

Audio-Visual Content Indexing, Filtering, and Adaptation

Shih-Fu Chang

Digital Video and Multimedia Group
ADVENT University-Industry Consortium
Columbia University

10/12/2001

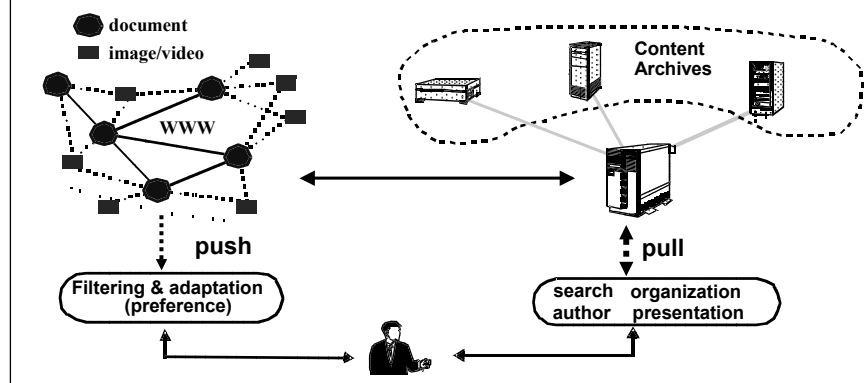
<http://www.ee.columbia.edu/dvmm>

10/2001

S.-F. Chang Columbia U.

1

Challenge: tools/systems for searching and filtering



■ Important Issues:

- Information Overload, Limited User Attention, Time-Sensitive
- Heterogeneous Platforms

■ Example Products:

- AltaVista, Google, DoCoMo/IBM, NewsTake, TiVo/Philips, PC DVR

10/2001

S.-F. Chang Columbia U.

2

MPEG-7 and issues

- ISO/IEC 15938

“Multimedia Content Description Interface”



- Issues:

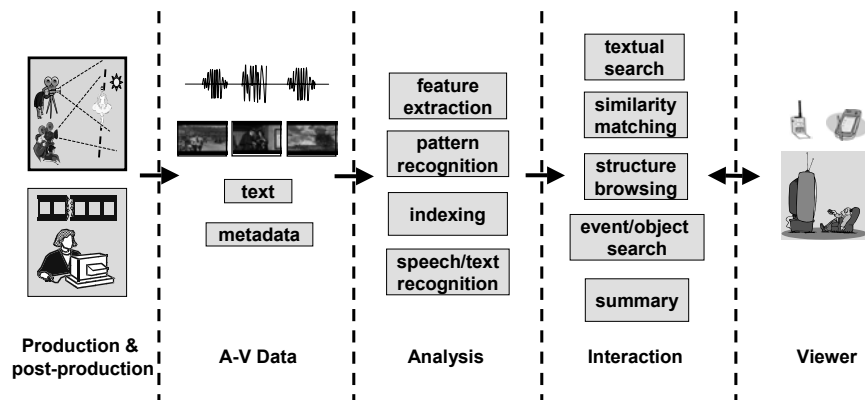
- Active research in analysis and indexing
- System issues: delivery, compression
- What applications? ...

10/2001

S.-F. Chang Columbia U.

3

Re-examine Content Production/Usage Chain



10/2001

S.-F. Chang Columbia U.

4

Types of Content Indexing

- **Reverse engineering**
 - Production structure parsing (shot, camera)
- **Metadata preservation**
 - Format, Production, Credits, etc.
- **Feature measurement**
 - audio-visual, spatio-temporal objects and features
- **Content recognition and organization**
 - Semantic meaning description
 - Scene grouping and skimming

10/2001

S.-F. Chang Columbia U.

5

When do we need automated tools

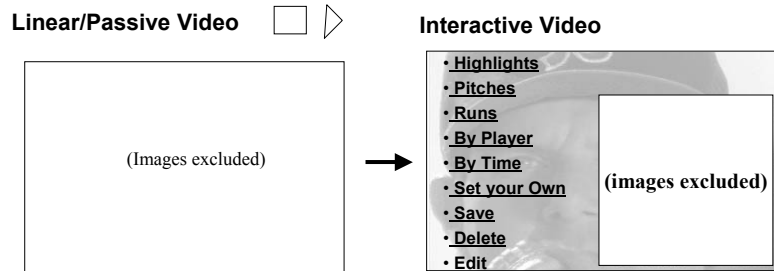
- **Conditions for high impact**
 - Metadata that's not available from production
 - Work that humans are not good at (e.g., low-level computation)
 - Large volume, low individual value
 - Time sensitive
 - Acceptable performance
- **Less Promising Areas**
 - Content from digital production tools with metadata
 - Prime content that can afford manual solutions

10/2001

S.-F. Chang Columbia U.

