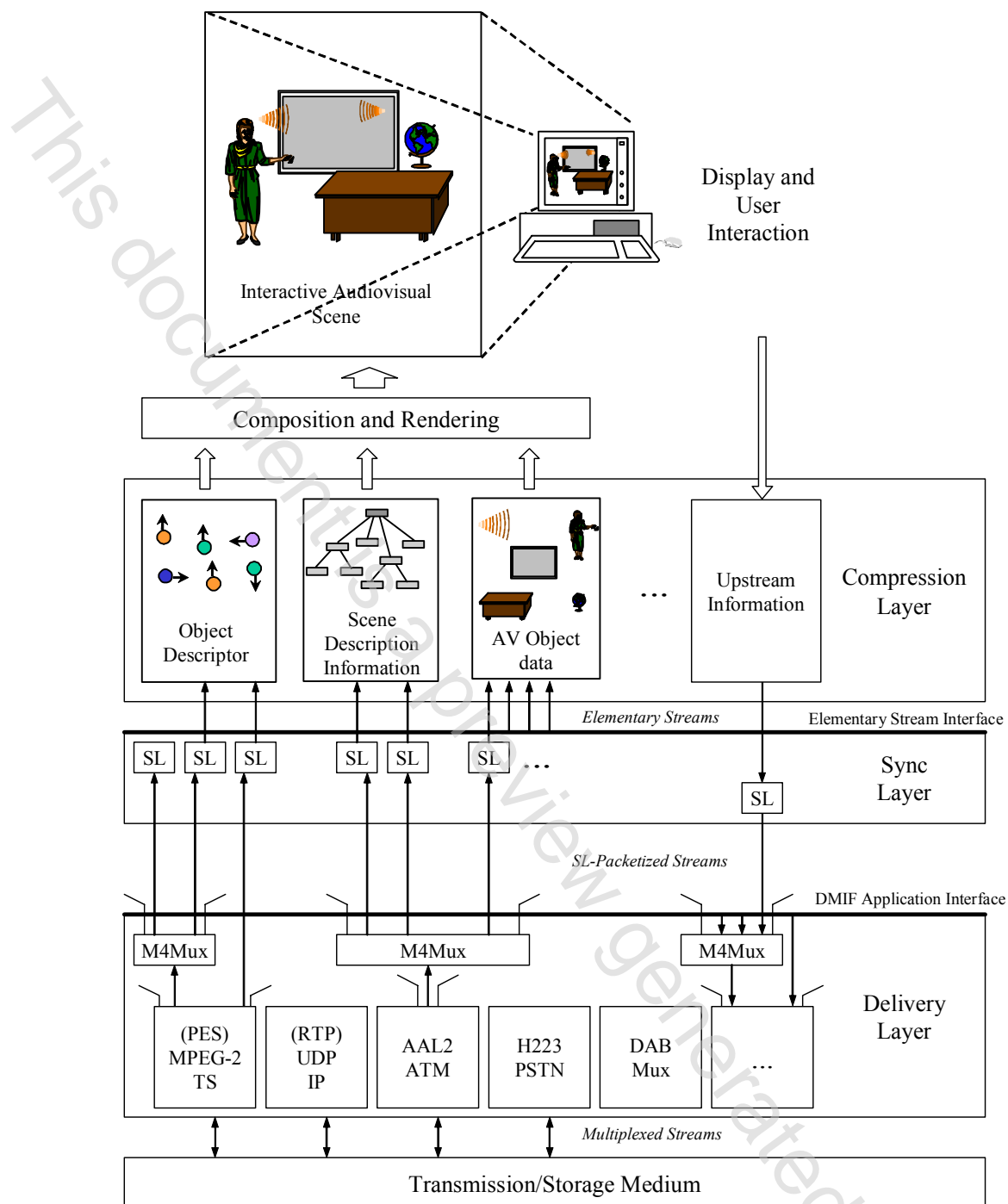
- a terminal model for time and buffer management;
- a coded representation of metadata for the identification, description and logical dependencies of the elementary streams (object descriptors and other descriptors);
- a coded representation of descriptive audio-visual content information (object content information – OCI);
- an interface to intellectual property management and protection (IPMP) systems;
- a coded representation of synchronization information (sync layer – SL); and
- a multiplexed representation of individual elementary streams in a single stream (M4Mux).

