# 5 Multimedia Content Description

Literature:

B.S. Manjunath et al. (eds.): Introduction to MPEG-7 - Multimedia Content Description Interface, Wiley 2002
MPEG-7 Overview,
http://www.chiariglione.org/mpeg/standards/mpeg-7/mpeg-7.htm

---

# MPEG-7 Description Terminology

- Feature:
  - Distinctive characteristic of the data which siginifies something to somebody
- Descriptor:
  - Representation of a feature
    » Defines syntax and semantics of feature representations
  - A feature may be represented by several descriptors
- Descriptor value:
  - Instantiation of a descriptor
- Description scheme:
  - Structured composition of descriptions and description schemes
- Description:
  - Instance of a description scheme with appropriate descriptor values
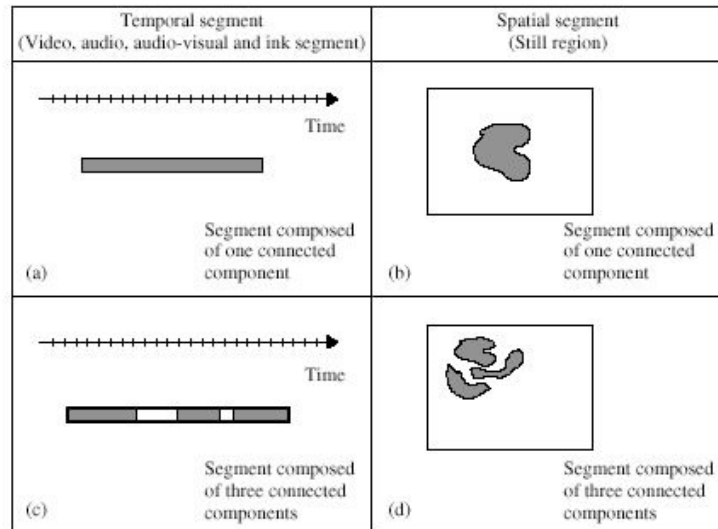
# Metadata Classification in MPEG-7

- Content Management
  - Media information (e.g. file name, format, resolution)
  - Creation information (e.g. creator, location, date)
  - Usage information (e.g. rights)
- Content Structure
  - Segments
  - Segment relations
- Content Semantics

Following slides: Details on Content Structure

# Structural Content Description: Segments

- A segment represents a section of an audio-visual content item.
- The Segment Description Scheme (DS) is an abstract class
  (in the sense of object-oriented programming).
- It has nine major subclasses:
  - Still Region DS (spatial)
    - » ImageText DS
  - Video Segment DS (temporal)
    - » Analytic edited video segment DSs
  - Moving Region DS (spatiotemporal)
    - » VideoText DS
  - Audio Segment DS (temporal)
  - AudioVisual Segment DS (temporal)
  - AudioVisual Region DS (spatiotemporal)
  - Still Region 3D DS (3D spatial)
  - Ink Segment DS (electronic ink from pen, smartboard etc. )
  - Multimedia Segment DS (composite of segments)
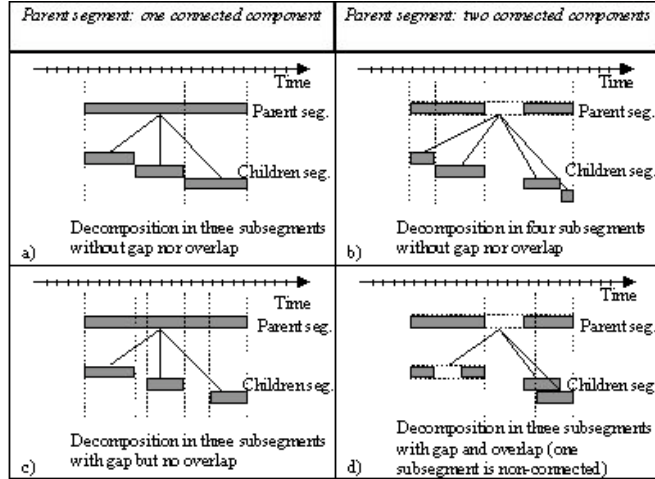
# Examples of Segments