6

Case 1: Live Sports Video Filtering and Navigation



- Time sensitive interest
- Massive production and audience
- Time compressibility
- Temporal structure and production rules

Beyond Ringer and Clip Download

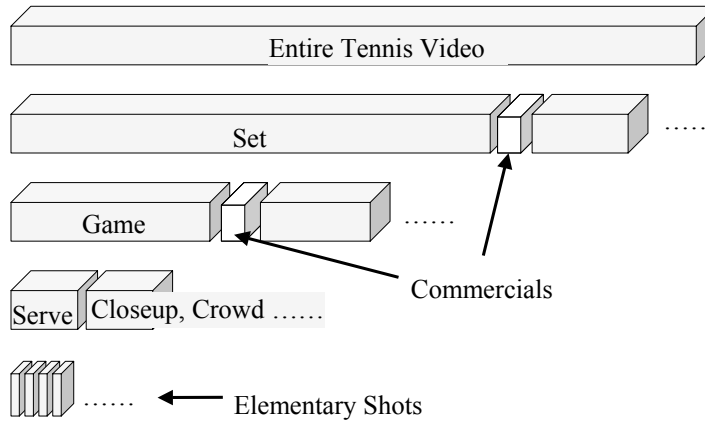
Services:

- Messaging
- Localized information
- Media/game download
- Multi-person games
- TV phone
- On-site purchase

Time-sensitive short video messages suiting personal needs

Image showing sports video on mobile handset excluded

Regular Structure and Views in Sports Video

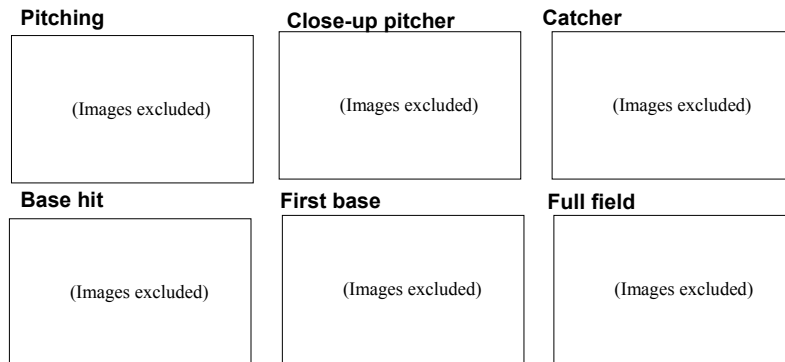


10/2001

S.-F. Chang Columbia U.

9

Regular set of views



- **Important issues:**
 canonical view \leftrightarrow beginning of recurrent semantic unit
 view transition pattern \leftrightarrow types of events

10/2001

S.-F. Chang Columbia U.

10

Real-Time Sports Video Parsing

- Detect and classify recurrent semantic units
- Real-time processing
 - simple features
 - compressed-domain processing
- Combine global-level and object-level classification

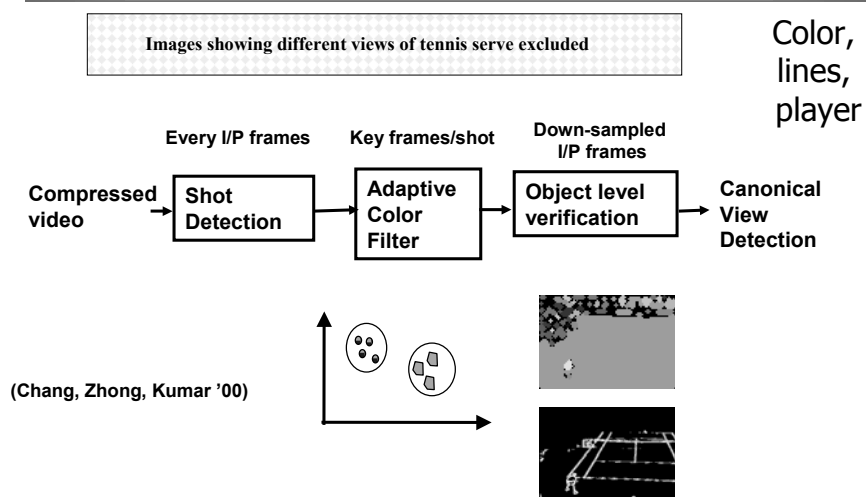
(Chang, Zhong, Kumar '01)

10/2001

S.-F. Chang Columbia U.

11

Detecting recurrent semantic units: Serve

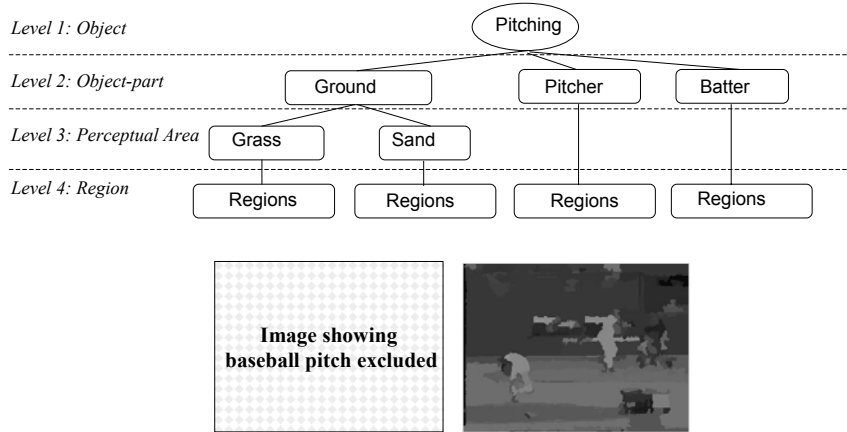


10/2001

S.-F. Chang Columbia U.

