# Video Processing & Communications

# Video Coding Standards: Part I

Yao Wang
Polytechnic University, Brooklyn, NY11201
http://eeweb.poly.edu/~yao

Based on:  Y. Wang, J. Ostermann, and Y.-Q. Zhang, Video Processing and Communications, Prentice Hall, 2002.

# **Outline**

- Overview of Standards and Their Applications
- What do video coding standards define?
- ITU-T Standards for Audio-Visual Communications
  - H.261
  - H.263
  - H.263+, H.263++
- ISO Standards for
  - MPEG-1
  - MPEG-2
  - MPEG-4

# Timeline and "needs" of the standards

- H.261 (1990): Video conferencing
- MPEG-1 (1992): Non-interactive applications, VCD
- MPEG-2 and H.262 (1996): TV broadcast, DVD
- H.263 (Nov. 1995; Sept. 1997, Nov. 2000): Video conferencing
- MPEG-4 video (part 2) (1999): object-oriented coding
- H.264 and MPEG-4 part 10 (AVC) (2003): compression efficiency
- AVS (Audio and Visual Coding Standard) (2006): avoiding high licensing fees
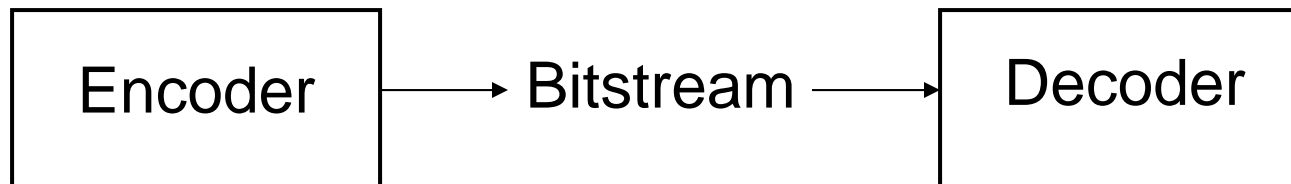- HEVC (will be H.265) (2013?): compression efficiency

From Amy Reibman

# Motivation for Standards

- ## Goal of standards:

  - *Ensuring interoperability:* Enabling communication between devices made by different manufacturers

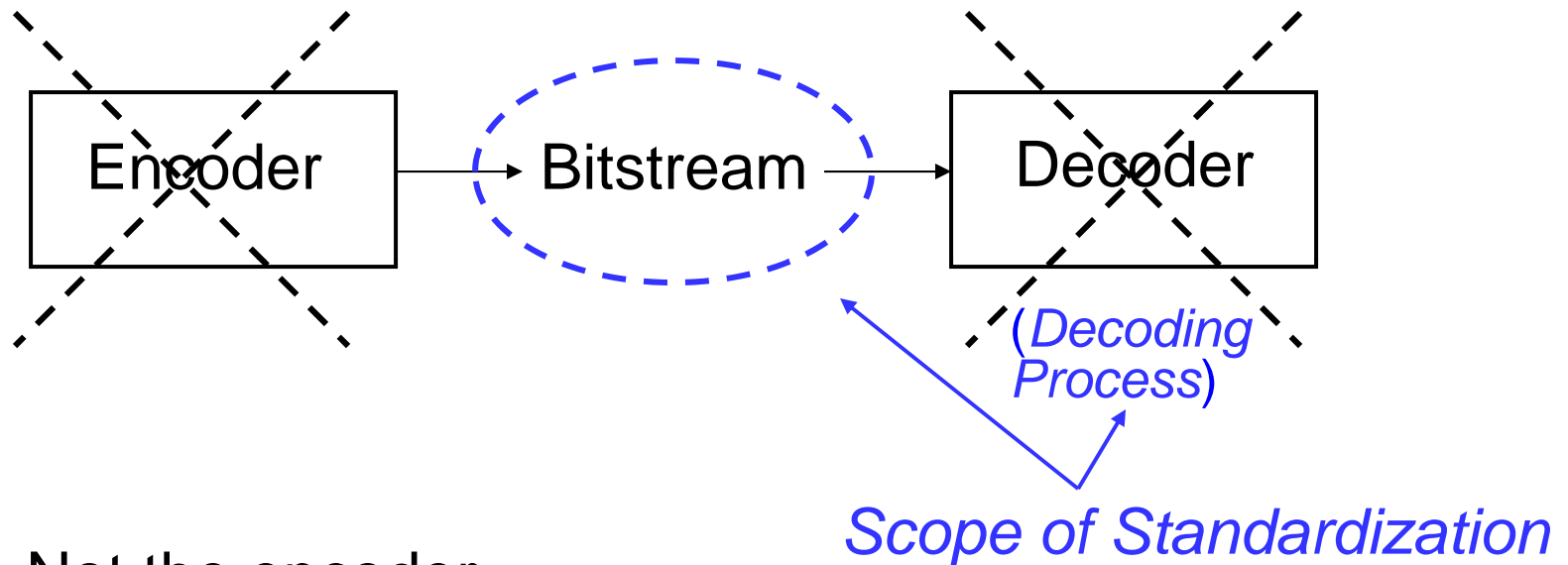  - Promoting a technology or industry

  - Reducing costs

From John Apostolopoulos

# What do the Standards Specify?

Encoder → Bitstream → Decoder

# What do the Standards Specify?



- Not the encoder
- Not the decoder
- Just the *bitstream syntax* and the *decoding process* (e.g., use IDCT, but not how to implement the IDCT)
  → Enables improved encoding & decoding strategies to be employed in a standard-compatible manner

From John Apostolopoulos

# Video coding standards

- Video coding standards define the operation of a decoder given a correct bitstream
- They do NOT describe an encoder

- Video coding standards typically define a toolkit
- Not all pieces of the toolkit need to be implemented to create a conforming bitstream

- Decoders must implement some subset of the toolkit to be declared "conforming"

# Current Image and Video Compression Standards

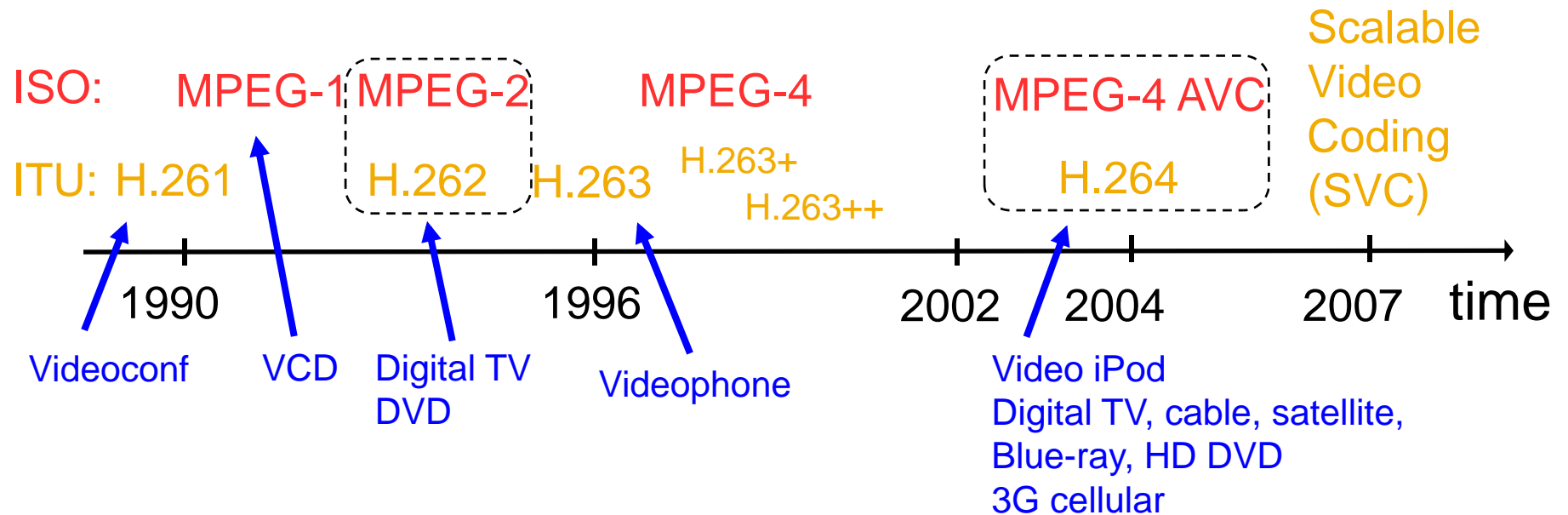| Standard | Application | Bit Rate |
|---|---|---|
| JPEG | Continuous-tone still-image compression | Variable |
| H.261 | Video telephony and teleconferencing over ISDN | p x 64 kb/s |
| MPEG-1 | Video on digital storage media (CD-ROM) | 1.5 Mb/s |
| MPEG-2 | Digital Television | 2-20 Mb/s |
| H.263 | Video telephony over PSTN | 33.6-? kb/s |
| MPEG-4 | Object-based coding, synthetic content, interactivity | Variable |
| JPEG-2000 | Improved still image compression | Variable |
| H.264 / MPEG-4 AVC | Improved video compression | 10's kb/s to Mb/s |

MPEG and JPEG: International Standards Organization (ISO)

H.26x family: International Telecommunications Union (ITU)

# History of Video Coding Standards

Scalable Video Coding (SVC)

ISO:   MPEG-1   MPEG-2        MPEG-4              MPEG-4 AVC

ITU: H.261   H.262   H.263   H.263+              H.264
                                    H.263++

1990              1996           2002  2004    2007   time

Videoconf   VCD   Digital TV      Video iPod
                  DVD             Digital TV, cable, satellite,
            Videophone            Blue-ray, HD DVD
                                  3G cellular

