


**Concept-Based Video Retrieval**  
Cees Snoek and Marcel Worring

with contributions by:  
many

Intelligent Systems Lab Amsterdam,  
University of Amsterdam, The Netherlands



3

**The science of labeling**

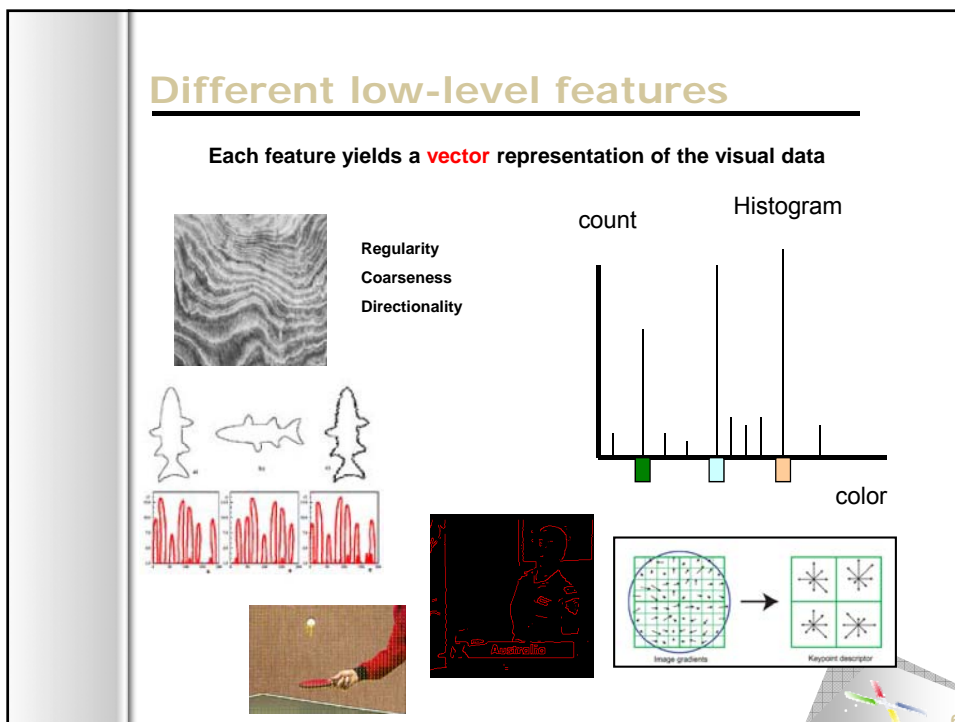
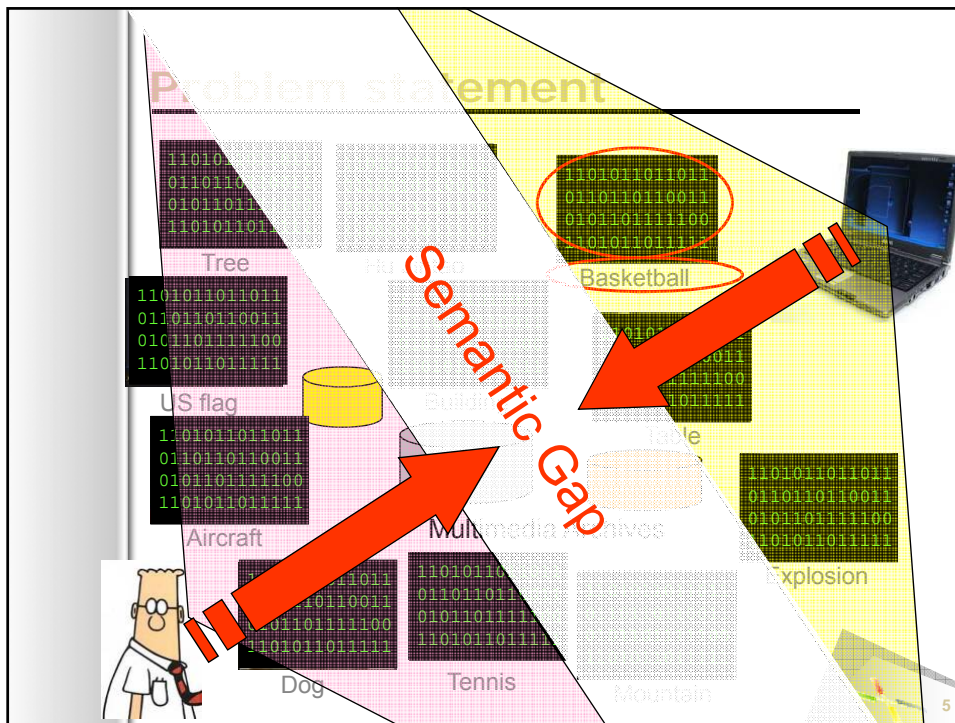
➤ To understand anything in science, things have to have a name that is recognized and is universal