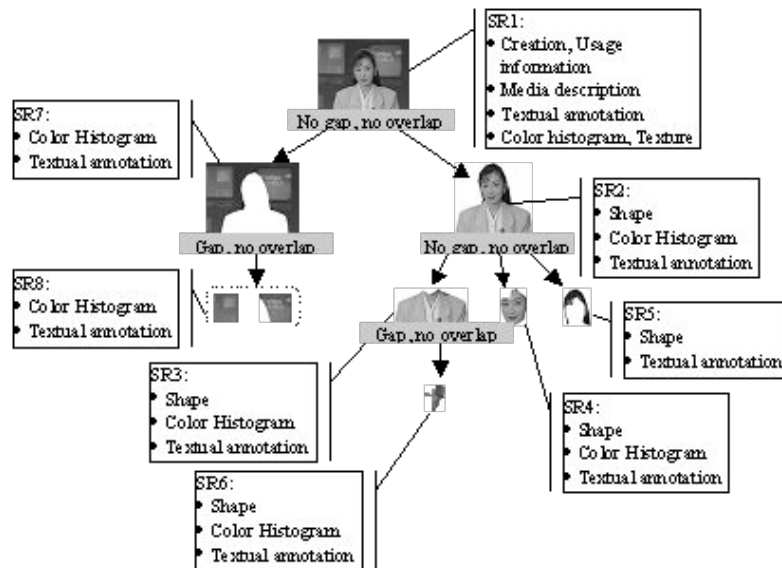| Temporal segment (Video, audio, audio-visual and ink segment) | Spatial segment (Still region) |
|---|---|
| Time<br><br>Segment composed of one connected component<br>(a) | Segment composed of one connected component<br>(b) |
| Time<br><br>Segment composed of three connected components<br>(c) | Segment composed of three connected components<br>(d) |

# Segment Attributes

- Generic features
  - (media information, creation information, usage information, annotations)
- Media type dependent features:
  - (visual features, audio features)
- Specific features for segments
  - Mask Descriptor
    - » Spatial mask, Temporal mask, Spatio-temporal mask
  - Importance of descriptors
    - » MatchingHint: relative importance of descriptors
    - » PointOfView: relative importance of segments for a specific point of view (PointOfView given as string, e.g. "Home team" for soccer game)
  - Ink segment descriptors
    - » Handwriting recognition information (recognizer, lexicon)
    - » Handwriting recognition result (quality, accuracy-scored results)

# Segment Decomposition

- Segments can be decomposed into subsegments
  - Subsegments may overlap in time/space
  - Subsegments may not cover the full extents of parent segment
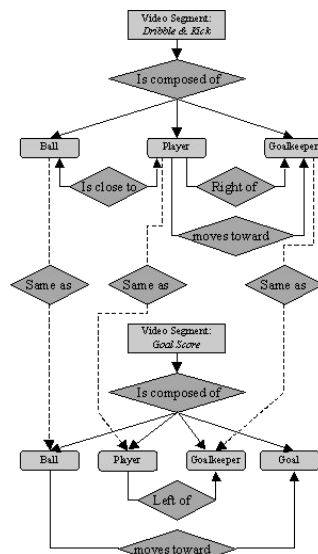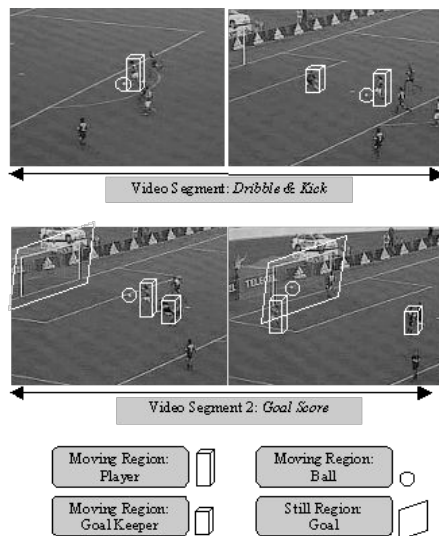  - Decomposition may result in segments of different nature



| Parent segment: one connected component | Parent segment: two connected components |

a) Decomposition in three subsegments without gap nor overlap

b) Decomposition in four subsegments without gap nor overlap

c) Decomposition in three subsegments with gap but no overlap

d) Decomposition in three subsegments with gap and overlap (one subsegment is non-connected)

---

# Example of Image Description



SR1:
- Creation, Usage information
- Media description
- Textual annotation
- Color histogram, Texture

