# Multimedia Information Retrieval

http://morpheus.micc.unifi.it/learning/
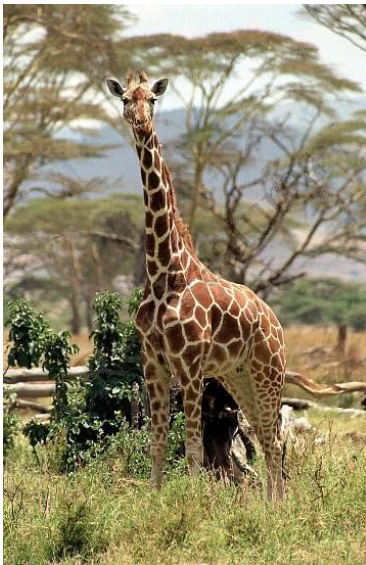
# State-of-the-art Multimedia Search Engines

☐ Work better for simple concepts,
   e.g. Two people kissing, A picture of a giraffe
☐ Don't work for complex queries
   e.g. A picture of a brick home with black shutters and
   white pillars, with a pickup truck in front of it (image)

# Examples

- Find the pictures of giraffe
  - Keyword: giraffe
  - http://images.google.it/images?svnum=10&hl=it&lr=&rls=GGGL%2CGGGL%3A2006-19%2CGGGL%3Ait&q=giraffe&btnG=Cerca

- A picture of a brick home with black shutters and white pillars, with a pickup truck in front of it (image)
  - brick home shutters
  - http://images.google.it/images?sourceid=navclient-ff&ie=UTF-8&rls=GGGL,GGGL:2006-19,GGGL:it&q=brick+home+shutters+

- # In Google Video try searching for "Bush" or "Bush speaking about Iraq"

# Why this happens?

- Most of these search engines are keyword based
  - "False" multi-media search engine
  - Have to represent your idea in keywords
  - These keywords are expected to appear in the filename, or corresponding webpage

- Therefore……
  - Unable to handle semantic meaning of images
  - Unable to handle visual position
  - Unable to handle time information
  - Unable to use images as query
  - ……….

*When I use a word," Humpty Dumpty said, in rather a scornful tone, "it means just what I choose it to mean—neither more nor less.*



- Try to search the image of the logo of Osama

# How Google does it?

- No image processing. Textual context!
  - In videos it uses closed captions and transcriptions
- File names, nearby words
- Distance from image to words
- "give me images with *flower* in the file name or near the image"

# Solution

- ……it would be great to have multimedia search engine intelligent enough to <u>associate its own keywords based on what's in the image</u>.

- Content-based information retrieval (CBIR)

- **Different from text IR:**
  - Structure of data is more complex. Efficiency is an issue
  - Using of metadata
  - Characteristics of multimedia data
  - Operations to be performed
- **Aspects:**
  - Data modeling: Extract and maintain the features of objects
  - Data retrieval: based not only on description but on content

# Retrieval process

- **Query specification**
  - ☐ fuzzy predicates: *similar to*
  - ☐ content predicates: *images containing an apple*
  - ☐ data type predicates: video, ...
- **Query processing and optimization**
  - ☐ Parsed, compiled, optimized for order of execution
  - ☐ Problem: many data types, different processing for each
- **Answer**
  - ☐ Relevance: similarity to query
- **Iteration**
  - ☐ Bad quality, so need to refine

- **Multimedia Information Retrieval is quite big in scope:**
- **Data examples:**
  - 2D/3D color/grayscale images: e.g., brain scans, scientific databases of vector fields
  - (2D) video,
  - (1D) voice/music; (1D) time series: e.g., financial/marketing time series; DNA/genomic databases
- **Query examples:**
  - find photographs with the same color distribution as this
  - find companies whose stock prices move as this one
  - find brain scans with a texture of a tumor
  - Find videos where something happens

# Some solutions

- Reduce the problem to search for multi-dimensional points (feature vectors, but vector space is not used)
- Define a distance measure
  - for time series: e.g., Euclidean distance between vectors
  - for images: e.g., color distribution (Euclidean distance); another approach: *mathematical morphology*
  - Other features as vectors
- Often, for search within distance, the vectors are organized in R(*/+)-trees or other spatial trees
- Clustering plays important role

# Query types

- **All within given distance**
  - Find all images that are within 0.05 distance from this one
- **Nearest-neighbor**
  - Find 5 stocks most similar to IBM
- **All pairs within given distance**
  - Further: clustering
- **Whole object vs. sub-pattern match**
  - Find parts of image that are...
  - E.g., in $512 \times 512$ brain scans, find pieces similar to the given $16 \times 16$ typical X-ray of a tumor
  - Like *passage retrieval* for text documents

# Open problems

- How similarity function can be defined?
- What features of images (video, sound) there are?
- How to better specify the importance of individual features? (*Give me similar houses*: similar = size? color? structure? Architectural style?)
- How to determine the objects in an image?
- Integration with DBMSs and SQL for fast access and rich semantics
  - Integration with XML
  - Ranking: by similarity, taking into account history, profile