 naming 'categories'	 naming chemical elements	 naming human genome
 naming living organisms	 naming rocks and minerals	 naming textual information

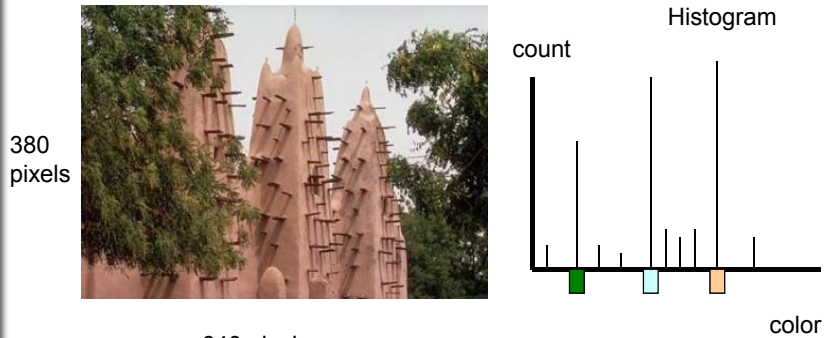
**What about naming video information?**



4



### Basic example: color histogram



380 pixels

640 pixels

Total 243200 pixels

count

Histogram

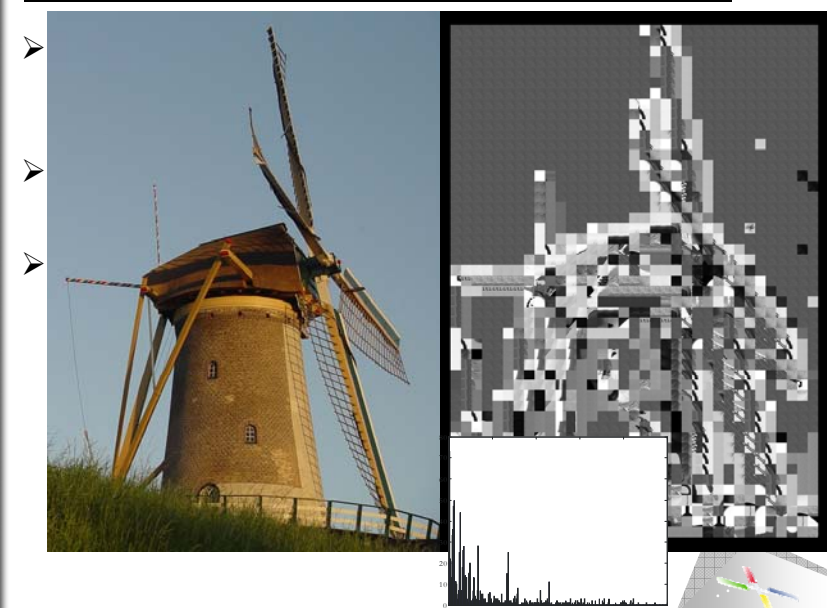
color

Histogram is a summary of the data summarizing in this case color characteristics

7

### Advanced example: codebook model

Leung and Malik. IJCV, 2001.  
Sivic and Zisserman. ICCV, 2003.  
van Gemert, PhD thesis, UvA, 2008.



7

## The goal: semantic video indexing

- Is the process of automatically detecting the presence of a semantic concept in a video stream

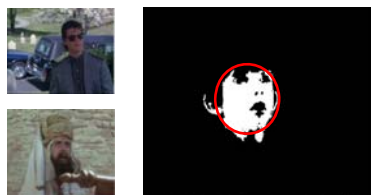


Airplane

9

## Semantic indexing

- The computer vision approach
  - ✓ Building detectors one-at-the-time



A face detector for frontal faces

3 years later



A face detector for non-frontal faces

One (or more) PhD for every new concept

10

### So how about these?

Animal Building Road Beach Boat

Graphic People Car Vegetation Overlaid Text

Studio Setting Outdoor

And the > 1000 others .....