12

Learning spatio-temporal rules

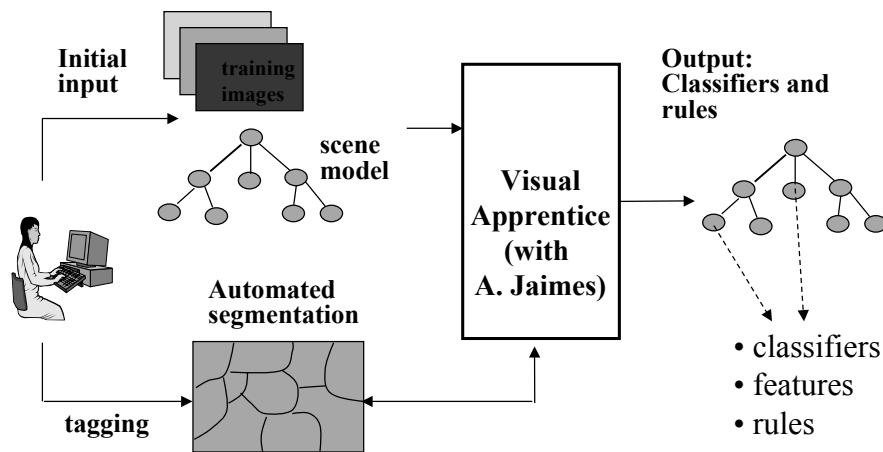


10/2001

S.-F. Chang Columbia U.

13

Systematic learning of object rules



(Jaimes and Chang '99)

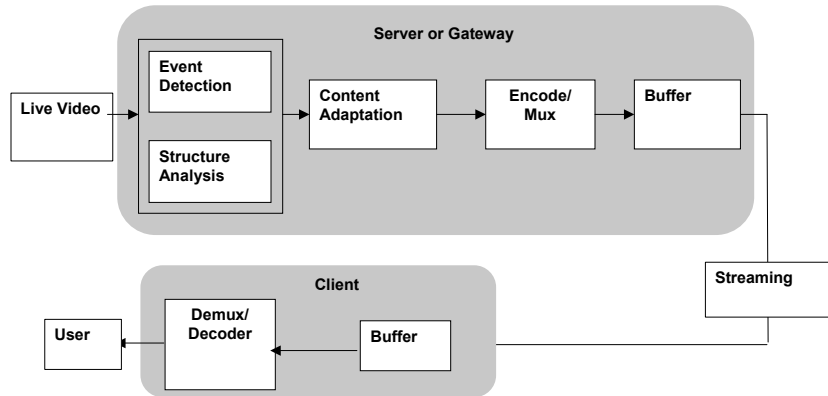
10/2001

S.-F. Chang Columbia U.

14

Application: Time-Sensitive Mobile Video

Content Adaptive Streaming

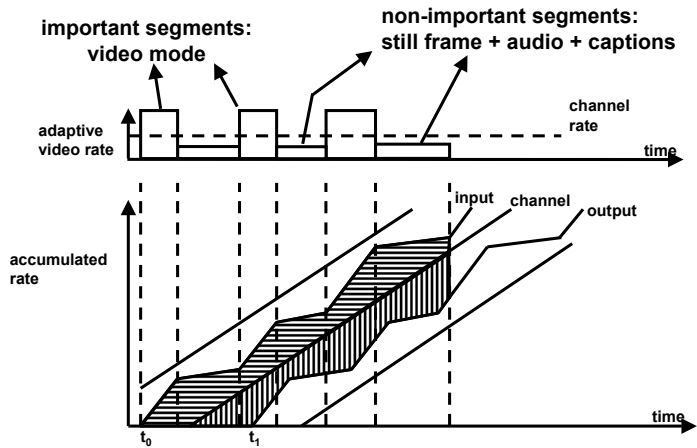


10/2001

S.-F. Chang Columbia U.

15

Content Adaptive Streaming

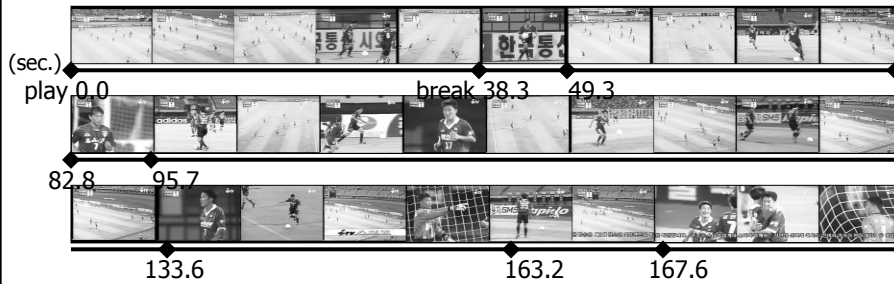


10/2001

S.-F. Chang Columbia U.

16

Recurrent semantic units in soccer video?