# Open problems

- Object/event detection (computer vision and pattern recognition)
- Automatic feature selection
- Spatial indexing data structures (more than 1D)
- New types of data.
    - What features to select? How to determine them?
- Mixed-type data (e.g., web pages, or images with sound and description)
- What clustering/IR methods are better suited for what features? (What features for what methods?)
- Similar methods in data mining, ...
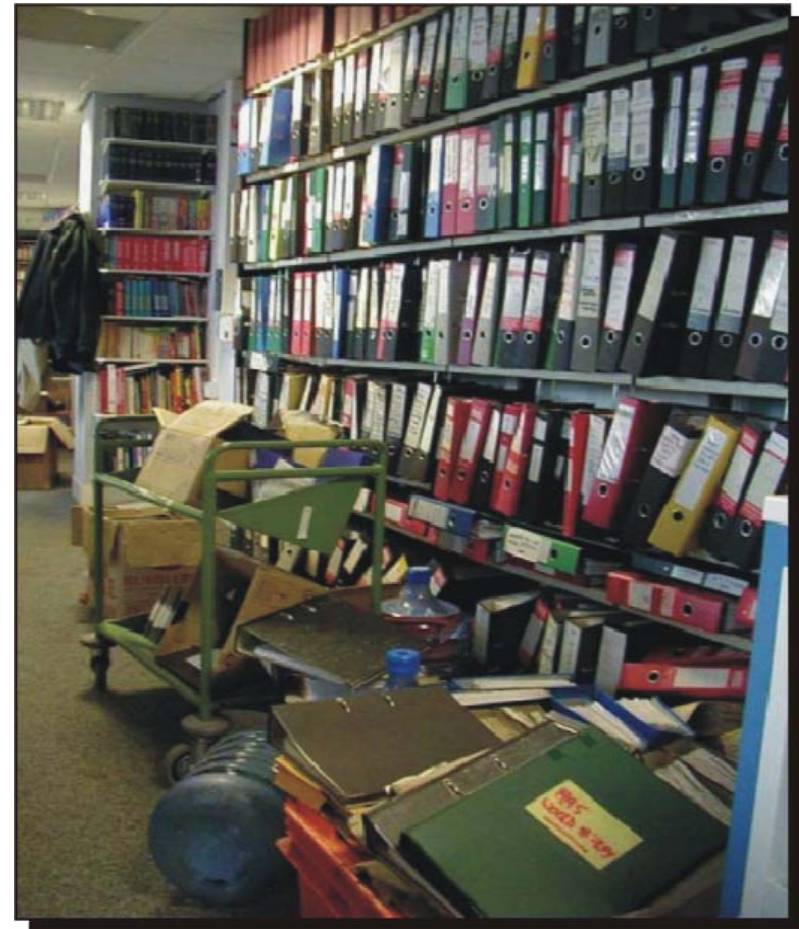
# Content-based Video Retrieval

- Application
- Implementation
- Experience from TREC video track
  - Feature Extraction Task (High-level Semantics Feature)
  - Manual Retrieval Task (One-run Retrieval)
  - Interactive Retrieval Task (Multiple-run with Feedback)
  - Results & Demo (CMU and IBM)
- Conclusion

# Application

- Increasing demand for visual information retrieval
  - Retrieve useful information from databases
  - Sharing and distributing video data through internet

- Example: BBC
  - BBC archive has +500k queries plus 1M new items … per year;
  - From the BBC …
    - Police car with blue light flashing
    - Government plan to improve reading standards
    - Two shot of  Kenneth Clarke and William Hague

# Application ( Cont. )

- Past project: ASSAVID in collaboration with BBC sports library:

- Develop automatic annotation systems for sports videos

# Application ( Cont. )

- Video Surveillance
  - Find where else the person appears
- Experience On-Demand
  - Help to remember previous events
- Provide useful information on traveling
  - Equipment on cars to retrieve useful multimedia information according to your location/preference
- ………
- Video content is plentiful … its now available digitally … we can work on it directly … so it follows

# Application ( Cont. )

# Application ( Cont. )

# Typical Retrieval Framework

- User : provide query information that represents his information needs
- Database: store a large collection of video data
- Goal: Find the most relevant shots from the database
  - Shots: "paragraph" in video, typically 20 – 40 seconds, which is the basic unit of video retrieval
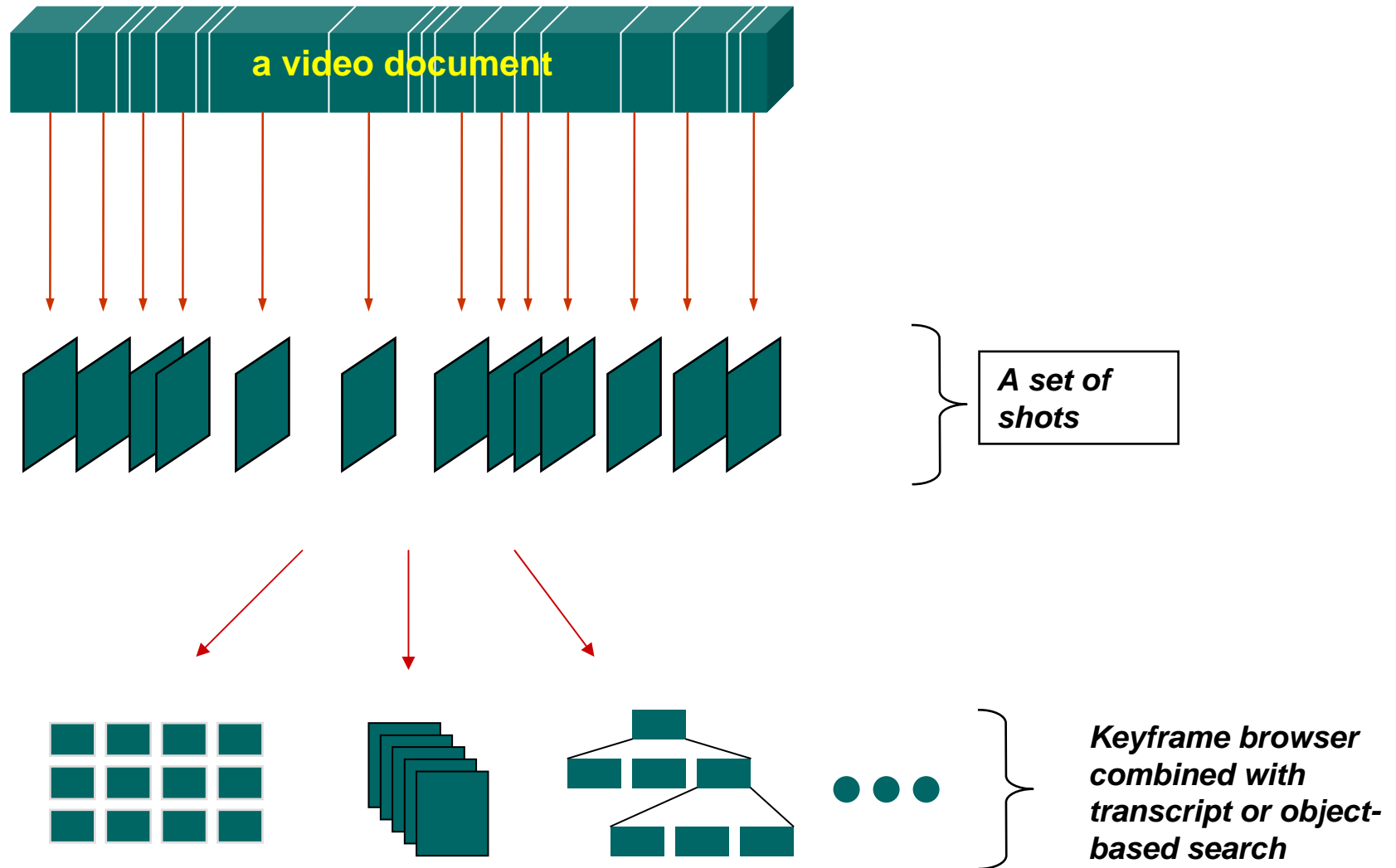
# Bridging the Gap