11

### Generic concept detection in a nutshell

Labeled examples

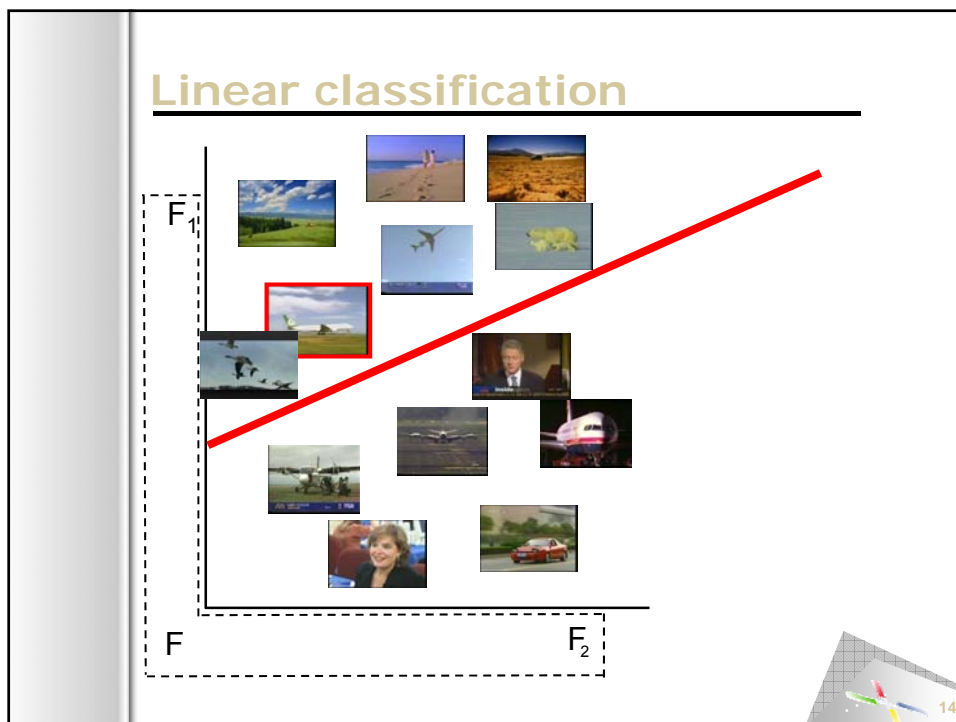
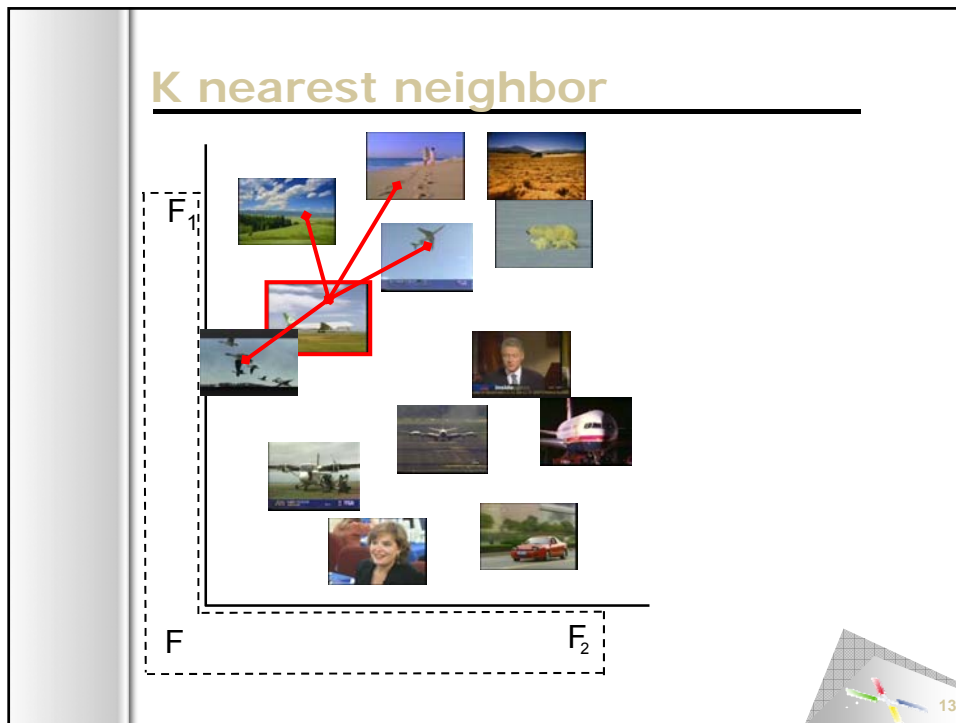
Video