SR7:
- Color Histogram
- Textual annotation

No gap, no overlap

Gap, no overlap

No gap, no overlap

SR2:
- Shape
- Color Histogram
- Textual annotation

SR8:
- Color Histogram
- Textual annotation

Gap, no overlap

SR5:
- Shape
- Textual annotation

SR3:
- Shape
- Color Histogram
- Textual annotation

SR4:
- Shape
- Color Histogram
- Textual annotation

SR6:
- Shape
- Color Histogram
- Textual annotation
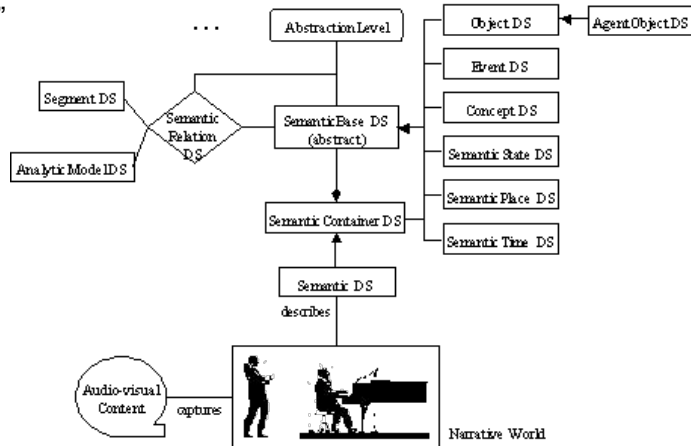
# Structural Relations of Segments

- Content structure:
  - Either hierarchical segment decomposition
  - Or general segment relationship graph
- Predefined structural relations in MPEG-7 (can be extended):
  - Spatial:
    - » South, north, west, east, northwest, northeast, southwest, southeast, left, right, below, above, over, under
  - Temporal:
    - » Precedes, follows, meets, metBy, overlaps, overlappedBy, contains, during, strictContains, strictDuring, starts, startedBy, finishes, finishedBy, coOccurs, contiguous, sequential, coBegin, coEnd, parallel, overlapping
  - Generic:
    - » Identical, union, disjoint
- For each relation, the inverse relation is implicitly defined.

---

# Video Segmentation with Moving Regions

# Content Semantics in MPEG-7

- Event: Occasion when something happens
  - Occurs at some time and place
  - Populated by objects and people
- "Narrative world" for a piece of content

# Content Semantics in MPEG-7: Example
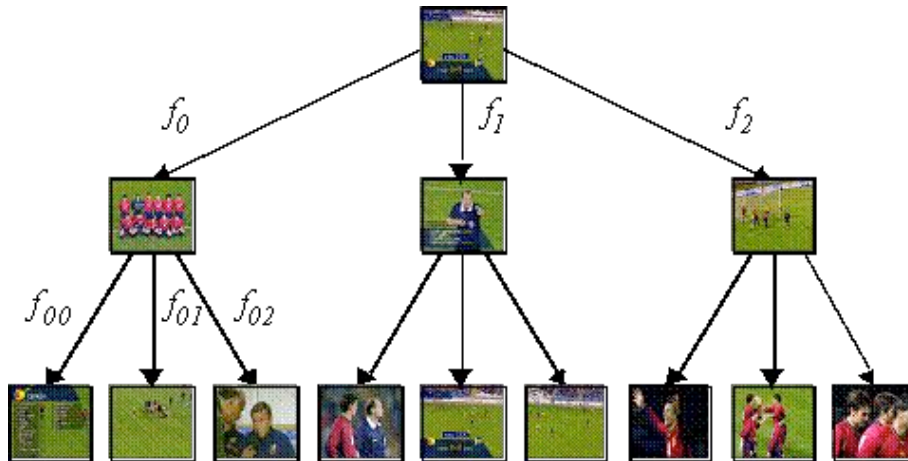
# Relating Structure and Semantics: Example
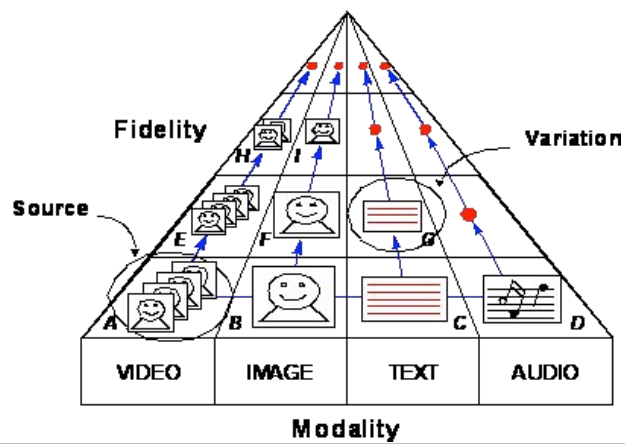
- Creating a semantic index for a video sequence