Video Database

User
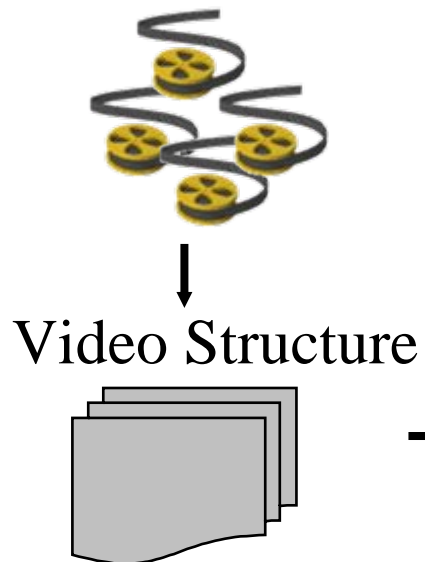
Result

# Automatically Structure Video Data

- The first step for video retrieval: Video "programmes" are structured into logical scenes, and physical shots
- If dealing with text, then the structure is obvious:
  - paragraph, section, topic, page, etc.

- All text-based indexing, retrieval, linking, etc. builds upon this structure;
- Automatic shot boundary detection and selection of representative keyframes is usually the first step;

# Typical automatic structuring of video
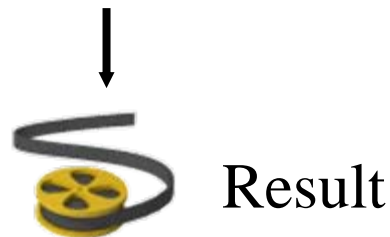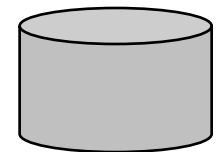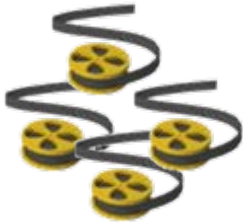
# Bridging the Gap

Video Database

User

Video Structure

Information Need

Result

# Ideal solution

Video Database

User

Video Structure

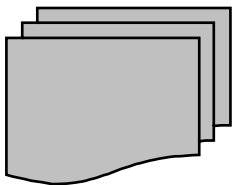Information Need
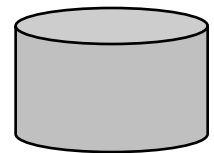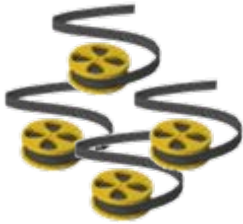
Understanding the semantic meaning and retrieve

Result

# Ideal solution

Video Database

Video Structure

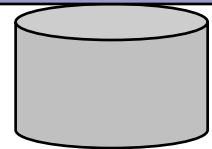Understanding semantic meaning and retrieve

However,
1. Hard to represent query in natural language and for computer to understand
2. Computers have no experience
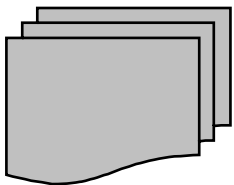3. Other representation restriction like position, time

Result

# Alternative Solution

Video Database

User

Video Structure

Provide evidence of relevant information ( text, image, audio)

Information Need

Match and combine

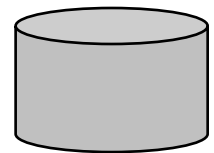Result

# Evidence-based Retrieval System

- General framework for current video retrieval system
- Video retrieval based on the evidence from both users and database, including
  - Text information
  - Image information
  - Motion information
  - Audio information
- Return a relevant score for each evidence
- Combination of the scores