Training

Testing

It is an outdoorcraft probability 0.95

12



## Support vector machine

**SVM usually is a good choice**

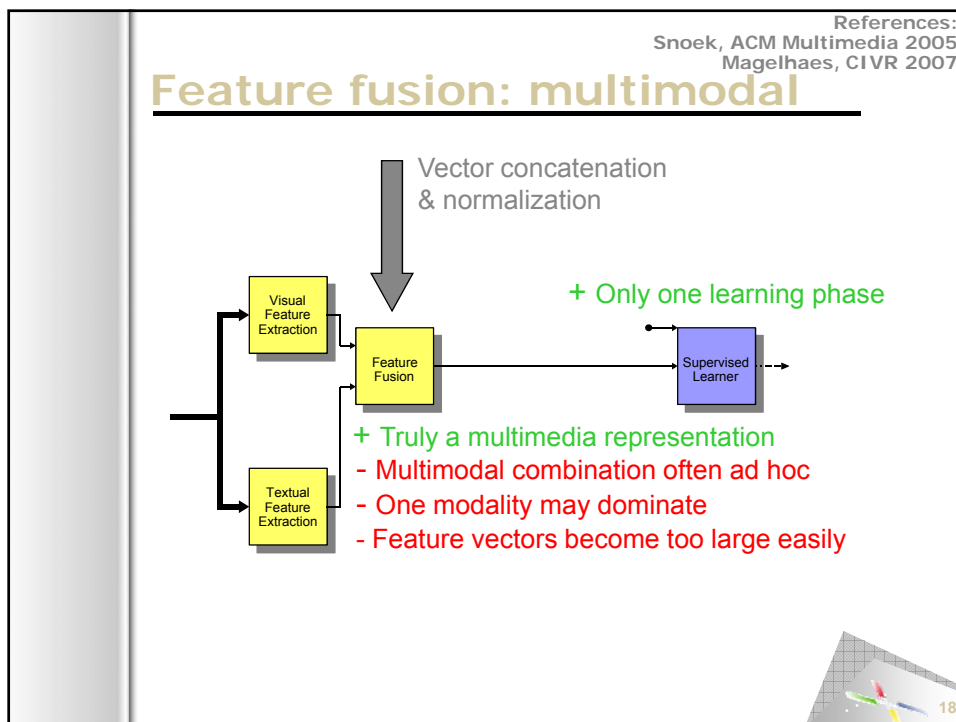
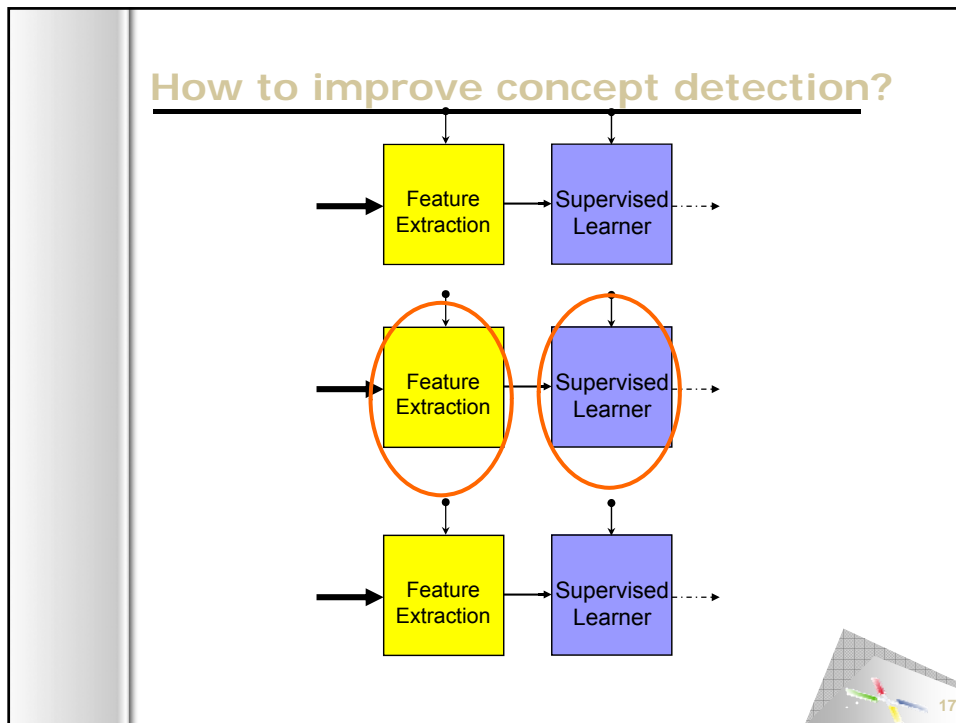
15

