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

**Part 2:
Visual**

*Technologies de l'information — Codage des objets audiovisuels —
Partie 2: Codage visuel*
**iTeh STANDARD PREVIEW
(standards.iteh.ai)**

[ISO/IEC 14496-2:1999](https://standards.iso.org/iso-iec/14496-2:1999)

<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

Contents

1	Scope.....	1
2	Normative references.....	1
3	Definitions.....	2
4	Abbreviations and symbols	8
4.1	Arithmetic operators	9
4.2	Logical operators	9
4.3	Relational operators.....	9
4.4	Bitwise operators	10
4.5	Conditional operators	10
4.6	Assignment.....	10
4.7	Mnemonics.....	10
4.8	Constants.....	10
5	Conventions.....	10
5.1	Method of describing bitstream syntax.....	10
5.2	Definition of functions	12
5.2.1	Definition of next_bits() function.....	12
5.2.2	Definition of bytealigned() function.....	12
5.2.3	Definition of nextbits_bytealigned() function.....	12
5.2.4	Definition of next_start_code() function.....	12
5.2.5	Definition of next_resync_marker() function.....	12
5.2.6	Definition of transparent_mb() function	13
5.2.7	Definition of transparent_block() function	13

iTeh STANDARD PREVIEW
(standards.iteh.ai)

ISO/IEC 14496-2:1999

<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-21362908bba5/iso-iec-14496-2-1999>

5.3	Reserved, forbidden and marker_bit.....	13
5.4	Arithmetic precision	13
6	Visual bitstream syntax and semantics	13
6.1	Structure of coded visual data.....	13
6.1.1	Visual object sequence	14
6.1.2	Visual object	14
6.1.3	Video object	14
6.1.4	Mesh object.....	19
6.1.5	Face object.....	20
6.2	Visual bitstream syntax	24
6.2.1	Start codes.....	24
6.2.2	Visual Object Sequence and Visual Object	27
6.2.3	Video Object Layer	29
6.2.4	Group of Video Object Plane	34
6.2.5	Video Object Plane and Video Plane with Short Header	34
6.2.6	Macroblock	48
6.2.7	Block.....	54
6.2.8	Still Texture Object	55
6.2.9	Mesh Object	64
6.2.10	Face Object.....	67
6.3	Visual bitstream semantics	77
6.3.1	Semantic rules for higher syntactic structures.....	77
6.3.2	Visual Object Sequence and Visual Object	77
6.3.3	Video Object Layer	83
6.3.4	Group of Video Object Plane	91
6.3.5	Video Object Plane and Video Plane with Short Header	91
6.3.6	Macroblock related.....	101
6.3.7	Block related.....	104
6.3.8	Still texture object	104
6.3.9	Mesh object.....	109
6.3.10	Face object.....	112

iTeH STANDARD PREVIEW
(standards.iteh.ai)

ISO/IEC 14496-2:1999
<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

7	The visual decoding process	117
7.1	Video decoding process.....	117
7.2	Higher syntactic structures.....	118
7.3	VOP reconstruction.....	118
7.4	Texture decoding	119
7.4.1	Variable length decoding.....	119
7.4.2	Inverse scan.....	120
7.4.3	Intra dc and ac prediction for intra macroblocks.....	121
7.4.4	Inverse quantisation	123
7.4.5	Inverse DCT	126
7.5	Shape decoding.....	126
7.5.1	Higher syntactic structures.....	127
7.5.2	Macroblock decoding	127
7.5.3	Arithmetic decoding.....	136
7.5.4	Grayscale Shape Decoding.....	138
7.6	Motion compensation decoding.....	140
7.6.1	Padding process	141
7.6.2	Half sample interpolation	144
7.6.3	General motion vector decoding process	144
7.6.4	Unrestricted motion compensation.....	146
7.6.5	Vector decoding processing and motion-compensation in progressive P-VOP	146
7.6.6	Overlapped motion compensation	148
7.6.7	Temporal prediction structure	150
7.6.8	Vector decoding process of non-scalable progressive B-VOPs	150
7.6.9	Motion compensation in non-scalable progressive B-VOPs.....	151
7.7	Interlaced video decoding.....	155
7.7.1	Field DCT and DC and AC Prediction.....	155
7.7.2	Motion compensation	155
7.8	Sprite decoding	162
7.8.1	Higher syntactic structures.....	163
7.8.2	Sprite Reconstruction.....	163

iTeH STANDARD PREVIEW
(standards.iteh.ai)