- Game → a sequence of play and break segments
- Play/break → No canonical views/events
- Sporadic events (start and end)
- Shot boundary \neq play boundary

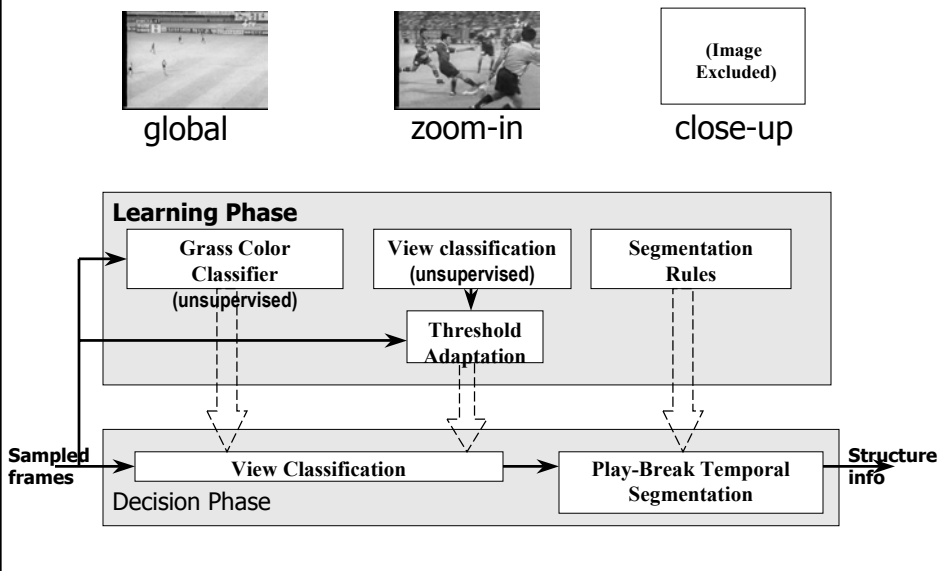


10/2001

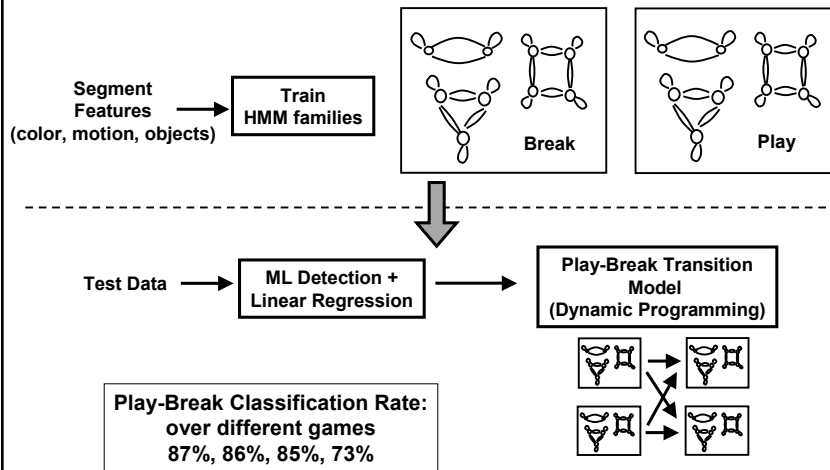
S.-F. Chang Columbia U.

17

Heuristic Model: Features → Views → Structure



Modeling Transitions with HMM

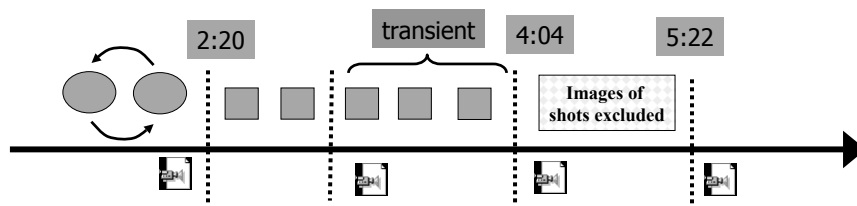


10/2001

S.-F. Chang Columbia U.

19

Case 2: Long Programs with Loose Structures



- **Goal:**
 - **Scene/Topic Segmentation, Browsing, Preview**
- **Challenge:**
 - **Diverse production styles and views**

10/2001

S.-F. Chang Columbia U.

20

Domain Consideration

- **Film:**
 - **Does not satisfy impact conditions**
 - **Not much value except**
 - archived collection in libraries
 - less popular selection in PVR/STB
 - **Challenging, rich styles and media**
 - **Fun**

10/2001

S.-F. Chang Columbia U.

21

Problems and Approaches

- **Explore production models**
- **Explore multimedia integration**
- **Incorporate structural and special cues**
- **Explore viewer's perceptual models**

- **Goal:**

Computable features/theories + models



Scene structures and summaries

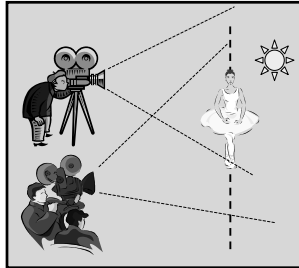
10/2001

S.-F. Chang Columbia U.

22

Production Constraints → Video Scenes

■ Camera placement [180° rule]



overlapping background
causes chromatic consistency

■ Lighting/chromaticity continuity

10/2001

S.-F. Chang Columbia U.