## Supervised Learner

- Support Vector Machine
  - ✓ Learns from provided examples
  - ✓ Maximizes margin between two classes
- Depends on many parameters
  - ✓ Select best of multiple parameter combinations
  - ✓ Using cross validation

$$p(\omega|\vec{s}) = \frac{1}{1 + \exp(\alpha\gamma(\vec{s}) + \beta)}$$

16





References:  
van de Sande, CIVR 2008

## Feature fusion: unimodal

- + Codebook model reduces dimensionality
- Combination still ad hoc
- One feature may dominate

References:  
Wu, ACM Multimedia 2004  
Snoek, ACM Multimedia 2005

## Classifier fusion: multimodal

- + Focus on modality strength
- + Fusion in semantic space

- Expensive in terms of learning effort
- Possible loss of feature space correlation

References:  
 Snoek, TRECVID 2006  
 Wang, ACM MIR 2007

## Classifier fusion: unimodal

+ Aggregation functions reduce learning effort  
 + Offers opportunity to use all available examples efficiently  
 - Linear function likely to be sub-optimal

21

References: IBM 2003  
 Naphade and Huang, TMM 3(1) 2001

## Modeling relations

- Exploitation of conceptual co-occurrence
  - ✓ Concepts do not occur in vacuum
  - ✓ In contrast, they are related