These various elements are described functionally in this subclause and specified in the normative clauses that follow.

### 0.2 Architecture

The information representation specified in ISO/IEC 14496 describes the means to create an interactive audio-visual scene in terms of coded audio-visual information and associated scene description information. The entity that composes and sends, or receives and presents such a coded representation of an interactive audio-visual scene is generically referred to as an "audio-visual terminal" or just "terminal". This terminal may correspond to a standalone application or be part of an application system.





**Figure 1 — The ISO/IEC 14496 Terminal Architecture**

The basic operations performed by such a receiver terminal are as follows. Information that allows access to content complying with ISO/IEC 14496 is provided as initial session set up information to the terminal. Part 6 of ISO/IEC 14496 defines the procedures for establishing such session contexts as well as the interface to the delivery layer that generically abstracts the storage or transport medium. The initial set up information allows, in a recursive manner, to locate one or more elementary streams that are part of the coded content representation. Some of these elementary streams may be grouped together using the multiplexing tool described in ISO/IEC 14496-1.

Elementary streams contain the coded representation of either audio or visual data or scene description information or user interaction data. Elementary streams may as well themselves convey information to identify streams, to describe logical dependencies between streams, or to describe information related to the content of the streams. Each elementary stream contains only one type of data.

Elementary streams are decoded using their respective stream-specific decoders. The audio-visual objects are composed according to the scene description information and presented by the terminal's presentation device(s). All these processes are synchronized according to the systems decoder model (SDM) using the synchronization information provided at the synchronization layer.

These basic operations are depicted in Figure 1, and are described in more detail below.

### 0.3 Terminal Model: Systems Decoder Model

The systems decoder model provides an abstract view of the behavior of a terminal complying with ISO/IEC 14496-1. Its purpose is to enable a sending terminal to predict how the receiving terminal will behave in terms of buffer management and synchronization when reconstructing the audio-visual information that comprises the presentation. The systems decoder model includes a systems timing model and a systems buffer model which are described briefly in the following subclauses.

#### 0.3.1 Timing Model

The timing model defines the mechanisms through which a receiving terminal establishes a notion of time that enables it to process time-dependent events. This model also allows the receiving terminal to establish mechanisms to maintain synchronization both across and within particular audio-visual objects as well as with user interaction events. In order to facilitate these functions at the receiving terminal, the timing model requires that the transmitted data streams contain implicit or explicit timing information. Two sets of timing information are defined in ISO/IEC 14496-1: clock references and time stamps. The former convey the sending terminal's time base to the receiving terminal, while the latter convey a notion of relative time for specific events such as the desired decoding or composition time for portions of the encoded audio-visual information.

#### 0.3.2 Buffer Model

The buffer model enables the sending terminal to monitor and control the buffer resources that are needed to decode each elementary stream in a presentation. The required buffer resources are conveyed to the receiving terminal by means of descriptors at the beginning of the presentation. The terminal can then decide whether or not it is capable of handling this particular presentation. The buffer model allows the sending terminal to specify when information may be removed from these buffers and enables it to schedule data transmission so that the appropriate buffers at the receiving terminal do not overflow or underflow.

### 0.4 Multiplexing of Streams: The Delivery Layer

The term delivery layer is used as a generic abstraction of any existing transport protocol stack that may be used to transmit and/or store content complying with ISO/IEC 14496. The functionality of this layer is not within the scope of ISO/IEC 14496-1, and only the interface to this layer is considered. This interface is the DMIF Application Interface (DAI) specified in ISO/IEC 14496-6. The DAI defines not only an interface for the delivery of streaming data, but also for signaling information required for session and channel set up as well as tear down. A wide variety of delivery mechanisms exist below this interface, with some of them indicated in Figure 1. These mechanisms serve for transmission as well as storage of streaming data, i.e., a file is considered to be a particular instance of a delivery layer. For applications where the desired transport facility does not fully address the needs of a service according to the specifications in ISO/IEC 14496, a simple multiplexing tool (M4Mux) with low delay and low overhead is defined in ISO/IEC 14496-1.

