

MPEG-4: An object-based multimedia coding standard supporting mobile applications

Atul Puri^a and Alexandros Eleftheriadis^b

^a *AT&T Laboratories, NSL 3-237, 100 Shultz Drive – Middletown, Red Bank, NJ 07701, USA*

^b *Department of Electrical Engineering, Columbia University, New York, NY 10027, USA*

The ISO MPEG committee, after successful completion of the MPEG-1 and the MPEG-2 standards is currently working on MPEG-4, the third MPEG standard. Originally, MPEG-4 was conceived to be a standard for coding of limited complexity audio-visual scenes at very low bit-rates; however, in July 1994, its scope was expanded to include coding of scenes as a collection of individual audio-visual objects and enabling a range of advanced functionalities not supported by other standards. One of the key functionalities supported by MPEG-4 is robustness in error prone environments. In general, the MPEG-4 standard provides solutions for coding of natural or synthetic video and audio, as well as a system for multiplex/demultiplex and description of scenes in a truly flexible manner. With focus on the mobile multimedia functionality, we present an overview of the current status as well as the details of the MPEG-4 coding standard. We also discuss profiles, a mechanism used for partitioning MPEG-4 into realizable subsets. Finally, plans for testing and verification of current MPEG-4 (Version 1) standard, ongoing work for MPEG-4 Version 2, as well as directions for MPEG-7, the next MPEG standard, are briefly discussed.

1. Introduction

The need for mobile communications is ever increasing due to the sense of timeliness and flexibilities it offers. Increasingly, in mobile communications, a diverse set of media such as speech, data, synthetic and natural images as well as synthetic and natural video are becoming necessary as mobile communications aims to supplement or replace traditional fixed communications. In general, multimedia is expensive in the sense of its bandwidth requirement, with video being highly bandwidth intensive. Efficient compression of video is therefore critical to making any form of multimedia feasible. Further, the feasibility of mobile multimedia of acceptable quality poses a significant additional challenge. This is so because wireless channels impose a fairly harsh environment for multimedia communications, and while the goal of compression is to squeeze redundancy out of signals to fit them on limited available bandwidth, the requirements for robust delivery necessitate some amount of redundancy. As in the case for wired or wireless environments, the success of multimedia terminals, products or services [8] depends on many factors, of particular significance is interworking which is facilitated by standardization.

Mobile multimedia applications can be classified into two primary classes, indoor and outdoor. Mobile indoor applications are characterized by lower mobility and higher bandwidth (about 1 Mbit/s or higher) while mobile outdoor applications typically tend to involve higher mobility (including higher speeds) and relatively lower bandwidths (a few kbit/s to a few tens of kbit/s or so). Of course, a number of other applications [43] in between these two extremes also exist. Considering the limitations of the existing standards when used in mobile environment, this is

an area of active research. However, the focus of this paper is to examine particular considerations for robustness in the state of the art standards being developed. As a passing reference, the ISO MPEG-1 video standard [18] was primarily optimized for coding of noninterlaced video at bit-rates of 1.2–1.5 Mbit/s and the ISO MPEG-2 video standard [15,19] was primarily optimized for coding of interlaced video at bit-rates of 4–9 Mbit/s. Furthermore, the MPEG-1 standard assumed a relatively error free channel and the MPEG-2 standard, due to its generic nature, only considered the very basic error resilience techniques such as slice synchronization, intra refresh, and a mechanism to facilitate error concealment, motion vectors for intra coded blocks.

The currently ongoing MPEG standard (MPEG-4) [3,23,40,42] was started in 1993 with intended completion by late 1998. Its original focus was modified in July 1994 from that of coding with high efficiency of videophone scenes at very low bit-rates, to flexible coding of generic scenes facilitating a number of important functionalities not supported by other standards. Among the functionalities [3] that were considered important for MPEG-4 were content-based coding, universal accessibility (which includes robustness to errors) and good coding efficiency. Further, MPEG-4 video is being optimized for bit-rates ranging from about 10 kbit/s to around 1.5 Mbit/s and is expected to be applicable to even higher bit-rates. Incidentally, the range of bit-rates discussed for MPEG-4 encompasses the bit-rates applicable to both indoor and outdoor mobile applications. It is worthwhile pointing out that the MPEG standards [18,19] are essentially decoding standards and thus only specify the bitstream representation and the semantics of the decoding process, in other words, the encoding algorithm is not standardized. Furthermore, unlike earlier standards,

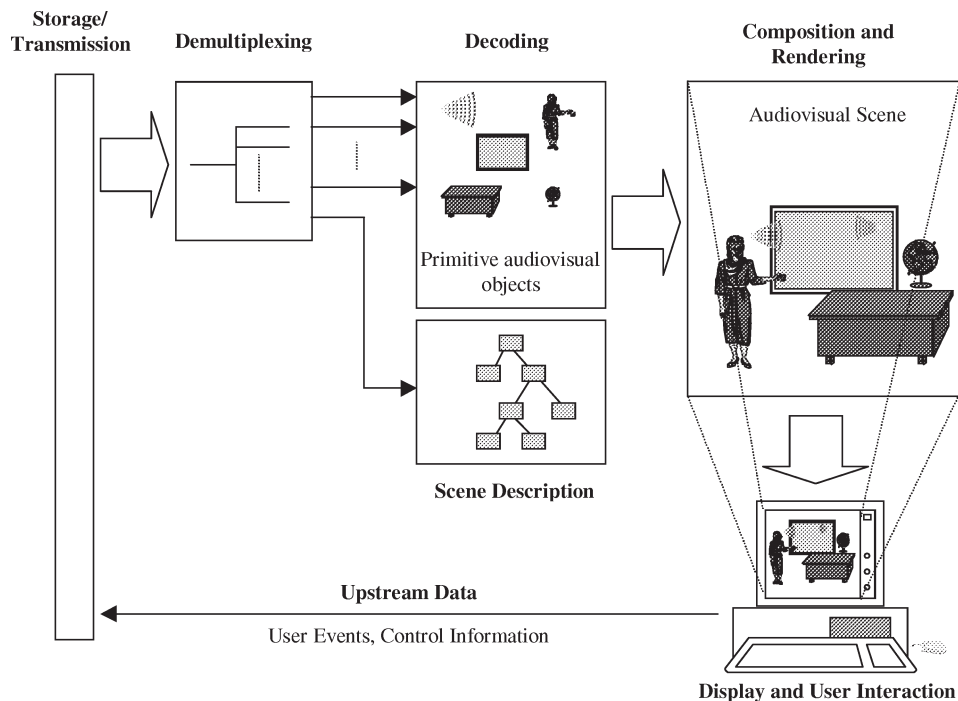


Figure 1. A high level view of an MPEG-4 terminal.

MPEG-4 is expected to consist of at least two versions – MPEG-4 Version 1 and Version 2; this paper mainly addresses the MPEG-4 Version 1 standard. Details of versioning of MPEG-4 are still evolving and the current status is discussed in section 9. The specification of MPEG-4 [28, 33,35], since MPEG-4 is designed to be a truly multimedia standard, goes much beyond that of previous MPEG standards and addresses not only audio coding [28], video coding [35] and multiplexing of coded data [33] but also coding of text/graphics and synthetic images [35] as well as flexible representation of audio-visual scene and composition [33].

Figure 1 shows a high level view of an MPEG-4 terminal [45,46,48]. We use the term “terminal” in a generic sense, including both standalone hardware as well as software running on general purpose computers. A set of individually coded audiovisual objects (natural or synthetic) are obtained multiplexed from a storage or transmission medium. They are accompanied with scene description information, which describes how these objects should be combined in space and time in order to form the scene intended by the content creator. The scene description is thus used during composition and rendering, which results in individual frames or audio samples being presented to the user. In addition, the user may have the option to interact with the content, either locally or with the source, using an upstream channel (if available).

The rest of the paper is organized as follows. In section 2, we present an overview of the ITU-T video and systems standards. In section 3, we review the applications, requirements, tests, and the organization of the MPEG-4 standard. Next, section 4 discusses MPEG-4 visual tools with emphasis on error resilience tools. In section 5, MPEG-4 audio standard is briefly discussed. In section 6, MPEG-4

systems is discussed in detail. In section 7, we discuss the issue of profiles, and in section 8 the plans for verification tests are presented. Section 9 discusses the work in progress for version 2 of MPEG-4, as well as the plans for MPEG-7, the next MPEG standard. Section 10 summarizes the key points presented in the paper.

2. Related ITU-T standards

Besides the ISO standards, the ITU-T has also developed video and audio coding as well as multiplex standards. The ITU-T H.263 standard [19] is aimed at coding of video at low bit-rates of 10–24 kbit/s and is based on the earlier ITU-T H.261 video standard [22] which was optimized at 64 kbit/s (although it allows a range of 64 kbit/s to 2 Mbit/s). In a general sense, the H.263 standard [23] uses the motion compensated DCT coding framework which is also common to the H.261, the MPEG-1 and the MPEG-2 standards. This consists of partitioning each picture into macroblocks, where a macroblock consists of 16×16 luminance (Y) block (composed of $4, 8 \times 8$ blocks) and the corresponding 8×8 chrominance blocks of Cb, and Cr. Each macroblock can be coded as intra (original signal) or as inter (prediction error signal). Spatial redundancy is exploited by DCT coding. Temporal redundancy is exploited by motion compensation which is used to determine the prediction error signal. Block DCT coefficients are quantized and entropy coded. Details of H.263 include motion compensation with accuracy of half-pixel (like that of MPEG-1, whereas H.261 supports only integer-pixel accuracy) as well as optional modes such as PB-frames (a substitute for B-pictures of MPEG-1), unrestricted motion

vector, advanced 8×8 block prediction, and syntax based arithmetic coding. Incidentally, these modes are options that are negotiated between a decoder and an encoder. The ITU-T has continued work on further embellishing H.263 by adding yet many more features and optional modes [26].

The ITU-T H.324 [27] is a multimedia communications standard consisting of component standards, such as V.34 modem, H.223 multiplexer, H.245 control protocol, G.723.1 audio decoder, and H.263 (or H.261) video decoder. H.223 is the multiplexer used to mix audio, video, data and control channels together for transmission on V.34 modem. The ITU-T H.223 Annex A multiplexer has been designed for error prone channels and therefore features a robust packet synchronization and constant packet length. ITU-T H.324 Annex C specifies features of multimedia terminals operating in mobile radio environments, in terms of differences with normal terminals.

3. MPEG-4 overview

MPEG-4 was originally intended for very high compression coding of audio-visual information at very low bit-rates of 64 kbit/s or under. When MPEG-4 video was started, it was anticipated that with continuing advances in advanced (non-block based) coding schemes, for example, in region based and model based coding, a scheme capable of achieving very high compression, mature for standardization would emerge. By mid 1994, two things became clear. First, video coding schemes that were likely to be mature within the time frame of MPEG were likely to offer only moderate increase in compression (say, by a factor of 1.5 or so) over then existing methods as compared to the original goal of MPEG-4. Second, a new class of multimedia applications were emerging that required increasing levels of functionality than that provided by any other video standard at bit-rates in the range of 10 kbit/s to 1024 kbit/s. This led to broadening of the original scope of MPEG-4 to a larger range of bit-rates and important new functionalities [3]. Basically, three important trends were identified, which are as follows:

- The trend towards wireless communications.
- The trend towards interactive computer applications.
- The trend towards integration of audio-visual data into a number of applications.

The focus and scope of MPEG-4 was redefined as the intersection of the traditionally separate industries of telecommunications, computer, and TV/film where audio-visual applications exist. The mission and the focus statement of MPEG-4 explaining the trends leading up to MPEG-4 and what can be expected in the future are documented in the MPEG-4 Proposal Package Description (PPD) document [3]. Figure 2 shows the application areas of interest to MPEG-4 arising at the intersection of the aforementioned industries.

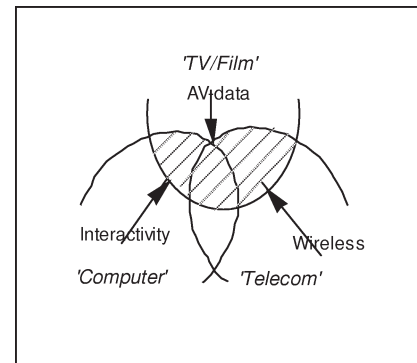


Figure 2. Applications areas addressed by MPEG-4 (shaded region).

To make the discussion a bit more concrete, we now provide a few examples of applications or application classes [45] that the MPEG-4 standard is aimed at:

- Internet and Intranet video.
- Wireless video.
- Video databases.
- Interactive home shopping.
- Video e-mail, home movies.
- Virtual reality games, simulation and training.

With this revised understanding of the goals of MPEG-4, the MPEG-4 work was subsequently reorganized and partitioned into the following subgroups:

- *Requirements* – develops requirements, application scenarios and meaningful clustering of coding tool combinations for interoperability (profiles).
- *Tests* – develops methods for subjective and objective assessment and conducts tests.
- *Video* – develops coded representation of moving pictures of natural origin.
- *Synthetic and Natural Hybrid Coding (SNHC)* – develops coded representation of synthetic audio, graphics and moving images.
- *Audio* – develops coded representation of audio of natural origin.
- *Systems* – develops techniques for multiplexing/demultiplexing and presentation of moving images, audio, graphics and data.
- *Digital Media Integration Framework (DMIF)* – develops standard interfaces between digital storage media, networks, servers and clients for delivery bitstreams in networked environments.
- *Implementation studies* – evaluates realizability of coding tools and techniques.

On the way to becoming an International Standard, MPEG-4 undergoes a sequence of interim steps as a Working Draft, a Committee Draft, a Final Committee Draft and a Draft International Standard; in table 1, we provide the schedule of these steps for MPEG-4 Version 1.

Table 1
MPEG-4 Version 1 workplan.

Working Draft	Committee Draft	Final Committee Draft	Draft International Standard	International Standard
November 1996	November 1997	July 1998	November 1998	January 1999

Table 2
Functionalities expected to be supported by MPEG-4 Version 1.

Content-based interactivity

Hybrid natural and synthetic data coding: The ability to code and manipulate natural and synthetic objects in a scene including decoder controllable methods of compositing of synthetic data with ordinary video and audio, allowing for interactivity.

Improved temporal random access: The ability to efficiently access randomly in a limited time and with fine resolution parts (frames or objects) within an audio-visual sequence. This also includes the requirement for conventional random access.

Content-based manipulation and bitstream editing: The ability to provide manipulation of contents and editing of audio-visual bitstreams without the requirement for transcoding.

Universal access

Robustness in error prone environments: The capability to allow robust access to applications over a variety of wireless and wired networks and storage media. Sufficient robustness is required, especially, for low bit-rate applications under severe error conditions.

Content-based scalability: The ability to achieve scalability with fine granularity in spatial, temporal or amplitude resolution, quality or complexity. Content based scaling of audio-visual information requires these scalabilities.

Compression

Improved coding efficiency: The ability to provide subjectively better audio-visual quality at bit-rates compared to existing or emerging video coding standards.

The MPEG-4 standard (ISO/IEC 14496) [46] is planned to consist of the following basic parts. Other parts may be added when the need is identified.

- ISO/IEC 14496-1: Systems.
- ISO/IEC 14496-2: Visual (Natural and Synthetic Video).
- ISO/IEC 14496-3: Audio (Natural and Synthetic Audio).
- ISO/IEC 14496-4: Conformance.
- ISO/IEC 14496-5: Software.
- ISO/IEC 14496-6: DMIF.

3.1. MPEG-4 functionalities and requirements

Now that we have some idea of the type of applications MPEG-4 is aimed for, we clarify the three basic functionality classes [3,15] that the MPEG-4 standard is addressing. They are as follows:

- *Content-based interactivity* allows the ability to interact with objects in a scene. Currently such interaction is typically only possible for synthetic objects; extending such interaction to natural and hybrid synthetic/natural objects is important to enable new audio-visual applications.
- *Universal accessibility* means the ability to access audio-visual data over a diverse range of storage and transmission media. Due to increasing trend toward mobile communications, it is important that access be available to applications via wireless networks. Thus acceptable

performance is needed over error-prone environments and at low bit-rates.

- *Improved compression* is needed to allow increase in efficiency in transmission or decrease in amount of storage required. For low bit-rate applications, high compression is very important to enable new applications.

Although we have looked at general classes of functionalities being addressed by MPEG-4 it is desirable to look at specific functionalities that MPEG-4 Version 1 expects to offer; in table 2 we now show a list of 6 such functionalities [3,6,15] and show their clustering into three functionality classes.