ISO/IEC 14496-2:1999
<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

7.8.3	Low-latency sprite reconstruction	164
7.8.4	Sprite reference point decoding	165
7.8.5	Warping	165
7.8.6	Sample reconstruction	167
7.9	Generalized scalable decoding	167
7.9.1	Temporal scalability	169
7.9.2	Spatial scalability	172
7.10	Still texture object decoding	175
7.10.1	Decoding of the DC subband	175
7.10.2	ZeroTree Decoding of the Higher Bands	176
7.10.3	Inverse Quantization	181
7.11	Mesh object decoding	188
7.11.1	Mesh geometry decoding	188
7.11.2	Decoding of mesh motion vectors	191
7.12	Face object decoding	193
7.12.1	Frame based face object decoding	193
7.12.2	DCT based face object decoding	194
7.12.3	Decoding of the viseme parameter fap 1	195
7.12.4	Decoding of the viseme parameter fap 2	196
7.12.5	Fap masking	196
7.13	Output of the decoding process	196
7.13.1	Video data	197
7.13.2	2D Mesh data	197
7.13.3	Face animation parameter data	197
8	Visual-Systems Composition Issues	197
8.1	Temporal Scalability Composition	197
8.2	Sprite Composition	198
8.3	Mesh Object Composition	199
9	Profiles and Levels	199
9.1	Visual Object Types	200
9.2	Visual Profiles	202
9.3	Visual Profiles@Levels	202

9.3.1	Natural Visual	202
9.3.2	Synthetic Visual.....	202
9.3.3	Synthetic/Natural Hybrid Visual.....	203
Annex A (normative) Coding transforms		205
A.1	Discrete cosine transform for video texture.....	205
A.2	Discrete wavelet transform for still texture	205
A.2.1	Adding the mean	205
A.2.2	Wavelet filter	206
A.2.3	Symmetric extension	206
A.2.4	Decomposition level	207
A.2.5	Shape adaptive wavelet filtering and symmetric extension	207
Annex B (normative) Variable length codes and arithmetic decoding		209
B.1	Variable length codes	209
B.1.1	Macroblock type.....	209
B.1.2	Macroblock pattern	210
B.1.3	Motion vector	212
B.1.4	DCT coefficients.....	214
B.1.5	Shape Coding	227
B.1.6	Sprite Coding.....	233
B.1.7	DCT based facial object decoding.....	234
B.2	Arithmetic Decoding	246
B.2.1	Aritmetic decoding for still texture object	246
B.2.2	Arithmetic decoding for shape decoding	251
B.2.3	Face Object Decoding.....	254
Annex C (normative) Face object decoding tables and definitions		256
Annex D (normative) Video buffering verifier		269
D.1	Introduction	269
D.2	Video Rate Buffer Model Definition	269
D.3	Comparison between ISO/IEC 14496-2 VBV and the ISO/IEC 13818-2 VBV (Informative).....	272
D.4	Video Complexity Model Definition	273
D.5	Video Reference Memory Model Definition	274

iTech STANDARD PREVIEW
(standards.iteh.ai)