### 0.5 Synchronization of Streams: The Sync Layer

Elementary streams are the basic abstraction for any streaming data source. Elementary streams are conveyed as sync layer-packetized (SL-packetized) streams at the DMIF Application Interface. This packetized representation additionally provides timing and synchronization information, as well as fragmentation and random access information. The sync layer (SL) extracts this timing information to enable synchronized decoding and, subsequently, composition of the elementary stream data.

### 0.6 The Compression Layer

The compression layer receives data in its encoded format and performs the necessary operations to decode this data. The decoded information is then used by the terminal's composition, rendering and presentation subsystems.

#### 0.6.1 Object Description Framework

The purpose of the object description framework is to identify and describe elementary streams and to associate them appropriately to an audio-visual scene description. Object descriptors serve to gain access to ISO/IEC 14496 content.

Object content information and the interface to intellectual property management and protection systems are also part of this framework.

An object descriptor is a collection of one or more elementary stream descriptors that provide the configuration and other information for the streams that relate to either an audio-visual object or a scene description. Object descriptors are themselves conveyed in elementary streams. Each object descriptor is assigned an identifier (object descriptor ID), which is unique within a defined name scope. This identifier is used to associate audio-visual objects in the scene description with a particular object descriptor, and thus the elementary streams related to that particular object.

Elementary stream descriptors include information about the source of the stream data, in form of a unique numeric identifier (the elementary stream ID) or a URL pointing to a remote source for the stream. Elementary stream descriptors also include information about the encoding format, configuration information for the decoding process and the sync layer packetization, as well as quality of service requirements for the transmission of the stream and intellectual property identification. Dependencies between streams can also be signaled within the elementary stream descriptors. This functionality may be used, for example, in scalable audio or visual object representations to indicate the logical dependency of a stream containing enhancement information, to a stream containing the base information. It can also be used to describe alternative representations for the same content (e.g. the same speech content in various languages).

### 0.6.1.1 Intellectual Property Management and Protection

The intellectual property management and protection (IPMP) framework for ISO/IEC 14496 content consists of a normative interface that permits an ISO/IEC 14496 terminal to host one or more IPMP Systems in the form of monolithic IPMP Systems or modular IPMP Tools. The IPMP interface consists of IPMP elementary streams and IPMP descriptors. IPMP descriptors are carried as part of an object descriptor stream. IPMP elementary streams carry time variant IPMP information that can be associated to multiple object descriptors.

The IPMP System, or IPMP Tools themselves are non-normative components that provides intellectual property management and protection functions for the terminal. The IPMP Systems or Tools uses the information carried by the IPMP elementary streams and descriptors to make protected ISO/IEC 14496 content available to the terminal.

The intellectual property management and protection (IPMP) framework for ISO/IEC 14496 content consists of a set of tools that permits an ISO/IEC 14496 terminal to support IPMP functionality. This functionality is provided by two different complementary technologies, supporting different levels of interoperability:

- The IPMP framework as defined in 7.2.3, consists of a normative interface that permits an ISO/IEC 14496 terminal to host one or more IPMP Systems. The IPMP interface consists of IPMP elementary streams and IPMP descriptors. IPMP descriptors are carried as part of an object descriptor stream. IPMP elementary streams carry time variant IPMP information that can be associated to multiple object descriptors. The IPMP System itself is a non-normative component that provides intellectual property management and protection functions for the terminal. The IPMP System uses the information carried by the IPMP elementary streams and descriptors to make protected ISO/IEC 14496 content available to the terminal.
- The IPMP framework extension, as specified in ISO/IEC 14496-13 allows, in addition to the functionality specified in ISO/IEC 14496-1, a finer granularity of governance. ISO/IEC 14496-13 provides normative support for individual IPMP components, referred to as IPMP Tools, to be normatively placed at identified points of control within the terminal systems model. Additionally ISO/IEC 14496-13 provides normative support for secure communications to be performed between IPMP Tools. ISO/IEC 14496-1 also specifies specific normative extensions at the Systems level to support the IPMP functionality described in ISO/IEC 14496-13.