23

Audition Psychology → Audio Scenes

A. Bregman '90: Auditory Scene Analysis

- Unrelated sounds seldom begin/end at the same time.
- Sounds from the same source change properties smoothly and gradually over time.
- Changes in acoustic states affect all components of the sound (e.g., walking away from a ringing bell)
- Human perception uses long term grouping (e.g., a series of footsteps).

- An audio scene change is said to occur when majority of the sound sources change.

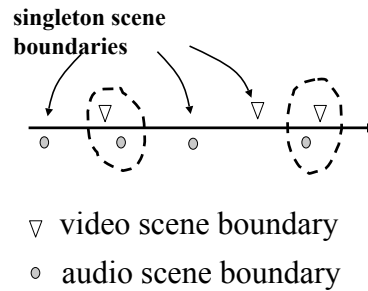
10/2001

S.-F. Chang Columbia U.

24

Explore Audio-visual association

- Complementary audio-visual structures
- first hour of Blade Runner
 - video: 28 scenes
 - audio only: 33 scenes
- audio/video agreement: 24
 - video change only: silent or transient scenes
 - audio change only: mood/state change



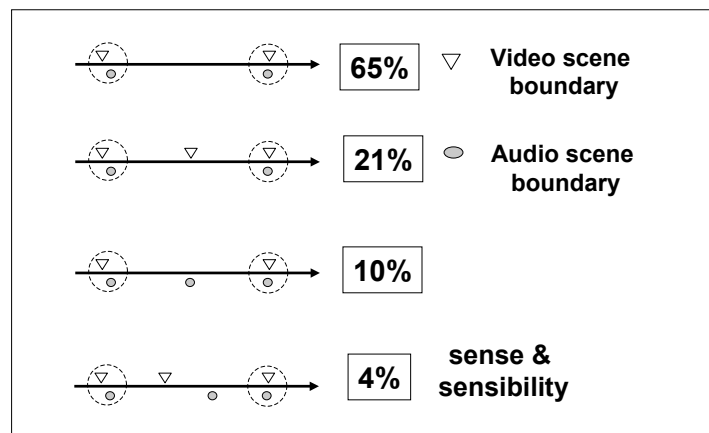
Sundaram & Chang ICASSP '00

10/2001

S.-F. Chang Columbia U.

25

Mixing audio-visual scenes



10/2001

S.-F. Chang Columbia U.

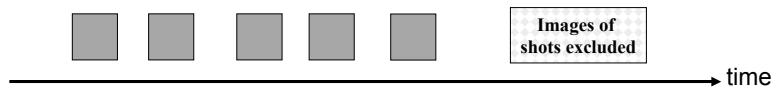
26

Film: Common Scene Structures

- 4 common scene types in film

- *Progressive or Coherent:*

- same location, similar camera takes, consistent audio



(Demo)

10/2001

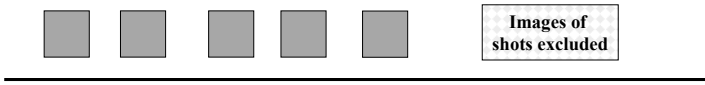
S.-F. Chang Columbia U.

27

Common Scene Types (2)

- *Transient :*

- different locations, dissimilar camera takes, evolving audio



- *MTV-type: fast, dissimilar shots, consistent audio*

(Demo)

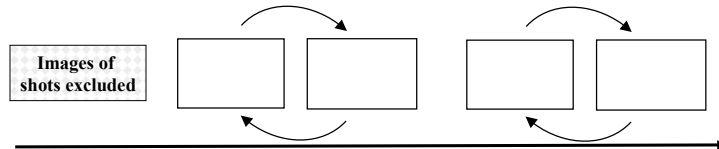
10/2001

S.-F. Chang Columbia U.

28

Common Scene Types (3)

- *Dialog*: repeating structures, consistent audio

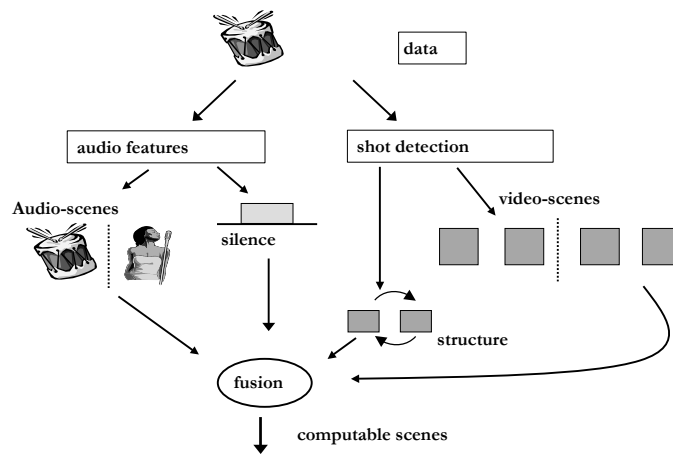


10/2001

S.-F. Chang Columbia U.

29

Computable Scene Detection Architecture

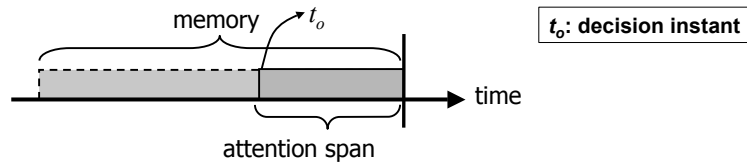


10/2001

S.-F. Chang Columbia U.