# Keyword-based System

Video Database

User

Video Structure

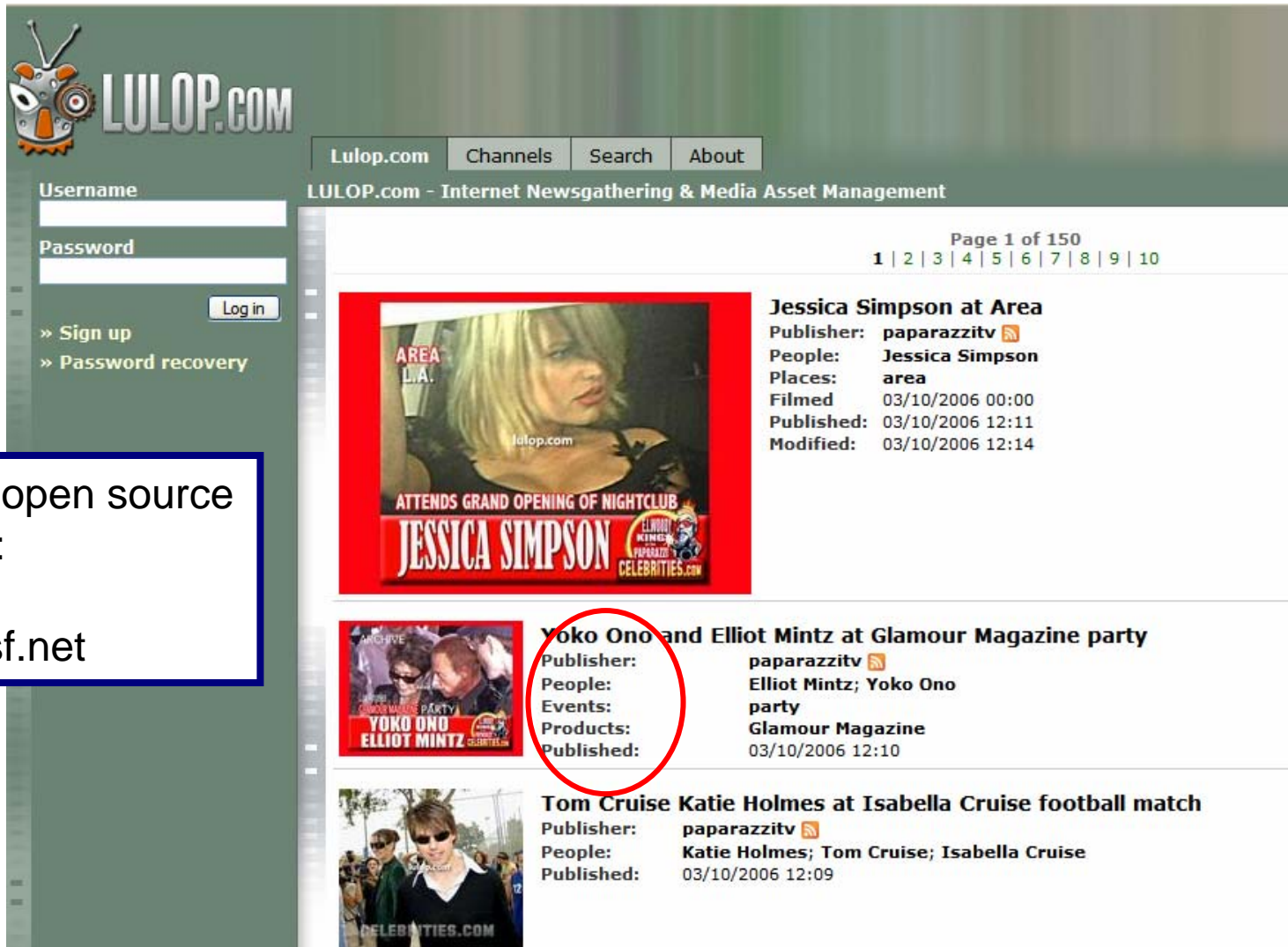Automatic Annotation ⟷ Keyword

Information Need

Including filename, video title, caption, related web page

# Keyword-based System

Video Database

User

Video Structure

| Automatic Annotation |
|---|

Keyword

Information Need

| Manual Annotation |
|---|

# Manual Annotation

- Manually creating annotation/keywords for image / video data
- Examples: Gettyimages.com (image retrieval)
- Pros:
  - Represent the semantic meaning of video
- Cons
  - Time-consuming, labor-intensive
  - Keyword is not enough to represent information need

# Manual annotation using metadata

Try the open source version:

lulop2.sf.net

# Tagging

Hello, world! This is a channel of Lulop.com built with the sole purpose of demonstrating the concept of Shot Tagging, which is indexing individual scenes and selctions within a video, in this particular case using a

Use tags on the left to select your vi
Play the video you have selected an
Your scene will be added to the othe

Tags.lulop.com is purely a demo pag

y as it is.

"The Color Purple" afterparty  Hot Ho

Lohan  Moms Day  Nicky's Mom's big

Valeria Marini  Venice Cinema Festival  W

"Derailed" premiere  "Lindsay Lohan is

"Match Point" Photocall  "Match Point" P

"Un posto al sole"  160 E 44th St, New Yo

pregnancy  407  60th Birthday of King Ca

Film Festival  Aaron  Aaron Carter  Acad

acapulco  ACE YOUNG  ad  Adam Brody

Adrian Brody  ADRIAN GRENIER  adrie

Restaurant  aircraft  Airport  Alain Favey
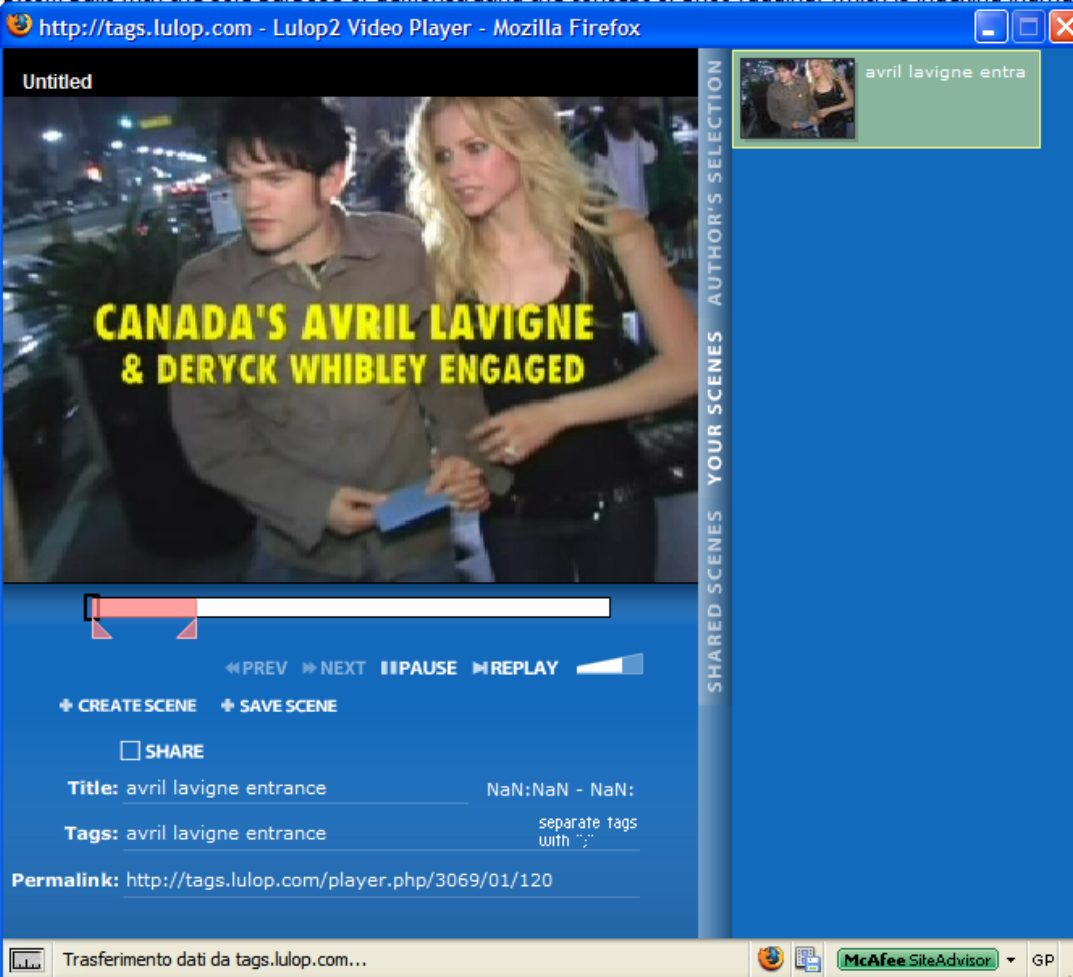
WHATEVER YOU BRING WE SING  Alcoholi