An application may choose not to use an IPMP System, thereby offering no management and protection features.

### 0.6.1.2 Object Content Information

Object content information (OCI) descriptors convey descriptive information about audio-visual objects. The main content descriptors are: content classification descriptors, keyword descriptors, rating descriptors, language descriptors, textual descriptors, and descriptors about the creation of the content. OCI descriptors can be included directly in the related object descriptor or elementary stream descriptor or, if it is time variant, it may be carried in an elementary stream by itself. An OCI stream is organized in a sequence of small, synchronized entities called events that contain a set of OCI descriptors. OCI streams can be associated to multiple object descriptors.

### 0.6.2 Scene Description Streams

Scene description addresses the organization of audio-visual objects in a scene, in terms of both spatial and temporal attributes. This information allows the composition and rendering of individual audio-visual objects after the respective decoders have reconstructed the streaming data for them. For visual data, ISO/IEC 14496-11 does not mandate particular composition algorithms. Hence, visual composition is implementation dependent. For audio data, the composition process is defined in a normative manner in ISO/IEC 14496-11 and ISO/IEC 14496-3.

The scene description is represented using a parametric approach (BIFS - Binary Format for Scenes). The description consists of an encoded hierarchy (tree) of nodes with attributes and other information (including event sources and targets). Leaf nodes in this tree correspond to elementary audio-visual data, whereas intermediate nodes group this material to form audio-visual objects, and perform grouping, transformation, and other such operations on audio-visual objects (scene description nodes). The scene description can evolve over time by using scene description updates.

In order to facilitate active user involvement with the presented audio-visual information, ISO/IEC 14496-11 provides support for user and object interactions. Interactivity mechanisms are integrated with the scene description information, in the form of linked event sources and targets (routes) as well as sensors (special nodes that can trigger events based on specific conditions). These event sources and targets are part of scene description nodes, and thus allow close coupling of dynamic and interactive behavior with the specific scene at hand. ISO/IEC 14496-11, however, does not specify a particular user interface or a mechanism that maps user actions (e.g., keyboard key presses or mouse movements) to such events.

Such an interactive environment may not need an upstream channel, but ISO/IEC 14496 also provides means for client-server interactive sessions with the ability to set up upstream elementary streams and associate them to specific downstream elementary streams.

### 0.6.3 Audio-visual Streams

The coded representation of audio and visual information are described in ISO/IEC 14496-3 (Audio) and ISO/IEC 14496-2 (Visual) and ISO/IEC 14496-10 (Advanced Video Coding) respectively. The reconstructed audio-visual data are made available to the composition process for potential use during the scene rendering.

### 0.6.4 Upchannel Streams

Downchannel elementary streams may require upchannel information to be transmitted from the receiving terminal to the sending terminal (e.g., to allow for client-server interactivity). Figure 1 indicates the flowpath for an elementary stream from the receiving terminal to the sending terminal. The content of upchannel streams is specified in the same part of the specification that defines the content of the downstream data. For example, upchannel control streams for video downchannel elementary streams are defined in ISO/IEC 14496-2.

### 0.6.5 Interaction Streams

The coded representation of user interaction information is not in the scope of ISO/IEC 14496. But this information shall be translated into scene modification and the modifications made available to the composition process for potential use during the scene rendering.

## 0.7 Application Engine

The MPEG-J is a programmatic system (as opposed to a conventional parametric system) which specifies API(s) for interoperation of MPEG-4 media players with Java code. By combining MPEG-4 media and safe executable code, content creators may embed complex control and data processing mechanisms with their media data to intelligently manage the operation of the audio-visual session. The parametric MPEG-4 System forms the Presentation Engine while the MPEG-J subsystem controlling the Presentation Engine forms the Application Engine.