Besides the new functionalities, MPEG-4 is also supporting the basic functionalities such as synchronization of audio and video, auxiliary data streams capability, multi-point capability, low delay mode, coding of a variety of audio types, interoperability with other audio-visual systems, support for interactivity, ability to efficiently operate in the 9.6 to 1024 kbit/s range, ability to operate in different media environments, and the ability to operate in low complexity mode.

To keep up with marketplace needs for practical timely standards and to follow the evolving trends, the requirements collection process for MPEG-4 is kept flexible. The major restructuring of MPEG-4 effort in July 1994 to expand its scope was a response to the evolving trends in the marketplace. Evaluating requirements for MPEG-4 is an ongoing exercise that uses both top down (common requirements of related applications) and bottom up approach (related functions provided by a tool, which may be needed

Table 3
List of MPEG-4 first evaluation formal tests and their explanation.

Compression
<p><i>Class A sequences at 10, 24 and 48 kbit/s:</i> Coding to achieve the highest compression efficiency. Input video resolution is CCIR-601 and although any spatial and temporal resolution can be used for coding, the display format is CIF on a windowed display. The test method employed is SS.</p> <p><i>Class B sequences at 24, 48 and 112 kbit/s:</i> Coding to achieve the highest compression efficiency. Input video resolution is CCIR-601 and although any combination of spatial and temporal resolutions can be used for coding, the display format is CIF on a windowed display. The test method employed is SS.</p> <p><i>Class C sequences at 320, 512 and 1024 kbit/s:</i> Coding to achieve the highest compression efficiency. Input video resolution is CCIR-601 and although any combination of spatial and temporal resolution can be used for coding, the display format is CCIR-601 on a full display. The test method employed is DSCQS.</p>
Error robustness
<p><i>Error resilience at 24 kbit/s for Class A, 48 kbit/s for Class B, and 512 kbit/s for Class C:</i> Test with high random bit error rate (BER) of 10^{-3}, multiple burst errors with 3 bursts of errors with 50% BER within a burst, and a combination of high random bit errors and multiple burst errors. The display format for Class A and Class B sequences is CIF on a windowed display and for Class C sequences is CCIR-601 on full display. The test method employed for Class A and Class B is SS and that for Class C is DSCQS.</p> <p><i>Error recovery at 24 kbit/s for Class A, 48 kbit/s for Class B and 512 kbit/s for Class C:</i> Test with long burst errors of 50% BER within a burst and a burst length of 1 to 2 seconds. Display format for Class A and Class B is CIF on a windowed display and Class C is CCIR-601 on full display. The test method employed for Class A and Class B is SS and that for Class C is DSCQS.</p>
Scalability
<p><i>Object scalability at 48 kbit/s for Class A, 320 kbit/s for Class E, and 1024 kbit/s for Class B/C sequences:</i> Coding to permit dropping of specified objects resulting in remaining scene at lower than total bit-rate; each object and the remaining scene is evaluated separately by experts. The display format for Class A is CIF on a windowed display and for Class B/C and Class E is CCIR-601 on a full display. The test method employed for Class A is SS, for Class B/C is DSCQS, and for Class E is DSIS.</p> <p><i>Spatial scalability at 48 kbit/s for Class A, and 1024 kbit/s for Class B/C/E sequences:</i> Coding of a scene as two spatial layers with each layer using half of the total bit-rate, however, full flexibility in choice of spatial resolution of objects in each layer is allowed. The display format for Class A is CIF on a windowed display and that for Class B/C/E is CCIR-601 on a full display. The test method employed for Class A is SS, and that for Class B/C/E is DSCQS.</p> <p><i>Temporal scalability at 48 kbit/s for Class A, and 1024 kbit/s for Class B/C/E sequences:</i> Coding of a scene as two temporal layers with each layer using half of the total bit-rate, however, full flexibility in choice of temporal resolution of objects in each layer is allowed. The display format for Class A is CIF on a windowed display and that for Class B/C/E is CCIR-601 on a full display.</p>

in various applications). The collected requirements are clustered and translated into general directions for coding methods/tools development for individual subgroups such as Video, Audio, SNHC, Systems and DMIF. Finally, taking into account requirements of similar applications, a clustering of tools into profiles takes place (see section 7).

3.2. Tests and evaluation

We now discuss the testing and evaluation that took place in the competitive phase to determine the potential of the proposed technologies for MPEG-4. We also indicate how the outcome of tests and evaluation was used to initiate the collaborative phase.

3.2.1. Video tests

The competitive phase of MPEG-4 video began with an open call for proposals in November 1994 (and subsequently revised [6]), inviting technical proposals for the first testing and evaluation [4] which took place in October 1995. A proposal package description (PPD) and a test/evaluation procedures document was finalized by July 1995. Some functionalities were formally tested while the others were informally evaluated by experts.

Video scenes are classified from relatively simple to more complex by categorizing them into three classes: Class A, Class B and Class C. Two other classes of scenes, Class D (stereoscopic) and Class E (hybrid of natural and synthetic) were additionally defined. Many of the test scenes were presegmented (semiautomatically) into objects and a segmentation mask was provided along with test scenes. Also, since a variety of functionalities had to be tested, three type of tests were devised: *Single Stimulus* (SS), *Double Stimulus Impairment Scale* (DSIS) and *Double Stimulus Continuous Quality Scale* (DSCQS).

Table 3 summarizes the list of formal subjective tests [4], an explanation of each test and the type of method employed for each test.

The results [30] of tests revealed that DCT based coding performed reasonably well. Furthermore, it seemed possible to code shape of objects efficiently allowing object based functionalities. A collaborative effort was begun by first defining a set of core experiments (in November 1995) and soon after a Verification Model, VM (in January 1996) as the basis for further experimentation. Further, tentative plans were also made for a second test for verification of adopted coding algorithms/tools by mid 1997.

3.2.2. Audio tests

The MPEG-4 Audio also underwent subjective testing, similar to video. Three classes of audio test sequences, Class A, B and C were identified:

- Class A: Single source sequences consisting of a clean recording of a solo instrument.
- Class B: Single source with background sequences consisting of a person speaking with background noise.
- Class C: Complex sequences consisting of an orchestral recording.

All sequences were originally sampled at 48 kHz with 16 bits/sample and were monophonic in nature. For generating reference formats, filters were specified to downsample them to 24, 16 and 8 kHz. A number of bit-rates such as 2, 6, 16, 24, 40 and 64 kbit/s were selected for testing of audio/speech. The first three bit-rates are obviously only suitable for speech material. The audio test procedures used were as defined in ITU-R Recommendation 814.

The proposals submitted for testing included variants of MPEG-2 Advanced Audio Coding (AAC), variations of MPEG-1 audio coding and new coding schemes. For specific bit-rates, some candidates outperformed the reference coding schemes, although for all combinations tested, no single scheme was the clear winner. After subjective testing, the collaborative work started and an initial MPEG-4 Audio VM was developed. The MPEG-4 Audio development underwent a core experiments process similar to that of the MPEG-4 Video development process.

3.2.3. SNHC tests

The SNHC group started its work much later than the video group. Its focus was primarily on coding for storage and communication of 2D and 3D scenes involving synthetic images, sounds, and animated geometry and its integration into scenes that contain coded natural images/video and sound. Further, it was sought that the coded representation should also facilitate various forms of interactions.

In the SNHC call for proposals [29] and PPD [30], proposals in the areas of efficient compression and simplification of synthetic data, parameterized animated models, new primitive operations for compositing of natural and hybrid objects, scalability, real-time interactivity with synthetic/hybrid environments, modeling of timing and synchronization, synthetic audio, etc., were sought.

In the competitive phase, for the purpose of standardized evaluation, a database of test data set was established. The actual evaluation of proposals by a group of experts took place in September 1996. The evaluation criteria was based on the functionality addressed such as coding efficiency, quality of decoded model, real-time interactivity, anticipated performance in future, and implementation cost. After the evaluation, the collaborative phase was begun by defining a verification model (VM).

3.3. Video development

Video development was started by identifying a number of needed tools and the options available for each tool as well as a reference framework. A total of about 40 core experiments were defined by January 1996 and were categorized for evaluation by the following 4 ad hoc groups:

- Coding efficiency – prediction, frame texture coding, quantization and rate control.
- Shape and object texture coding – binary and grey scale shape coding, object texture coding.
- Robust coding – error resilience and error concealment.
- Multifunctional coding – bandwidth and complexity scalability, object manipulation, post processing.

A reference coding framework known as the first Verification Model (VM1) was released on 24th January 1996. It supported the following features:

- Coding of arbitrary shaped objects using Video Object Planes (VOPs).
- Coding of binary and grey scale shape of arbitrary shaped objects.
- Padding of pixels to fill the region outside of an object to full blocks for motion compensation and DCT.
- Macroblock based motion-texture (motion compensated DCT) coding derived from H.263.
- A mode allowing separation of motion and texture data for increased error resilience.

During the March 1996 MPEG meeting, a number of additional features were added to VM1 and thus VM2 [29] was released on 29th March 1996. The additional features were as follows:

- Bidirectional VOPs derived from combination of H.263 PB-frames mode and MPEG-1/2 B-pictures.
- DC coefficients prediction for intra macroblocks as per MPEG-1/2.
- Extended motion vector range.
- Quantization visibility matrices as per MPEG-1/2.

The process of iterative development and refinement of video VM's via core experiments has continued and there have been seven iterations on the first VM. At the October 1997 meeting, a number of mature tools from VM8 were accepted for MPEG-4 Video Version 1 for the Committee Draft. The remaining tools as well as some new tools are being considered for MPEG-4 Video Version 2. In section 4, we describe the basic coding methods formed by the tools accepted for the MPEG-4 Video Version 1 standard.

3.4. Audio development

The MPEG-4 Audio coding effort occurred in parallel with the MPEG-2 AAC (formerly, Non-Backward Compatible (NBC)) coding effort. The MPEG-2 standard origi-

nally had an audio coding mode called backward compatible (BC) mode which as the name suggests was backward compatible with MPEG-1 audio coding. However, at a late stage in MPEG-2 it was discovered that the BC audio coding was rather inefficient compared to non compatible solutions and thus work on NBC mode was begun and overlapped with the MPEG-4 schedule. The NBC mode was renamed AAC and became a new part of MPEG-2 achieving International Standard status in April 1997 (although it had reached a mature status in mid 1996).

Towards the very low bit-rate end a valid question to ask is why not use the existing ITU-T coders? As an answer to this question, the ITU-T speech coders currently operate at 6.3/5.3 kbit/s (G.723), 8 kbit/s (G.729), 16 kbit/s (G.728), 32 kbit/s (G.721) 48/56/64 kbit/s (G.722). In comparison, MPEG-4 speech coding is being designed to operate at bit rates between 2–24 kbit/s for the 8 kHz mode and 14–24 kbit/s for the 16 kHz mode, whereas ITU-T coders do not operate at bit-rates as low 2 kbit/s for the 8 kHz mode, or 14–24 kbit/s for the 16 kHz mode. Furthermore, MPEG-4 speech coders are being designed for bit-rate scalability, complexity scalability and multi-bit-rate operation from 2–24 kbit/s. The coding quality of the coder is comparable to that of the ITU coder at corresponding bit-rates but the MPEG-4 speech coder can operate down to 2 kbit/s. The quality at 2 kbit/s is “communication quality” and could be used for usual conversation, performing better than the FS1016 4.8 kbit/s coder.

Therefore, the MPEG-4 Audio VMs have targeted bit-rates from 2 kbit/s to 64 kbit/s; a number of coding schemes are used to cover portions of this range. Besides coding efficiency, content based coding of audio objects and scalability are being investigated. There have been a total of four iterations of audio VM, from VM1 to VM4; the last VM was released in July 1997. In fact, the more mature tools of Audio VM3 have been accepted for the audio part [28] of the MPEG-4 Version 1 standard.

In section 5, we briefly discuss the basic coding techniques accepted for part 3 of the MPEG-4 Version 1 standard.

3.5. SNHC development

There have been a total of four iterations of SNHC VM, from VM1 to VM4; the last VM was released in July 1997. In fact, the more mature tools of SNHC VM3 have been accepted for the visual part of the MPEG-4 Version 1 standard; the remaining tools have been left in VM4 for consideration for the next version of MPEG-4.

In section 4 we describe the SNHC tools expected to be included in the visual part of the MPEG-4 Version 1 standard. In section 5, we discuss the SNHC tools expected to be included in the audio part of the MPEG-4 Version 1 standard.

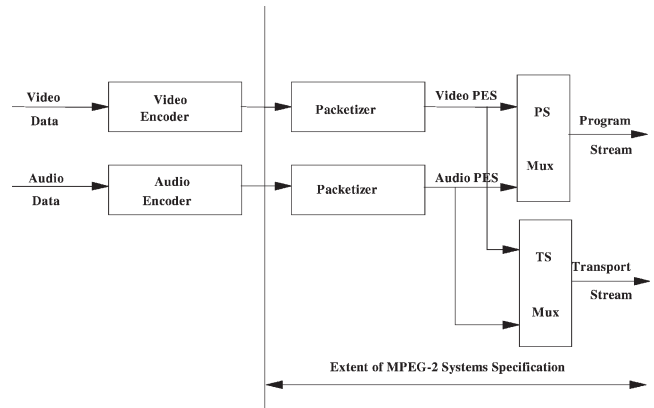


Figure 3. MPEG-2 Systems.

3.6. Systems development

The Systems layer in MPEG has been traditionally responsible for integrating media components into a single system, providing multiplexing and synchronization services for audio and video streams.

In MPEG-2 [3,15], for example, these are the primary functionalities, and were designed for two types of transport facilities. The first, Program Stream, is intended for reliable media such as storage devices, and can only carry a single program (combinations of synchronized audio and video streams). It also provided backwards compatibility with MPEG-1 [23]. The second, Transport Stream, is intended for potentially unreliable media and can carry multiple programs. This is shown in figure 3. The object-based nature of MPEG-4 necessitates a much more complex Systems layer since, in addition to still addressing multiplexing and synchronization, it must also provide for ways to combine simple audio or visual objects into meaningful scenes.

Considering the object-based nature of MPEG-4, a key requirement from the System part is the capability to combine individual audiovisual objects in scenes. In late 1995, this was accomplished by using Java [14]. Issues of performance and compliance soon arose. Clearly, it is essential for a content creator to be assured that the content generated will be shown in an identical way (or nearly so) regardless of the terminal used, if both such terminals comply to the standard. A three-step approach was adopted, involving three different flexibility levels, as shown in figure 4. In level 0, no programmability was allowed. In level 1, facilities were provided to combine different tools into algorithms, while in level 2 even individual tools were considered as targets for programmable behavior. The group was also renamed to MPEG-4 System and Description Languages (MSDL), separating system and syntactic description [5]. After further examination, in late 1996 it was decided that any meaningful operation of level 1 would require the complexity of implementing a level 2 system, and hence this intermediate level was eliminated.

The group subsequently focused on a parametric (bit-stream oriented, non-programmable) solution for describing how objects should be combined together. Using the

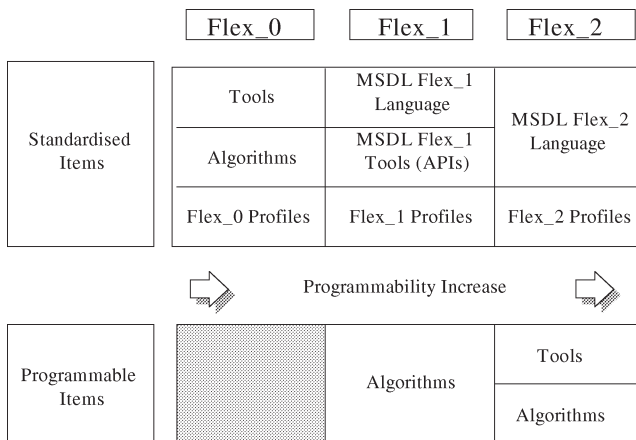


Figure 4. Evolution of MPEG-4 systems architecture (1996).

above figure, that would be a level 0 design. The group also reverted to the use of the traditional term “Systems”, reflecting the varied components that it addresses. A programmable approach is still being considered and is discussed in more detail in section 6.2.4.

In addition to the overall architectural issues, the issue of multiplexing in MPEG-4 also underwent several stages of evolution. The H.223 Annex A multiplexer was used as a basis, including error protection tools (interleaving and ARQ). A key requirement [48], however, for MPEG-4 was the need to be transport-independent. As a result, services that belong to a transport layer were subsequently removed from the set of specified tools, so that efficient implementation of MPEG-4 systems could be performed in a very broad range of environment (broadcast, ATM, IP, and wireless).