Alejandro Fernandez  Alessandro Botturi

Alessandro Zanni  ALISON MELNICK  Allis

NOLASCO  anastacia  André Heller  And

Andy Roddick  Ang Lee  Angela Bassett

l lavigne  bandana  bar  being le

Lohan fan  Brandor

brown  Bruce Willis  bruce willis a

uilera  Citro, Rallye di Montecarlo  cit

Cuthbert  finger  Firecrotch  foo  fo

ng story  John Travolta  kiss  leni  len

m Anderson  panties flash  paris hi

lo  Tag 1  Tag 2  Tag 3  test, avril or

ssi, intervista  wedding  white bag

# Manual annotation using taxonomy
(sort of…)

# Manual annotation using taxonomy



Taxonomy is exploited for retrieval

# Manual annotation using 4 Ws



Who

What

Where

When

Frameline uses
MPEG-47: i.e.
MPEG-4 with
MPEG-7 stream

# Speech and OCR transcription

Video Database

User

Video Structure

Annotation

Keyword

Information Need

Speech Transcription

OCR Transcription

# Query using speech/OCR information



**Query:**
Find pictures of Harry Hertz, Director of the National Quality Program, NIST



**Speech**:
We're looking for people that have a broad range of expertise that have business knowledge that have knowledge on quality management on quality improvement and in particular …

**OCR:**

```
H,arry Hertz a Director aro 7 wa-
,i,,ty Program
,Harry Hertz a Director
```

# Automatic face and OCR recognition

# Closed captions

# DVD subtitle ripping and OCR

# What we lack?

# Image-based Retrieval

Video Database

User

Video Structure

Text Information

Keyword

Information Need

Image Feature

Query Images

# Global Low-level Image Feature

- **Color-based Feature**
  - ☐ Color Histogram
  - ☐ Color Percentage
  - ☐ Color Correlogram
  - ☐ Color Moments
- **Texture-based Feature**
  - ☐ Gabor Filter
  - ☐ Wavelet
- **Shape/Structure Feature**

# Regional Low-level Image Feature

- Segmentation into objects – hard problem !
- Extract low-level features from each regions

# Image Search

- **Feature Representation**
  - Image: represented as a series of real number, or a vector of features, *(f1, …., fn)*
  - Distance Function: The distance between two vectors, typically Euclidean Distance
  - Probably "Nearest is relevant"
    - The nearest images in the database is relevant to the query images.

# Finding Similar Images



botanic1     CnScenery9     botanic3     raffles1

exar     CnScenery7     shangrila     foothills

# But…..

- Low-level feature doesn't work in all the cases

# Find similar objects

# High-level Image Feature

- Objects: Persons, Roads, Cars, Skies…
- Scenes: Indoors, Outdoors, Cityscape, Landscape, Water, Office, Factory…
- Event: Parade, Explosion, Picnic, Playing Soccer…
- Generated from low-level features

# Image-based Retrieval

Video Database

User

Video Structure

Text Information ←→ Keyword

Information Need

Image Feature

Low-level Feature

High-level Feature

Query Images

# More Evidence in Video Retrieval

Video Database

User

Video Structure

| Text Information | ↔ | Keyword |

| Image Information | ↔ | Query Images |

| Motion Information | ↔ | Motion |

| Audio Information | ↔ | Audio |

Information Need

# Combination of multi-modal results

- Difference characteristics between multi-modal information
  - □ Text-based Information: better for middle and high level queries
    - e.g. Find the video clip of dancing women wearing dresses
  - □ Image-based Information: better for low and middle level queries
    - e.g. Find the video clip of green trees
- Combination of multi-modal information

# Other Useful Technique

- Query Expansion

- Cross-Modal Relation

- Relevance Feedback

# Recap

- Video Retrieval is to bridge the gap between user information need and video database
- Multi-modal evidence
    - Text-based (most popular)
    - Image-based
    - Motion-based
    - Audio-based
- Combination of the evidence

# Introduction to TREC Video Retrieval Track

- Full Name: Text REtrieval Conference
- TREC Video Track web site: http://www-nlpir.nist.gov/projects/trecvid/
- TREC series sponsored by the National Institute of Standards and Technology (NIST) with additional support from other U.S. government agencies
  - Goal is to encourage research in information retrieval

# Introduction to TREC Video Retrieval Track

- **Video Retrieval Track started in 2001**
  - □ Goal is investigation of content-based retrieval from digital video
  - □ Focus on the *shot* as the unit of information retrieval rather than the scene or story/segment/clip
- **Current state-of-the-art Video Retrieval Competition**
  - □ 17 active participants, including groups from CMU, IBM Research, Microsoft Research Asia, MediaMill, LIMSI, Dublin City University.