Salembier 2002

---

# Navigation and Access

- Description schemes to facilitate navigation and access of audio-visual content:
    - Summaries
        - » Browsing, navigation, discovery, visualization, sonification
    - Views and partitions
        - » Representations in multiple domains, resolutions
    - Variations
        - » Different versions adapted to delivery conditions

# Example: Summary as Hierarchy of Key Frames



$f_0$      $f_1$      $f_2$

$f_{00}$      $f_{01}$   $f_{02}$

---

# Variations

- Components of a complex multimedia object may exist in various variations (different resolutions, languages, etc.)
  - Server or proxy server should be able to select the appropriate variation

## MPEG-7 Visual Description Tools

- Descriptors for the following basic visual features:
    - Color, Texture, Shape, Motion, Localization, and Face recognition
    - Each category consists of elementary and sophisticated Descriptors
- Basic structures for composing visual features:
    - Grid layout
    - Time series
    - Multiple (2D/3D) view
    - Spatial 2D coordinates
    - Temporal interpolation

## Principles of Automatic Feature Extraction

- Colour:
    - Histogram, colour clusters
- Texture:
    - Spectral distribution, energy
- Motion:
    - Vector histogram, parametric models
- Contours:
    - Moments, wavelet coefficients
- Faces:
    - Vector basis and similarity matching

- Usage of compressed data formats:
    - E.g. frequency space transformation (JPEG), motion estimation (MPEG-2) can be re-used for feature extraction

# Shape Descriptors

- Region shapes
  - Pixel distribution, using both boundary and internal pixels
  - Can describe complex objects with multiple disconnected regions
  - Shape analysis based on moments
    - » Angular Radial Transformation (ART)
- Contour shapes
  - Based on Curvature Scale-Space (CSS) representation of contour
  - Recognized characteristic contour shapes
  - Similar to human perception
- Desirable properties of extraction methods
  - Able to handle complex shapes
  - Robust to minor deformations, perspective transformations, movement, splits, occlusions etc.
  - Compact and efficient

---

# Examples for Shape Descriptors
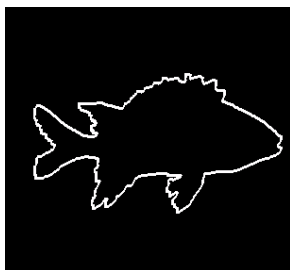
Region shapes:



Contour shapes:

# Angular Radial Transformation (ART)

- Convert image information into angular and radial parts
- Represent image as coefficients of basis functions
- First 36 basis functions:
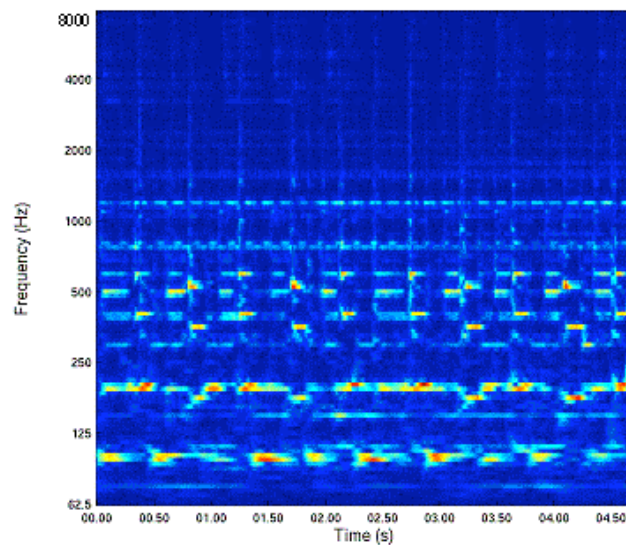
---