3.7. DMIF development

At a recent MPEG meeting, the significance of the DMIF activity has been recognized and DMIF has been given the status of a new group [7]. Previously, DMIF was an ad hoc group operating under the systems group. The charter of the DMIF group is to develop standards for interfaces between Digital Storage Media (DSM), networks, servers and clients for the purpose of managing DSM resources and controlling the delivery of MPEG bitstreams and associated data. The ongoing work of this group is expected to result in part 6 of the MPEG-4 standard.

4. MPEG-4 visual

The ongoing work on MPEG-4 visual standard specification [35] consists of tools and methods from two major areas – coding of (natural) video and coding of synthetic video (visual part of the SNHC work). We address both these areas; in sections 4.1–4.4 we discuss tools and techniques relevant to natural video coding and in sections 4.5–4.7 we discuss tools and techniques relevant to synthetic video coding.

4.1. MPEG-4 video coding basics

In this and the next section, we describe the coding methods and tools of MPEG-4 video; the encoding description is borrowed from Video VM7 [37], the decoding description follows [35]. An input video sequence can be defined as a sequence of related snapshots or pictures, separated in time. In MPEG-4, each picture is considered as consisting of temporal instances of objects that undergo a variety of changes such as translations, rotations, scaling, brightness and color variations etc. Moreover, new objects enter a scene and/or existing objects depart, leading to the presence of temporal instances of certain objects only in certain pictures. Sometimes, scene change occurs, and thus the entire scene may either get reorganized or replaced by a new scene. Many of MPEG-4 functionalities require access not only to entire sequence of pictures, but to an entire object, and further, not only to individual pictures, but also to temporal instances of these objects within a picture. A temporal instance of a video object can be thought of as a snapshot of arbitrary shaped object that occurs within a picture, such that like a picture, it is intended to be an access unit, and, unlike a picture, it is expected to have a semantic meaning.

The concept of Video Objects (VOs) and their temporal instances, Video Object Planes (VOPs) is central to MPEG-4 video. A VOP can be fully described by texture variations (a set of luminance and chrominance values) and (explicit or implicit) shape representation. In natural scenes, VOPs are obtained by semi-automatic or automatic segmentation, and the resulting shape information can be represented as a *binary shape* mask. On the other hand, for hybrid (of natural and synthetic) scenes generated by blue screen composition, shape information is represented by an 8-bit component, referred to as *grey scale shape*. In figure 5, we show the decomposition of a picture into a number of separate VOPs. The scene consists of two objects (head and shoulders view of a human, and a logo) and the background. The objects are segmented by semi-automatic or automatic means and are referred to as VOP1 and VOP2, while the background without these objects is referred to as VOP0. Each picture in the sequence is segmented into VOPs in this manner. Thus, a segmented sequence contains a set of VOP0s, a set of VOP1s and a set of VOP2s, in other words, in our example, a segmented sequence consists of VO0, VO1 and VO2.

Each VO is encoded separately and multiplexed to form a bitstream that users can access and manipulate (cut, paste, etc.). The encoder sends together with VOs, information about scene composition to indicate where and when VOPs of a VO are to be displayed. This information is however optional and may be ignored at the decoder which may use user-specified information about composition.

In figure 6 we show a high level logical structure of a VO based coder. Its main components are VO Segmenter/Formatter, VO Encoders, Systems Multiplexer Systems Demultiplexer, VO Decoders and VO Compositor. VO Segmenter segments the input scene into VOs for en-

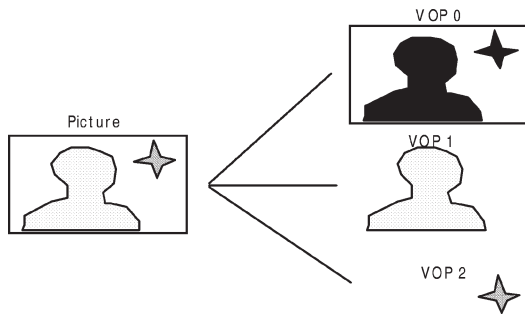


Figure 5. Semantic segmentation of a picture in to VOPs.

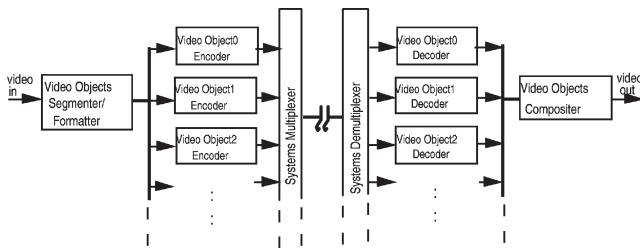


Figure 6. Logical structure of Video Object based codec of MPEG-4 video.

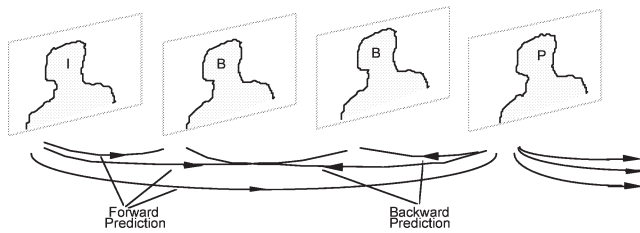


Figure 7. An example prediction structure when using I-, P- and B-VOPs.

coding by VO Encoders. The coded data of various VOs is multiplexed for storage or transmission, following which it is demultiplexed and decoded by VO decoders and offered to composer, which renders the decoded scene.

To see how coding takes place in a video object encoder, consider a sequence of VOPs. Extending the concept of intra (I-) pictures, predictive (P-) and bidirectionally predictive (B-) pictures of MPEG-1/2 to VOPs, I-VOP, P-VOP and B-VOP result. If two consecutive B-VOPs are used between a pair of reference VOPs (I- or a P-VOPs), the resulting coding structure is as shown in figure 7.

In figure 8 we show the internal structure of the VM based encoder which codes a number of VOs of a scene. Its main components are: Motion Coder, Texture and Shape Coder. The Motion Coder uses macroblock and block motion estimation and compensation, similar to H.263 and MPEG-1/2 but modified to work with arbitrary shapes. The Texture Coder uses block DCT coding based on H.263 and MPEG-1/2 but much better optimized; further it is also adapted to work with arbitrary shapes. An entirely new component is the Shape Coder. The partial data of VOs (such as VOPs) is buffered and sent to the Systems Multiplexer.

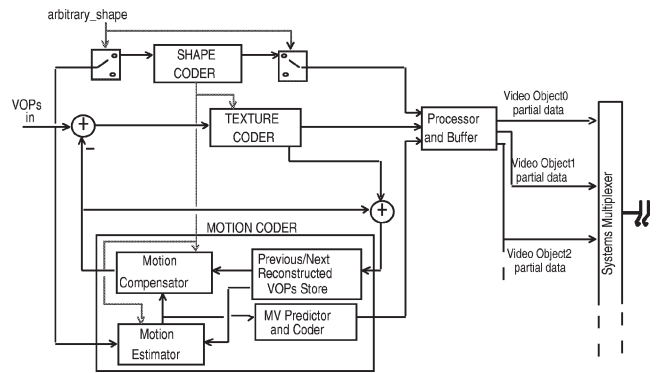


Figure 8. Detailed structure of video objects encoder.

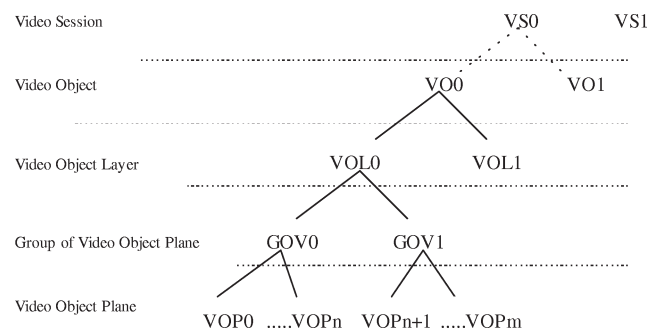


Figure 9. Class hierarchy for structuring coded video data.

From a top-down perspective, the organization of coded MPEG-4 Video data can be described by the following class hierarchy:

- **VideoSession:** A Video Session represents the highest level in the class hierarchy and simply consists of an ordered collection of Video Objects. This class has only been a place holder for video VM and core experiments work and, since composition of objects is now handled by systems, it is not needed.
- **VideoObject:** A Video Object (2D + time) represents a complete scene or a portion of a scene with a semantic meaning.
- **VideoObjectLayer (VOL):** A Video Object Layer (2D+time) represents various instantiations of an Video Object. For instance, different VOLs may correspond to different layers, such as in the case of scalability.
- **GroupOfVideoObjectPlanes (GOV):** Group of Video Object Planes are optional entities and are essentially access units for editing, tune-in or synchronization.
- **VideoObjectPlane (VOP):** A Video Object Plane represents a snap shot in time of a Video Object. A simple example may be an entire frame or a portion of a frame. Different coding methods from MPEG-1/2 such as intra (I-) coding, predictive (P-) coding and bidirectionally predictive (B-) coding can now be applied to VOPs.

The class hierarchy used for representation of coded bit-stream described above is shown by the tree structure of figure 9.

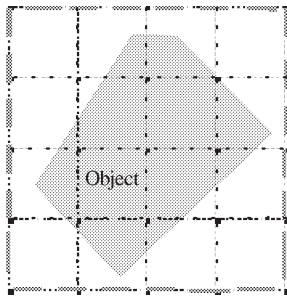


Figure 10. A VOP in a bounding box.

4.2. Video coding details

4.2.1. Binary shape coder

Compared to other standards, the ability to represent arbitrary shapes is an important capability of the MPEG-4 video standard. In general, shape representation can be either implicit (based on chroma-key and texture coding) or explicit (boundary coding separate from texture coding). Implicit shape representation, although it offers less encoding flexibility, can result in quite usable shapes while being relatively simple and computationally inexpensive. Explicit shape representation although it can offer flexible encoding and somewhat better quality shapes, it is more complex and computationally expensive. Regardless of the implications, the explicit shape representation was chosen in MPEG-4 video; we now briefly describe the essence of this method [35,37] without its many details.

For each VO given as a sequence of VOPs of arbitrary shapes, the corresponding sequence of binary alpha planes is assumed to be known (generated via segmentation or via chroma-key). For the binary alpha plane, a rectangular bounding box enclosing the shape to coded is formed such that its horizontal and vertical dimensions are multiples of 16 pixels (macroblock size). For efficient coding, it is important to minimize the number of macroblocks contained in the bounding box. The pixels on the boundaries or inside the object are assigned a value of 255 and are considered opaque while the pixels outside the object but inside the bounding box are considered transparent and are assigned a value of 0. If a 16×16 block structure is overlaid on the bounding box, three types of binary alpha blocks exist: completely transparent, completely opaque, and partially transparent (or partially opaque). Figure 10 shows an example of an arbitrary shape VOP with a bounding box and the overlaid 16×16 block structure; the opaque area is shown shaded whereas the transparent area is shown unshaded.

Coding of each 16×16 binary alpha block representing shape can be performed either lossy or losslessly. The degree of lossiness in coding the shape of a video object is controlled by a threshold which can take values of 0, 16, 32, 64, ..., 256. The higher the value of this threshold, the more lossy the shape representation; a zero value implies lossless shape coding. Within the global bound of specified lossiness, local control, if needed, can be exerted by selecting a maximum subsampling factor on a 16×16 binary alpha that results in just acceptable distortion. The

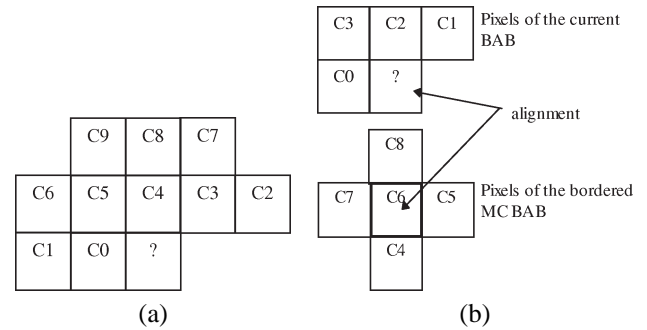


Figure 11. Pixel templates used for (a) intra and (b) inter context determination of a binary alpha block (BAB). Pixel to be coded is marked with '?'. To

estimation of this factor is iterative and consists of using the same subsampling factor in both dimensions and determining the acceptability of resulting shape quality. To be specific, a 4:1 downsampled binary alpha block is used first and if the shape errors are higher than acceptable, a 2:1 downsampled binary alpha block is used next, again if it is found unacceptable, an unsampled binary alpha block is used.

Further, each binary alpha block can be coded in intra mode or in inter mode, similar to coding of texture macroblocks. In intra mode, no explicit prediction is performed. In inter mode, shape information is differenced with respect to the prediction obtained using a motion vector, the resulting binary shape prediction error may or may not be coded. The motion vector of a binary alpha block is estimated at the encoder by first finding a suitable initial candidate from among the motion vectors of 3 previously decoded surrounding texture macroblocks as well as the 3 previously decoded surrounding shape binary alpha blocks. Next, the initial candidate is either accepted as the shape motion vector, or is used as the starting basis for a new motion vector search, depending on if the resulting prediction errors of the initial motion vectors are below a threshold. The motion vector is coded differentially and included in the bitstream. Following this procedure, a binary alpha block is assigned a mode from among the following choices:

1. Zero differential motion vector and no inter shape update.
2. Nonzero differential motion vector and no inter shape update.
3. Transparent.
4. Opaque.
5. Intra shape.
6. Zero differential motion vector and inter shape update.
7. Nonzero differential motion vector and inter shape update.

Depending on the coding mode and whether it is an I-, P- or B-VOP, a variable length codeword is assigned identifying the coding type of the binary alpha block. The

entropy coding of shape data is performed by using the context information determined on a pixel basis to drive an adaptive arithmetic coder. The pixels of binary alpha block are raster scanned, however, a binary alpha block maybe transposed. Next, a context number is determined and is used to index a probability table, and further, the indexed probability is used to drive the arithmetic coder. To determine the context, different templates of surrounding pixels are used for intra and inter coded binary alpha blocks, as shown in figure 11.

The decoding of binary alpha block follows the inverse sequence of operations with the exception of encoder specific tasks such as motion estimation, subsampling factor determination, mode decision, etc., which are readily extracted from the coded bitstream.

4.2.2. Motion Coder

The Motion Coder [35,37] consists of a Motion Estimator, Motion Compensator, Previous/Next VOPs Store and Motion Vector (MV) Predictor and Coder. In case of P-VOPs, Motion Estimator computes motion vectors using the current VOP and temporally previous reconstructed VOP available from the Previous Reconstructed VOPs Store. In case of B-VOPs, Motion Estimator computes motion vectors using the current VOP and temporally previous reconstructed VOP from the Previous Reconstructed VOP Store, as well as, the current VOP and temporally next VOP from the Next Reconstructed VOP Store. The Motion Compensator uses these motion vectors to compute motion compensated prediction signal using the temporally previous reconstructed version of the same VOP (reference VOP). The MV Predictor and Coder generates prediction for the MV to be coded. We now discuss the details of padding needed for motion compensation of arbitrary shaped VOPs, as well as the various modes of motion compensation allowed.

In the reference VOP, based on its shape information, two types of macroblocks require padding, those that lie on the boundary and (depending on encoding choice, some or all of) the other that lie outside of the VOP. Macroblocks that lie on the VOP boundary are padded by first replicating the boundary pixels in the horizontal direction, followed by replicating the boundary pixels in the vertical direction making sure that if a pixel can be assigned a value by both horizontal and vertical padding, it is assigned an average value. Next, the macroblocks that lie outside of the VOP are padded by extending the boundary macroblock pixels, up, down, left and right, and averaging wherever a pixel is assigned a value from more than one directions. The processing for the previous step can be reduced by only padding macroblocks that are outside of the VOP but right next to the boundary pixels.

The basic motion estimation and compensation is performed on 16×16 luminance block of a macroblock. The motion vector is specified to half-pixel accuracy. The motion estimation is performed by full search to integer pixel accuracy vector and using it as the initial estimate, a half

pixel search is performed around it. The luminance block motion vector is scaled by a factor of 2 for each component and rounded for use on 8×8 chrominance blocks. MPEG-4 video, like H.263, supports an unrestricted range for motion estimation and compensation. Basically, motion vectors are allowed to point out of the VOP bounding box, by extending the reference VOP bounding box in all four directions. Further, a larger range of motion vectors is supported for motion vector coding in MPEG-4 as compared to H.263.