30

The Viewer/Listener Model



- **Memory (M):** Net amount of data remaining in viewer's memory, e.g., 32 sec.
- **Attention span (AS):** The most recent data occupying user's attention, e.g., 16 sec.
- **Measure the group, long-term coherence between M and AS**

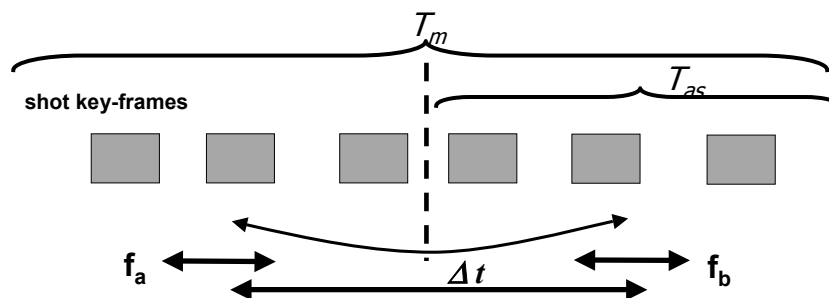
(Kender and Yeo 98, Sundaram and Chang 00)

10/2001

S.-F. Chang Columbia U.

31

Video Scene Segmentation

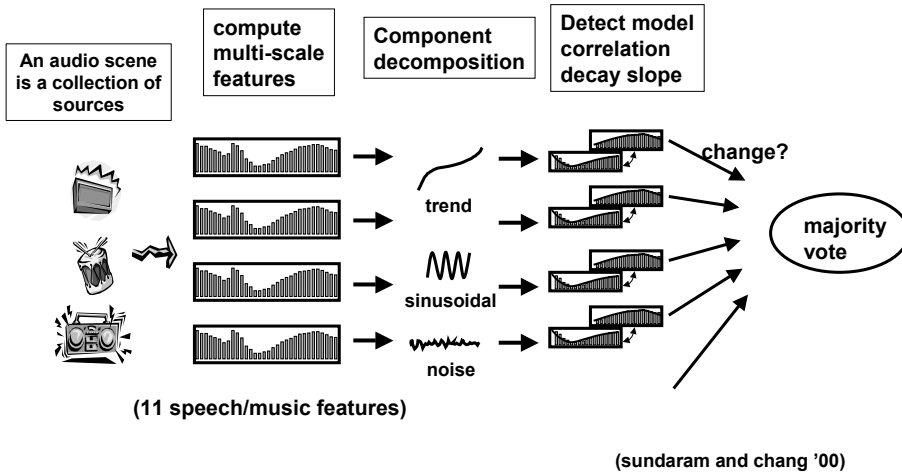


- **Group coherence based on viewer memory model**

$$R(a, b) = (1 - d(a, b)) \cdot f_a \cdot f_b \cdot (1 - \Delta t / T_m)$$

$$C(t_o) = \left(\sum_{a \in T_{as}} \sum_{b \in \{T_m \setminus T_{as}\}} R(a, b) \right) / C_{\max}(t_o)$$

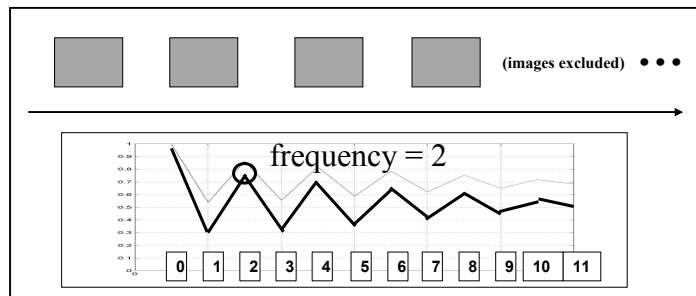
Audio Scene Segmentation



Detecting Dialog Structure

cyclic autocorrelation

$$\Delta(n) \equiv \frac{1}{N} \sum_{i=0}^{N-1} d(o_i, o_{\text{mod}(i+n, N)})$$



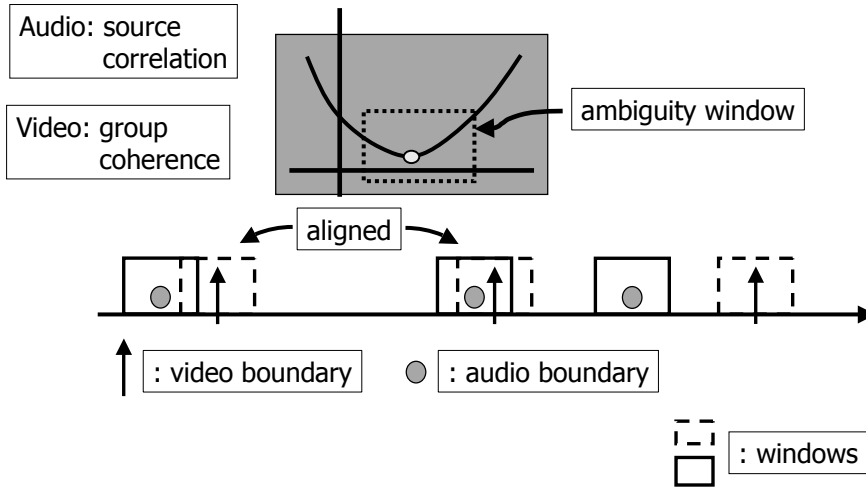
Accuracy: 80%/94% - 91%/100% precision recall

10/2001

S.-F. Chang Columbia U.