ISO/IEC 14496-2:1999
<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

D.6	Interaction between VBV, VCV and VMV (informative).....	274
D.7	Video Presentation Model Definition (informative).....	275
Annex E (informative) Features supported by the algorithm.....		277
E.1	Error resilience.....	277
E.1.1	Resynchronization.....	277
E.1.2	Data Partitioning.....	278
E.1.3	Reversible VLC.....	278
E.1.4	Decoder Operation.....	279
E.1.5	Adaptive Intra Refresh (AIR) Method.....	282
E.2	Complexity Estimation.....	284
E.3	Resynchronization in Case of Unknown Video Header Format.....	284
Annex F (informative) Preprocessing and postprocessing.....		286
F.1	Segmentation for VOP Generation.....	286
F.1.1	Introduction.....	286
F.1.2	Description of a combined temporal and spatial segmentation framework.....	286
F.1.3	References.....	288
F.2	Bounding Rectangle of VOP Formation.....	289
F.3	Postprocessing for Coding Noise Reduction.....	290
F.3.1	Deblocking filter.....	290
F.3.2	Deringing filter.....	292
F.3.3	Further issues.....	294
F.4	Chrominance Decimation and Interpolation Filtering for Interlaced Object Coding.....	294
Annex G (normative) Profile and level indication and restrictions.....		296
Annex H (informative) Patent statements.....		298
H.1	Patent statements.....	298
Annex I (informative) Bibliography.....		300
Annex J (normative) View dependent object scalability.....		301
J.1	Introduction.....	301
J.2	Decoding Process of a View-Dependent Object.....	301
J.2.1	General Decoding Scheme.....	301
J.2.2	Computation of the View-Dependent Scalability parameters.....	303
J.2.3	VD mask computation.....	304

J.2.4	Differential mask computation.....	305
J.2.5	DCT coefficients decoding.....	305
J.2.6	Texture update.....	305
J.2.7	IDCT.....	306
Annex K (normative) Decoder configuration information.....		307
K.1	Introduction.....	307
K.2	Description of the set up of a visual decoder (informative).....	307
K.2.1	Processing of decoder configuration information.....	308
K.3	Specification of decoder configuration information.....	309
K.3.1	VideoObject.....	309
K.3.2	StillTextureObject.....	309
K.3.3	MeshObject.....	309
K.3.4	FaceObject.....	310
Annex L (informative) Rate control.....		311
L.1	Frame Rate Control.....	311
L.1.1	Introduction.....	311
L.1.2	Description.....	311
L.1.3	Summary.....	314
L.2	Multiple Video Object Rate Control.....	314
L.2.1	Initialization.....	315
L.2.2	Quantization Level Calculation for I-frame and first P-frame.....	315
L.2.3	Update Rate-Distortion Model.....	317
L.2.4	Post-Frameskip Control.....	317
L.3	Macroblock Rate Control.....	319
L.3.1	Rate-Distortion Model.....	319
L.3.2	Target Number of Bits for Each Macroblock.....	319
L.3.3	Macroblock Rate Control.....	320
Annex M (informative) Binary shape coding.....		322
M.1	Introduction.....	322
M.2	Context-Based Arithmetic Shape Coding.....	322
M.2.1	Intra Mode.....	322

iTeh STANDARD PREVIEW
(standards.iteh.ai)

ISO/IEC 14496-2:1999
<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

M.2.2	Inter Mode	323
M.3	Texture Coding of Boundary Blocks	324
M.4	Encoder Architecture	324
M.5	Encoding Guidelines	325
M.5.1	Lossy Shape Coding	325
M.5.2	Coding Mode Selection	326
M.6	Conclusions	326
M.7	References	326
Annex N	(normative) Visual profiles@levels	328

iTeh STANDARD PREVIEW
(standards.iteh.ai)

[ISO/IEC 14496-2:1999](https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999)

<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

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.

International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 3.

In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1. 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.

International Standard ISO/IEC 14496-2 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 29, *Coding of audio, picture, multimedia and hypermedia information*.

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 testing

— Part 6: Delivery Multimedia Integration Framework (DMIF)

iTeh STANDARD PREVIEW
(standards.iteh.ai)

[ISO/IEC 14496-2:1999](https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999)

<https://standards.iteh.ai/catalog/standards/sist/abd25b6a-6026-4a3d-b31d-213b2908bba5/iso-iec-14496-2-1999>

Annexes A to D, G, J, K and N form a normative part of this part of ISO/IEC 14496. Annexes E, F, H, I, L and M are for information only.