# Curvature-Scale Space Computation

- Curvature is a local measure of how fast a curvature is turning
  - Curvature zero crossing points are essential for contours
  - Contour is sampled with increasing precision and smoothed stepwise to retrieve curvature zero-cronssings of various scales
- http://www.ee.surrey.ac.uk/Research/VSSP/imagedb/affine.html
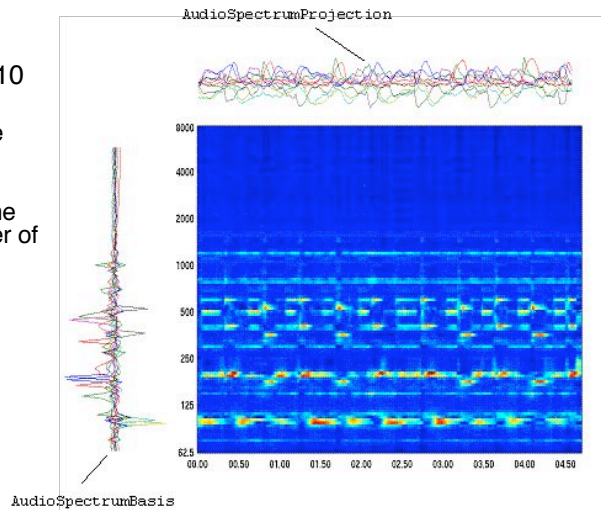
# MPEG-7 Audio Description Tools

- Low-level audio descriptors:
  - Basic: Instantaneous waveform and power values
  - Basic spectral: Log-frequency power spectrum and spectral features (centroid, spread, flatness)
    - » AudioSpectrumEnvelope: Spectrogram of the signal
  - Signal parameters: Fundamental frequency
  - Temporal timbral: Log attack time and temporal centroid
  - Spectral timbral: Specialized spectral features
- High-levell audio descriptors:
  - Sound recognition and indexing
  - Musical instrument timbre description
  - Melody description tools
  - Spoken language recognition

# Spectral Analysis with AudioSpectrumEnvelope

## Data-Reduced Spectral Representation

- Reconstruction of sonogram using a compact representation of 10 vectors
    - required storage space 10(M+N) values

    M number of time points, N number of spectrum bins
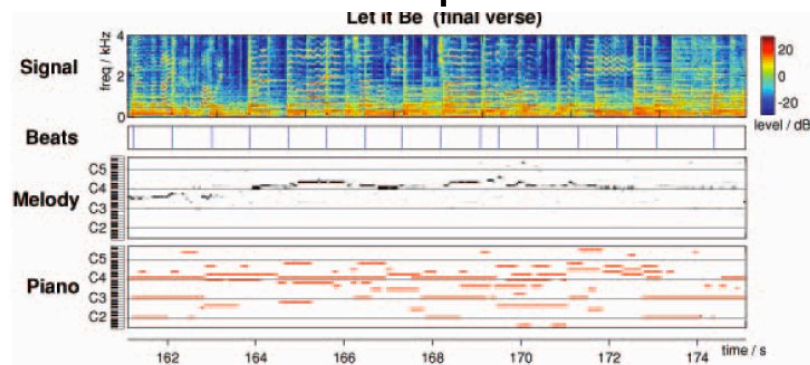
---

# 5 Multimedia Content Description

Literature:
*Communications of the ACM* 49(8), August 2006,
Special section on Music Information Retrieval, pp. 28-60

# Timescales of Musical Information

- Individual music note events
  - Extraction of the music score
  - Identification of instrument playing
- Chords (simultaneous notes)
  - Identification of chords
- Phrase level
  - Tempo extraction
  - Identification of phrases (based on repetition/alternation of segments)
    e.g. identificaton of chorus
- Piece level
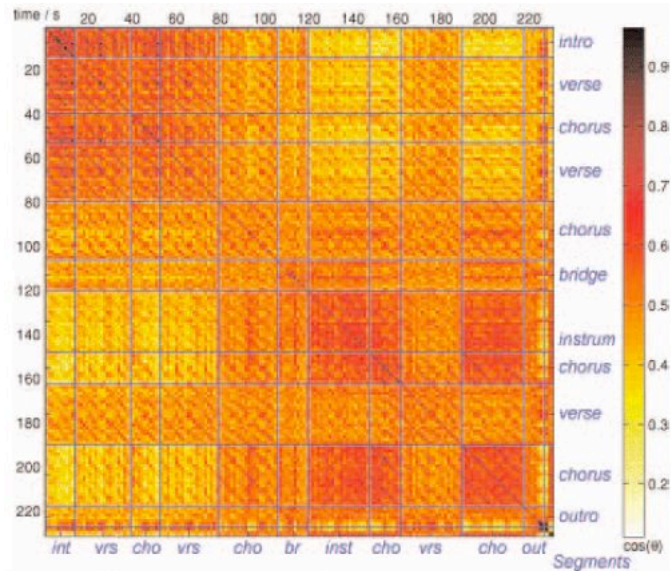  - Genre identification ("rock", "jazz", "classical")