34

Aligning Scene Boundaries

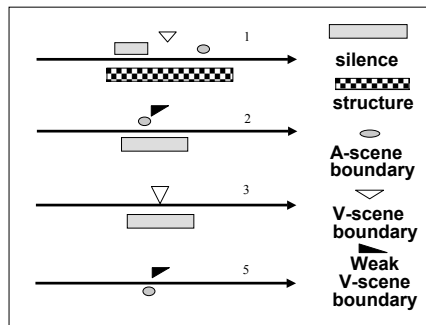


Fusion Rules

■ Synchronized a-scene and v-scene → scene

■ Exceptions:

- Ignore scenes within structures
- Ignore (weak v-scene, weak a-scene)
- Add (strong v-scene, silence)
- Require tighter synchronization if (weak v, normal a)



Open Issues

■ Video

■ Use higher-level cognitive models

- Currently focus on group coherence between groups
- Syntax structure within scenes not considered
- Potential use of information theory and temporal modeling

■ Audio

■ Source separation (music, speech, noise, etc)

■ Scene-level organization

■ Visualization, summarization, matching

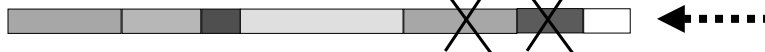
10/2001

S.-F. Chang Columbia U.

37

Scene Skimming

Reduce to a short-time preview



syntax
preserved

Find the skim with the highest utility,
satisfying the
production syntax constraints

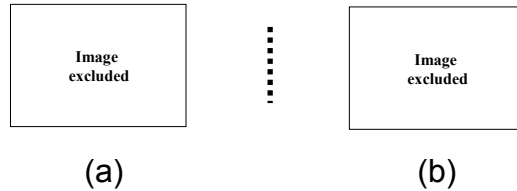
10/2001

S.-F. Chang Columbia U.

38

Utility of a shot

- Explore Perceptual Model and Scene Syntax



- Is comprehension time related to the visual spatio-temporal complexity of the shot ?
- The presence of detail robs a shot of its screen time.

10/2001

S.-F. Chang Columbia U.

39

Shot complexity and comprehension time

- Represent shot by its key-frame
- A shot is selected at random
- The subject was asked to *correctly* answer four questions in minimum time:
 - Who
 - What
 - When
 - Where

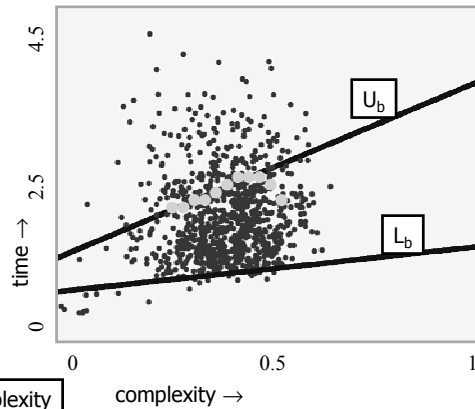
"Why" was not asked

Time-complexity relationship

- Plot of average time vs. complexity shows two bounds

$$U_b(c) = 2.40c + 1.11$$

$$L_b(c) = 0.61c + 0.68$$



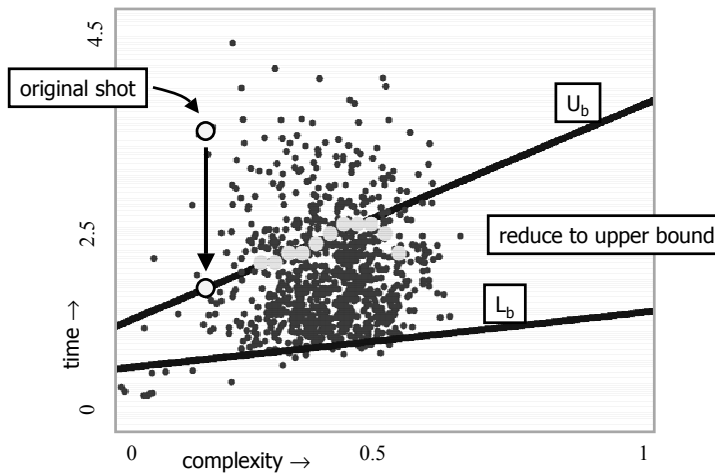
comprehension time increases with complexity

10/2001

S.-F. Chang Columbia U.

41

Shot condensation



10/2001

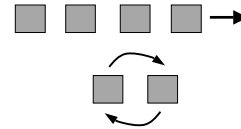
S.-F. Chang Columbia U.

42

Film syntax

The specific arrangement of shots so as to bring out their mutual relationship. [sharff 82].

- Minimum number of shots in a scene
- The particular ordering of the shots
- The specific duration of the shots, to direct viewer attention
- Changing the scale of the shots



Film makers think in terms of phrases of shots and not individual shots.