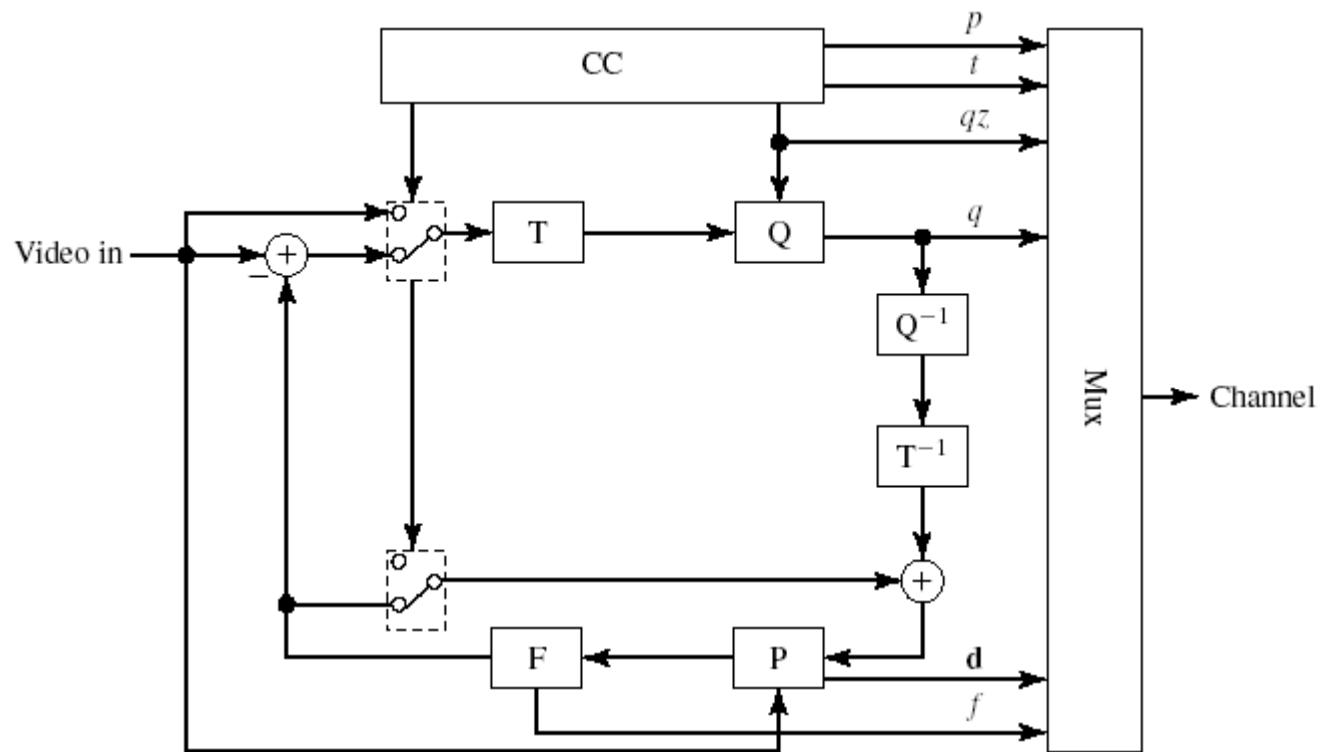- Recent development:
  – HEVC, 2012

# H.261 Video Coding Standard

- For video-conferencing/video phone
  - Video coding standard in H.320
  - Low delay (real-time, interactive)
  - Slow motion in general
- For transmission over ISDN
  - Fixed bandwidth: px64 Kbps, p=1,2,…,30
- Video Format:
  - CIF (352x288, above 128 Kbps)
  - QCIF (176x144, 64-128 Kbps)
  - 4:2:0 color format, progressive scan
- Published in 1990
- Each macroblock can be coded in intra- or inter-mode
- Periodic insertion of intra-mode to eliminate error propagation due to network impairments
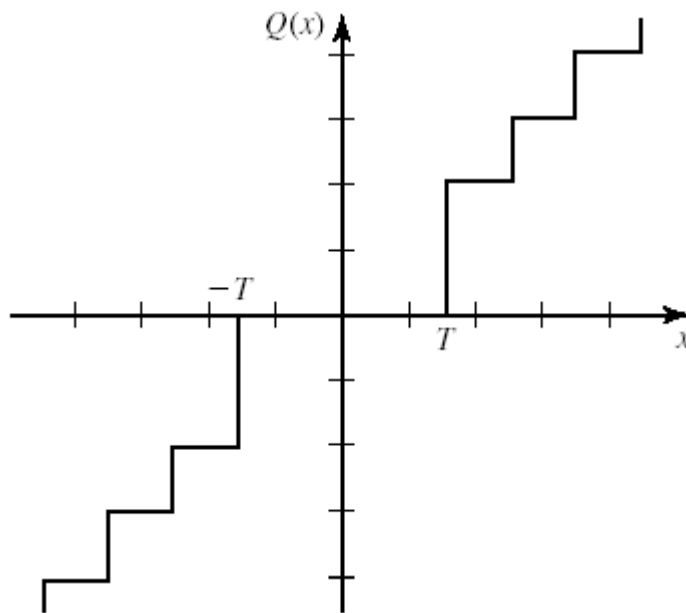- Integer-pel accuracy motion estimation in inter-mode

# H.261 Encoder



F: Loop filter;  P: motion estimation and compensation

# DCT Coefficient Quantization



**DC Coefficient in Intra-mode:**
Uniform, stepsize=8

**Others:**
Uniform with deadzone,
stepsize=2~64 (MQUANT)

**Deadzone:**
To avoid too many small
coefficients being coded, which
are typically due to noise

# Motion Estimation and Compensation

- Integer-pel accuracy in the range [-16,16]
- Methods for generating the MVs are not specified in the standard
  - Standards only define the bitstream syntax, or the decoder operation)
- MVs coded differentially (DMV)
- Encoder and decoder uses the decoded MVs to perform motion compensation
- Loop-filtering can be applied to suppress propagation of coding noise temporally
  - Separable filter  [1/4,1/2,1/4]
  - Loop filter can be turned on or off

# Variable Length Coding

- DCT coefficients are converted into runlength representations and then coded using VLC (Huffman coding for each pair of symbols)
    - Symbol: (Zero run-length, non-zero value range)
- Other information are also coded using VLC (Huffman coding)

# **Parameter Selection and Rate Control**

- MTYPE (intra vs. inter, zero vs. non-zero MV in inter)
- CBP (which blocks in a MB have non-zero DCT coefficients)
- MQUANT (allow the changes of the quantizer stepsize at the MB level)
  - should be varied to satisfy the rate constraint
- MV (ideally should be determined not only by prediction error but also the total bits used for coding MV and DCT coefficients of prediction error)
- Loop Filter on/off

# H.263 Video Coding Standard

- H.263 is the video coding standard in H.323/H.324, targeted for visual telephone over PSTN or Internet

- Developed later than H.261, can accommodate computationally more intensive options
  - Initial version (H.263 baseline): 1995
  - H.263+: 1997
  - H.263++: 2000

- Goal: Improved quality at lower rates

- Result: Significantly better quality at lower rates
  - Better video at 18-24 Kbps than H.261 at 64 Kbps
  - Enable video phone over regular phone lines (28.8 Kbps) or wireless modem
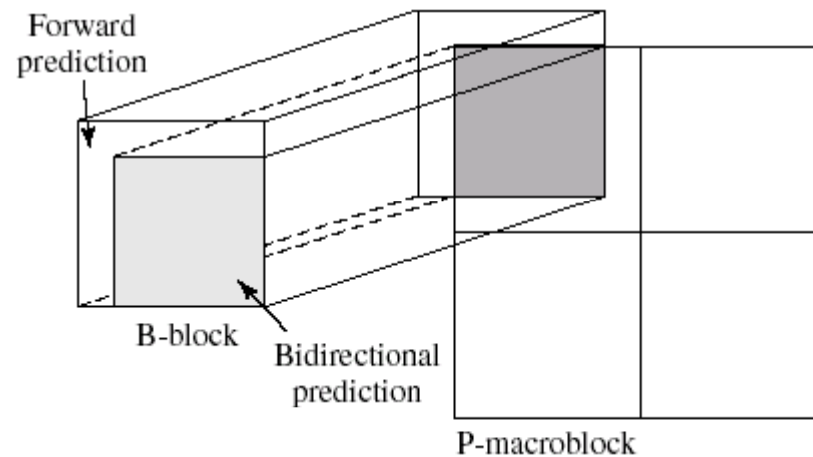
# Improvements over H.261

- Better motion estimation
  - half-pel accuracy motion estimation with bilinear interpolation filter
  - Larger motion search range [-31.5,31], and unrestricted MV at boundary blocks
  - More efficient predictive coding for MVs (median prediction using three neighbors)
  - overlapping block motion compensation (option)
  - variable block size: 16x16 -> 8x8, 4 MVs per MB (option)
  - use bidirectional temporal prediction (PB picture) (option)
- 3-D VLC for DCT coefficients
  - (runlength, value, EOB)
- Syntax-based arithmetic coding (option)
  - 4% savings in bit rate for P-mode, 10% saving for I-mode, at 50% more computations
- The options, when chosen properly, can improve the PSNR 0.5-1.5 dB over default at 20-70 kbps range.

# PB-Picture Mode



Forward prediction
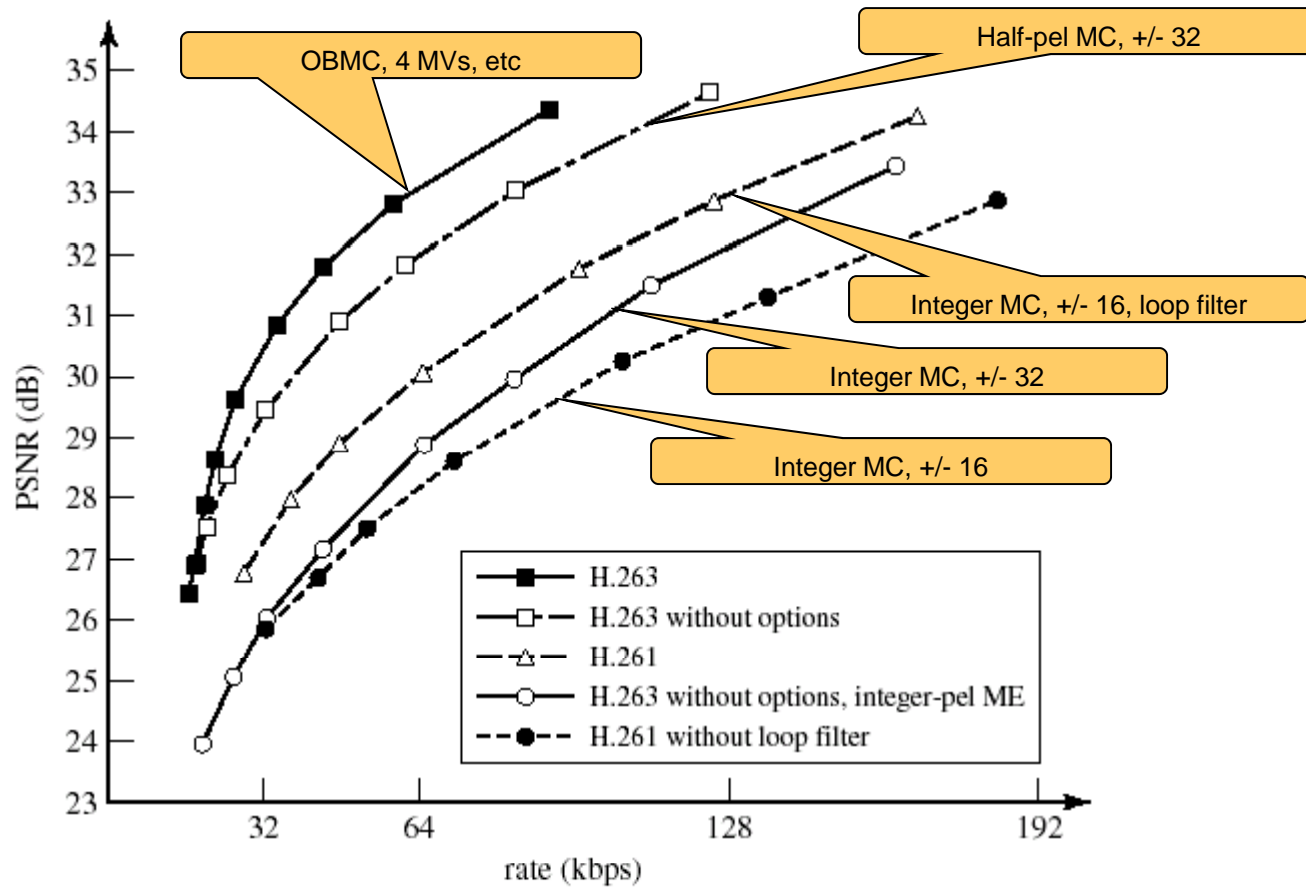
B-block — Bidirectional prediction

P-macroblock

PB-picture mode codes two pictures as a group. The second picture (P) is coded first, then the first picture (B) is coded using both the P-picture and the previously coded picture. This is to avoid the reordering of pictures required in the normal B-mode. But it still requires additional coding delay than P-frames only.

In a B-block, forward prediction (predicted from the previous frame) can be used for all pixels; backward prediction (from the future frame) is only used for those pels that the backward motion vector aligns with pels of the current MB. Pixels in the "white area" use only forward prediction.