Often a single motion vector for a 16×16 luminance block does not reduce the prediction errors sufficiently or when dealing with boundary macroblocks, motion vectors can be sent for individual 8×8 blocks. Further, the 8×8 block motion vectors are used to generate overlapped block motion compensated prediction. Both the 8×8 block motion compensation and overlapped motion compensated prediction are referred to as advanced prediction in H.263 and are adapted in MPEG-4 to work with arbitrary shaped VOPs.

An intra versus inter coding decision is performed to determine if motion vector(s) need to be sent for the macroblock being coded; further, a decision is also performed to determine if 16×16 or 8×8 block motion vectors will be sent for the macroblock being coded. All motion vectors are coded differentially using median of neighboring three decoded macroblock (or block in case of 8×8 coding) motion vectors as the prediction.

As mentioned earlier, a B-VOP is a VOP which is coded bidirectionally. For example, macroblocks in a B-VOP can be predicted using the forward, the backward or both using the forward and backward motion vectors; this has similarities to MPEG-1/2 in which B-pictures can use such motion vectors. However, MPEG-4 video also supports an H.263 based mode for motion compensation, referred to as the direct mode. In direct mode, the motion vector for a macroblock in a B-VOP is obtained by scaling of the P-VOP motion vector, and further correcting it by a small (delta) motion vector. The actual motion compensation mode to be used for a macroblock is decided taking into account the motion compensated prediction errors produced by various choices and the coding overhead of any additional motion vectors. All motion vectors (except delta) are coded differentially with respect to motion vectors of the same type.

MPEG-4 also supports efficient coding of interlaced video. It combines the macroblock based frame/field motion compensation of MPEG-2 with the normal motion compensation of MPEG-4, resulting in overall improved motion compensation. Furthermore, it allows motion compensation of arbitrary shaped VOPs of interlaced video whereas MPEG-2 only supports rectangular pictures of interlaced video.

4.2.3. Video Texture Coder

The Texture Coder [35,37] codes the luminance and chrominance variations of blocks forming macroblocks

Table 4
Nonlinear scaler for DC coefficients of DCT blocks.

Component	DC_scaler for Quantizer (Q_p) range			
	1-4	5-8	9-24	25-31
Luminance	8	$2Q_p$	$Q_p + 8$	$2Q_p - 16$
Chrominance	8	$(Q_p + 13)/2$		$Q_p - 6$

within a VOP. Two types of macroblocks exist, those that lie inside the VOP and those that lie on the boundary of the VOP. The blocks that lie inside the VOP are coded using DCT coding similar to that used in H.263 but optimized in MPEG-4. The blocks that lie on the VOP boundary are first padded and then coded similar to the block that lie inside the VOP. The remaining blocks are transparent (they lie inside the bounding box but outside of the coded VOP shape) and are not coded at all.

The texture coder uses block DCT coding and codes blocks of size 8×8 similar to H.263 and MPEG-1/2, with the difference that since VOP shapes can be arbitrary, the blocks on the VOP boundary require padding prior to texture coding. The general operations in the texture encoder are: DCT on original or prediction error blocks of size 8×8 , quantization of 8×8 block DCT coefficients, scanning of quantized coefficients and variable length coding of quantized coefficients. For inter (prediction error block) coding, the texture coding details are similar to that of H.263 and MPEG-1/2. However, for intra coding of texture data, a number of improvements are included. We now discuss the quantization for intra and inter macroblocks, followed by coefficient prediction, scanning and entropy of intra macroblocks, and finally the entropy coding of inter blocks.

Typically, the DC coefficients of DCT of blocks belonging to an intra macroblock, are scaled by a constant scaling factor of 8, however, in MPEG-4 video, a nonlinear scaler [44] as per table 4 is used to provide a higher coding efficiency while keeping the blockiness artifacts under the visibility threshold. The characteristics of nonlinear scaling are different between the luminance and chrominance blocks and further depends on the quantizer used for the block.

MPEG-4 video supports two techniques of quantization, one referred to as the H.263 quantization method (with deadzone for intra and inter), and the other, the MPEG quantization method (no deadzone for intra but uses deadzone for inter, and intra and inter quantization matrices). Further, the quantization matrices are downloadable like in MPEG-1/2, but with the difference that it is possible to update matrices partially.

Unlike H.263, the quantized intra DC coefficients are predicted [44] with respect to 3 previous decoded DC coefficients, for example, quantized DC coefficients of blocks A, B and C when predicting quantized DC value for block X in figure 12. Although MPEG-1/2 also allows prediction of DC coefficients, however the gradient based prediction of MPEG-4 is more effective. In computing the prediction for block X, if the absolute value of a horizontal gradient

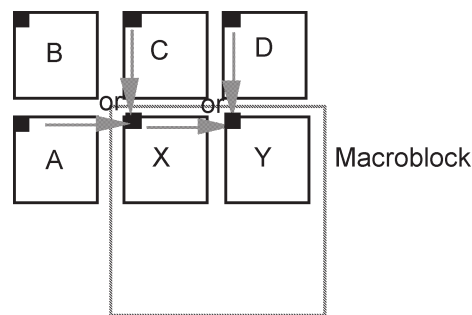


Figure 12. Prediction of DC coefficients of blocks in an intra macroblock.

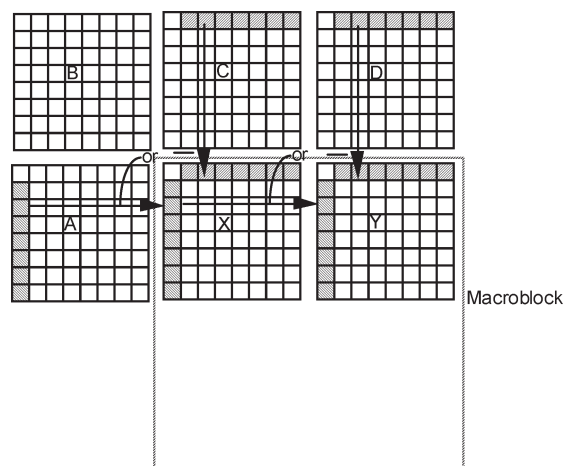


Figure 13. Prediction of AC coefficients of blocks in an intra macroblock.

$|QDC_A - QDC_B|$ is less than the absolute value of a vertical gradient $|QDC_B - QDC_C|$, then the QDC value of block C is the prediction, else QDC value of block A is used as prediction. This process is independently repeated for every block of an intra macroblock using horizontally and vertically adjacent blocks. Further, the procedure is identical for luminance and chrominance blocks.

Not only the DC coefficients of intra blocks are predicted and coded differentially, so are some of the AC coefficients [44]. In particular, on a block basis, either the first row or the first column of AC coefficients of DCT blocks of each intra macroblock are predicted. The direction (horizontal or vertical) used for prediction of DC coefficient of a block is also used for predicting the corresponding first column or row of AC coefficients. The prediction direction can differ from block to block within each intra macroblock. Further, the AC coefficient prediction can be disabled for a macroblock when it does not work well. Figure 13 shows the prediction of quantized AC coefficients belonging to the first column or the first row of block X from the corresponding quantized AC coefficients of block A or block C.

The predicted DC and AC coefficients (as well as the unpredicted AC coefficients) of DCT blocks of intra macroblocks are scanned by one of the three scans [44]: alternate-horizontal, alternate-vertical (MPEG-2 interlace scan) and the zigzag scan (normal scan used in H.263 and MPEG-1). The actual scan used depends on the coefficient

0	1	2	3	10	11	12	13
4	5	8	9	17	16	15	14
6	7	19	18	26	27	28	29
20	21	24	25	30	31	32	33
22	23	34	35	42	43	44	45
36	37	40	41	46	47	48	49
38	39	50	51	56	57	58	59
52	53	54	55	60	61	62	63

(a)

0	4	6	20	22	36	38	52
1	5	7	21	23	37	39	53
2	8	19	24	34	40	50	54
3	9	18	25	35	41	51	55
10	17	26	30	42	46	56	60
11	16	27	31	43	47	57	61
12	15	28	32	44	48	58	62
13	14	29	33	45	49	59	63

(b)

0	1	5	6	14	15	27	28
2	4	7	13	16	26	29	42
3	8	12	17	25	30	41	43
9	11	18	24	31	40	44	53
10	19	23	32	39	45	52	54
20	22	33	38	46	51	55	60
21	34	37	47	50	56	59	61
35	36	48	49	57	58	62	63

(c)

Figure 14. Scans for intra blocks: (a) alternate-horizontal, (b) alternate-vertical, (c) zigzag.

predictions used. For instance, if AC coefficient prediction is disabled for a intra macroblock, all blocks in that macroblock are zigzag-scanned. If AC coefficient prediction is enabled and DC coefficient prediction was selected from the horizontally adjacent block, alternate-vertical scan is used, likewise, if AC coefficient prediction is enabled and DC coefficient prediction was selected from the vertically adjacent block, alternate-horizontal scan is used. Figure 14 shows the three scans used.

A three dimensional variable length code is used to code the scanned DCT events of intra blocks. An event is a combination of three items (last, run, level). The ‘last’ indicates if a coefficient is the last nonzero coefficient of a block or not, the ‘run’ indicates number of zero coefficients preceding the current nonzero coefficient and level indicates the amplitude of the quantized coefficient.

The DCT coefficients of inter blocks, unlike DCT coefficients of intra blocks do not undergo any prediction or adaptive scanning, in fact they use the fixed zigzag scan of figure 14(c). The scanned coefficients of inter blocks are also coded by a three dimensional variable length code table with similar structure as the intra variable length code table but with code entries optimized for inter statistics.

Finally, as mentioned earlier, MPEG-4 also supports efficient coding of interlaced video. It combines the macroblock based frame/field DCT coding of MPEG-2 with the improved DC coefficient coding, quantization, scanning and variable length coding of normal MPEG-4 video coding resulting in improved coding efficiency. Furthermore, it allows DCT coding of arbitrary shaped VOPs of interlaced video where as MPEG-2 only supports rectangular pictures of interlaced video.

4.2.4. Sprite coding

In computer games, a sprite refers to an synthetic object that undergoes some form of transformation (including animation). Also, in the literature, and in connection with highly efficient representation of natural video, the term ‘mosaic’ or ‘world image’ is used to describe a large image built by integration of many frames of a sequence spatially and/or many frames of a sequence temporally; in MPEG-4 terminology, such an image is referred to as a static sprite. Static sprites can improve the overall coding efficiency, for example, by coding the background only once and warping it to generate the rendition required at a specific time instance.

A static sprite [35,37] is usually built offline and can be used to represent synthetic or natural objects. It is quite suitable for natural objects that undergo rigid motion and where a wall paper like rendering is sufficient. One of the main components in coding using natural sprites is generation of the sprite itself. For generating a static sprite, the entire video object is assumed to be available. For each VOP in the VO, the global motion is estimated according to a transformation model (say, perspective transformation) using which a VOP is then registered with the sprite by warping the VOP to sprite coordinate system. Finally, the warped VOP is blended with the sprite which is used for estimation of motion of the subsequent VOP.

A number of choices regarding the transformations models exist such as stationary, translation, magnification–rotation–translation, affine, and perspective transformation. Each transformation can be defined as either a set of coefficients or the motion trajectories of some reference points; the former, is convenient for performing the transformations whereas the later for encoding the transformations. If four reference points are used, perspective transformation can be employed for warping and is defined by the following:

$$\begin{aligned}x' &= (ax + by + c)/(gx + hy + l), \\y' &= (dx + ey + f)/(gx + hy + l),\end{aligned}$$

where $\{a, b, c, d, e, f, g, h, l\}$ are the coefficients of the transformation, (x, y) is one of the reference points of interest in current VOP which corresponds to point (x', y') in the sprite, expressed in sprite coordinate system.

Once the sprite is available, global motion between the current VOP and the sprite is estimated, using the perspective transform, for example. The reconstructed VOPs are generated from the sprite by directly warping the quantized sprite using specified motion parameters. Residual error between the original VOP and the warped sprite is not sent.

4.3. Scalable video coding

Scalability of video is the property that allows a video decoder to decode portions of the coded bitstreams to generate decoded video of quality commensurate with the amount of data decoded. In other words, scalability allows a simple

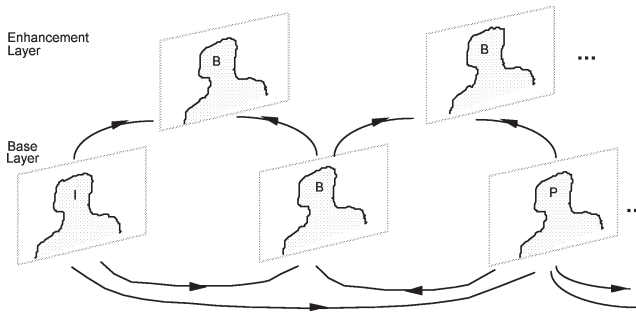


Figure 15. Temporal scalability.

video decoder to decode and produce basic quality video while an enhanced decoder may decode and produce enhanced quality video, all from the same coded video bitstream. This is possible because scalable video encoding ensures that input video data is coded as two or more layers, an independently coded base layer and one or more enhancement layers coded dependently, thus producing scalable video bitstreams. The first enhancement layer is coded with respect to the base layer, the second enhancement layer with respect to the first enhancement layer and so forth.

MPEG-4 video offers a generalized scalability [35,37, 44] framework supporting both the Temporal and the Spatial scalabilities, the primary type of scalabilities. Temporally scalable encoding offers decoders a means to increase temporal resolution of decoded video using decoded enhancement layer VOPs in conjunction with decoded base layer VOPs. Spatial scalability encoding on the other hand offers decoders a means to decode and display either the base layer or the enhancement layer output; typically, since base layer uses one-quarter resolution of the enhancement layer, the enhancement layer output provides the better quality, albeit requiring increased decoding complexity. The MPEG-4 generalized scalability framework employs modified B-VOPs that only exist in enhancement layer to achieve both temporal and spatial scalability; the modified enhancement layer B-VOPs use the same syntax as normal B-VOPs but for modified semantics which allows them to utilize a number of interlayer prediction structures needed for scalable coding.

Figure 15 shows an example of the prediction structure used in temporally scalable coding. The base layer is shown to have one-half of the total temporal resolution to be coded, the remaining one-half is carried by the enhancement layer. Base layer is coded independently as in normal video coding whereas the enhancement layer uses B-VOPs that use both, an immediate temporally previous decoded base layer VOP as well as an immediate temporally following decoded base layer VOP for prediction.

Next, figure 16 shows an example of prediction structure used in spatially scalable coding. The base layer is shown to have one-quarter resolution of the enhancement layer. Base layer is coded independently as in normal video coding whereas the enhancement layer mainly uses B-VOPs that use both, an immediate previous decoded enhancement

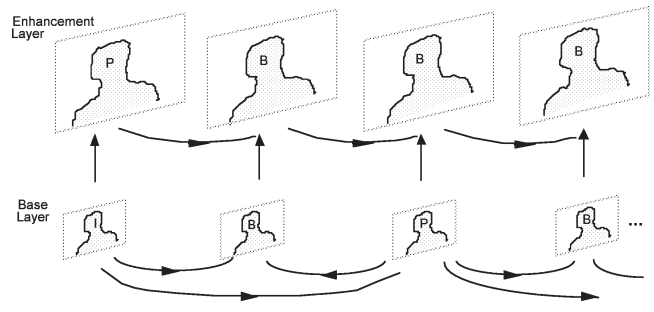


Figure 16. Spatial scalability.

layer VOP as well as a coincident decoded base layer VOP for prediction.

In reality, some flexibility is allowed in choice of spatial and temporal resolutions for base and enhancement layers as well as the prediction structures allowed for the enhancement layer to cope with a variety of conditions in which scalable coding may be needed. Further, both spatial and temporal scalability with rectangular VOPs and temporal scalability of arbitrary shape VOPs is expected to be supported in MPEG-4 Version 1. Figures 15 and 16 are applicable not only to rectangular VOP scalability but also to arbitrary shape VOP scalability (in this case only the shaded region depicting the head and shoulder view is used for predictions in scalable coding, and the rectangle represents the bounding box). MPEG-4 Video Version 1 supports the spatial and temporal scalability with rectangular VOPs as well as temporal scalability with arbitrary shape VOPs, however, spatial scalability with arbitrary shape VOPs will be supported in Version 2.

4.4. Robust video coding

Truly robust video coding requires a diversity of strategies. MPEG-4 video offers a number of tools which an encoder operating in the error resilient mode [37] can employ. MPEG-4 video also offers other tools (not specific to error resilience) that can be used to provide robust video coding. We now discuss the various available tools and how they can be used by themselves or in conjunction to provide robust video coding.