- What is sports?
  - ✓ Answer: a combination of various individual sports

22

References: IBM 2003  
 Qi, ACM Multimedia 2007  
 Liu, IEEE TMM 2008

## Modeling relations

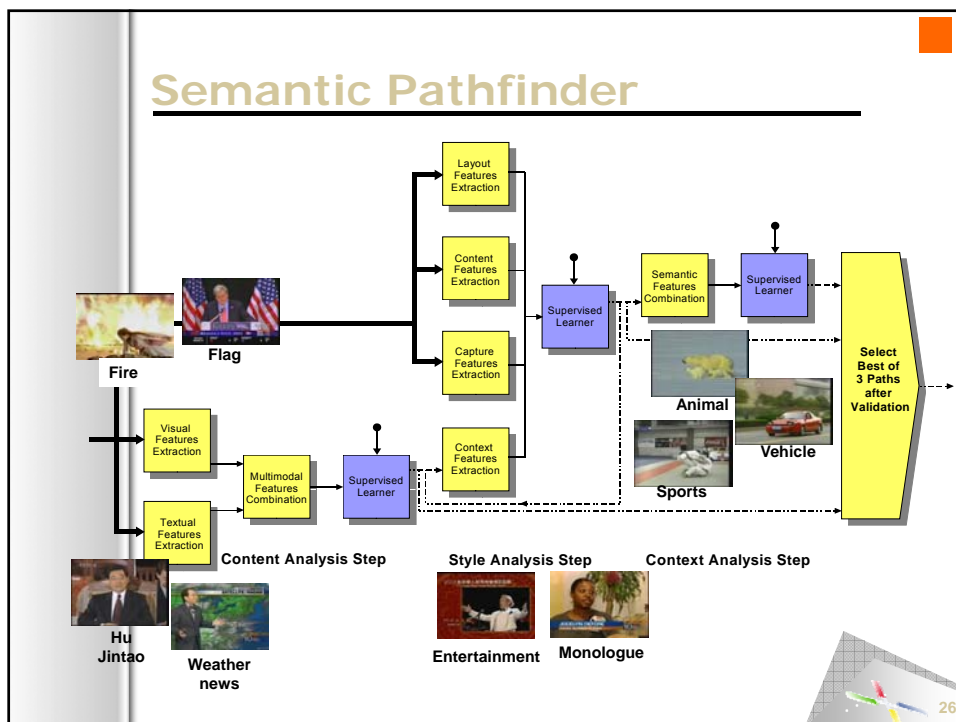
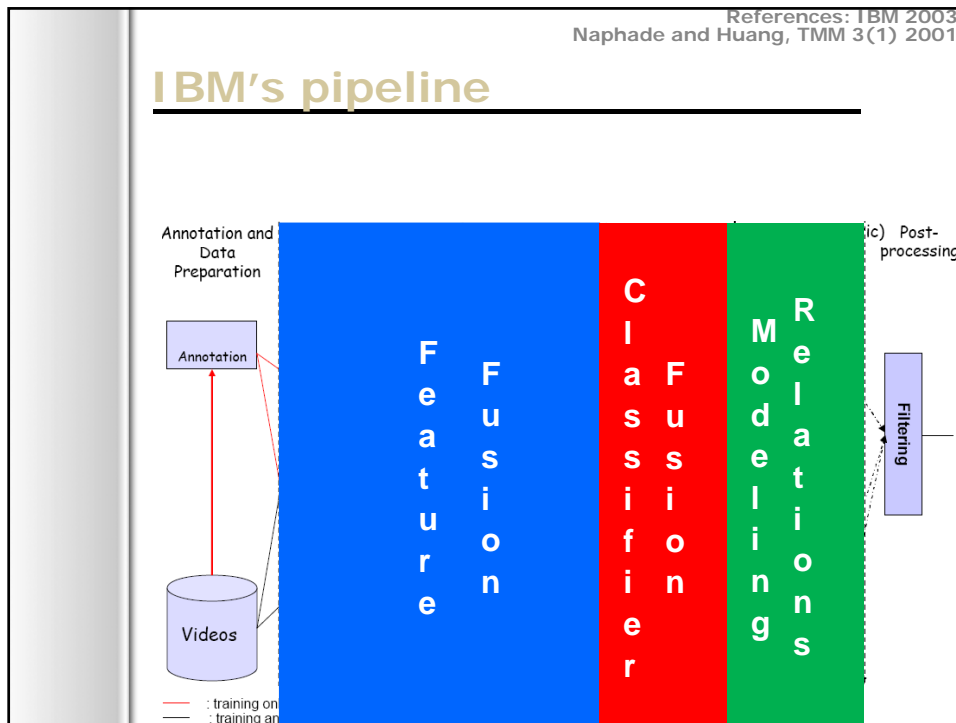
- Learning co-occurrence
  - ✓ Explicitly model relations: using graphical models
    - Computationally complex
    - Limited scalability
  - ✓ Implicitly learn relations: using SVM, or data mining tools
    - Assumes classifier learns relations
    - Suffers from error propagation