An improved PB-frame mode was defined in H.263+, that removes the previous restriction.

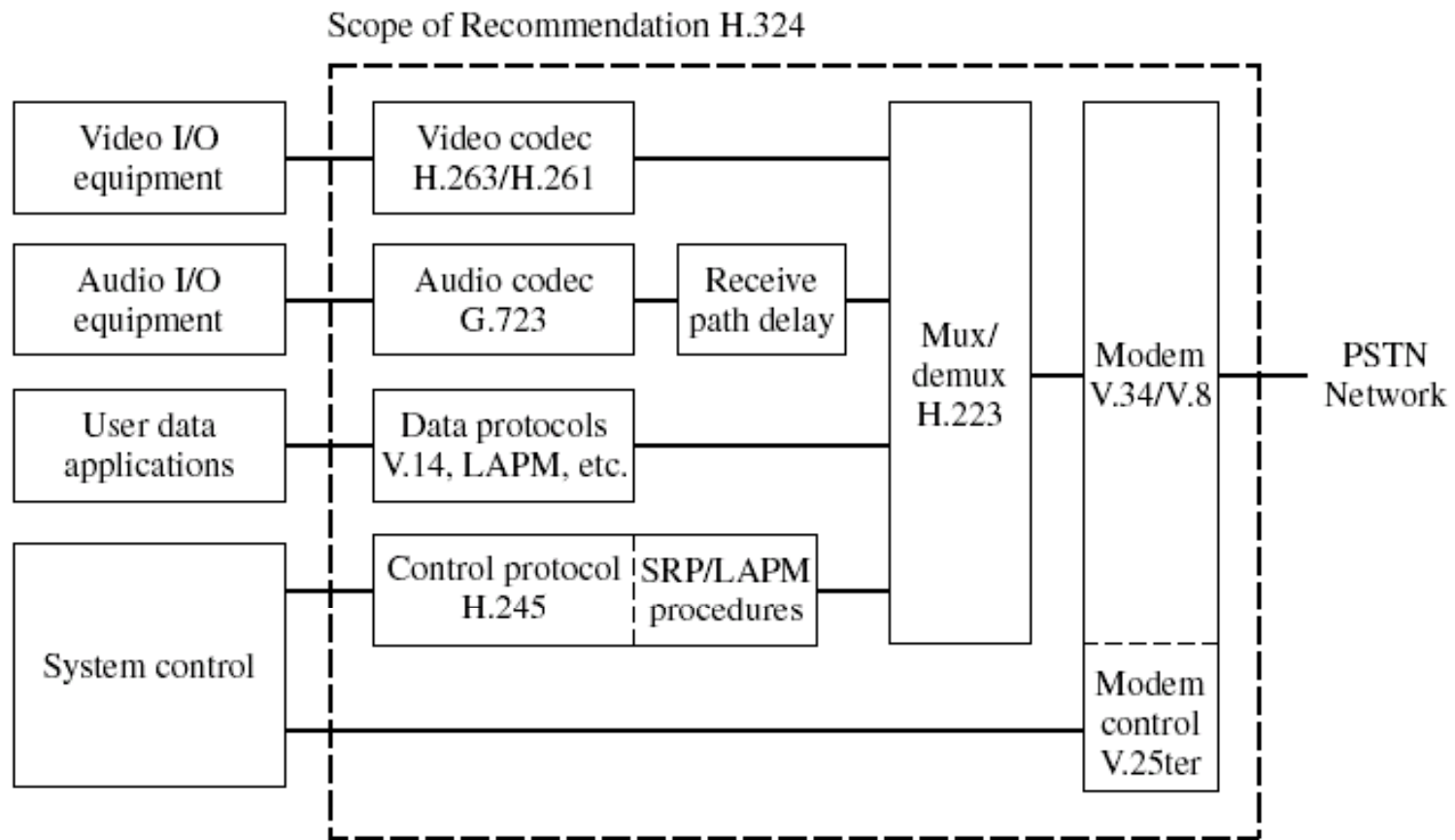# Performance of H.261 and H.263



Forman, QCIF, 12.5 Hz

# ITU-T Multimedia Communications Standards

| Network | System | Video | Audio | Mux | Control |
|---------|--------|-------|-------|-----|---------|
| PSTN | H.324 | H.261/3 | G.723.1 | H.223 | H.245 |
| N-ISDN | H.320 | H.261 | G.7xx | H.221 | H.242 |
| B-ISDN/ATM | H.321 | H.261 | G.7xx | H.221 | Q.2931 |
| | H.310 | H.261/2 | G.7xx/MPEG | H.222.0/1 | H.245 |
| QoS LAN | H.322 | H.261/3 | G.7xx | H.221 | H.242 |
| Non-QoS LAN | H.323 | H.261/3 | G.7xx | H.225.0 | H.245 |

# H.324 Terminal
## (multimedia communication over PSTN)



Scope of Recommendation H.324

| Video I/O equipment | Video codec H.263/H.261 | | Mux/ demux H.223 | Modem V.34/V.8 | PSTN Network |
| Audio I/O equipment | Audio codec G.723 | Receive path delay | | | |
| User data applications | Data protocols V.14, LAPM, etc. | | | | |
| System control | Control protocol H.245 | SRP/LAPM procedures | | Modem control V.25ter | |

# MPEG-1 Overview

- Audio/video on CD-ROM  (1.5 Mbps, SIF: 352x240).
  - Maximum: 1.856 mbps, 768x576 pels
- Start late 1988, test in 10/89, Committee Draft 9/90
- ISO/IEC 11172-1~5 (Systems, video, audio, compliance, software).
- Prompted explosion of digital video applications: MPEG1 video CD and downloadable video over Internet
- Software only decoding, made possible by the introduction of Pentium chips, key to the success in the commercial market
- MPEG-1 Audio
  - Offers 3 coding options (3 layers), higher layer have higher coding efficiency with more computations
  - MP3 = MPEG1 layer 3 audio

# MPEG-1 video vs H.261

- Developed at about the same time
- Must enable random access (Fast forward/rewind)
  - Using GOP structure with periodic I-picture and P-picture
- Not for interactive applications
  - Do not have as stringent delay requirement
- Fixed rate (1.5 Mbps), good quality (VHS equivalent)
  - SIF video format (similar to CIF)
    - CIF: 352x288, SIF: 352x240
  - Using more advanced motion compensation
    - Half-pel accuracy motion estimation, range up to +/- 64
  - Using bi-directional temporal prediction
    - Important for handling uncovered regions
  - Using perceptual-based quantization matrix for I-blocks (same as JPEG)
    - DC coefficients coded predictively

# Group of Picture Structure in MPEG



1 GOP

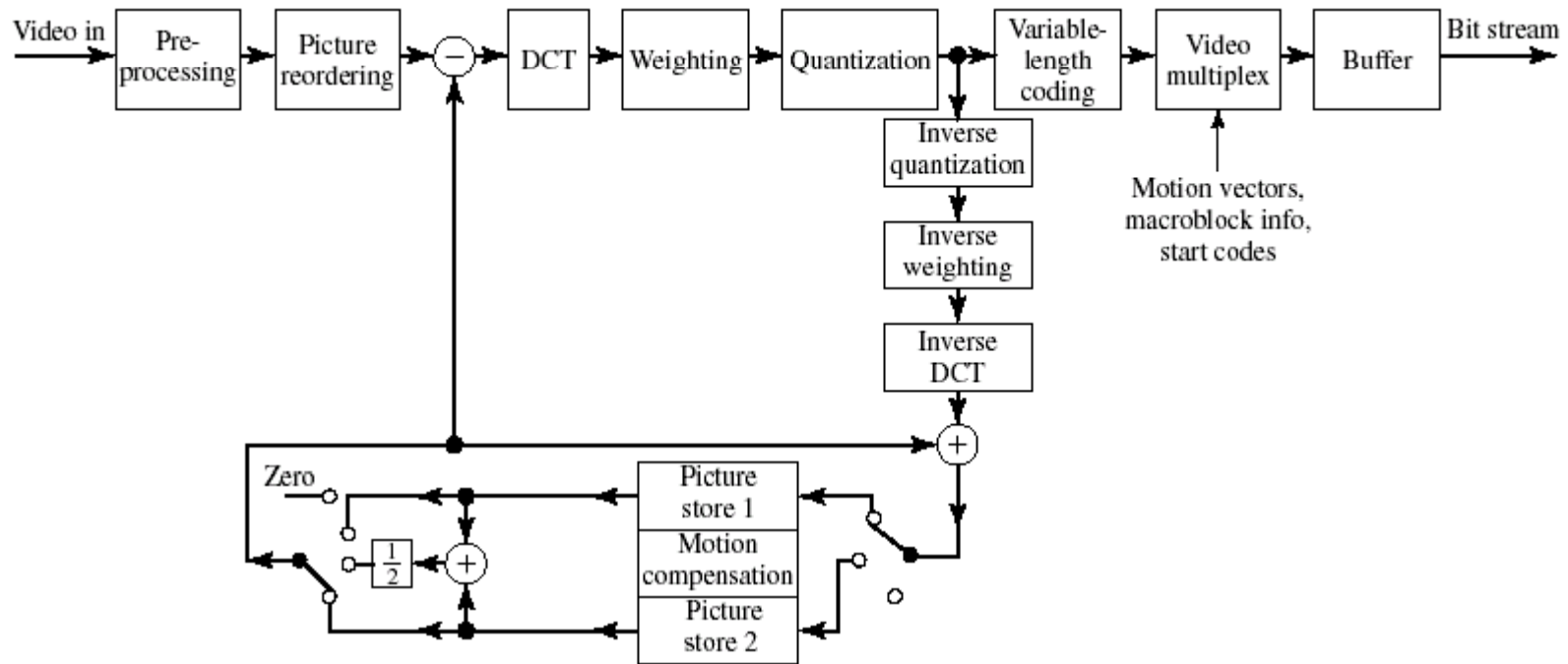| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|
| I | B | B | P | B | B | B | I |

Encoding order:   4   2   3   8   5   6   7

# MPEG-1 Video Encoder

# MPEG2 Overview

- A/V broadcast (TV, HDTV, Terrestrial, Cable, Satellite, High Speed Inter/Intranet) as well as DVD video
- 4~8 Mbps for TV quality, 10-15 for better quality at SDTV resolutions (BT.601)
- 18-45 Mbps for HDTV applications
  - MPEG-2 video high profile at high level is the video coding standard used in HDTV
- Test in 11/91, Committee Draft 11/93
- ISO/IEC 13818-1~6 (Systems, video, audio, compliance, software, DSM-CC)
- Consist of various profiles and levels
- Backward compatible with MPEG1
- MPEG-2 Audio
  - Support 5.1 channel
  - MPEG2 AAC: requires 30% fewer bits than MPEG1 layer 3

# MPEG2 vs. MPEG1 Video

- MPEG1 only handles progressive sequences (SIF).
- MPEG2 is targeted primarily at interlaced sequences and at higher resolution (BT.601 = 4CIF).
- More sophisticated motion estimation methods (*frame/field prediction mode*) are developed to improve estimation accuracy for interlaced sequences.
- Different DCT modes and scanning methods are developed for interlaced sequences.
- MPEG2 has various scalability modes.
- MPEG2 has various profiles and levels, each combination targeted for different application

# Frame vs. Field Picture



Frame pictures      Field pictures

Scan lines viewed edgewise

| Fields | T1 & B1 | T2 & B2 | Field | T1 | B1 | T2 | B2 |
| Frame | 1 | 2 | Frame | 1 | | 2 | |