# Main tasks in TREC

- Shot boundary detection
- Semantic Feature Extraction Task
- Video Retrieval Task
  - Manual Retrieval: Human formulate a query and then automatically retrieve from collection
  - Interactive Retrieval: Full human access and feedback

# Where are they?

Video Database

Retrieval Task

User

Text Information

Keyword

Video Structure

Information Need

Shot Boundary Detection

Image Feature

Query Images

Low-level Feature

High-level Feature

Feature Extraction

# Video Data

- Difficult to get video data for use in TREC because ©

- Used mainly Internet Archive
  - advertising, educational, industrial, amateur films 1930-1970
  - produced by corporations, non-profit organisations, trade groups, etc.
  - Noisy, strange color, but real archive data
  - 73.3 hours partitioned as follows:

# Shot Boundary Detection

- Fundamental primitive of most/all work in content-based video retrieval

a video document

A set of video shots

# Feature Extraction

- Extracted high-level semantic feature from video
- Assign a video clip to one or more of several categories of video



High-level features: Cityscape, Lake, Trees, Water, Sky

# Feature Extraction

- Interesting itself but when it serves to help video navigation and search then its importance increases

- Benefits:

  - Retrieval - Find video from a particular class

  - Filtering - Remove irrelevant and distracting information categories from summaries and visualizations

# The Features

- Face
  - Clip contains at least one human face with the nose, mouth, and both eyes visible. Pictures of a face meeting the above conditions count
- People
  - Clip contains a group of two more humans, each of which is at least partially visible and is recognizable as a human
- On-screen Text
  - Clip contains superimposed text large enough to be read

# The Features

Indoor

    Clip contains a recognizably indoor location, i.e., inside a building

Outdoor

    Clip contains a recognizably outdoor location, i.e., one outside of
        buildings

Cityscape

    Clip contains a recognizably city/urban/suburban setting

Landscape

    Clip contains a predominantly natural inland setting, i.e., one with little
        or no evidence of development by humans. Scenes with bodies of
        water that are clearly inland may be included

# Non-Video (Audio) Features

Speech

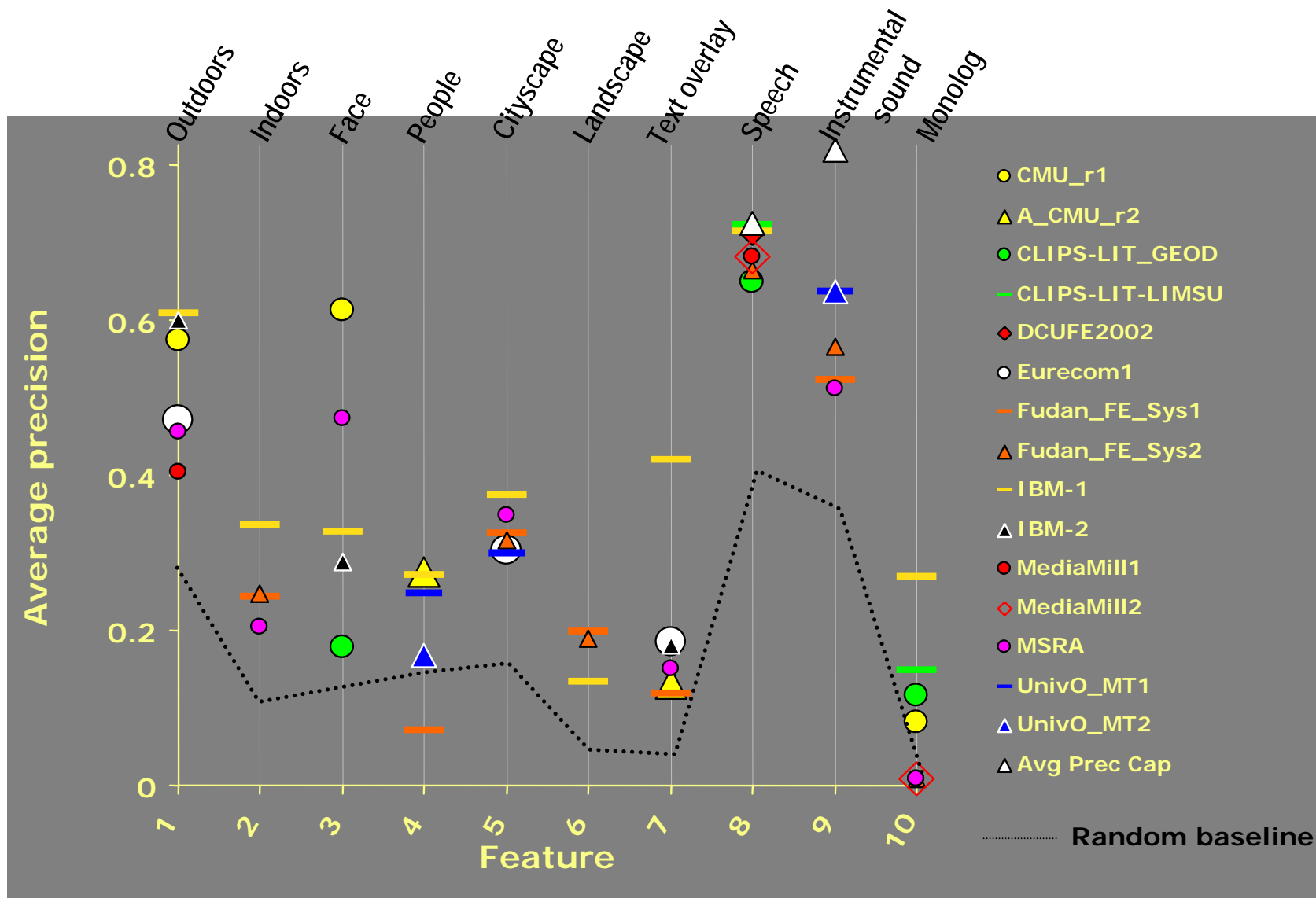A human voice uttering words is recognizable as such in this segment