---

# Automatic Score Transcription



- Beats determined by tempo-smoothed event detector
- Melody recognized by general-purpose support-vector classifier
  - Trained to recognize spectral slices to be labelled with pitch values

## Automatic Phrase Detection

- Self-similarity matrix
  - Looking for diagonal ridges off the main diagonal
  - Blue lines are manually inserted for comparison

---

## Example: Shazam Music Tagging

- Commercial service for mobile phones: Identify music from a short audio sample *(query by example)*
  - See http://www.shazam.com
- Challenges:
  - Distinguishing music from noise
  - Dealing with distortions
  - Keeping fingerprints small (in order to deal with millions of songs)
- Basic idea:
  - Spectrogram peaks (energy distribution in time and frequency)
  - Few "anchor" peaks are combined with peaks in a certain surrounding zone (time and frequency offsets)
    - » Combinatorial hashing creates 32b fingerprint hash token
  - Temporal alignment greatly accelerates matching process
- Real system:
  - "a few dozen" x86-based servers with 8GB of RAM, Linux
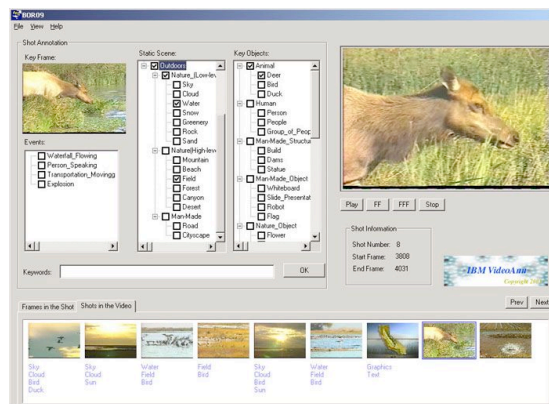  - June 2006: More than three million tracks

# 5 Multimedia Content Description

Literature:
www.virage.com

---

# IBM VideoAnnEx (1)

- Support tool for manual annotation of video sequences with MPEG-7 metadata
    - Experimental tool 2001-2003, no longer supported
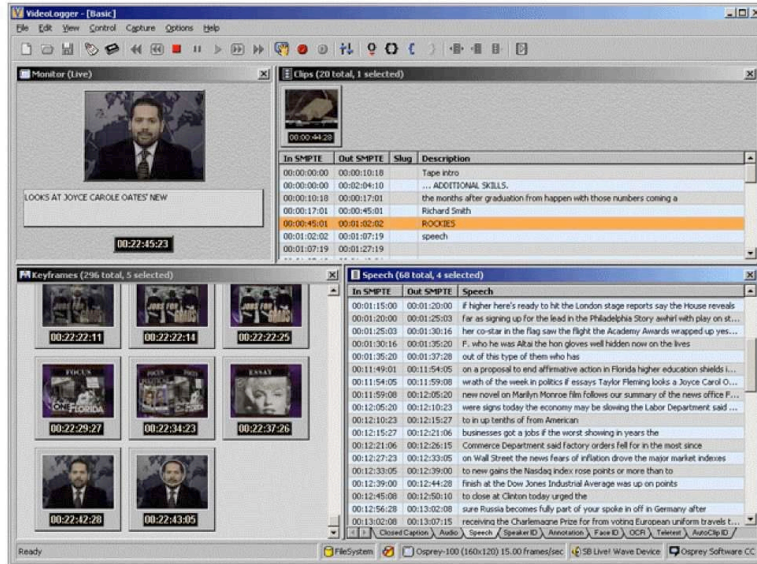    - Requires a basic lexicon of description items in addition to video file

# IBM VideoAnnEx (2)

# IBM VideoAnnEx (3)

# Virage VideoLogger

---

# Techniques used by Virage VideoLogger

- Signal analysis algorithms to generate key frames for visual overview
- Speech-to-text transcription
- Sound identification
- Speaker identification
  - voice identification and face identification
- Analysis of embedded textual information:
  - close captioning, teletext
- External metadata:
  - PowerPoint presentations
  - EDLs
  - GPS data
  - transcripts
- Manual annotation:
  - Effective user interface (hot keys etc.)