Two fields are combined into one frame before being coded
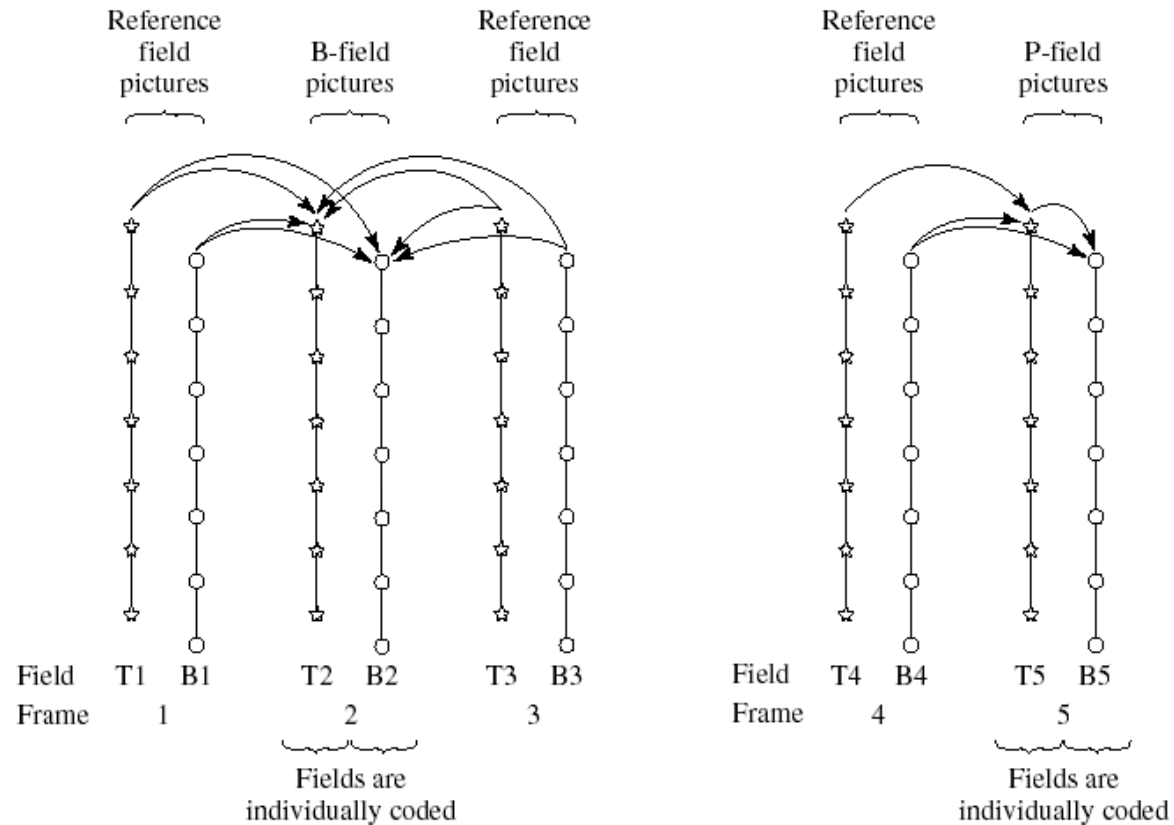
Fields are coded individually

# Motion Compensation for Interlaced Video

- Field prediction for field pictures
- Field prediction for frame pictures
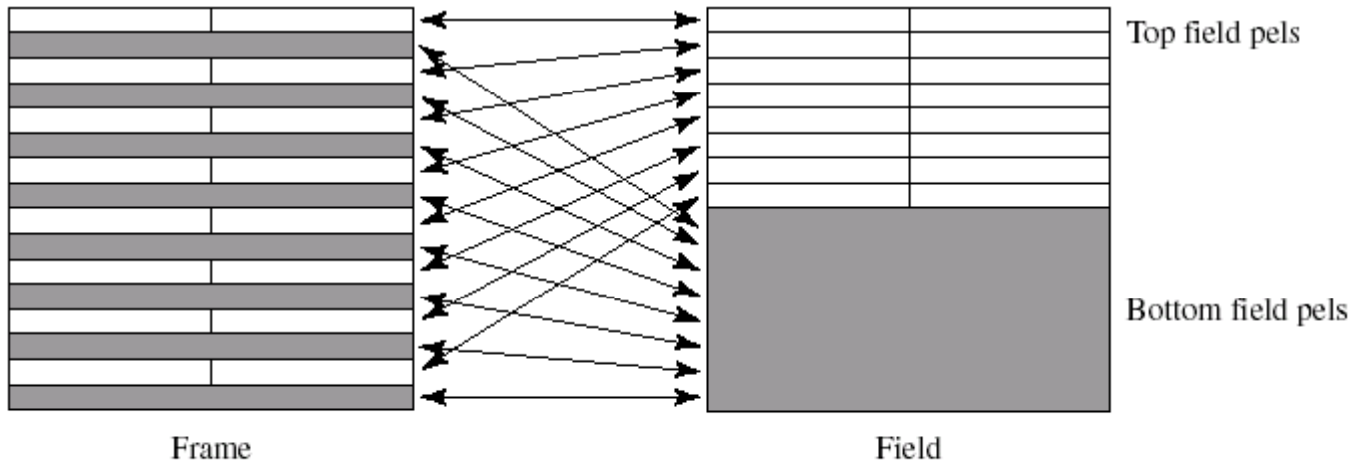- Dual prime for P-pictures
- 16x8 MC for field pictures

# Field prediction for field pictures

- Each field is predicted individually from the reference fields
    - A P-field is predicted from one previous field
    - A B-field is predicted from two fields chosen from two reference pictures
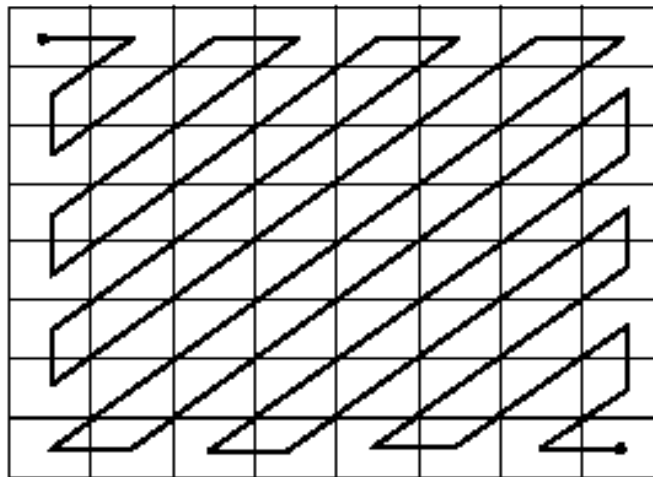
**Figure 13.18** Field prediction for frame pictures: the MB to be predicted is split into top field pels and bottom field pels. Each $16 \times 8$ field block is predicted separately with its own motion vector (P-frame) or two motion vectors (B-frame).
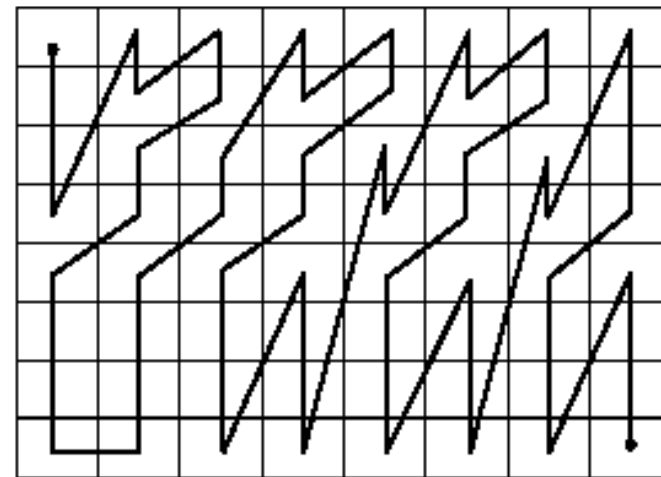
# DCT Modes

**Two types of DCT and two types of scan pattern:**
- **Frame DCT**: divides an MB into 4 blocks for Lum, as usual
- **Field DCT**: reorder pixels in an MB into top and bottom fields.
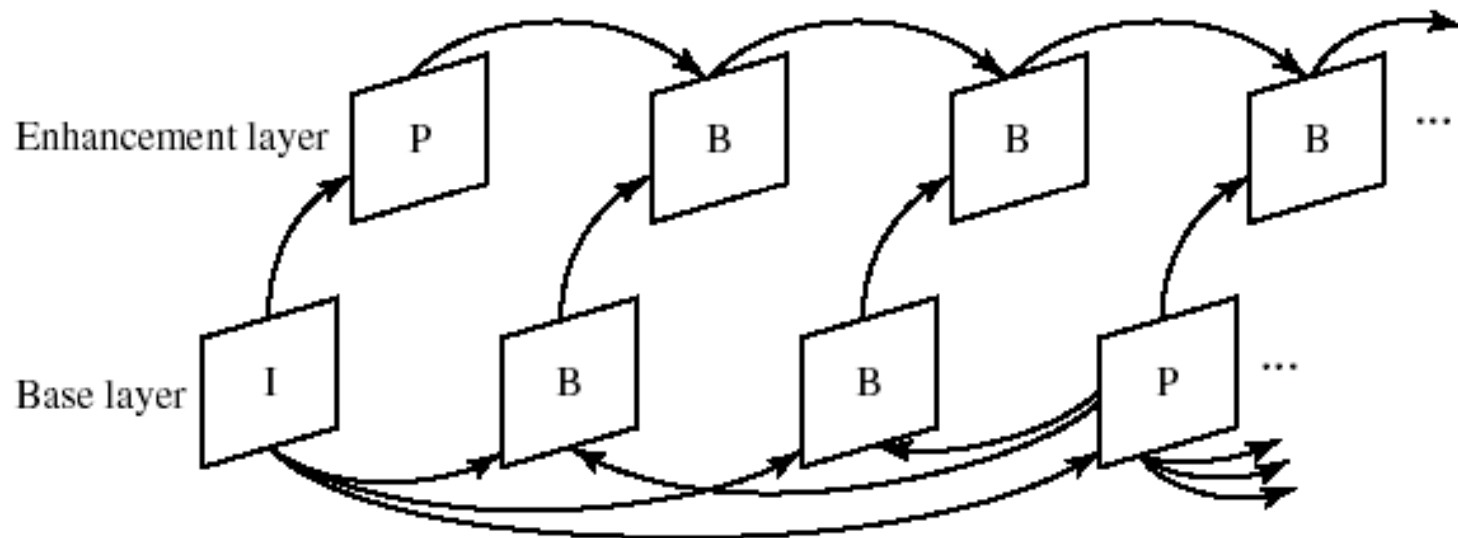


Zigzag scan                     Alternate scan

**Figure 13.19** The zigzag scan as known from H.261, H.263, and MPEG-1 is augmented by the alternate scan in MPEG-2, in order to code interlaced blocks that have more correlation in the horizontal than in the vertical direction.

# Temporal Scalability: Option 2



**Figure 13.24** A configuration in which temporal scalability enhancement layer may use the base layer and the enhancement layer for prediction. This arrangement is especially useful for coding of stereoscopic video.