4.4.1. Object priorities

The object based organization of MPEG-4 video potentially makes it easier to achieve a higher degree of error robustness due to the possibility of prioritizing each semantic object based on its relevance. Further, MPEG-4 systems, since it offers scene description and composition flexibilities can ensure that the reconstructed scenes are meaningful even if low priority objects are only partially available or become unavailable (say due to data loss or corruption). Currently, MPEG-4 systems offers a way of providing priorities to each stream; if these are found insufficient, priorities may also be assigned to VOs and VOLs in the video bitstream (this has not yet been resolved). Further, VOP types themselves lend to a form of automatic prioritization

Table 5
Recommended spacings for resync markers.

Bit-rate (kbit/s)	Spacing (bits)
≤ 24	480
25–48	736

since, B-VOPs are noncausal and do not contribute to error propagation and thus can be assigned a lower priority and perhaps even be discarded in case of severe errors.

Although in principle, coding of scenes as arbitrary shaped objects can be advantageous, the down side can be shape overhead, increase in decoding complexity and the sensitivity of shape coding (which can be interframe and context dependent) to errors.

4.4.2. Resynchronization

VOPs already offer a means of resynchronization to prevent accumulation of errors. Further, it is possible for an encoder to offer increased error resilience by placing resynchronization (resync) markers in the bitstream with approximately constant spacing. The resync marker is a unique 17 bit code that normally can not be emulated by any valid combination of VLC codewords that may precede it. The VM7 error resilient encoding recommends the spacings of table 5, in bits as a function of the coding bit-rate.

In fact, to enable recovery from errors, a video packet header is used which in addition to resync marker contains macroblock number, quantizer scale and an optional extension header (when present, it includes VOP time and coding type information). The timing information ensures that the decoder can determine the VOP to which the video packet header belongs.

4.4.3. Data partitioning

Data partitioning provides a mechanism to increase error resilience by separating the normal motion and texture data of all macroblocks in a video packet and send all of the motion data followed by a motion marker, followed by all of the texture data. The motion marker is a unique 17 bit code that cannot be emulated by any valid combination of VLC codewords that may precede it.

The motion data per macroblock is arranged to contain coded/not coded information followed by combined macroblock type and coded block pattern information followed by motion vector(s); the motion data of the next macroblock follows that of the previous macroblock till the motion data for all macroblocks in the video packet can be sent. The texture data per macroblock is arranged in two parts. The first part contains coded block information of luminance blocks in a macroblock followed by optional differential quantizer information, this is repeated for all macroblocks in the video packet. The second part contains coded DCT coefficients of a macroblock, followed by that of the next macroblock till DCT coefficients for all macroblocks in the video packet can be sent.

4.4.4. Reversible VLCs

The reversible VLCs offer a mechanism for a decoder to recover additional texture data in the presence of errors since the special design of reversible VLCs enables decoding of codewords in both the forward (normal) and the reverse direction. The encoder decides whether for coding of DCT coefficients, to use the reversible VLCs or normal VLCs (depending on the coding efficiency versus error resilience tradeoffs needed) by signalling this information as part of the bitstream. It is possible to invoke the error resilience mode independent of whether reversible VLCs are used or not.

The process of additional texture recovery in a corrupted bitstream starts by first detecting the error and searching forward in the bitstream to locate the next resync marker. Once the next resync marker is located, from that point, due to use of reversible VLCs for texture coding, the texture data can be decoded in the reverse direction until an error is detected. Further, when errors are detected in texture data, the decoder can use correctly decoded motion vector information to perform motion compensation and conceal these errors.

4.4.5. Other tools: intra update and scalable coding

Typically, intra coding of macroblocks although it can provide refresh from coding errors, is expensive when used very frequently due to higher coding cost when video data is coded in intra mode. In MPEG-4, considerable effort has been placed in improving the efficiency of intra coding and the resulting scheme offers higher efficiency than H.263 or MPEG-1 based intra coding. Thus, encoders requiring higher error resilience can choose to code increased number of macroblocks in intra coding mode than with previous standards, providing an improved refresh from coding errors.

Scalable coding can offer a means of graceful degradation in quality when packet errors due to noisy conditions or packet losses due to congestion on the network are likely. Since scalable coding involves independent coding of the base layer, the base layer data can be assigned a higher priority and be better protected. Since the enhancement layers only offer improvement in spatial or temporal resolution, the enhancement layer data can be assigned a lower priority.

4.4.6. Correction and concealment strategies

Due to the channel specific nature of the degree and type of error correction needed, MPEG-4 is not likely to recommend a specific error correction method, but leaves it up to the chosen data transport layer to implement the needed technique. Further, error concealment strategies although encouraged are not standardized by MPEG-4; perhaps the work done on this topic in MPEG-2 can be useful.

Table 6
Viseme number 1 to 14, its related phoneme set and examples.

1	2	3	4	5	6	7	8	9	10	11	12	13	14
(p,b,m)	(f,v)	(T,D)	(t,d)	(k,g)	(tS,dZ,S)	(s,z)	(n,l)	(r)	(A:)	(e)	(I)	(Q)	(U)
put	far	think	tip	call	chair	sir	lot	red	car	bed	tip	top	book
bed	voice	that	doll	gas	join	zeal	not						
mill					she								

4.5. Facial animation coding

The Facial Animation Parameters (FAPs) and the Facial Definition Parameters (FDPs) [35] are sets of parameters designed to allow animation of faces reproducing expressions, emotions and speech pronunciation, as well as, definition of facial shape and texture. The same set of FAPs when applied to different facial models result in reasonably similar expressions and speech pronunciation without the need to initialize or calibrate the model. The FDPs, on the other hand, allow the definition of a precise facial shape and texture in the setup phase. If the FDPs are used in the setup phase, it is also possible to precisely produce the movements of particular facial features. Using a phoneme to FAP conversion it is possible to control facial models accepting FAPs via text to speech (TTS) systems; this conversion is not standardized. Since it is assumed that every decoder has a default face model with default parameters the set up stage is not necessary to create face animation but for customizing the face at the decoder.

The FAP set contains two high level parameters *visemes* and *expressions*. A viseme is a visual correlate to a phoneme. The viseme parameter allows viseme rendering (without having to express them in terms of other parameters) and enhances the result of other parameters, ensuring the correct rendering of visemes. Only static visemes which are clearly distinguished are included in the standard set; examples of such visemes are shown in table 6. The expression parameter similarly allows the definition of high level facial expressions. The facial expression parameter values are defined by textual descriptions such as, joy, sadness, anger, fear, disgust, surprise.

All the parameters involving translational movement are expressed in terms of the Facial Animation Parameter Units (FAPU). These units are defined in order to allow interpretation of the FAPs on any facial model in a consistent way, producing reasonable results in terms of expression and speech pronunciation. The measurement units are shown in figure 17 and are defined as follows:

IRISD0: Iris diameter (by definition it is equal to the distance between upper and lower eyelid) in neutral face; $IRISD = IRISD0/1024$;

ES0: eye separation; $ES = ES0/1024$;

ENS0: eye-nose separation; $ENS = ENS0/1024$;

MNS0: mouth-nose separation; $MNS = MNS0/1024$;

MW0: mouth width; $MW = MW0/1024$;

AU: angle unit = 10^{-5} rad.

The FDPs (figure 18) are used to customize the proprietary face model of the decoder to a particular face or to

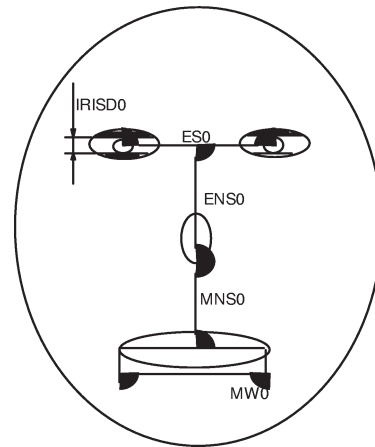


Figure 17. Facial animation parameter units.

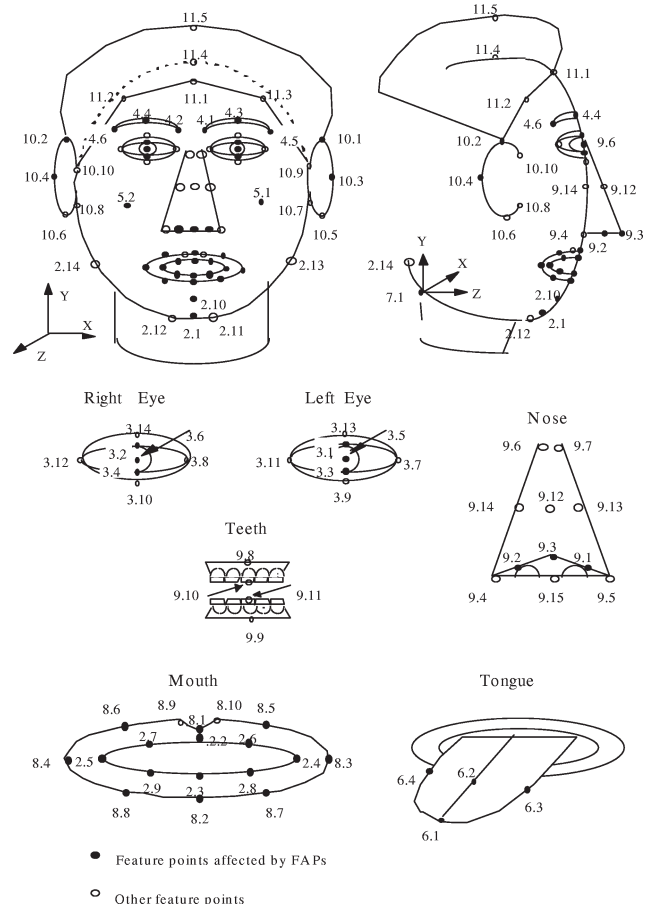


Figure 18. Facial definition feature set.

download a face model along with the information about how to animate it. The FDPs are normally transmitted once per session, followed by a stream of compressed FAPs. However, if the decoder does not receive the FDPs, the use of FAPUs ensures that it can still interpret the FAP stream. This ensures minimal operation in broadcast or teleconferencing applications.

The FDP set is specified using the FDP node (in MPEG-4 systems) which defines the face model to be used at the receiver. Two options are supported:

- Calibration information is downloaded so that the proprietary face of the receiver can be configured using facial feature points and optionally a 3D mesh or texture.
- A face model is downloaded with the animation definition of the Facial Animation Parameters. This face model replaces the proprietary face model in the receiver.

4.6. Object mesh coding

For general natural or synthetic visual objects, mesh based representation [35] can be useful for enabling a number of functions such temporal rate conversion, content manipulation, animation, augmentation (overlay), transfiguration (merging or replacing natural video with synthetic) and others. MPEG-4 includes a tool for triangular mesh based representation of general purpose objects.

A visual object of interest, when it first appears (as a 2D VOP) in the scene, is tassellated into triangular patches resulting in a 2D triangular mesh. The vertices of the triangular patches forming the mesh are referred to as the *node points*. The node points of the initial mesh are then tracked as the VOP moves within the scene. The 2D motion of a Video Object can thus be compactly represented by the motion vectors of the node points in the mesh. Motion compensation can then be achieved by texture mapping the patches from VOP to VOP according to affine transforms. Coding of video texture or still texture (to be discussed next) of object is performed by the normal texture coding tools of MPEG-4. Thus, efficient storage and transmission of the mesh representation of a moving object (dynamic mesh) requires compression of its geometry and motion.

The initial 2D triangular mesh is either a uniform mesh or a Delaunay mesh, the mesh triangular topology (links between node points) is not coded; only the 2D node point coordinates $\vec{p}_n = (x_n, y_n)$ are coded. A uniform mesh can be completely specified using five parameters such as the number of nodes horizontally and the number of nodes vertically, the horizontal and the vertical dimensions of each quadrangle consisting of two triangles, and the type of splitting applied on each quadrangle to obtain triangles. For a Delaunay mesh, the node point coordinates are coded by first coding the boundary node points and then the interior node points of the mesh. To encode the interior node positions, the nodes are traversed one by one using a *nearest*

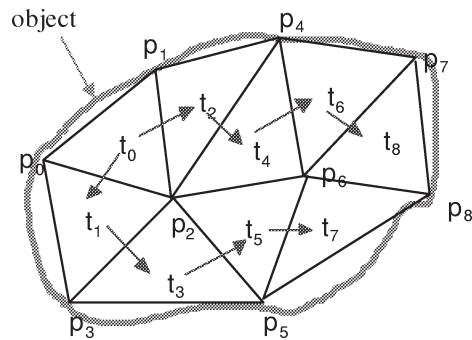


Figure 19. 2D Mesh representation of an object, and coding of mesh geometry.

neighbor strategy. A linear ordering of the node points is computed such that each node is visited only once. When a node is visited, its position is differentially coded with respect to the position of previous coded node used as the predictor. By sending the total number of node points and the number of boundary node points, the decoder knows how many node points will follow, and how many of those are boundary nodes; thus it is able to reconstruct the polygonal boundary and the locations of all nodes. The mesh based representation of an object and the traversal of nodes for mesh geometry coding is illustrated in figure 19 by an example.

First, the total number of nodes and the number of boundary nodes is encoded. The top-left node \vec{p}_0 is coded without prediction. Then, the next clockwise boundary node \vec{p}_1 is found and the difference between \vec{p}_0 and \vec{p}_1 is encoded; then all other boundary nodes are encoded in a similar fashion. Then, the not previously encoded interior node that is nearest to the last boundary node is found and the difference between these is encoded; this process is repeated until all the interior nodes are covered. The mesh geometry is only encoded when a new mesh needs to be initialized with respect to a particular VOP of the corresponding visual object; it consists of the initial positions of the mesh nodes.

The mesh motion is encoded [52] at subsequent time instants to describe the motion of the corresponding video object; it consists of a motion vector for each mesh node such that the motion vector points from a node point of the previous mesh in the sequence to a node point of the current mesh. The mesh bitstream syntax consists of the two parts: mesh geometry and mesh motion. The node coordinates and node motion vectors are specified to one-half pixel accuracy.

4.7. Still texture coding

The Discrete Wavelet Transform (DWT) [35,37] is used to code still image data employed for texture mapping. Besides coding efficiency, an important requirement for coding texture map data is that it should be coded in a manner facilitating continuous scalability, thus allowing many resolution/qualities to be derived from the same coded bitstream.

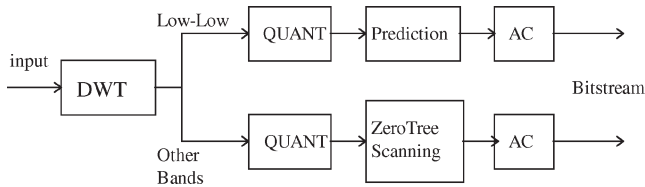


Figure 20. Block diagram of the DWT encoder of still image texture.

While DCT based coding is able to provide comparable coding efficiency as well as a few scalability layers, DWT based coding offers flexibility in organization and number of scalability layers.

The principle of DWT encoding is shown in figure 20.

The basic modules of a zero-tree wavelet based coding scheme are as follows:

1. Decomposition of the texture using discrete wavelet transform (DWT).
2. Quantization of the wavelet coefficients.
3. Coding of the lowest frequency subband using a predictive scheme.
4. Zero-tree scanning of the higher order subband wavelet coefficients.
5. Entropy coding of the scanned quantized wavelet coefficients and the significance map.

A 2D separable wavelet decomposition is applied to the still texture to be coded. The wavelet decomposition is performed using a Daubechies (9,3) tap biorthogonal filter which has been shown to provide good compression performance. The filter coefficients are:

Lowpass =

$$\begin{bmatrix} 0.033 & 145 & 630 & 368 & 12 & -0.066 & 291 & 260 & 736 & 24 & -0.176 & 776 & 695 & 296 & 65 \\ 0.419 & 844 & 651 & 329 & 52 & 0.994 & 368 & 911 & 043 & 60 & 0.419 & 844 & 651 & 329 & 52 \\ 0.176 & 776 & 695 & 296 & 65 & -0.066 & 291 & 260 & 736 & 24 & 0.033 & 145 & 630 & 368 & 12 \end{bmatrix}$$

Highpass =

$$[-0.353 & 553 & 390 & 593 & 27 & 0.707 & 106 & 781 & 186 & 55 & -0.353 & 553 & 390 & 593 & 27]$$

A group delay is applied to each filter to avoid the phase shift on both of the image domain and the wavelet domain.