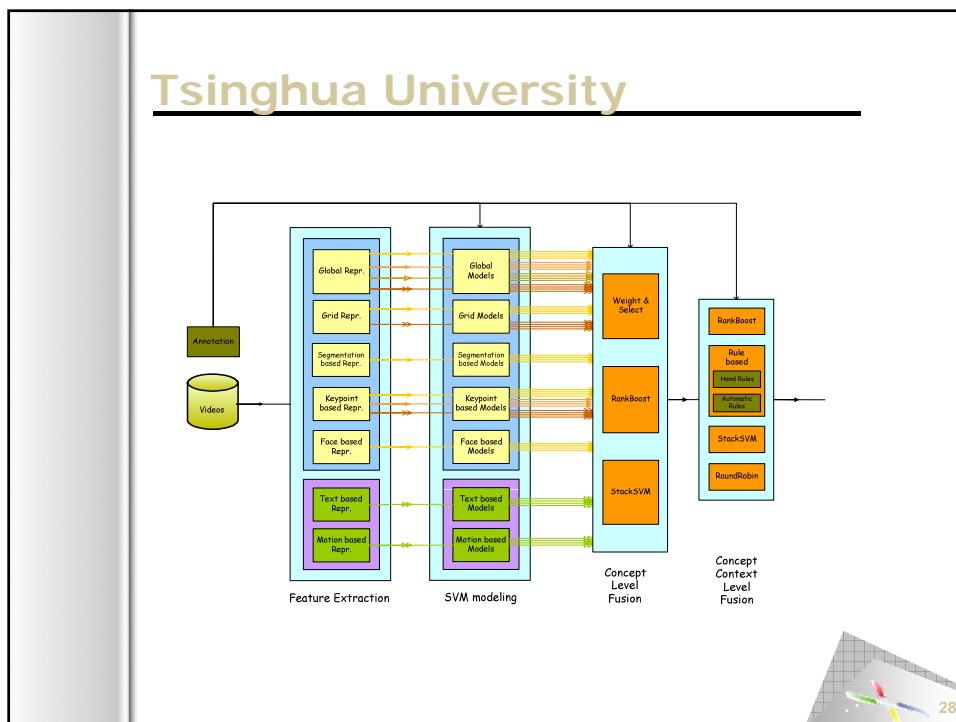
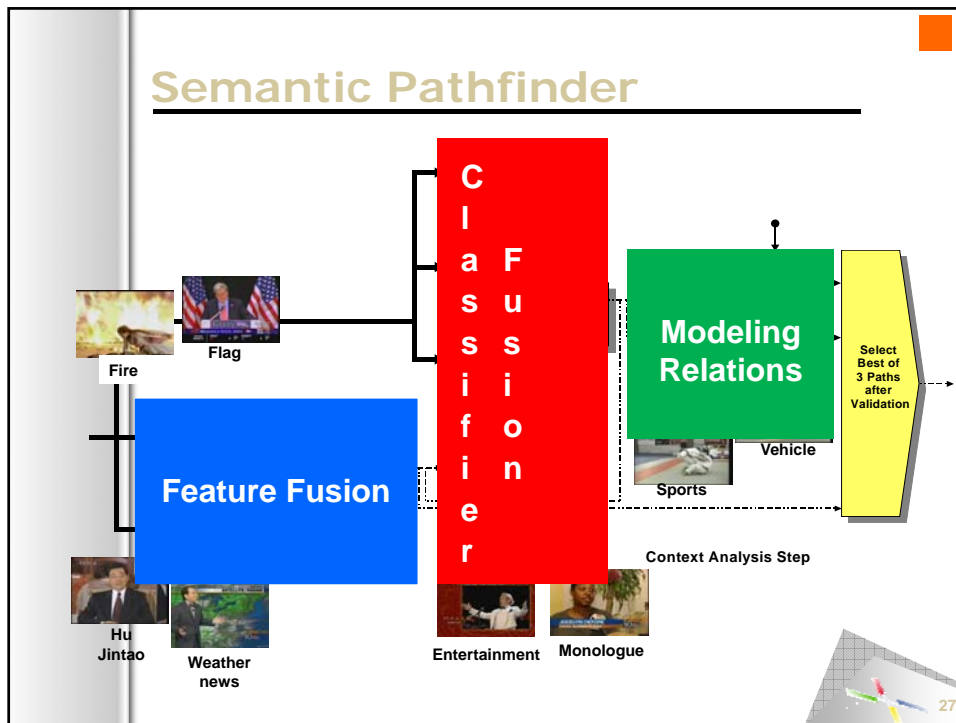
23