# Profiles and Levels in MPEG-2

| | | | Simple (I, P) (4:2:0) | Main (I, P, B) (4:2:0) | SNR (I, P, B) (4:2:0) | Spatial (I, P, B) (4:2:0) | High (I, P, B) (4:2:0; 4:2:2) | Multiview (I, P, B) (4:2:0) | 4:2:2 (I, P, B) (4:2:0; 4:2:2) |
|---|---|---|---|---|---|---|---|---|---|
| | Low | Pels/line | | | 352 | 352 | | 352 | |
| | | Lines/frame | | | 288 | 288 | | 288 | |
| | | fps | | | 30 | 30 | | 30 | |
| | | mbps | | | 4 | 4 | | 8 | |
| | Main | Pels/line | 720 | 720 | 720 | | 720 | 720 | 720 |
| | | Lines/frame | 576 | 576 | 576 | | 576 | 576 | 512/608 |
| | | fps | 30 | 30 | 30 | | 30 | 30 | 30 |
| | | mbps | 15 | 15 | 15 | | 20 | 25 | 50 |
| Level | High-1440 | Pels/line | 1440 | | 1440 | 1440 | 1440 | | |
| | | Lines/frame | 1152 | | 1152 | 1152 | 1152 | | |
| | | fps | 60 | | 60 | 60 | 60 | | |
| | | mbps | 60 | | 60 | 80 | 100 | | |
| | High | Pels/line | 1920 | | | 1920 | 1920 | 1920 | |
| | | Lines/frame | 1152 | | | 1152 | 1152 | 1152 | |
| | | fps | 60 | | | 60 | 60 | 60 | |
| | | mbps | 80 | | | 100 | 130 | 300 | |

I, P, B: allowable picture types. Maximum bit rates include all layers in case of scalable bit streams.

**Profiles:** tools

**Levels:** parameter range for a given profile

**Main profile at main level** (mp@ml) is the most popular, used for digital TV

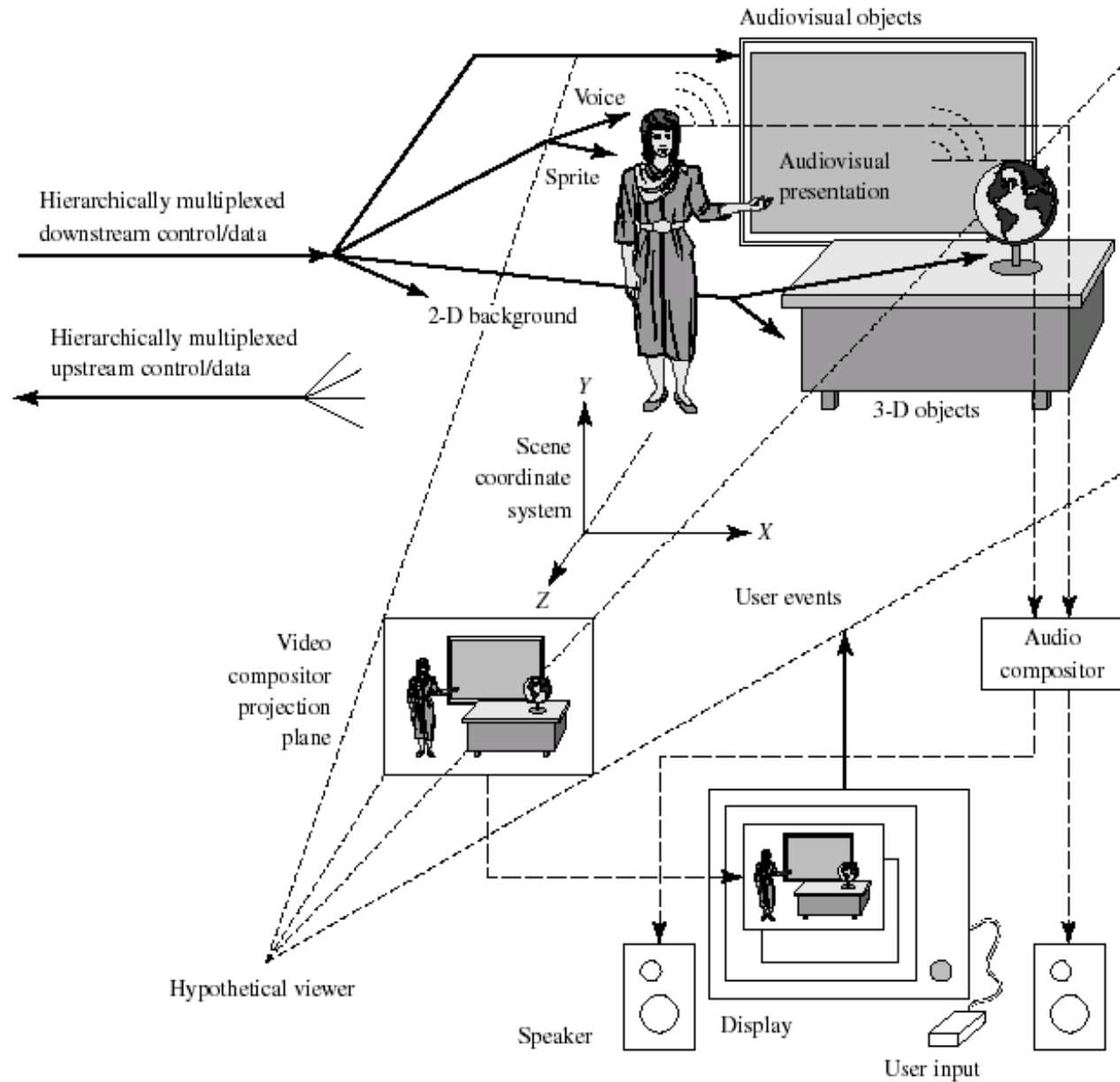**Main profile at high level (mp@hl):** HDTV

**4:2:2 at main level** (4:2:2@ml) is used for studio production

# MPEG-4 Overview

- ## Functionalities beyond MPEG-1/2

  - Interaction with individual objects
    - The displayed scene can be composed by the receiver from coded objects
  - Scalability of contents
  - Error resilience
  - Coding of both natural and synthetic audio and video

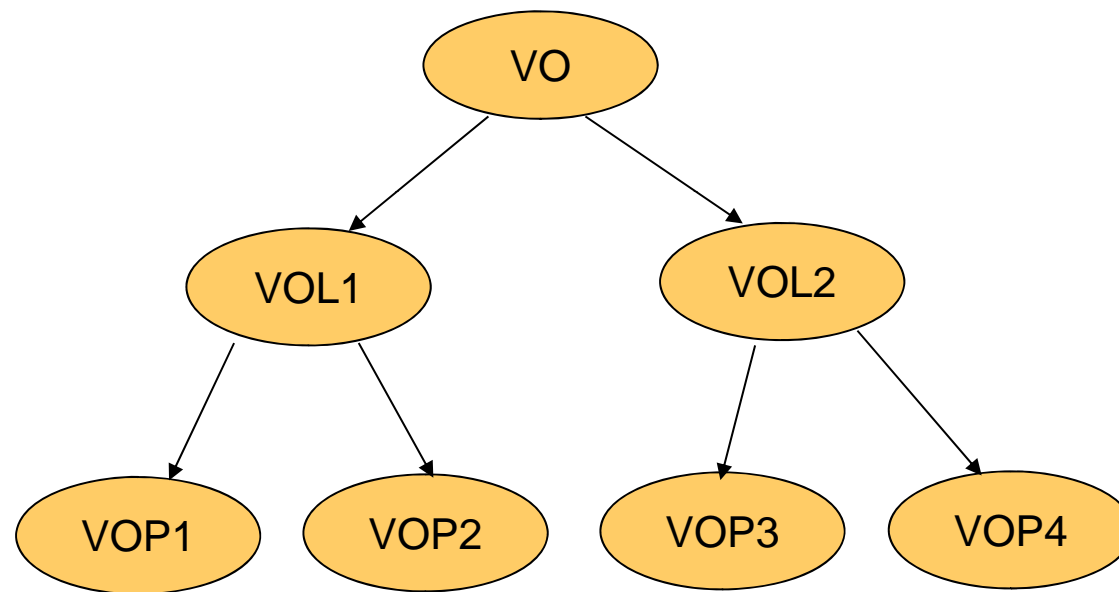The displayed scene is composed by the receiver based on desired view angle and objects of interests

# Object-Based Coding

- Entire scene is decomposed into multiple objects
  - Object segmentation is the most difficult task!
  - But this does not need to be standardized ☺

- Each object is specified by its shape, motion, and texture (color)
  - Shape and texture both changes in time (specified by motion)

- MPEG-4 assumes the encoder has a segmentation map available, specifies how to code (actually decode!) shape, motion and texture

# Object Description Hierarchy in MPEG-4
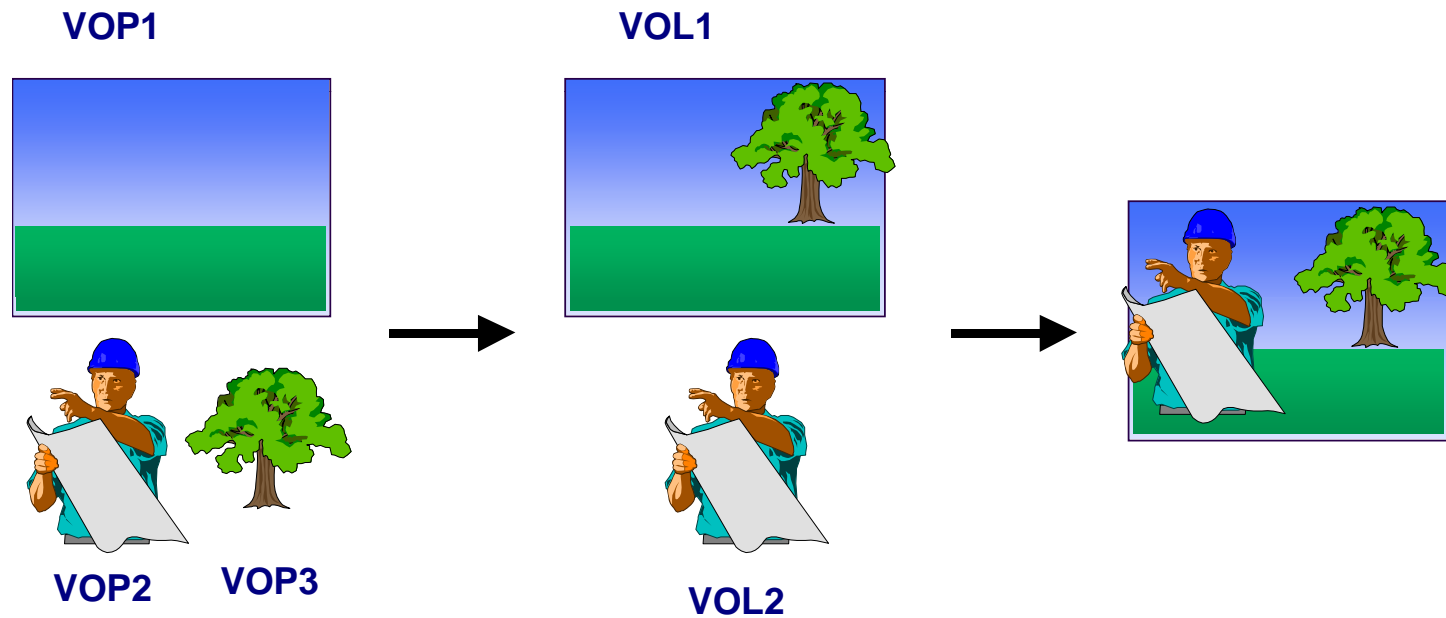


VO: video object
VOL: video object layer
      (can be different parts of a VO or different rate/resolution representation of a VOL)
VOP: video object plane

# Example of Scene Composition

VOP1

VOL1

VOP2    VOP3

VOL2

The decoder can compose a scene by including different VOPs in a VOL

# Object-Based Coding Basics
## (Chap 10)

- Entire scene is decomposed into multiple objects

  – Object segmentation is the most difficult task!