The wavelet coefficients of the lowest band are coded independently from the other bands. These coefficients are quantized using a uniform midrise quantizer. After quantization of the lowest subband coefficients, an implicit prediction (same as that used for DC prediction in intra DCT coding) is applied to compute the prediction error which is then encoded using an adaptive arithmetic coder which uses min-max coding.

The wavelet coefficients of the higher bands are first quantized by multilevel quantization which provides the flexibility needed to tradeoff number of levels, type of scalability (spatial or SNR), complexity and coding efficiency. Different quantization step sizes (one for luminance and one for chrominance) can be specified for each level of scalability. All the quantizers of the higher bands are uniform

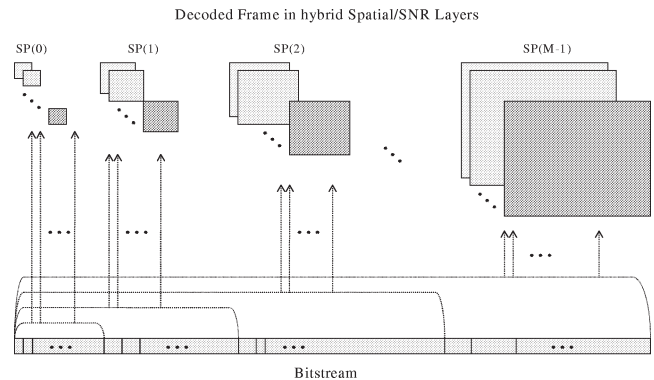


Figure 21. Scalable decoding of still object texture.

mid-rise quantizer with a dead zone 2 times the quantization step size. The quantization step sizes are specified by the encoder in the bitstream. In order to achieve the fine granularity of SNR scalability, a bi-level quantization scheme is used for all the multiple quantizers. This quantizer is also a uniform mid rise quantizer with a dead zone 2 times the quantization step size. The coefficients that lie outside the dead zone (in the current and previous pass) are quantized with a 1 bit accuracy. The number of quantizers is equal to the maximum number of bitplanes in the wavelet transform representation. In this bi-level case, instead of the quantization step sizes, the maximum number of the bitplanes is specified in the bitstream.

After quantization, each wavelet coefficient is either zero or nonzero. The coefficients of all bands (except the lowest) are scanned by zero-tree scanning. Zero-tree scanning is based on the observation that strong correlation exists in the amplitudes of the wavelet coefficients across scales, and on the idea of partial ordering of the coefficients. The coefficient at the coarse scale is called the parent, and all coefficients at the same spatial location, and of similar orientation, at the next finer scale are that parent's children. Since the lowest frequency subband is coded separately, the wavelet trees start from the adjacent higher bands. In order to achieve a wide range of scalability levels efficiently as needed by the application, a multiscale zerotree coding scheme is employed. The zero-tree symbols and the quantized values are coded using an adaptive arithmetic coder which uses a three-symbol alphabet.

The process of scalable decoding the various spatial/SNR layers from a single DWT coded bitstream is shown in figure 21.

5. MPEG-4 audio

The benefits of the MPEG-4 speech coder can be exploited in a number of applications. As an example, the MPEG-4-based Internet-phone system offers robustness against packet loss or change in transmission bit-rates. Furthermore, low bit-rate is useful for "Party talk". Besides speech coding MPEG-4 also offers multichannel audio coding based on optimized MPEG-2 AAC coding. MPEG-4

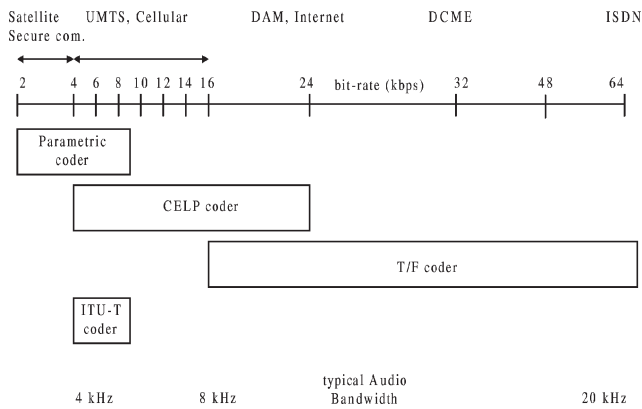


Figure 22. Natural audio coding in MPEG-4.

also offers solutions for medium bit-rate audio coding. Furthermore, it supports the concept of audio objects. Just as video scenes are made from visual objects, audio scenes may be usefully described as the spatiotemporal combination of audio objects. An “audio object” is a single audio stream coded using one of the MPEG-4 coding tools, like CELP or Structured Audio. Audio objects are related to each other by mixing, effects processing, switching, and delaying them, and may be spatialized to a particular 3D location. The effects processing is described abstractly in terms of a signal-processing language (the same language used for Structured Audio), so content providers may design their own empirically, and include them in the bitstream.

5.1. Natural audio

As mentioned earlier, Natural Audio coding [16,28] in MPEG-4 consists of the following:

- The lowest bit-rate range between 2 and 6 kbit/s is covered by *Parametric Coding* (mostly for speech coding).
- The medium bit-rates between 6 and 24 kbit/s use *Code Excited Linear Predictive (CELP)* with two sampling rates, 8 and 16 kHz, for a broader range of audio signals.
- For higher bit-rates starting at about 16 kbit/s, frequency domain coding techniques are applied, e.g., optimized version of MPEG-2 *Advanced Audio Coding (AAC)*. The audio signals in this region typically have bandwidths starting at 8 kHz.

Figure 22 provides a composite picture of the applications of MPEG-4 audio and speech coding, the signal bandwidth and the type of coders used.

5.1.1. Parametric coder

The parametric coder core provides two sets of tools. The HVXC coding tools (Harmonic Vector eXcitation Coding) allow coding of speech signals at 2 kbit/s while the Individual Line coding tools allow coding of non-speech signals like music at bit rates of 4 kbit/s and higher. Both sets of tools allow independent change of speed and pitch during the decoding and can be combined to handle a wider range of signals and bit rates.

5.1.2. CELP coder

The CELP coder is designed for speech coding at two different sampling frequencies, namely, 8 kHz and 16 kHz. The speech coders using 8 kHz sampling rate are referred to as narrowband coders while those using 16 kHz sampling rate are wideband coders. The CELP coder includes tools offering a variety of functions including bit rate control, bit rate scalability, speed control, complexity scalability and speech enhancement. Using the narrowband and the wideband CELP coders, it is possible to span a wide range of bit rates (4 kbps to 24 kbps). Real-time bit-rate control in small steps can be provided. A common structure of tools have been defined for both the narrowband and wideband coders; many tools and processes have been designed to be commonly usable for both narrowband and wideband speech coders.

5.1.3. Time/frequency coder

The high-end audio coding in MPEG-4 [28] is based on MPEG-2 AAC coding. MPEG-2 AAC is a state-of-the-art audio compression algorithm that provides compression superior to that provided by older algorithms. AAC is a transform coder and uses a filterbank with a finer frequency resolution that enables superior signal compression. AAC also uses a number of new tools such as temporal noise shaping, backward adaptive linear prediction, joint stereo coding techniques and Huffman coding of quantized components, each of which provide additional audio compression capability. Furthermore, AAC supports a wide range of sampling rates and bit-rates, from one to 48 audio channels, up to 15 low frequency enhancement channels, multi-language capability and up to 15 embedded data streams. MPEG-2 AAC provides a 5-channel audio coding capability, while being a factor of two better in coding efficiency relative to MPEG-2 BC; since AAC has no such backward compatibility requirement and it incorporates the recent advances, in MPEG formal listening tests for 5-channel audio signals, it provided slightly better audio quality at 320 kbit/s than MPEG-2 BC can provide at 640 kbit/s.

5.2. Text to Speech

The Text-to-Speech (TTS) conversion system synthesizes speech as its output when a text is accessed as its input. In other words, when the text is accessed, the TTS changes the text into a string of phonetic symbols and the corresponding basic synthetic units are retrieved from a pre-prepared database. Then the TTS concatenates the synthetic units to synthesize the output speech with the rule-generated prosody. Some application areas for MPEG-4 TTS are as follows:

- artificial story teller (or story teller on demand);
- synthesized speech output synchronized with Facial Animation (FA);
- speech synthesizer for avatars in various Virtual Reality (VR) applications;

- voice news paper;
- dubbing tools for animated pictures;
- voice Internet;
- transportation timetables.

The MPEG-4 TTS [28] can not only synthesize speech according to the input speech with a rule-generated prosody, but also executes several other functions. They are as follows:

- (1) Speech synthesis with the original prosody from the original speech.
- (2) Synchronized speech synthesis with Facial Animation (FA) tools.
- (3) Synchronized dubbing with moving pictures not by recorded sound but by text and some lip shape information.
- (4) Trick mode functions such as stop, resume, forward, backward without breaking the prosody even in the applications with Facial Animation (FA)/ Motion Pictures (MP).
- (5) Users can change the replaying speed, tone, volume, speaker's sex, and age.

MPEG-4 TTS can be used for many languages in the world since it adopts the concept of the language code such as the country code for an international call. Presently, only 25 countries, i.e., the current ISO members, have their own code numbers, to identify that their own language has to be synthesized, except the International Phonetic Alphabet (IPA) code assigned as 0. However, 8 bits have been assigned for the language code to ensure that all countries can be assigned language code when asked in the future. IPA could be used to transmit all languages.

For MPEG-4 TTS, only the interface bitstream profiles are the subject of standardization. Because there are already many different types of TTS and each country has several or a few tens of different TTSs synthesizing its own language, it is impossible to standardize all the things related to TTS. However, it is believed that almost all TTSs can be modified to accept MPEG-4 TTS interface very quickly by a TTS expert because of the rather simple structure of the MPEG-4 TTS interface bitstream profiles.

5.3. Structured audio

Structured audio formats use ultra-low bit-rate algorithmic sound models to code and transmit sound. MPEG-4 standardizes an algorithmic sound language and several related tools for the structured coding of audio objects. Using these tools, algorithms which represent the exact specification of a sound scene are created by the content designer, transmitted over a channel, and executed to produce sound at the terminal. Structured audio techniques in MPEG-4 [28] allow the transmission of synthetic music and sound

effects at bit-rates from 0.01 to 10 kbit/s, and the concise description of parametric sound post-production for mixing multiple streams and adding effects processing to audio scenes

MPEG-4 does not standardize a synthesis method, but a signal-processing language for describing synthesis methods. SAOL, pronounced "sail", stands for "Structured Audio Orchestra Language" and is the signal-processing language enabling music-synthesis and effects post-production in MPEG-4. It falls into the music-synthesis category of "Music V" languages; that is, its fundamental processing model is based on the interaction of oscillators running at various rates. However, SAOL has added many new capabilities to the Music V language model which allow for more powerful and flexible synthesis description. Using this language, any current or future synthesis method may be described by a content provider and included in the bitstream. This language is entirely normative and standardized, so that every piece of synthetic music will sound exactly the same on every compliant MPEG-4 decoder, which is an improvement over the great variety in MIDI-based synthesis systems.

The techniques required for automatically producing a Structured Audio bitstream from an arbitrary sound are beyond today's state of the art and are referred to as "automatic source separation" or "automatic transcription". In the mean time, content authors will use special content creation tools to directly create Structured Audio bitstreams. This is not a fundamental obstacle to the use of MPEG-4 Structured Audio, because these tools are very similar to the ones that content authors use already; all that is required is to make them capable of producing MPEG-4 output bitstreams. There is no fixed complexity which is adequate for decoding every conceivable Structured Audio bitstream. Simple synthesis methods are very low-complexity, and complex synthesis methods require more computing power and memory. Since the description of the synthesis methods is under the control of the content provider, the content provider is responsible for understanding the complexity needs of the bitstreams. Past versions of structured audio systems with similar capability have been optimized to provide multitimbral, highly polyphonic music and post-production effects in real-time on a 150 MHz Pentium computer or simple DSP chip.

6. MPEG-4 systems

The Systems [33] part of MPEG-4 perhaps represents the most radical departure from previous MPEG specifications. Indeed, the object-based nature of MPEG-4 necessitates a totally new approach on what the Systems layer is required to provide. Issues of synchronization and multiplexing are, of course, still very essential. Note, however, that an MPEG-4 scene may be composed of several objects, and hence synchronization between a large number of streams is required. In addition, the spatio-temporal po-

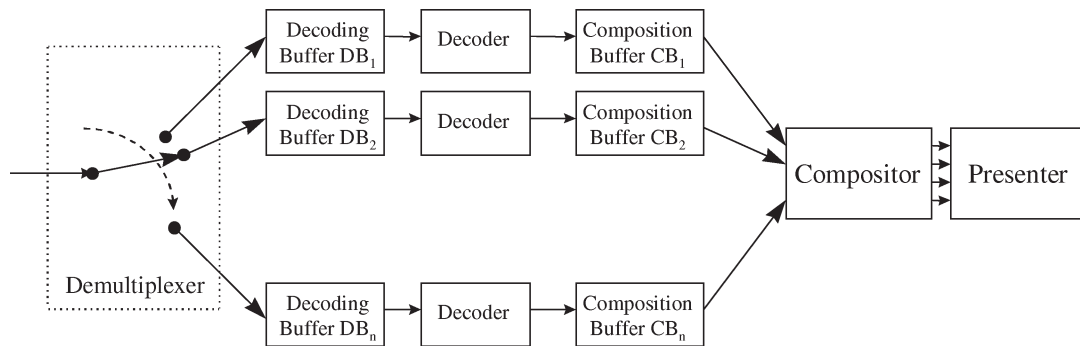


Figure 23. System decoder model.

sitioning of such objects forming a scene (or scene description) is a key component. Finally, issues of interactivity are also quite new in MPEG as well.

6.1. System decoder model

A key problem in designing an audiovisual communication system is ensuring that time is properly represented and reconstructed by the terminal. This serves two purposes: (1) it ensures that “events” occur at designated times as indicated by the content creator, and (2) that the sender can properly control the behavior of the receiver. The latter is essentially providing an open-loop flow control mechanism, a requirement for a specification that covers broadcast channels (without an upstream, or feedback, channel). Assuming a finite set of buffer resources at the receiver, by proper clock recovery and timestamping of events the source can always ensure that these resources are not exhausted.

Clock recovery is typically performed using clock references. The receiving system has a local system clock which is controlled by a PLL, driven by the differences in received clock references and the local clock references at the time of their arrival. This way the receiver’s clock speed is increased or decreased, matching that of the sender. In addition, coded units are associated with decoding time stamps, indicating the time instance in which a unit is removed from the receiving system’s decoding buffer. The combination of clock references and time stamps is sufficient for full control of the receiver.

In order to properly address specification issues without making unnecessary assumption about system implementation, MPEG-4 Systems defines a System Decoder Model. This represents an idealized unit, in which operations can be unambiguously controlled and characterized. The MPEG-4 System Decoder Model exposes resources available at the receiving terminal, and defines how they can be controlled by the sender or content creator. The model is shown in figure 23. It is composed of a set of decoders (for the various audio or visual object types), provided with two types of buffers: decoding and composition.

The decoding buffers have the same functionality as in previous MPEG specifications, and are controlled by clock references and decoding timestamps. In MPEG-2, each

program had its own clock. Proper synchronization was ensured by using the same clock for coding and transmitting the audio and video components. In MPEG-4, each individual object is assumed to have its own clock, or Object Time Base (OTB). Of course, several objects may share the same clock. In addition, coded units of individual objects (Access Units – AUs, corresponding to an instance of a video object or a set of audio samples) are associated with Decoding Timestamps (DTSs). Note that the decoding operation at DTS is considered (in this ideal model) to be instantaneous.

The composition buffers which are present at the decoder outputs form a second set of buffers. Their use is related to object persistence. In some situations, a content creator may want to reuse a particular object after it has been presented. By exposing a composition buffer, the content creator can control the lifetime of data in this buffer for later use. This lifetime is controlled by an Expiration Timestamp (ETS). A decoded object is assumed to remain in this buffer until its ETS is reached, and can be used repeatedly until that time. This feature may be particularly useful in low bandwidth wireless environments. MPEG-4 defines an additional timestamp, the Composition Timestamp (CTS), which defines the time at which data is taken from the composition buffer for (instantaneous) composition and presentation.