10/2001

S.-F. Chang Columbia U.

43

The progressive phrase

“Two well chosen shots will create expectations of the development of narrative; the third well-chosen shot will resolve those expectations.”
[sharff 82].

Hence, a phrase (a group of shots) must at least have three shots.



Maximal compression:
eliminate all the dark shots.

10/2001

S.-F. Chang Columbia U.

44

The dialog

"Depicting a conversation between m people requires $3m$ shots." [sharff 82].

Hence, a dialog must at least have six shots



Maximal compression:
eliminate all the dark shots.

10/2001

S.-F. Chang Columbia U.

45

Utility Optimization Framework

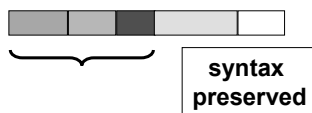
Objective function for each skim

$$O(\vec{t}, \vec{c}, \phi) \equiv 1 - U(\vec{t}, \vec{c}, \phi) + \gamma_0 R(\vec{t}, \phi) + \gamma_1 P(\phi)$$

Utility of Skim

Rhythm Penalty

Dropped shot penalty



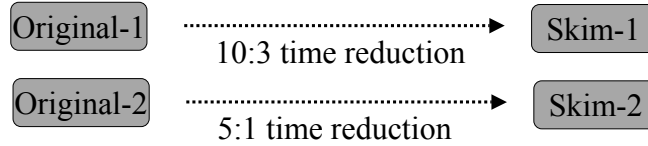
Find the skim with the minimal cost, satisfying the syntax constraints

10/2001

S.-F. Chang Columbia U.

46

Scene Skimming

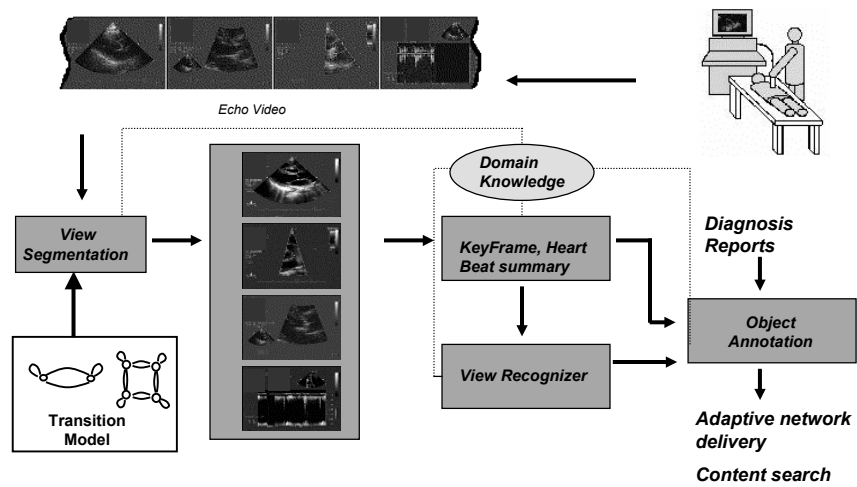


10/2001

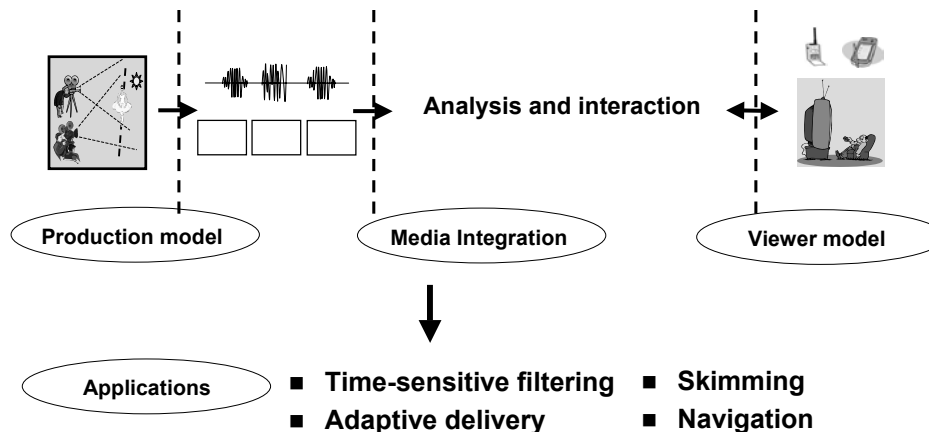
S.-F. Chang Columbia U.

47

Echo Video Digital Library & Remote Medicine



Conclusions



10/2001

S.-F. Chang Columbia U.

49

Acknowledgements

- **Live Sports Filtering:**
Di Zhong and Raj Kumar
- **Soccer Video Parsing:**
Lexing Xie, Peng Xu, Ajay Divakaran, Anthony Vetro, Huifang Sun
- **Scene Segmentation and Skimming:**
Hari Sundaram
- **Medical Video:**
Shahram Ebadollahi

10/2001

S.-F. Chang Columbia U.

50

More Information

- **Columbia Digital Video/Multimedia Group**

<http://www.ee.columbia.edu/dvmm>

- **ADVENT Industry/University Consortium**

<http://www.ee.columbia.edu/advent>

- <http://www.ee.columbia.edu/~sfchang>