  – But this does not need to be standardized ☺

- Each object is specified by its shape, motion, and texture (color)

  – Shape and texture both changes in time (specified by motion)

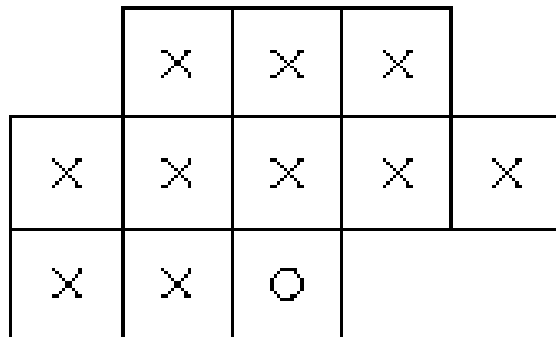# Generic Block Diagram for Object-Based Coding

# Shape Coding Methods

- Shape is specified by alpha maps
  - Binary alpha map: specifies whether a pel belongs to an object
  - Gray scale alpha map: a pel belong to the object can have a transparency value in the range (0-255)
- Bitmap coding
  - Run-length coding
  - Pel-wise coding using context-based arithmetic coding
  - Quadtree coding
- Contour coding
  - Chain coding
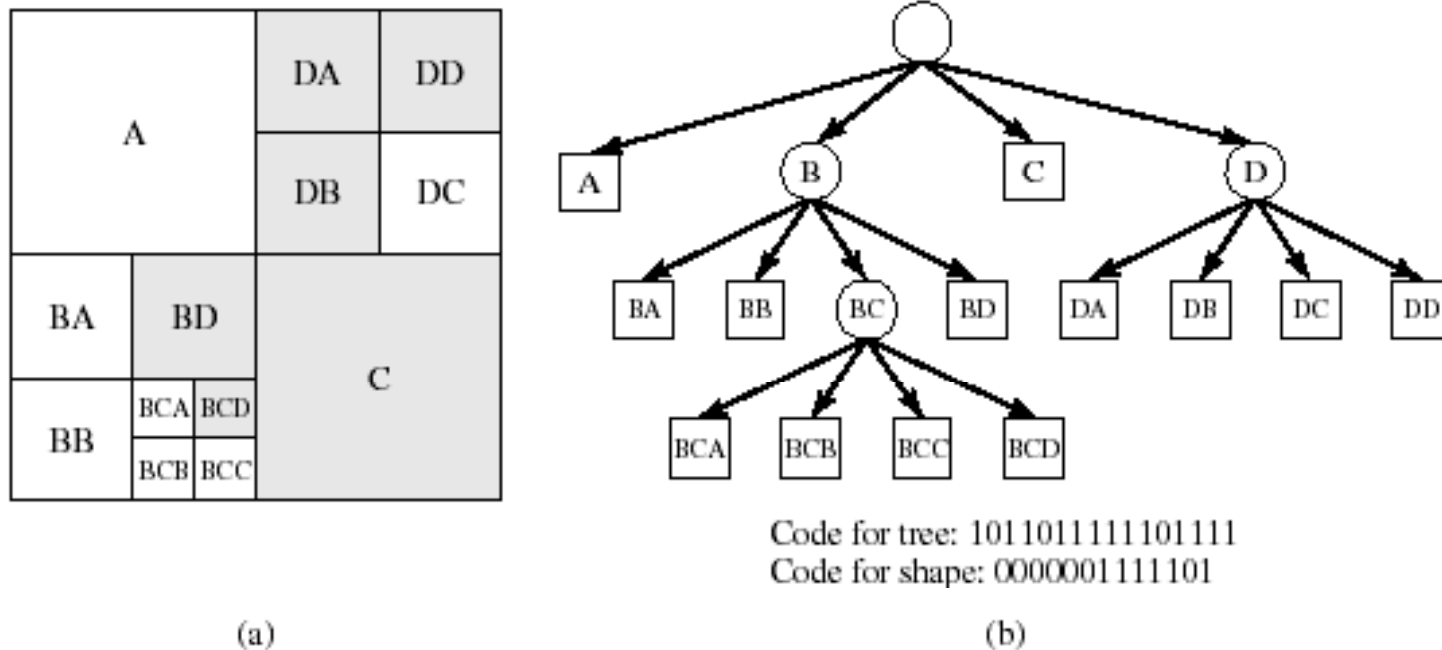  - Fourier descriptors
  - Polygon approximation

# Context-based Arithmetic Coding



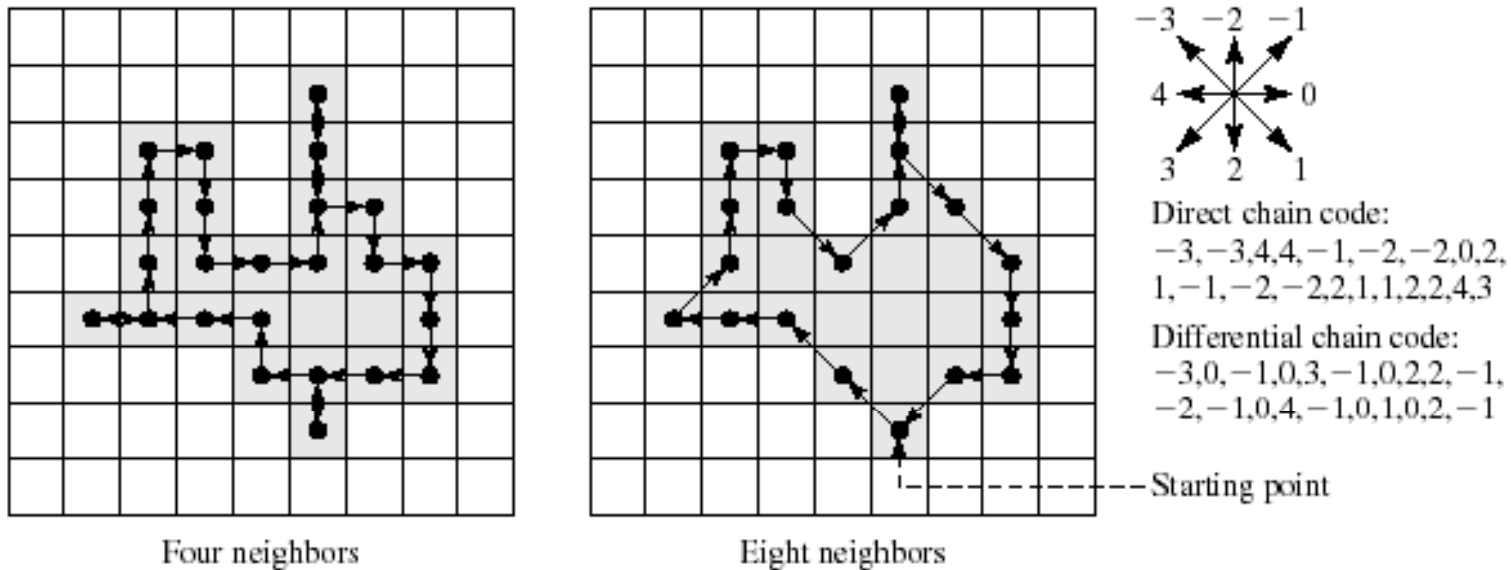**Figure 10.2** Template for defining the context of the pel to be coded (designated by "O").

# Quadtree Shape Coding



Code for tree: 101 1011 111 101 111
Code for shape: 0000001 111 101

(a)                                                    (b)

**Figure 10.3**   (a) An image is divided into a quad-tree that describes the object shape. (b) Related quad-tree—with code that describes the tree and labels for each square, such that the shape can be recovered.

−3  −2  −1

4  ←  →  0

3   2   1

Direct chain code:
−3,−3,4,4,−1,−2,−2,0,2,
1,−1,−2,−2,2,1,1,2,2,4,3

Differential chain code:
−3,0,−1,0,3,−1,0,2,2,−1,
−2,−1,0,4,−1,0,1,0,2,−1

Starting point
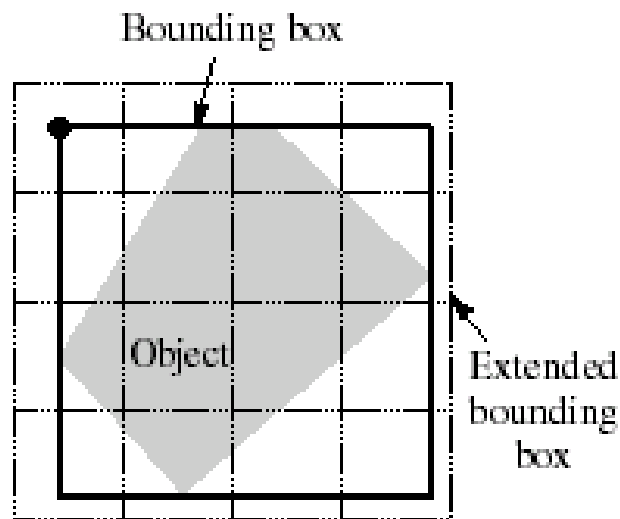
Four neighbors

Eight neighbors

**Figure 10.4** Chain code for pels with four and eight neighbors. We give examples for a direct chain code and for a differential code of the eight-connected chain. The first symbols of the two codes are identical and define the starting direction. The following symbols of the differential chain code are created by aligning direction 0 of the direction star with the direction of the last coded symbol.

# Coding of Texture with Arbitrary Shape

- Texture extrapolation through padding
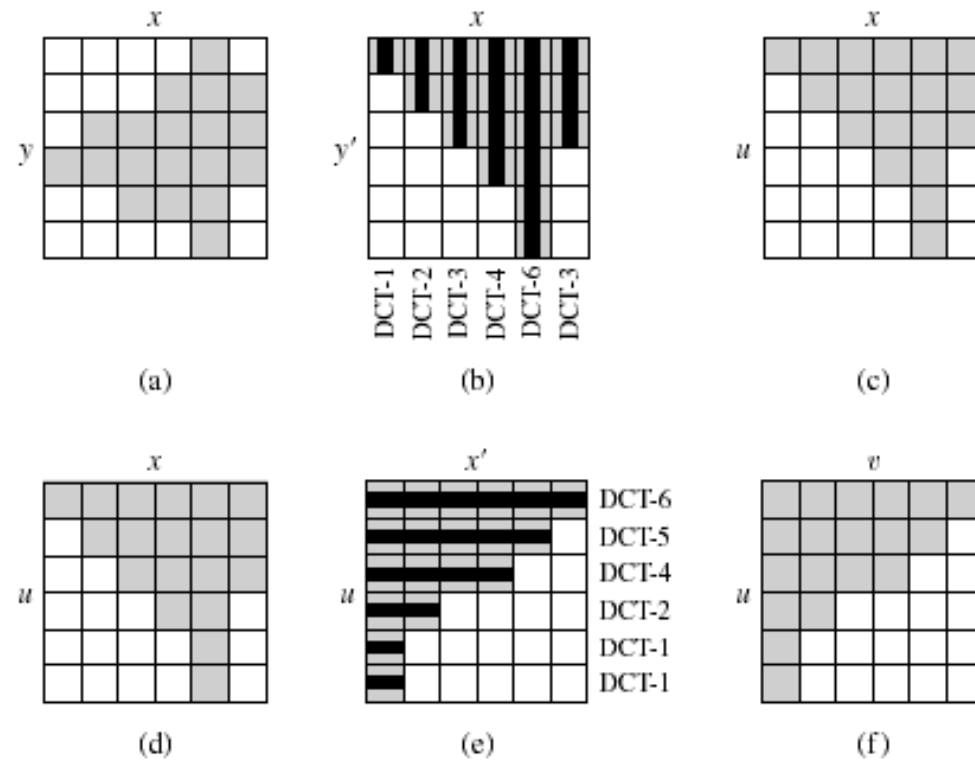- Shape Adaptive DCT