In order to coordinate the various objects, a single System Time Base (STB) is assumed to be present at the receiving system. All object time bases are subsequently mapped into the system time base, so that a single notion of time exists in the terminal. For clock recovery purposes, a single stream must be designated as the master. The current specification does not indicate the stream that has this role, but a plausible candidate is the one which contains the scene description. Note also that, in contrast with MPEG-2, the resolution of both the STB and the Object Clock References (OCRs) is not mandated by the specification. In fact, the size of the OCR fields for individual Access Units is fully configurable, as discussed later.

In other designs (e.g., IETF’s RTP), the assumption of a globally known clock can be made (provided by other network services such as NTP). There is work underway to provide a unified methodology so that mapping of MPEG-4 timing architecture can be seamlessly performed in such

an environment as well. This is facilitated by the flexible multiplexing methodology adopted in MPEG-4, and is discussed later.

6.2. Scene description

Scene description refers to the specification of the spatio-temporal positioning and behaviour of individual objects. It is a totally new component in the MPEG specifications, and allows the easy creation of compelling audiovisual content. Scene description involves an architectural component, i.e., the proper way to conceptually organize audiovisual information, and a syntactic component which is the mapping of such architecture into a bitstream. Note that the scene description is transmitted in a separate stream from the individual media objects. This allows one to change the scene description without operating on any of the constituent objects themselves.

6.2.1. Architecture

The architecture of MPEG-4's scene description is based on the Virtual Reality Modeling Language (VRML) [2,20]. Scenes are described as hierarchies of nodes forming a tree. Leafs of the tree correspond to media objects (audio or visual, natural or synthetic), while intermediate nodes perform operations on their underlying nodes (grouping, transformation, etc.). Nodes also expose attributes through which their behaviour can be controlled. This hierarchical design is shown in figure 24.

There are however several key differences between the domains that VRML and MPEG-4 address, which necessitate the adoption of slightly different approaches. MPEG-4 is concerned with the description of highly dynamic scenes, where temporal evolution is more dominant than navigation. VRML, on the other hand, allows the definition of a static 3D world (static in the sense that all the objects in that world are predefined and cannot be changed) in which a user is allowed to navigate. As a result, scene descriptions in MPEG-4 actually have their own time base, and can be updated at any time. Updates can take the form of a node's replacement, elimination, insertion, or attribute value modification. The presence of a time base and decoding time stamps ensures the application of such updates at the correct time instances.

In addition, MPEG-4 also needs to address pure 2D composition of objects. The complexity and cost of 3D graphics versus the possibility of developing low-cost or low-power systems makes it desirable to partition the space of graphics capabilities. In figure 24 we see an example where both 2D and 3D components coexist.

As a result of the close relationship of the MPEG-4 scene description with VRML, the types of scene description capabilities provided by MPEG-4 are essentially those provided by VRML nodes. It is currently examined if MPEG-4 can be cast as an extension of VRML, with appropriately defined subsets as MPEG-4 specific profiles.

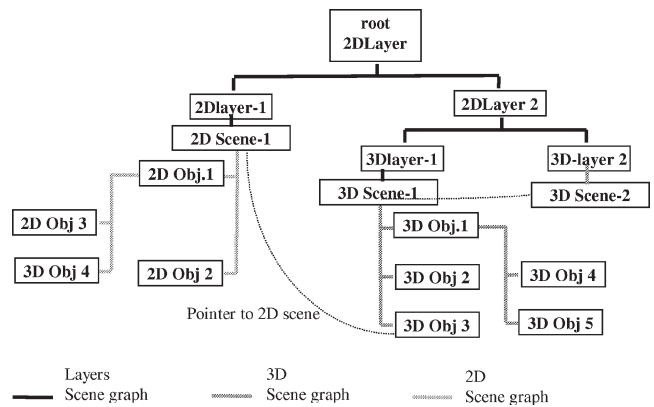


Figure 24. Scene description.

6.2.2. Binary format for scenes

The mapping of the scene description into a parametric form, suitable for low-overhead transmission has resulted in a scene description format called BINARY Format for Scenes (BIFS). This representation format associates each node with a node type. Nodes are then represented by their node type and a set of attribute value specifications. To avoid the specification of all attributes of a node, default values are used, and attributes are individually addressable within a node. This way one can, for example, only specify a non-default value for attribute X of a node of type Y.

Furthermore, nodes can optionally be reused. In this case, they are associated with an identifier which can appear in place of a node of the same type. In figure 24, for example, the node "3D Obj 3" contains the child node "2D Scene-1" which is used elsewhere in the scene as well. As mentioned earlier, scene descriptions can be updated: nodes can be inserted, deleted, replaced, or their attributes can be modified.

In addition to the VRML set of nodes, MPEG-4 defines its own set of nodes, particularly to handle media objects including natural audio and video, still images, face animation, basic MIDI, text-to-speech synthesizer, streaming text, 2-D composition operators, layout control, etc.

6.2.3. Interactivity

Interactivity in MPEG-4 can take two forms, client-based and server-based. In the client-based case, user operations affect the local scene description. The VRML model of events and routes is used within the BIFS format, thus providing direct support for a large variety of circumstances. In addition, and noting that user events can be transformed to scene updates, one can also have a form of interactivity that does not necessitate normative support by the specification: user events can form a secondary source of scene updates.

Server-based interaction requires the presence of an upstream channel. The details of such a mode of interactivity are still under investigation, as they are slightly complicated by the needs for network-independence.

6.2.4. Adaptive audio-visual session

An additional, flexible mechanism for describing scenes is also being developed, called Adaptive Audio-Visual Session (AAVS). This is based on the use of the Java language for constructing scenes, but not for composition, rendering, or decoding. By taking the programmable aspect of the system outside of the main data flow of decoders and the rendering engine (which have to operate extremely fast), overall performance can be kept high. Note that the AAVS approach is being designed as an extension of BIFS. In that sense, a BIFS scene can be a subset of an AAVS scene (but not vice versa). The AAVS work although started for MPEG-4 Version 1, will only be included in Version 2.

The benefits of the AAVS approach are in three major categories. First, because of its flexible nature, the scene description is capable of adapting to terminal capabilities hence providing a form of graceful degradation. Secondly, when scene description operations can be expressed concisely in a flexible way (e.g., a spline trajectory), using a flexible approach may result in increased compression. Finally, by allowing programmability, user interaction can extend beyond the modes defined in a parametric scene description and become as rich as the content creator wishes. The use of programmability in an audiovisual terminal certainly opens up a very broad spectrum of opportunities for a new generation multimedia applications.

6.3. Associating scene description with elementary streams

6.3.1. Object descriptors

Individual object data as well as scene description information is carried in separate Elementary Streams (ES). As a result, BIFS media nodes need a mechanism to associate themselves with the ESs that carry their data (coded natural video object data etc.). A direct mechanism would necessitate the inclusion of transport-related information into the scene description. As we mentioned earlier, an important requirement in MPEG-4 is transport-independence. As a result, an indirect way was adopted, using Object Descriptors (ODs).

Each media node is associated with an object identifier, which in turn uniquely identifies an OD. Within an OD, there is information on how many ESs are associated with this particular object (may be more than one for scalable video/audio coding, or multi-channel audio coding), and information describing each of those streams. The latter information includes the type of the stream, as well as how to locate it within the particular networking environment used. This approach simplifies remultiplexing (e.g., going through a wired-wireless interface), as there is only one entity that may need to be modified.

6.3.2. Stream map table

The object descriptor allows unique reference of an elementary stream by an id; this id may be assigned by an application layer when the content is created. The transport channel in which this stream is carried may only be assigned

at a later time by a transport entity; it is identified by a channel association tag associated to an elementary stream id by a stream map table. In interactive applications, the receiving terminal may select the desired elementary streams, send a request and receive the stream map table in return. In broadcast and storage applications, the complete stream map table must be included in the applications signalling channel.

6.4. Multiplexing

The key underlying concept in the design of the MPEG-4 multiplexer is network independence. MPEG-4 content may be delivered across a wide variety of channels, from very low bit rate wireless, to high-speed ATM and broadcast systems, to DVDs. A critical design question was what should be the tools included in the specification for mandatory implementation. Clearly, the broad spectrum of channels could not allow a single solution to be used. At the same time, inclusion of a large number of different tools and configuration would make implementations extremely complex and – through excessive fragmentation through profiles – make interoperability extremely hard to achieve in practice. Consequently, the assumption was made that MPEG-4 would not provide specific transport-layer features but would instead make sure that it could be easily mapped to existing such layers. This is accomplished by allowing several components of the multiplexer to be configurable, thus allowing designers to achieve the desired trade-off between functionality and efficiency. In addition, it is assumed that Quality of Service (QoS) guarantees may be made available by the underlying transport service if so desired.

The overall multiplexing architecture of MPEG-4 is shown in figure 25. At the top-most level, we have the Access Unit Layer (AL) which provides the basic conveyor of timing and framing information. It is at this level where timestamps and clock references are provided. The AL header, however, is very flexible: the presence of OCR/DTS is optional, and their resolution (number of bits) is configurable. In addition, the header contains information about framing (start or end of a coded unit, random access indicator), as well as sequence numbering. The latter is particularly useful in error-prone wireless or broadcast environments, where preemptive retransmission of critical information (e.g., scene description) may be performed. Using a sequence number a receiver that has already accurately received the information can ignore the duplicate. In order to be able to “bootstrap” the demultiplexing process, the AL header configuration information is carried in a channel that has a predefined configuration.

Immediately below the AL we have the optional Flexible Multiplexer or “FlexMux” layer. This is a very simple design, intended for systems that may not provide native multiplexing services. An example is the data channel available in GSM cellular telephones. Its use, however, is entirely optional, and does not affect the operation of the

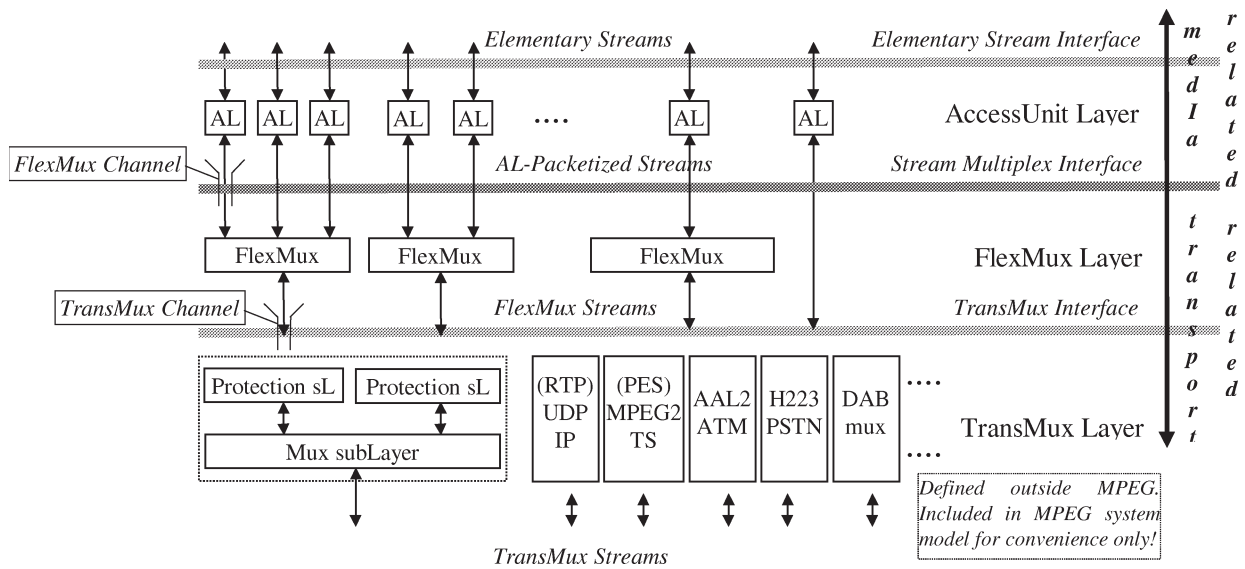


Figure 25. MPEG-4 multiplexing architecture.

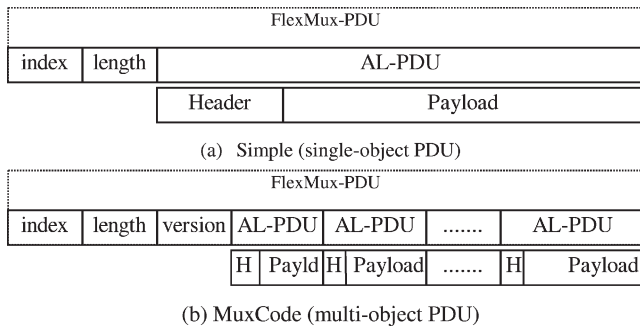


Figure 26. FlexMux modes.

rest of the system. As shown in the right side of figure 25 there can be a “null” connection directly from the Access Unit Layer to the lower layers. The FlexMux provides two modes of operation, a “simple” and a “muxcode” mode, as shown in figure 26.

In the simple mode, data from a single object (or scene description) is present in the FlexMux payload (appropriately encapsulated as an AL PDU). In the MuxCode mode, data from multiple objects are placed in predefined positions within the FlexMux payload. The ‘index’ field of the FlexMux header indicates which of the modes is used, depending on its value.

Finally, at the lowest level, we have the Transport Multiplexing or “TransMux” layer which is not specified by MPEG-4. This can be any current or future transport layer facility, including MPEG-2 Transport Stream, RTP, H.223, etc. For a mobile system, this layer must provide its own protection sublayer to ensure the desired QoS characteristics.

6.5. File format

The need for a file format for storage-based delivery (e.g., CD-ROM or DVD) of MPEG-4 content was recognized quite early during the MPEG-4 Systems work. In

addition to being a standardized format for interchange of coded MPEG-4 data, it was envisaged that the file format would support a host of functionalities including random access to individual objects or portions of objects and editability of coded objects. An early solution addressing these capabilities was adopted and consisted of a simple layer that substituted the TransMux layer and where random access information is provided in the form of directories; this would potentially allow a user to rapidly move back and forth in an MPEG-4 file having immediate access to Access Units of data. However, the anticipated role of file format for MPEG-4 was revised to include streaming of coded MPEG-4 among others, and with a new set of requirements in place, an open call was made inviting proposals for file format.

A number of proposals, some backed by large segments of the multimedia industry, were received. Among the proposals received were Quicktime (Apple, IBM, Oracle, Netscape, SGI and Sun), Advanced Streaming Format (Microsoft and Intel) and Integrated Intermedia Format (AT&T and Columbia Univ.). A review of the proposals revealed that for an MPEG-4 file format to be widely accepted in the multimedia community, it needed to focus on providing a rich set of functionality while supporting existing multimedia (legacy) content. One of the proposals (Quicktime) was selected as the starting basis for MPEG-4 file format. This is expected to be embellished by incorporating key tools from the other proposals to satisfy the unique requirements of MPEG-4. It has also been decided that MPEG-4 will use ‘.mp4’ as the extension for its file name. The file format will be included in MPEG-4 Version 2.

6.6. Syntactic description

Source coding, with its bit-oriented nature, directly conflicts with the byte-oriented structure of modern microprocessors and makes the task of handling coded audio-

visual information more difficult. A simple example is fast decoding of variable length codes; every programmer that wishes to use information using entropy coding must hand-code the tables so that optimized execution can be achieved. General-purpose programming languages such as C++ and Java do not provide native facilities for coping with such data. Even though other facilities already exist for representing syntax (e.g., ASN.1 – ISO International Standards 8824 and 8825), they cannot cope with the intricate complexities of source coding operations (variable length coding etc.).

MPEG-4 Systems has adopted an object-based syntactic description language for the definition of its bitstream syntax (Flavor – Formal Language for Audio-Visual Object Representation [5,10–12]). It is designed as an extension of C++ and Java in which the type system is extended to incorporate bitstream representation semantics (hence forming a syntactic description language). This allows the description, in a single place, of both the in-memory representation of data as well as their bitstream-level (compressed) representation as well. Also, Flavor is a declarative language, and does not include methods or functions. By building on languages widely used in multimedia application development, one can facilitate integration with an application's structure.

Figure 27 shows a simple example of a Flavor representation. Note the presence of bitstream representation information right after the type within the class declaration. The map declaration is the mechanism used in Flavor to introduce constant or variable length code tables (1-to- n mappings); in this case, binary codewords (denoted using the 'Ob' construct) are mapped to values of type unsigned char. Flavor also has a full complement of object-oriented features pertaining to bitstream representation (e.g., "bitstream polymorphism") as well as flow control instructions (if, for, do-while, etc.). The latter are placed within the declaration part of a class, as they control the serialization of the class' variables into a bitstream.