Introduction

Purpose

This part of ISO/IEC 14496 was developed in response to the growing need for a coding method that can facilitate access to visual objects in natural and synthetic moving pictures and associated natural or synthetic sound for various applications such as digital storage media, internet, various forms of wired or wireless communication etc. The use of ISO/IEC 14496 means that motion video can be manipulated as a form of computer data and can be stored on various storage media, transmitted and received over existing and future networks and distributed on existing and future broadcast channels.

Application

The applications of ISO/IEC 14496 cover, but are not limited to, such areas as listed below:

IMM	Internet Multimedia
IVG	Interactive Video Games
IPC	Interpersonal Communications (videoconferencing, videophone, etc.)
ISM	Interactive Storage Media (optical disks, etc.)
MMM	Multimedia Mailing
NDB	Networked Database Services (via ATM, etc.)
RES	Remote Emergency Systems
RVS	Remote Video Surveillance
WMM	Wireless Multimedia

Profiles and levels

ISO/IEC 14496 is intended to be generic in the sense that it serves a wide range of applications, bitrates, resolutions, qualities and services. Furthermore, it allows a number of modes of coding of both natural and synthetic video in a manner facilitating access to individual objects in images or video, referred to as content based access. Applications should cover, among other things, digital storage media, content based image and video databases, internet video, interpersonal video communications, wireless video etc. In the course of creating ISO/IEC 14496, various requirements from typical applications have been considered, necessary algorithmic elements have been developed, and they have been integrated into a single syntax. Hence ISO/IEC 14496 will facilitate the bitstream interchange among different applications.

This part of ISO/IEC 14496 includes one or more complete decoding algorithms as well as a set of decoding tools. Moreover, the various tools of this part of ISO/IEC 14496 as well as that derived from ISO/IEC 13818-2 can be combined to form other decoding algorithms. Considering the practicality of implementing the full syntax of ISO/IEC 14496-2, however, a limited number of subsets of the syntax are also stipulated by means of “profile” and “level”.

A “profile” is a defined subset of the entire bitstream syntax that is defined by this part of ISO/IEC 14496. Within the bounds imposed by the syntax of a given profile it is still possible to require a very large variation in the performance of encoders and decoders depending upon the values taken by parameters in the bitstream.

In order to deal with this problem “levels” are defined within each profile. A level is a defined set of constraints imposed on parameters in the bitstream. These constraints may be simple limits on numbers. Alternatively they may take the form of constraints on arithmetic combinations of the parameters.

Object based coding syntax

Video object

A *video object* in a scene is an entity that a user is allowed to access (seek, browse) and manipulate (cut and paste). The instances of video objects at a given time are called *video object planes* (VOPs). The encoding process generates a coded representation of a VOP as well as composition information necessary for display. Further, at the decoder, a user may interact with and modify the composition process as needed.

The full syntax allows coding of rectangular as well as arbitrarily shaped video objects in a scene. Furthermore, the syntax supports both non-scalable coding and scalable coding. Thus it becomes possible to handle normal scalabilities as well as object based scalabilities. The scalability syntax enables the reconstruction of useful video from pieces of a total bitstream. This is achieved by structuring the total bitstream in two or more layers, starting from a standalone base layer and adding a number of enhancement layers. The base layer can be coded using a non-scalable syntax, or in the case of picture based coding, even using a syntax of a different video coding standard.

To ensure the ability to access individual objects, it is necessary to achieve a coded representation of its shape. A natural video object consists of a sequence of 2D representations (at different points in time) referred to here as VOPs. For efficient coding of VOPs, both temporal redundancies as well as spatial redundancies are exploited. Thus a coded representation of a VOP includes representation of its shape, its motion and its texture.

Face object

A 3D (or 2D) *face object* is a representation of the human face that is structured for portraying the visual manifestations of speech and facial expressions adequate to achieve visual speech intelligibility and the recognition of the mood of the speaker. A face object is animated by a stream of *face animation parameters* (FAP) encoded for low-bandwidth transmission in broadcast (one-to-many) or dedicated interactive (point-to-point) communications. The FAPs manipulate key feature control points in a mesh model of the face to produce animated visemes for the mouth (lips, tongue, teeth), as well as animation of the head and facial features like the eyes. FAPs are quantized with careful consideration for the limited movements of facial features, and then prediction errors are calculated and coded arithmetically. The remote manipulation of a face model in a terminal with FAPs can accomplish lifelike visual scenes of the speaker in real-time without sending pictorial or video details of face imagery every frame.