Special color is assigned to pels not belonging to the object

Bounding box can be extended to multiples of 8x8 if the resulting box is to be coded using 8x8 DCT
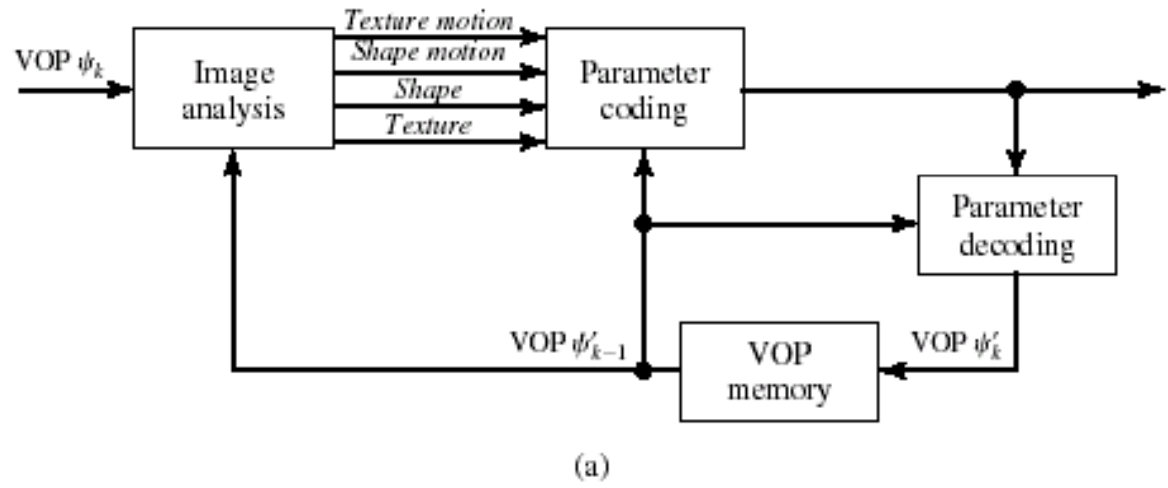
# Shape Adaptive DCT



**Figure 10.9** A shape-adaptive DCT requires transforms of length *n*: (a) original image segment, (b) pels shifted vertically, (c) location of DCT coefficients after vertical 1-D DCT, (d) location of DCT coefficients prior to horizontal 1-D DCT, (e) DCT coefficients shifted horizontally, (f) location of DCT coefficients after horizontal DCT.
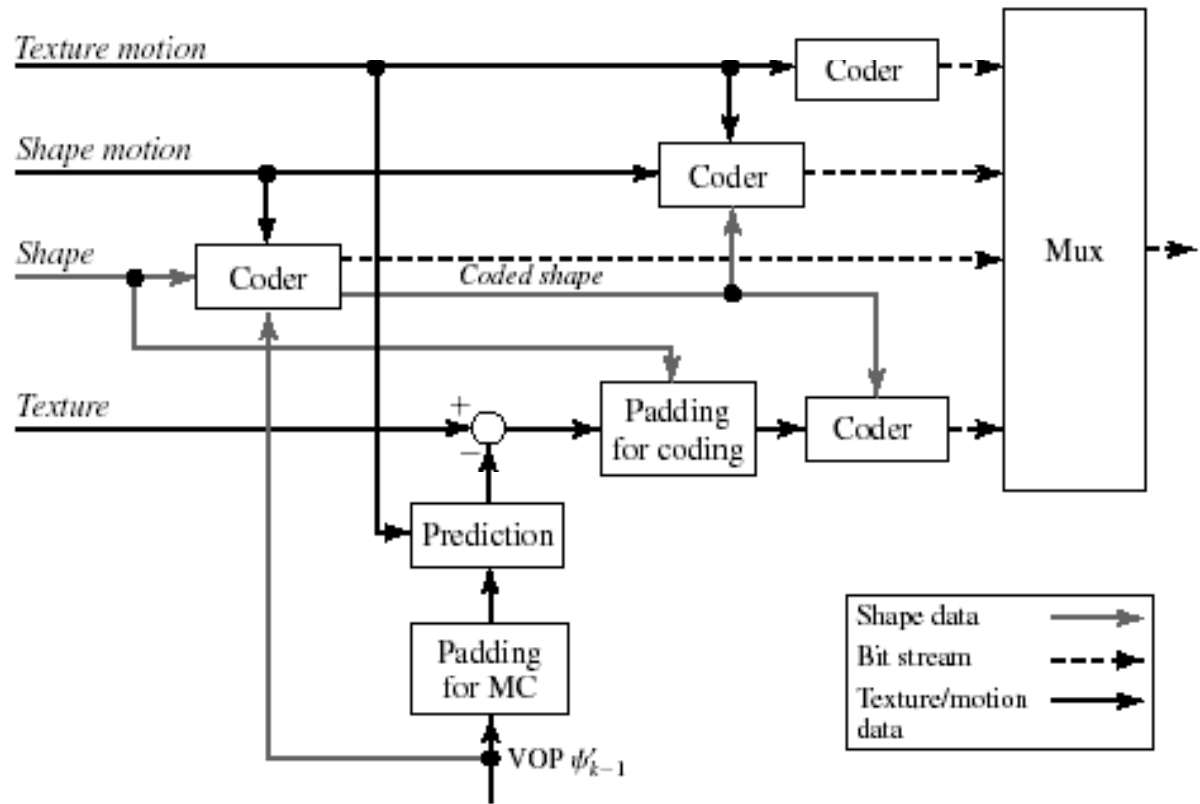
# MPEG-4 Shape Coding

- Uses block-based approach (block=MB)
  - Boundary blocks (blocks containing both the object and background)
  - Non-boundary blocks: either belong to the object or background
- Boundary block's binary alpha map (binary alpha block) is coded using context-based arithmetic coding
  - Intra-mode: context pels within the same frame
  - Inter-mode: context pels include previous frame, displaced by MV
    - Shape MV separate from texture MV
    - Shape MV predictively coded using texture MV
- Grayscale alpha maps are coded using DCT
- Texture in boundary blocks coded using
    - padding followed by conventional DCT
    - Or shape-adaptive DCT

MPEG4
video coder
overview

Texture motion
Shape motion
Shape
Texture

VOP $\psi_k$

Image
analysis

Parameter
coding

Parameter
decoding

VOP $\psi'_{k-1}$

VOP
memory

VOP $\psi'_k$

(a)

Details of
parameter
coding

Texture motion

Coder

Shape motion

Coder

Shape

Coder

Coded shape

Mux

Texture

Padding
for coding

Coder

Prediction

Padding
for MC

VOP $\psi'_{k-1}$

Shape data
Bit stream
Texture/motion
data

# Still Texture Coding

- MPEG-4 defines still texture coding method for intra frame, sprite, or texture map of an mesh object
- Use wavelet based coding method

# Mesh Animation

- An object can be described by an initial mesh and MVs of the nodes in the following frames
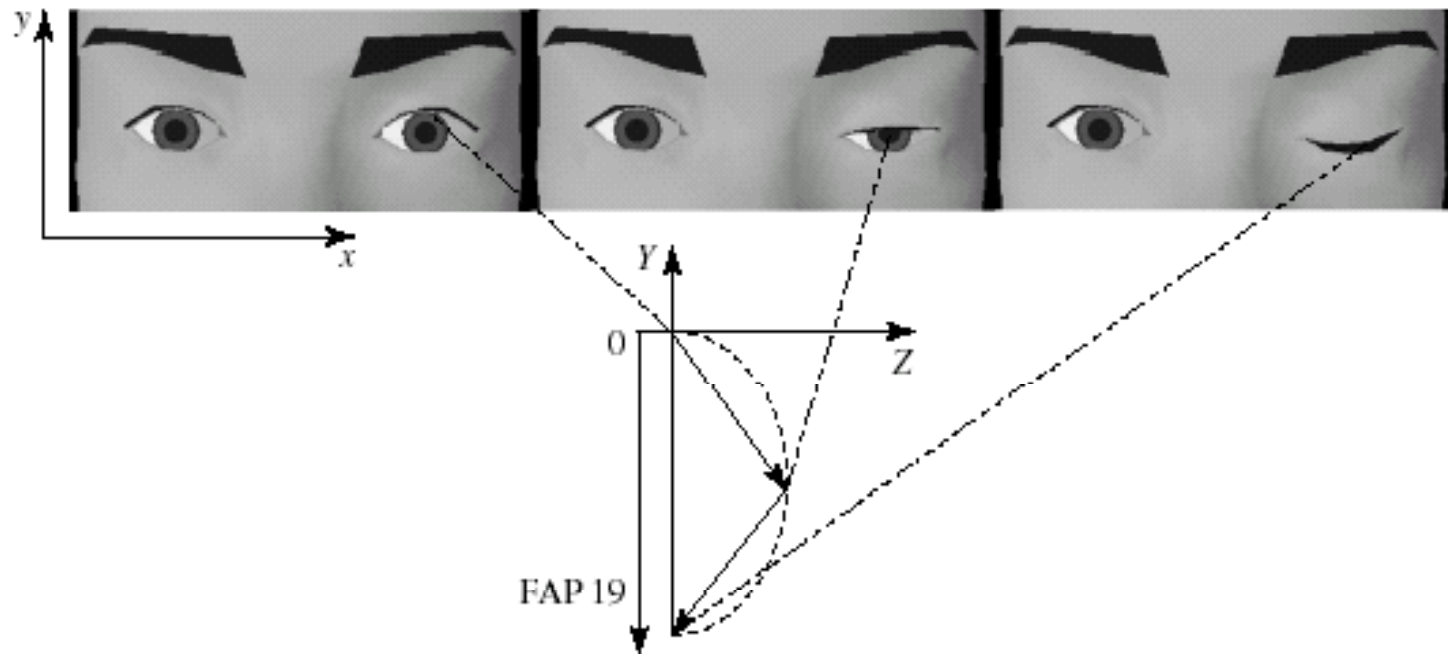- MPEG-4 defines coding of mesh geometry, but not mesh generation

# Body and Face Animation

- MPEG-4 defines a default 3-D body model (including its geometry and possible motion) through body definition table (BDP)

- The body can be animated using the body animation parameters (BAP)

- Similarly, face definition table (FDP) and face animation parameters (FAP) are specified for a face model and its animation

# Face Animation Through FAP



**Figure 13.36** Neutral state of the left eye (left) and two deformed animation phases for the eye blink (FAP 19). The FAP defines the motion of the eyelid in the negative y direction; the faceDefTable defines the motion in one of the vertices of the eyelid in the x and z directions.