The Java application is delivered as a separate elementary stream to the MPEG-4 terminal. There it will be directed to the MPEG-J run time environment, from where the MPEG-J program will have access to the various components and required data of the MPEG-4 player to control it.

In addition to the basic packages of the language (java.lang, java.io, java.util) a few categories of APIs have been defined for different scopes. For the Scene graph API the objective is to provide access to the scene graph: to inspect the graph, to alter nodes and their fields, and to add and remove nodes within the graph. The Resource API is used for regulation of performance: it provides a centralized facility for managing resources. This is used when the program execution is contingent upon the terminal configuration and its capabilities, both static (that do not change during execution) and dynamic. Decoder API allows the control of the decoders that are present in the terminal. The Net API provides a way to interact with the network, being compliant to the MPEG-4 DMIF Application Interface. Complex applications and enhanced interactivity are possible with these basic packages. The architecture of MPEG-J is presented in more detail in ISO/IEC 14496-11.

## 0.8 Extensible MPEG-4 Textual Format (XMT)

The Extensible MPEG-4 Textual (XMT) format is a textual representation of the multimedia content described in ISO/IEC 14496 using the Extensible Markup Language (XML). XMT is designed to facilitate the creation and maintenance of MPEG-4 multimedia content, whether by human authors or by automated machine programs. XMT is specified in ISO/IEC 14496-11.

The textual representation of MPEG-4 content has high-level abstractions, XMT-O, that allow authors to exchange their content easily with other authors or authoring tools, while at the same time preserving semantic intent. XMT also has low-level textual representations, XMT-A, covering the full scope and function of MPEG-4. The high-level XMT-O is designed to facilitate interoperability with the Synchronized Multimedia Integration Language (SMIL) 2.0, a recommendation from the W3C consortium, and also with Extensible 3D specification, X3D, developed by the Web3D consortium as the next generation of Virtual Reality Modeling Language (VRML).

The XMT language has grammars that are specified using the W3C XML Schema language. The grammars contain rules for element placement and attribute values, etc. These rules for XMT, defined using the Schema language, follow the binary coding rules defined in ISO/IEC 14496-11 and help ensure that the textual representation can be coded into correct binary according to ISO/IEC 14496-11 coding rules.

All constructs in the ISO/IEC 14496 specification have their parallel in the XMT textual format. For the Visual and Audio parts, XMT provides a means to reference external media streams of either pre-encoded or raw audiovisual binary content. While XMT does not contain a textual format for audiovisual media, it does contain hints in a textual format that allow an XMT tool to encode and embed the audiovisual media into a complete MPEG-4 presentation.



# Information technology — Coding of audio-visual objects —

## Part 1: Systems

### 1 Scope

This part of ISO/IEC 14496 specifies system level functionalities for the communication of interactive audio-visual scenes, i.e., the coded representation of information related to the management of data streams (synchronization, identification, description and association of stream content).

### 2 Normative references

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

ISO 639-2:1998, *Codes for the representation of names of languages — Part 2: Alpha-3 code*

ISO 3166-1:1997, *Codes for the representation of names of countries and their subdivisions — Part 1: Country codes*

ISO 9613-1:1993, *Acoustics — Attenuation of sound during propagation outdoors — Part 1: Calculation of the absorption of sound by the atmosphere*

ISO/IEC 10646-1:2000, *Information technology — Universal Multiple-Octet Coded Character Set (UCS) — Part 1: Architecture and Basic Multilingual Plane*

ISO/IEC 11172-2:1993, *Information technology — Coding of moving pictures and associated audio for digital storage media at up to about 1,5 Mbit/s — Part 2: Video*

ISO/IEC 11172-3:1993, *Information technology — Coding of moving pictures and associated audio for digital storage media at up to about 1,5 Mbit/s — Part 3: Audio*