A simple streaming connection can be made to a decoding terminal that animates a default face model. A more complex session can initialize a custom face in a more capable terminal by downloading *face definition parameters* (FDP) from the encoder. Thus specific background images, facial textures, and head geometry can be portrayed. The composition of specific backgrounds, face 2D/3D meshes, texture attribution of the mesh, etc. is described in ISO/IEC 14496-1. The FAP stream for a given user can be generated at the user's terminal from video/audio, or from text-to-speech. FAPs can be encoded at bitrates up to 2-3kbit/s at necessary speech rates. Optional temporal DCT coding provides further compression efficiency in exchange for delay. Using the facilities of ISO/IEC 14496-1, a composition of the animated face model and synchronized, coded speech audio (low-bitrate speech coder or text-to-speech) can provide an integrated low-bandwidth audio/visual speaker for broadcast applications or interactive conversation.

Limited scalability is supported. Face animation achieves its efficiency by employing very concise motion animation controls in the channel, while relying on a suitably equipped terminal for rendering of moving 2D/3D faces with non-normative models held in local memory. Models stored and updated for rendering in the terminal can be simple or complex. To support speech intelligibility, the normative specification of FAPs intends for their selective or complete use as signaled by the encoder. A masking scheme provides for selective transmission of FAPs according to what parts of the face are naturally active from moment to moment. A further control in the FAP stream allows face animation to be suspended while leaving face features in the terminal in a defined quiescent state for higher overall efficiency during multi-point connections.

The Face Animation specification is defined in ISO/IEC 14496-1 and this part of ISO/IEC 14496. This clause is intended to facilitate finding various parts of specification. As a rule of thumb, FAP specification is found in the part 2, and FDP specification in the part 1. However, this is not a strict rule. For an overview of FAPs and their interpretation, read subclauses "6.1.5.2 Facial animation parameter set", "6.1.5.3 Facial animation parameter units", "6.1.5.4 Description of a neutral face" as well as the Table C-1. The viseme parameter is documented in subclause "7.12.3 Decoding of the viseme parameter fap 1" and the Table C-5 in annex C. The expression parameter is

documented in subclause “7.12.4 Decoding of the expression parameter fap 2” and the Table C-3. FAP bitstream syntax is found in subclauses “6.2.10 Face Object”, semantics in “6.3.10 Face Object”, and subclause “7.12 Face object decoding” explains in more detail the FAP decoding process. FAP masking and interpolation is explained in subclauses “6.3.11.1 Face Object Plane”, “7.12.1.1 Decoding of faps”, “7.12.5 Fap masking”. The FIT interpolation scheme is documented in subclause “7.2.5.3.2.4 FIT” of ISO/IEC 14496-1. The FDPs and their interpretation are documented in subclause “7.2.5.3.2.6 FDP” of ISO/IEC 14496-1. In particular, the FDP feature points are documented in the Figure C-1.

Mesh object

A 2D *mesh object* is a representation of a 2D deformable geometric shape, with which synthetic video objects may be created during a composition process at the decoder, by spatially piece-wise warping of existing video object planes or still texture objects. The instances of mesh objects at a given time are called *mesh object planes* (mops). The geometry of mesh object planes is coded losslessly. Temporally and spatially predictive techniques and variable length coding are used to compress 2D mesh geometry. The coded representation of a 2D mesh object includes representation of its geometry and motion.

Overview of the object based non-scalable syntax

The coded representation defined in the non-scalable syntax achieves a high compression ratio while preserving good image quality. Further, when access to individual objects is desired, the shape of objects also needs to be coded, and depending on the bandwidth available, the shape information can be coded lossy or losslessly.

The compression algorithm employed for texture data is not lossless as the exact sample values are not preserved during coding. Obtaining good image quality at the bitrates of interest demands very high compression, which is not achievable with intra coding alone. The need for random access, however, is best satisfied with pure intra coding. The choice of the techniques is based on the need to balance a high image quality and compression ratio with the requirement to make random access to the coded bitstream.