Instrumental Sound

Sound produced by one or more musical instruments is recognizable as such in this segment
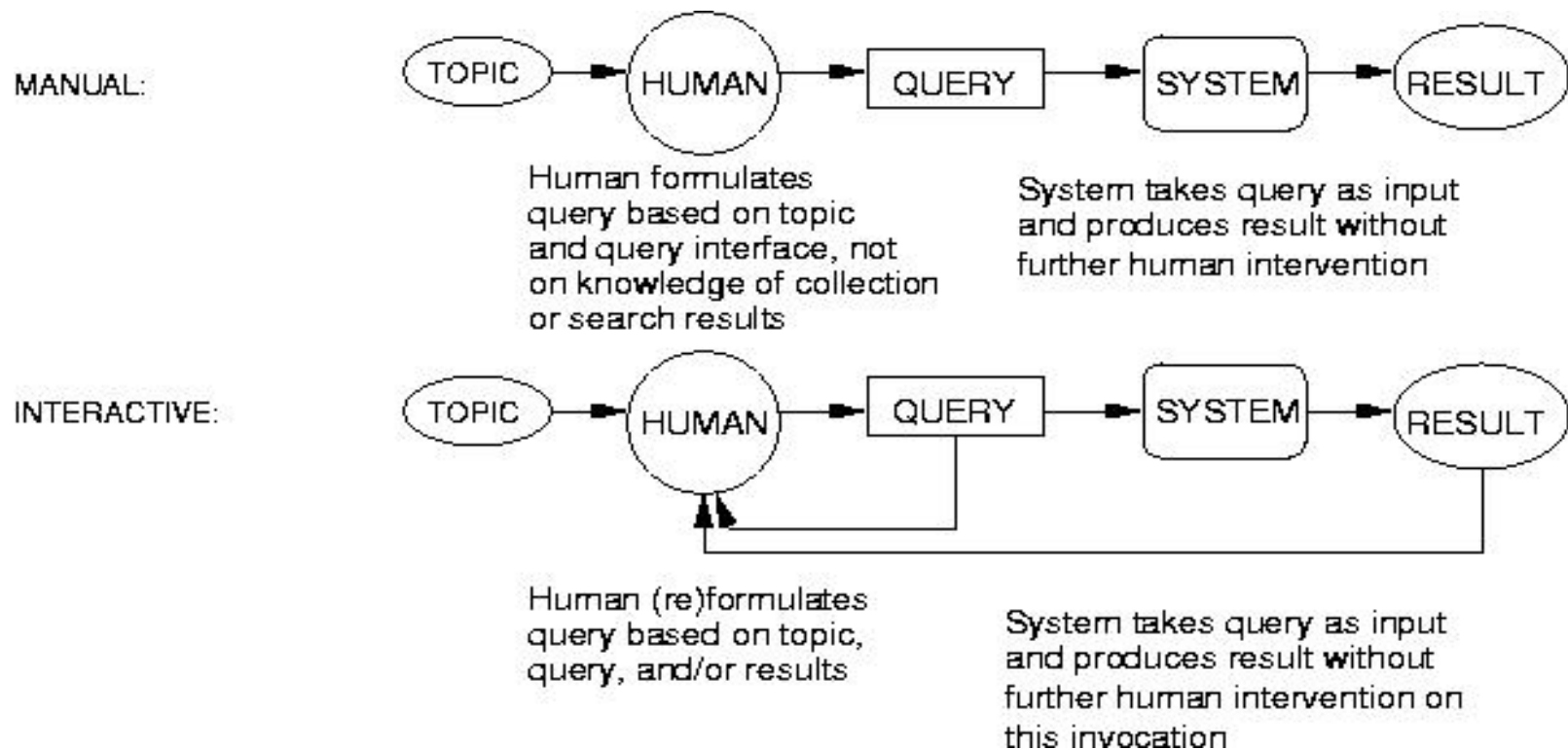
Monologues

Segment contains an event in which a single person is at least partially visible and speaks for a long time without interruption by another speaker. Pauses are ok if short

# TREC02 Results

# Video Search Task

- The most important task and final goal
- Manual & Interactive Search Task



MANUAL:

TOPIC → HUMAN → QUERY → SYSTEM → RESULT
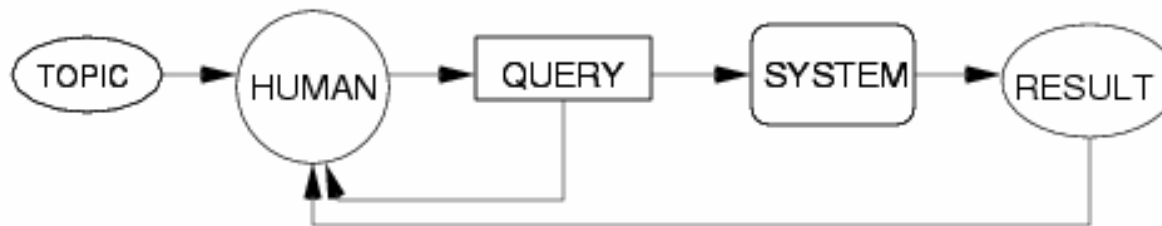
Human formulates query based on topic and query interface, not on knowledge of collection or search results

System takes query as input and produces result without further human intervention

INTERACTIVE:

TOPIC → HUMAN → QUERY → SYSTEM → RESULT

Human (re)formulates query based on topic, query, and/or results

System takes query as input and produces result without further human intervention on this invocation

# Queries for 2002 TREC Video Track

- ## Specific item or person
  - □ Eddie Rickenbacker, James Chandler, George Washington, Golden Gate Bridge, Price Tower in Bartlesville, OK

- ## Specific fact
  - □ Arch in Washington Square Park in NYC, map of continental US

- ## Instances of a category
  - □ football players, overhead views of cities, one or more women standing in long dresses

- ## Instances of events/activities
  - □ people spending leisure time at the beach, one or more musicians with audible music, crowd walking in an urban environment, locomotive approaching the viewer