A translator has been developed that automatically generates standard C++ and Java code from the Flavor source code [5], so that direct access to, and generation of, compressed information by application developers can be achieved with essentially zero programming. This way, a significant part of the work in developing a multimedia

```
map SampleVLC(unsigned char) {
    0b0, 2,
    0b10, 5,
    0b11, 7
}

class HelloBits {
    int(8) size;
    int(size) value1;
    unsigned char(SampleVLC) value2;
}
```

Figure 27. A simple example of syntactic description.

application (including encoders, decoders, content creation and editing suites, indexing and search engines) is eliminated.

7. Profiles and levels

Although there are many tools in the MPEG-4 standard, not every MPEG-4 decoder will have to implement all of them. In fact, only a few classes of MPEG-4 decoders are expected to exist with each class addressing clusters of applications with similar requirements; this is accomplished via the concept of profiles. Similar to that in MPEG-2, a profile is a defined sub-set of the entire bitstream syntax of all the tools. A level is a defined set of constraints imposed on parameters in the bitstream. Conformance tests will be carried out against defined profiles at defined levels. The purpose of defining conformance points in the form of profiles and levels is to facilitate bitstream interchange among different applications. Implementors of MPEG-4 are encouraged to produce decoders and bitstreams which correspond to those defined conformance regions. The discretely defined profiles and levels are the means of bitstream interchange between different applications of MPEG-4. The concept of profiles and levels of MPEG-2 Video [15,19] has been extended to MPEG-4 Video, Audio and Systems specifications.

Currently, MPEG is defining profiles for video objects, audio objects and systems for MPEG-4. In addition, it is likely that profiles may also be defined for composition. The work on profiles is ongoing and is thus likely to evolve; we now present the current structure of profiles envisaged and their key requirements [31,47,48]. Currently, three profiles each are being considered for each of the three parts, video, audio and system objects; the requirements for these profiles are being established.

The profiles currently being considered for video are as follows:

- Video simple object profile.
- Video random access object profile.
- Video main object profile.

In addition, recently, a profile known as ultra simple profile addressing mobile and other applications has been proposed and is likely to be accepted.

Similarly, for audio, a number of profiles are being discussed and are listed as follows:

- Audio speech object profile.
- Audio low delay object profile.
- Audio main object profile.

The three profiles being considered for systems are as follows:

- Systems simple profile.
- Systems interactive profile.
- Systems main profile.

Table 7
Wireless network.

	PCS	IMT2000 (1)	IMT2000 (2)
Data rate (total)	32 kbit/s	128 kbit/s	384 kbit/s
Error conditions	Random and burst	Random and burst	Random and burst

8. Verification tests

MPEG-4 is planning to hold verification tests to confirm the performance of its various combination of tools that will form profiles. The first of the series of such tests [1] will take place in March 1998 and is aimed to verify error resilience tools.

Version 1 of the MPEG-4 video codec is already implemented in software. It will be combined with simulations of the Universal Access Layer, a component of the MPEG-4 System Layer, and a TransMux. The particular TransMux to be implemented depends upon the application and corresponding network. For wireless applications the TransMux will be ITU's mobile multiplexing standard, H.223/Mobile. The overall encoding system to be simulated for the demonstration of MPEG-4 over a wireless network consists of an application layer consisting of an MPEG-4 error resilient video encoder and simulated audio data, an access unit layer consisting of adaptation layers for audio and video data, H.223/Mobile layer consisting of adaptation layers and a multiplexer. The multiplexed data is to be stored on PC hard drive and passes through physical layer simulation to the decoding system that performs the inverse operations such as H.223/Mobile demultiplexing and adaptation layers for audio and video, access unit layer consisting of adaptation layer for audio and video data, and back to application layer which consists of simulated audio data and MPEG-4 error resilient video decoder. The test conditions being considered are shown in table 7.

9. Beyond current MPEG-4 work

As mentioned earlier, the MPEG-4 work has been recently divided into two versions. The tools that were mature for standardization have been incorporated in Version 1 and are complete while the work on the remaining tools continues for Version 2. Version 2 tools are expected to either provide new functionalities not supported by Version 1 or provide some of the same functionality but a lot more efficiently (or with higher quality or lower implementation cost). Generally, Version 2 is expected to maintain backward compatibility with Version 1.

Besides MPEG-4, the MPEG committee has recently accepted a new work item on "Multimedia content description interface", and is dubbed as MPEG-7, the next MPEG standard. MPEG-7 has the primary goal of addressing the limited capabilities and interoperability problems of proprietary systems for content search and retrieval. Thus MPEG-7 will specify [49] a standard set of descriptors that can be used to

identify various types of multimedia information. This description shall be associated with the content itself, to allow fast and efficient searching for material of a user's interest. AV material that has MPEG-7 data associated with it can thus be indexed and searched. However, MPEG-7 is not expected to specify the actual search and retrieval engine.

10. Summary

In this paper we have introduced the MPEG-4 standard currently in progress. The necessary background information leading up to the MPEG-4 standard and the multiple facets of this standard were discussed including the directions for the future. The success of MPEG-4 will eventually depend on many factors such as market needs, competing standards, software versus hardware paradigms, complexity versus functionality tradeoffs, timing, profiles, etc. Technically, MPEG-4 appears to have a significant potential due to the integration of natural and synthetic worlds, computers and communication applications, and the functionalities and flexibilities it offers. Initially, perhaps only the very basic functionalities will be useful. As the demand for sophisticated multimedia grows, the advanced functionalities may be useful. Up-to-date information on progress on various topics in MPEG-4 can be found by visiting MPEG related websites [13,21,32,34,38].

References

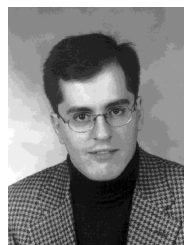
- [1] Ad Hoc Group on Error Resilience Core Experiments, Plan for March 1998 error resilience verification test, ISO/IEC JTC1/SC29/WG11 N1829, Stockholm (July 1997).
- [2] A.L. Ames, D.R. Nadeau and J.L. Moreland, *The VRML Sourcebook* (Wiley, New York, 1996).
- [3] AOE Group, MPEG-4 proposal package description (PPD) – Rev. 3, ISO/IEC JTC1/SC29/WG11 N0998, Tokyo (July 1995).
- [4] AOE Group, MPEG-4 testing and evaluation procedures document, ISO/IEC JTC1/SC29/WG11 N0999, Tokyo (July 1995).
- [5] O. Avaro, P. Chou, A. Eleftheriadis, C. Herpel and C. Reader, The MPEG-4 system and description languages, *Signal Processing: Image Communication* (Special issue on MPEG-4), to appear in 1997.
- [6] L. Chiariglione, MPEG-4 call for proposals, ISO/IEC JTC1/SC29/WG11 N0997, Tokyo (July 1995).
- [7] L. Chiariglione, Resolutions of 40th WG11 meeting, ISO/IEC JTC1/SC29/WG11 N1716, Stockholm (July 1997).
- [8] R. Cox, B. Haskell, Y. LeCun, B. Shahraray and L. Rabiner, On the applications of multimedia processing to communications, AT&T internal technical memo; to appear in *IEEE Transactions*.
- [9] T. Ebrahimi, Report of ad hoc group on definition of VMs for content based video representation, ISO/IEC JTC1/SC29/WG11 MPEG 96/0642, Munich (January 1996).
- [10] A. Eleftheriadis, The MPEG-4 system description language: From practice to theory, in: *Proceedings of 1997 IEEE International Conference on Circuits and Systems ISCAS '97*, Hong Kong (June 1997).
- [11] A. Eleftheriadis, Flavor: A language for media representation, *Proceedings, ACM Multimedia '97 Conference* (November 1997) (to appear).
- [12] Y. Fang and A. Eleftheriadis, A syntactic framework for bitstream-level representation of audio-visual objects, in: *Proceedings of 3rd IEEE International Conference on Image Processing ICIP '96*, Lausanne, Switzerland (September 1996).

- [13] Flavor Web Site: <http://www.ee.columbia.edu/flavor>.
- [14] J. Gosling, B. Joy and G. Steele, *The Java Language Specification* (Addison-Wesley, Reading, MA, 1996).
- [15] B.G. Haskell, A. Puri and A.N. Netravali, *Digital Video: An Introduction to MPEG-2* (Chapman and Hall, London, 1997).
- [16] B.G. Haskell, A. Puri and J. Osterman, Happenings in ISO MPEG: An introduction to MPEG-4, invited presentation at Data Compression Conference, Snow Bird (March 1997).
- [17] IEEE Transactions on Circuits and Systems for Video Technology (Special issue on MPEG-4) 7(1) (February 1997).
- [18] ISO/IEC 11172 International Standard (MPEG-1), Information technology – Coding of moving pictures and associated audio for digital storage media at up to about 1.5 Mbit/s (1993).
- [19] ISO/IEC 13818 International Standard (MPEG-2), Information technology – Generic coding of moving pictures and associated audio (also ITU-T Rec. H.262) (1995).
- [20] ISO/IEC 14472 draft international standard: Virtual reality modeling language (1997).
- [21] ISO/IEC JTC1/SC29/WG11 (MPEG) Web Site: <http://drogo.csel.it/mpeg>.
- [22] ITU-T Recommendation H.261, Video codec for audio-visual services at $p \times 64$ kbit/s (1990).
- [23] ITU-T, Draft ITU-T Recommendation H.263: Video coding for low bit-rate communication (December 1995).
- [24] ITU-T Recommendation H.223: Multiplexing protocol for low bit-rate multimedia communication (1995).
- [25] ITU-T Recommendation H.223 Annex A: Multiplexing protocol for low bit-rate mobile multimedia communication (1996).
- [26] ITU-T H.263+ Video Group, Draft 12 of ITU-T Recommendation H.263+ (May 1997).
- [27] D. Lindbergh, The H.324 multimedia communication standard, IEEE Communications Magazine 34(12) (December 1996) 46–51.
- [28] MPEG-4 Audio Group, MPEG-4 audio working draft version 4.0, ISO/IEC JTC1/SC29/WG11 N1745, Stockholm (July 1997).
- [29] MPEG-4 Integration Group, MPEG-4 SNHC call for proposals, ISO/IEC JTC1/SC29/WG11 N1195, Florence (March 1996).
- [30] MPEG-4 Integration Group, MPEG-4 SNHC proposal package description, ISO/IEC JTC1/SC29/WG11 N1199, Florence (March 1996).
- [31] MPEG-4 Requirements Ad Hoc Group, Draft of MPEG-4 requirements, ISO/IEC JTC1/SC29/WG11 N1238, Florence (March 1996).
- [32] MPEG-4 SNHC Web Site: <http://www.es.com/mpeg4-snhc>.
- [33] MPEG-4 Systems Group, MPEG-4 systems working draft version 5.0, ISO/IEC JTC1/SC29/WG11 N1825, Stockholm (July 1997).
- [34] MPEG-4 Systems Web Site: <http://garuda.imag.fr/MPEG4>.
- [35] MPEG-4 Video and SNHC Groups, MPEG-4 visual working draft version 4.0, ISO/IEC JTC1/SC29/WG11 N1797, Stockholm (July 1997).
- [36] MPEG-4 Video Group, MPEG-4 video verification model version 2.0, ISO/IEC JTC1/SC29/WG11 N1260, Florence (March 1996).
- [37] MPEG-4 Video Group, MPEG-4 video verification model version 7.0, ISO/IEC JTC1/SC29/WG11 N1642, Bristol (April 1997).
- [38] MPEG-4 Video Web Site: <http://wwwam.hhi.de/mpeg-video>.
- [39] H. Peterson, Report of the ad hoc group on MPEG-4 video testing logistics, ISO/IEC JTC1/SC29/WG11 Doc. MPEG95/0532, (November 1995).
- [40] A. Puri, Status and direction of the MPEG-4 standard, in: *International Symposium on Multimedia and Video Coding*, New York (October 1995); also published in a book by Plenum Press.
- [41] A. Puri, Report of ad hoc group on coordination of future core experiments in MPEG-4 video, ISO/IEC JTC1/SC29/WG11 MPEG 96/0669, Munich (January 1996).
- [42] A. Puri, MPEG-4: A flexible and extensible multimedia coding standard in progress, invited paper in IEEE Multimedia book (1996).
- [43] A. Puri, A.R. Reibman, R.L. Schmidt and B.G. Haskell, Robustness considerations in ISO MPEG-4 and ITU-T mobile video standards, in: *Proceedings MoMuC-3*, Princeton (September 1996) (Plenum Press, 1997).
- [44] A. Puri, R.L. Schmidt and B.G. Haskell, Improvements in DCT-based video coding, in: *Proc. SPIE Visual Communications and Image Processing*, San Jose (February 1997).
- [45] Requirements Group, MPEG-4 applications document, ISO/IEC JTC1/SC29/WG11 N1729, Stockholm (July 1997).
- [46] Requirements Group, MPEG-4 overview, ISO/IEC JTC1/SC29/WG11 N1730 (July 1997).
- [47] Requirements Group, MPEG-4 profile requirements version 4, ISO/IEC JTC1/SC29/WG11 N1728, Stockholm (July 1997).
- [48] Requirements Group, MPEG-4 requirements version 4, ISO/IEC JTC1/SC29/WG11 N1727, Stockholm (July 1997).
- [49] Requirements Group, MPEG-7 context and objectives version 4, ISO/IEC JTC1/SC29/WG11 N1733, Stockholm (July 1997).
- [50] Signal Processing: Image Communication (Special issue on MPEG-4, Part 1: Invited papers) 10(1–3) (May 1997).
- [51] Signal Processing: Image Communication (Special issue on MPEG-4, Part 2: Submitted papers) 10(4) (July 1997).
- [52] P.J.L. van Beek, A.M. Tekalp and A. Puri, 2D mesh geometry and motion compression for efficient object-based video compression, *IEEE Int. Conf. on Image Processing* (October 1997), to appear.



Atul Puri received his B.S. in electrical engineering from India in 1980, his M.S. in electrical engineering from City College of New York in 1982, and his Ph.D., also in electrical engineering, from the City University of New York in 1988. While working on his dissertation, he was a consultant in Visual Communications Research Department of Bell Labs and gained experience in developing algorithms, software and hardware for video communications. In 1988 he joined the same department at Bell Labs as a Member of Technical Staff. Since 1996, Dr. Puri has been a Principal Member of Technical Staff in the Image Processing Research Department of AT&T Labs and is presently located at Red Bank, NJ. Dr. Puri has represented AT&T at the Moving Pictures Expert Group standards of the International Standards Organization for the past 9 years and has actively contributed towards development of the MPEG-1, the MPEG-2 and the MPEG-4 audio-visual coding standards. Currently he is participating in video and systems parts of the MPEG-4 standard and is also a technical editor of the standard. He has been involved in research in video coding algorithms for a number of diverse applications such as videoconferencing, video on digital storage media, HDTV, 3D-TV and wireless video. His current interests are in the area of multimedia services and systems for Web/Internet. Dr. Puri holds over 14 patents and has applied for another 8 patents. He has also published over 30 technical papers in conferences and journals, including several invited papers. He is the co-author of a book entitled "Digital Video: An Introduction to MPEG-2". He is currently coediting another book. He has been the recipient of exceptional contribution and individual performance merit awards of AT&T. Furthermore, he has also received awards from the Communication Services business unit, and the AT&T Technical Journal. He has taught graduate courses on Digital Image and Video Coding at Columbia University. He is also an active member of IEEE, its Communications, and Signal Processing societies.

E-mail: apuri@research.att.com



Alexandros Eleftheriadis was born in Athens, Greece, in 1967. He received the Diploma in electrical engineering and computer science from the National Technical University of Athens, Greece, in 1990, and the M.S., M.Ph. and Ph.D. degrees in electrical engineering from Columbia University, New York, in 1992, 1994 and 1995, respectively. Since 1995 he has been an Assistant Professor in the Department of Electrical Engineering at Columbia University, where he is leading a re-

search team working in the areas of visual information representation and compression, video communication systems (including video-on-demand and Internet video), distributed multimedia systems, and the fundamentals of compression. During the summers of 1993 and 1994, he was with AT&T Bell Laboratories, Murray Hill, NJ, developing low bit rate model-assisted video coding techniques for videoconferencing applications. From 1990 until 1995, he was a Graduate Research Assistant in the Depart-

ment of Electrical Engineering at Columbia University. Dr. Eleftheriadis is a member of the ANSI NCITS L3.1 Committee and the ISO/IEC JTC1/SC29/WG11 (MPEG) group that develop national and international standards for audio-visual content representation and distribution. He is a member of IEEE, ACM, and the Technical Chamber of Greece.
E-mail: eleft@ctr.columbia.edu