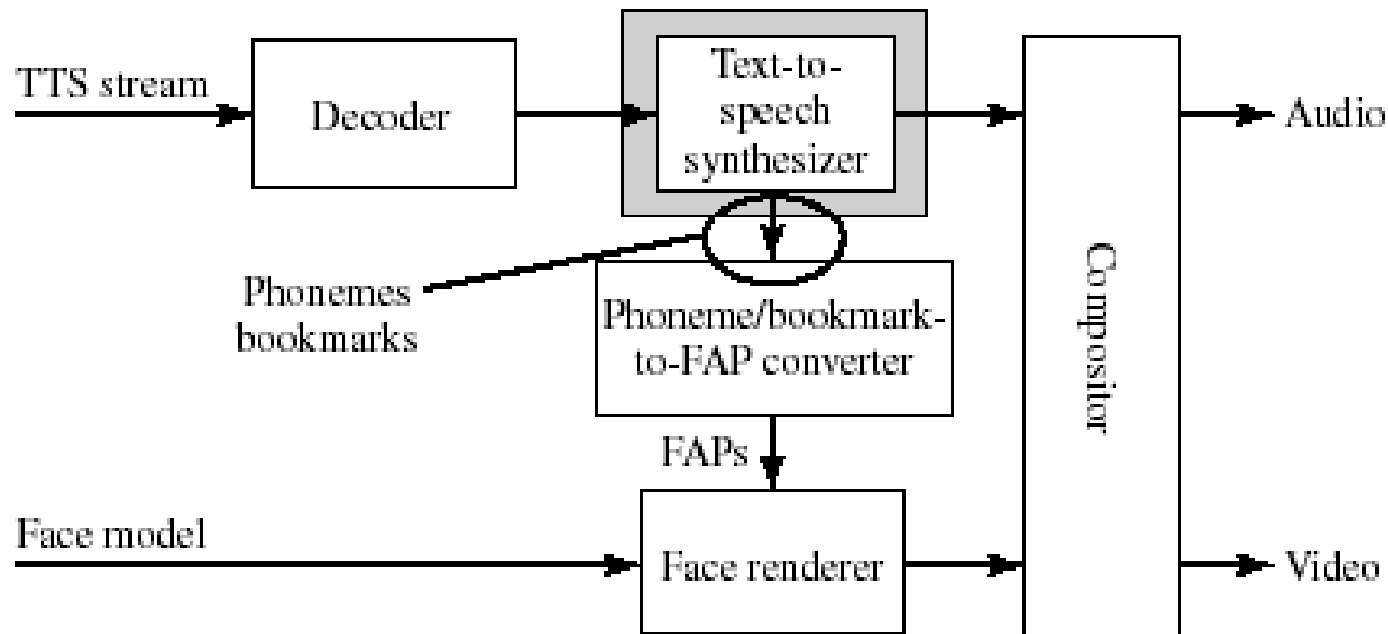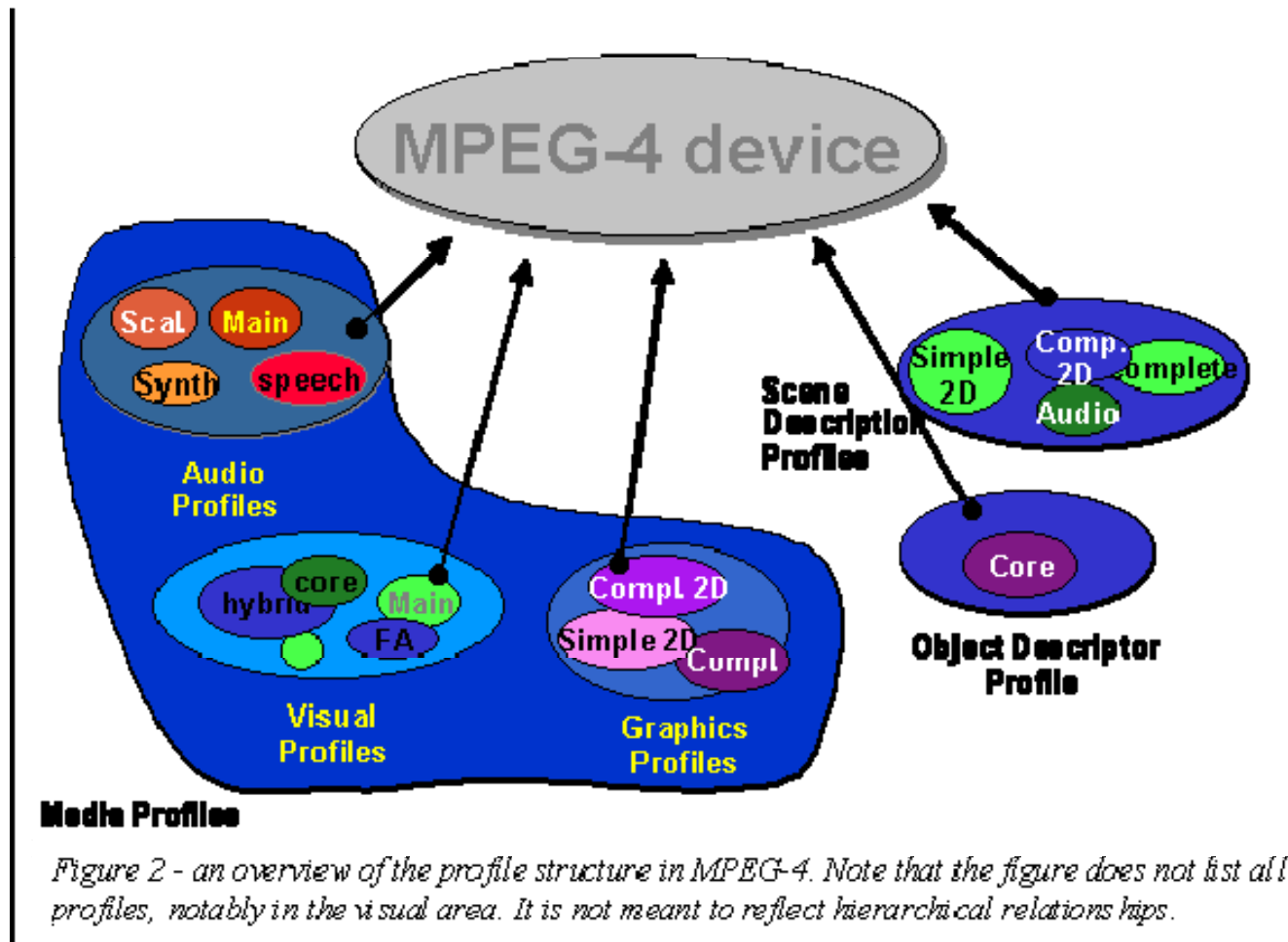# Text-to-Speech Synthesis with Face Animation



**Figure 13.37** MPEG-4 architecture for facial animation, allowing synchronization of facial expressions and speech generated by a proprietary text-to-speech synthesizer.

# MPEG-4 Profiles



Figure 2 - an overview of the profile structure in MPEG-4. Note that the figure does not list all profiles, notably in the visual area. It is not meant to reflect hierarchical relationships.
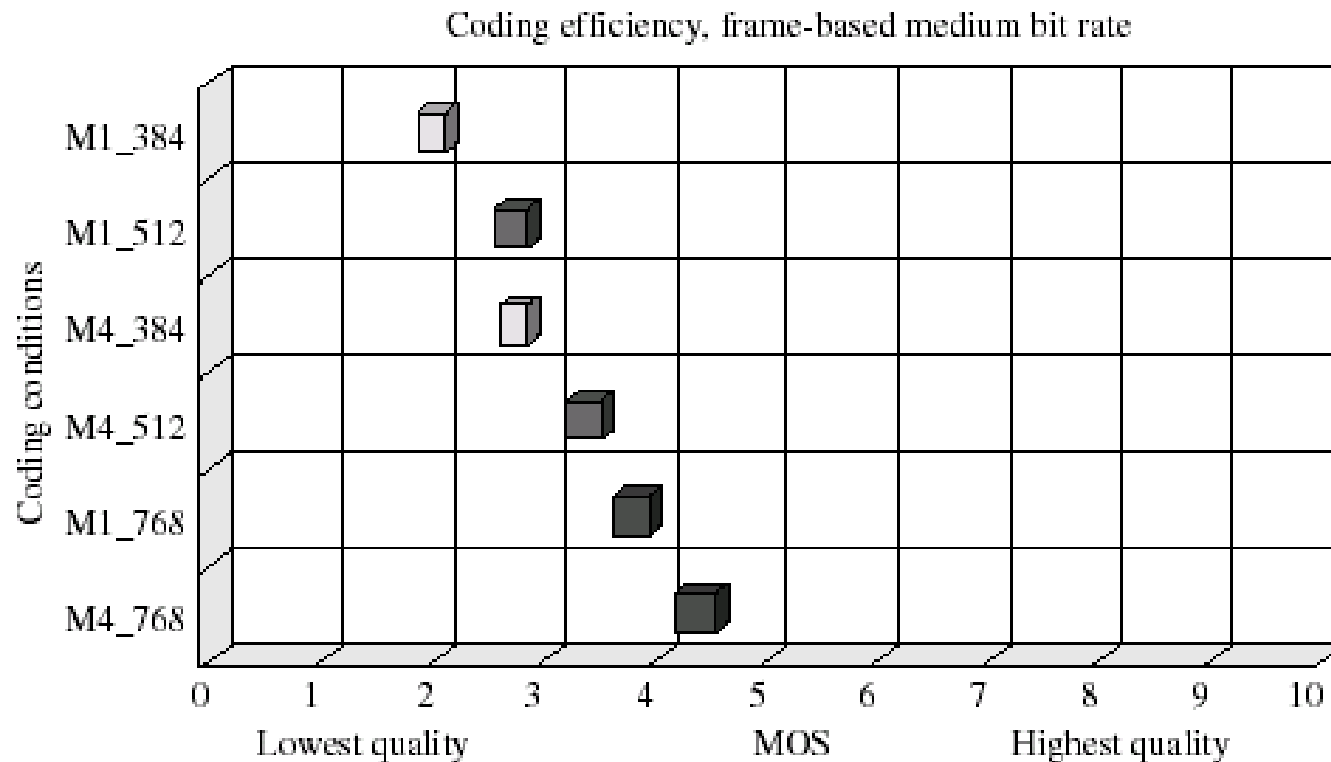
# Video Coding Efficiency Tools

- Sprite
  - Code a large background in the beginning of the sequence, plus affine mappings, which map parts of the background to the displayed scene at different time instances
  - Decoder can vary the mapping to zoom in/out, pan left/right

- Global motion compensation
  - Using 8-parameter projective mapping
  - Effective for sequences with large global motion

- DC and AC prediction: can predict DC and part of AC from either the previous and block above

- Quarter-pel motion estimation

- Similar to H.263
  - 3D VLC
  - Four MVs and Unrestricted MVs
  - OBMC not required

# MPEG-4 vs. MPEG-1 Coding Efficiency



Figure 13.39    Subjective quality of MPEG-4 Main profile versus MPEG-1. M4_x is an MPEG-4 coder operating at the rate of x kbps; M1_x is an MPEG-1 encoder operating at the given rate [27].

# Summary

- What does video coding standard define?
- H.261:
    - First video coding standard, targeted for video conferencing over ISDN
    - Uses block-based hybrid coding framework with integer-pel MC
- H.263:
    - Improved quality at lower bit rate, to enable video conferencing/telephony below 54 bkps (modems or internet access, desktop conferencing)
    - Half-pel MC and other improvement
- MPEG-1 video
    - Video on CD and video on the Internet (good quality at 1.5 mbps)
    - Half-pel MC and bidirectional MC
- MPEG-2 video
    - TV/HDTV/DVD (4-15 mbps)
    - Extended from MPEG-1, considering interlaced video

# References

- Chap. 13