A number of techniques are used to achieve high compression. The algorithm first uses block-based motion compensation to reduce the temporal redundancy. Motion compensation is used both for causal prediction of the current VOP from a previous VOP, and for non-causal, interpolative prediction from past and future VOPs. Motion vectors are defined for each 16-sample by 16-line region of a VOP or 8-sample by 8-line region of a VOP as required. The prediction error, is further compressed using the discrete cosine transform (DCT) to remove spatial correlation before it is quantised in an irreversible process that discards the less important information. Finally, the shape information, motion vectors and the quantised DCT information, are encoded using variable length codes.

Temporal processing

Because of the conflicting requirements of random access to and highly efficient compression, three main VOP types are defined. Intra coded VOPs (I-VOPs) are coded without reference to other pictures. They provide access points to the coded sequence where decoding can begin, but are coded with only moderate compression. Predictive coded VOPs (P-VOPs) are coded more efficiently using motion compensated prediction from a past intra or predictive coded VOPs and are generally used as a reference for further prediction. Bidirectionally-predictive coded VOPs (B-VOPs) provide the highest degree of compression but require both past and future reference VOPs for motion compensation. Bidirectionally-predictive coded VOPs are never used as references for prediction (except in the case that the resulting VOP is used as a reference for scalable enhancement layer). The organisation of the three VOP types in a sequence is very flexible. The choice is left to the encoder and will depend on the requirements of the application.

Coding of Shapes

In natural video scenes, VOPs are generated by segmentation of the scene according to some semantic meaning. For such scenes, the shape information is thus binary (binary shape). Shape information is also referred to as alpha plane. The binary alpha plane is coded on a macroblock basis by a coder which uses the context information, motion compensation and arithmetic coding.

For coding of shape of a VOP, a bounding rectangle is first created and is extended to multiples of 16×16 blocks with extended alpha samples set to zero. Shape coding is then initiated on a 16×16 block basis; these blocks are also referred to as binary alpha blocks.

Motion representation - macroblocks

The choice of 16×16 blocks (referred to as macroblocks) for the motion-compensation unit is a result of the trade-off between the coding gain provided by using motion information and the overhead needed to represent it. Each macroblock can further be subdivided to 8×8 blocks for motion estimation and compensation depending on the overhead that can be afforded.

Depending on the type of the macroblock, motion vector information and other side information is encoded with the compressed prediction error in each macroblock. The motion vectors are differenced with respect to a prediction value and coded using variable length codes. The maximum length of the motion vectors allowed is decided at the encoder. It is the responsibility of the encoder to calculate appropriate motion vectors. The specification does not specify how this should be done.

Spatial redundancy reduction

Both source VOPs and prediction errors VOPs have significant spatial redundancy. This part of ISO/IEC 14496 uses a block-based DCT method with optional visually weighted quantisation, and run-length coding. After motion compensated prediction or interpolation, the resulting prediction error is split into 8×8 blocks. These are transformed into the DCT domain where they can be weighted before being quantised. After quantisation many of the DCT coefficients are zero in value and so two-dimensional run-length and variable length coding is used to encode the remaining DCT coefficients efficiently.

Chrominance formats

This part of ISO/IEC 14496 currently supports the 4:2:0 chrominance format.

Pixel depth

This part of ISO/IEC 14496 supports pixel depths between 4 and 12 bits in luminance and chrominance planes.

Generalized scalability

The scalability tools in this part of ISO/IEC 14496 are designed to support applications beyond that supported by single layer video. The major applications of scalability include internet video, wireless video, multi-quality video services, video database browsing etc. In some of these applications, either normal scalabilities on picture basis such as those in ISO/IEC 13818-2 may be employed or object based scalabilities may be necessary; both categories of scalability are enabled by this part of ISO/IEC 14496.