ISO/IEC 13818-3:1998, *Information technology — Generic coding of moving pictures and associated audio information — Part 3: Audio*

ISO/IEC 13818-7:2003, *Information technology — Generic coding of moving pictures and associated audio information — Part 7: Advanced Audio Coding (AAC)*

ISO/IEC 14496-2:2004, *Information technology — Coding of audio-visual objects — Part 2: Visual*

ISO/IEC 14496-10:2003, *Information technology — Coding of audio-visual objects — Part 10: Advanced Video Coding*

ISO/IEC 14496-15:2004, *Information technology — Coding of audio-visual objects — Part 15: Advanced Video Coding (AVC) file format*

ISO/IEC 14772-1:1997, *Information technology — Computer graphics and image processing — The Virtual Reality Modeling Language — Part 1: Functional specification and UTF-8 encoding*

ISO/IEC 16262: 2002, *Information technology — ECMAScript language specification*

ITU-T Rec. H.262 (2000) | ISO/IEC 13818-2:2000, *Information technology — Generic coding of moving pictures and associated audio information — Part 2: Video*

ITU-T Rec. T.81 (1992) | ISO/IEC 10918-1:1994, *Information technology — Digital compression and coding of continuous-tone still images — Part 1: Requirements and guidelines*

IEEE Std 754-1985, *Standard for Binary Floating-Point Arithmetic*

Addison-Wesley: September 1996, *The Java Language Specification*, by James Gosling, Bill Joy and Guy Steele, ISBN 0-201-63451-1

Addison-Wesley: September 1996, *The Java Virtual Machine Specification*, by T. Lindholm and F. Yellin, ISBN 0-201-63452-X

Addison-Wesley: July 1998, *Java Class Libraries Vol. 1 The Java Class Libraries*, Second Edition Volume 1, by Patrick Chan, Rosanna Lee and Douglas Kramer, ISBN 0-201-31002-3

Addison-Wesley: July 1998, *Java Class Libraries Vol. 2 The Java Class Libraries*, Second Edition Volume 2, by Patrick Chan and Rosanna Lee, ISBN 0-201-31003-1

Addison-Wesley, May 1996, *Java API, The Java Application Programming Interface, Volume 1: Core Packages*, by J. Gosling, F. Yellin and the Java Team, ISBN 0-201-63453-8

DAVIC 1.4.1 specification, *Part 9: Information Representation*

ANSI/SMPTE 291M-1996, *Television — Ancillary Data Packet and Space Formatting*

SMPTE 315M -1999, *Television — Camera Positioning Information Conveyed by Ancillary Data Packets*

W3C Recommendation: August 2001 — *Synchronized Multimedia Integration Language (SMIL 2.0)*, <http://www.w3.org/TR/smil20/>

W3C Recommendation: May 2001 — *XML Schema*, <http://www.w3.org/TR/xmlschema-0/>

### 3 Additional reference

ISO/IEC 13522-6:1998, *Information technology — Coding of multimedia and hypermedia information — Part 6: Support for enhanced interactive applications*. This reference contains the full normative references to Java APIs and the Java Virtual Machine as described in the normative references above.

### 4 Terms and definitions

For the purposes of this part of ISO/IEC 14496, the following terms and definitions apply.

#### 4.1

##### **Access Unit (AU)**

An individually accessible portion of data within an *elementary stream*. An access unit is the smallest data entity to which timing information can be attributed.

#### 4.2

##### **Alpha Map**

The representation of the transparency parameters associated with a texture map.

#### 4.3

##### **Atom**

An object-oriented building block defined by a unique type identifier and length.

#### 4.4

##### **Audio-visual Object**

A representation of a natural or synthetic object that has an audio and/or visual manifestation. The representation corresponds to a node or a group of nodes in the BIFS scene description. Each audio-visual object is associated with zero or more *elementary streams* using one or more *object descriptors*.

#### 4.5

##### **Audio-visual Scene (AV Scene)**

A set of audio-visual objects together with scene description information that defines their spatial and temporal attributes including behaviors resulting from object and user interactions.