# TRECVid 2005 search topics



TRECVID 2005 INTERACTIVE VIDEO RETRIEVAL RESULTS

# TRECvid 2006



FULLY AUTOMATIC:

TOPIC → SYSTEM → RESULT

System takes query as input and produces result without any human intervention

MANUALLY–ASSISTED:

TOPIC → HUMAN → QUERY → SYSTEM → RESULT

Human formulates query based on topic and query interface, not on knowledge of collection or search results

System takes query as input and produces result without further human intervention

INTERACTIVE:

TOPIC → HUMAN → QUERY → SYSTEM → RESULT

Human (re)formulates query based on topic, query, and/or results

System takes query as input and produces result without further human intervention on this invocation

# Queries for 2006 TREC Video Track

- Example types of video needs

I'm interested in video material containing:

- □ a specific person
- □ one or more instances of a category of people
- □ a specific thing
- □ one or more instances of a category of things
- □ a specific event/activity
- □ one or more instances of a category of events/activities
- □ a specific location
- □ one or more instances of a category of locations
- □ combinations of the above

- Topics may target commercials as well as news content.

# Some TRECVid 2006 high level features

- **Sports**: Shots depicting any sport in action
- **Entertainment**: Shots depicting any entertainment segment in action
- **Weather**: Shots depicting any weather related news or bulletin
- **Court**: Shots of the interior of a court-room location
- **Office**: Shots of the interior of an office setting
- **Meeting**: Shots of a Meeting taking place indoors
- **Studio**: Shots of the studio setting including anchors, interviews and all events that happen in a news room
- **Outdoor**: Shots of Outdoor locations
- **Building**: Shots of an exterior of a building
- **Desert**: Shots with the desert in the background
- **Vegetation**: Shots depicting natural or artificial greenery, vegetation woods, etc.
- **Mountain**: Shots depicting a mountain or mountain range with the slopes visible
- **Road**: Shots depicting a road
- …

# Sample Query

- ■ XML Representation

<!DOCTYPE videoTopic SYSTEM "videoTopics.dtd">
<videoTopic num="077">

<textDescription text="Find pictures of George Washington" />

<imageExample
src="http://www.cia.gov/csi/monograph/firstln/955pres2.gif"
desc="face" />

<videoExample src="01681.mpg"  start="09m25.938s"
stop="09m29.308s" desc="face" />

</videoTopic>

# Evaluation Metric

- Goal:  Maximize the Mean Average Precision
  - Result set limited to 100 shots
  - Precision = (# relevant shots retrieved)/(total # shots retrieved)
  - Average precision:  compute precision after each retrieved relevant shot and then average these precisions over the total number of retrieved relevant shots in the collection for that topic
  - Submitting the maximum number of shots per result set can never lower the average precision for that submission
  - Mean Average Precision = average of the average precision measures for each topic

# Demo

- CMU Interactive Search System
- IBM Video Retrieval System
  http://mp7.watson.ibm.com/marvel/
- UvA MediaMill

# CMU Manual Retrieval System

Query → Text, Image

Text → Movie Info, Text Score

Image → Image Score, PRF Score

Retrieval Agents

Movie Info, Text Score, Image Score, PRF Score → Final Score

# Snapshot of the CMU system (2002)

# Snapshot of the CMU system (2002)

# CMU Filter Interface for using Image Features (2002)

# Mean AvgP vs. mean elapsed time



Wide variation in elapsed time.
Not the dominant factor in effectiveness

Mean average precision

Mean elapsed time (mins.)

# IBM Marvel system



Search using color similarity

# IBM Marvel: search using semantics from previous results
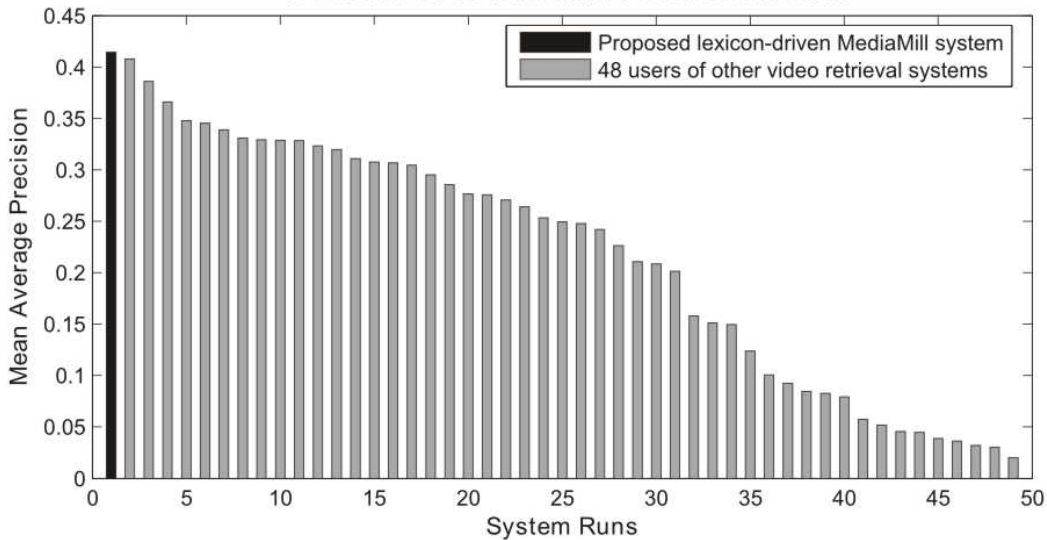
# UvA MediaMill

TRECVID 2004 Interactive Search Results



32 concept detectors

Interactive video search

TRECVID 2005 Interactive Search Results



101 concept detectors

# UvA MediaMill – cross browser

# Conclusion

- The goal of content-based video retrieval is to build more intelligent video retrieval engine via semantic meaning
- Many applications in daily life
- Combine evidence from different aspects
- Hot research topic, few business system
  - Check Techcrunch.com for info on business ventures
- State-of-the-art performance is still unacceptable for normal users, space to improve

# Virage: a business that has survived

# Reply.it – Multimedia asset management

# Credits

- Rong Yan – Carnegie Mellon
- Alexander Gelbukh