Although a simple solution to scalable video is the simulcast technique that is based on transmission/storage of multiple independently coded reproductions of video, a more efficient alternative is scalable video coding, in which the bandwidth allocated to a given reproduction of video can be partially re-utilised in coding of the next reproduction of video. In scalable video coding, it is assumed that given a coded bitstream, decoders of various complexities can decode and display appropriate reproductions of coded video. A scalable video encoder is likely to have increased complexity when compared to a single layer encoder. However, this part of ISO/IEC 14496 provides several different forms of scalabilities that address non-overlapping applications with corresponding complexities.

The basic scalability tools offered are temporal scalability and spatial scalability. Moreover, combinations of these basic scalability tools are also supported and are referred to as hybrid scalability. In the case of basic scalability, two layers of video referred to as the lower layer and the enhancement layer are allowed, whereas in hybrid scalability up to four layers are supported.

Object based Temporal scalability

Temporal scalability is a tool intended for use in a range of diverse video applications from video databases, internet video, wireless video and multiview/stereoscopic coding of video. Furthermore, it may also provide a migration path from current lower temporal resolution video systems to higher temporal resolution systems of the future.

Temporal scalability involves partitioning of VOPs into layers, where the lower layer is coded by itself to provide the basic temporal rate and the enhancement layer is coded with temporal prediction with respect to the lower layer. These layers when decoded and temporally multiplexed yield full temporal resolution. The lower temporal resolution systems may only decode the lower layer to provide basic temporal resolution whereas enhanced systems of the

future may support both layers. Furthermore, temporal scalability has use in bandwidth constrained networked applications where adaptation to frequent changes in allowed throughput are necessary. An additional advantage of temporal scalability is its ability to provide resilience to transmission errors as the more important data of the lower layer can be sent over a channel with better error performance, whereas the less critical enhancement layer can be sent over a channel with poor error performance. Object based temporal scalability can also be employed to allow graceful control of picture quality by controlling the temporal rate of each video object under the constraint of a given bit-budget.

Spatial scalability

Spatial scalability is a tool intended for use in video applications involving multi quality video services, video database browsing, internet video and wireless video, i.e., video systems with the primary common feature that a minimum of two layers of spatial resolution are necessary. Spatial scalability involves generating two spatial resolution video layers from a single video source such that the lower layer is coded by itself to provide the basic spatial resolution and the enhancement layer employs the spatially interpolated lower layer and carries the full spatial resolution of the input video source.

An additional advantage of spatial scalability is its ability to provide resilience to transmission errors as the more important data of the lower layer can be sent over a channel with better error performance, whereas the less critical enhancement layer data can be sent over a channel with poor error performance. Further, it can also allow interoperability between various standards.

Hybrid scalability

There are a number of applications where neither the temporal scalability nor the spatial scalability may offer the necessary flexibility and control. This may necessitate use of temporal and spatial scalability simultaneously and is referred to as the hybrid scalability. Among the applications of hybrid scalability are wireless video, internet video, multiviewpoint/stereoscopic coding etc.

Error Resilience

This part of ISO/IEC 14496 provides error robustness and resilience to allow accessing of image or video information over a wide range of storage and transmission media. The error resilience tools developed for this part of ISO/IEC 14496 can be divided into three major categories. These categories include synchronization, data recovery, and error concealment. It should be noted that these categories are not unique to this part of ISO/IEC 14496, and have been used elsewhere in general research in this area. It is, however, the tools contained in these categories that are of interest, and where this part of ISO/IEC 14496 makes its contribution to the problem of error resilience.

Patents

The International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC) draw attention to the fact that it is claimed that compliance with this part of ISO/IEC 14496 may involve the use of patents concerning the coded representation of picture information given in Annex H.

ISO and IEC take no position concerning the evidence, validity and scope of these patent rights.

The holders of these patent rights have assured ISO and IEC that they are willing to negotiate licences under reasonable and non-discriminatory terms and conditions with applicants throughout the world. In this respect, the statements of the holders of these patent rights are registered with ISO and IEC. Information may be obtained from the patent offices of the organizations listed in Annex H.

Attention is drawn to the possibility that some of the elements of this part of ISO/IEC 14496 may be the subject of patent rights other than those identified above. ISO and IEC shall not be held responsible for identifying any or all such patent rights.