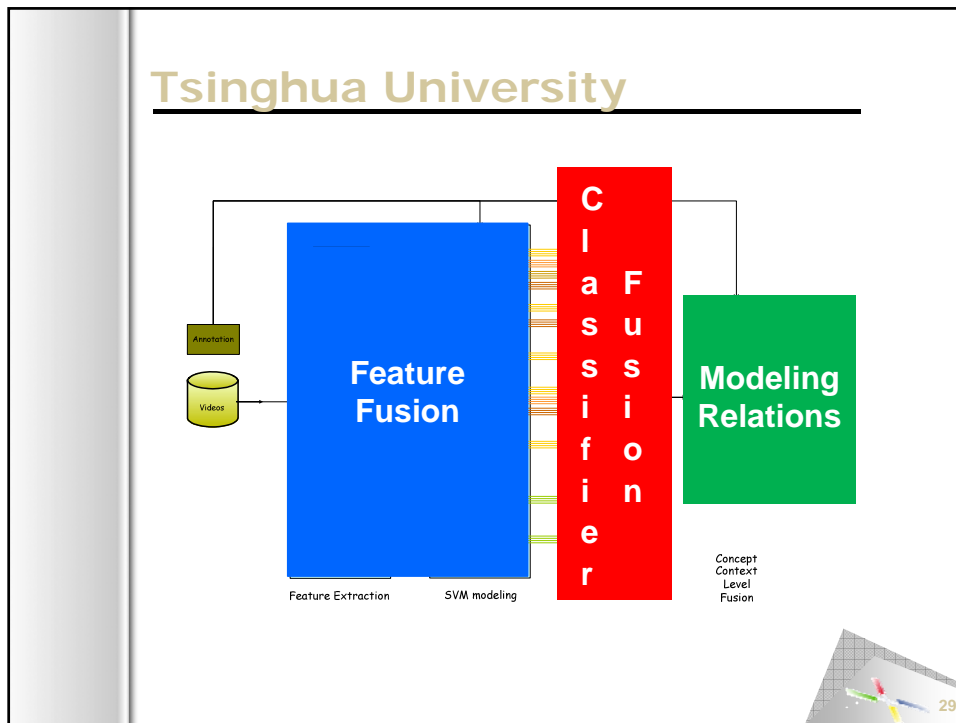
References: IBM 2003  
 Naphade and Huang, TMM 3(1) 2001

## IBM's pipeline

— : training only  
— : training and testing







### Fragmented research efforts...

Video analysis researchers

- ✓ Until 2001 everybody defined her or his own concepts
- ✓ Using **specific** and **small** data sets
- ✓ Hard to compare methodologies

Since 2001 worldwide evaluation by NIST





NIST




30

## NIST TRECVID benchmark

anno 2001

- Benchmark objectives
  - ✓ Promote progress in video retrieval research
  - ✓ Provide common dataset (shots, recognized speech, key frames)
  - ✓ Use open, metrics-based evaluation
  
- Large international field of participants

✓ and the 70 others...

- Currently the **de facto standard** for evaluation

<http://trecvid.nist.gov/>


## TRECVID Evolution: data, tasks, participants,...

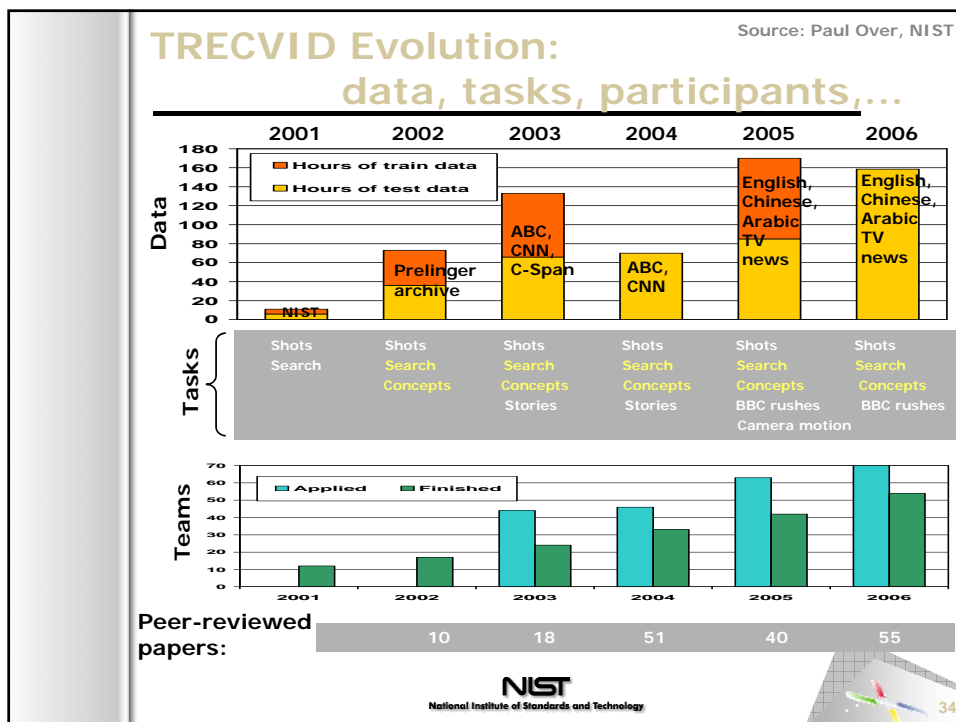
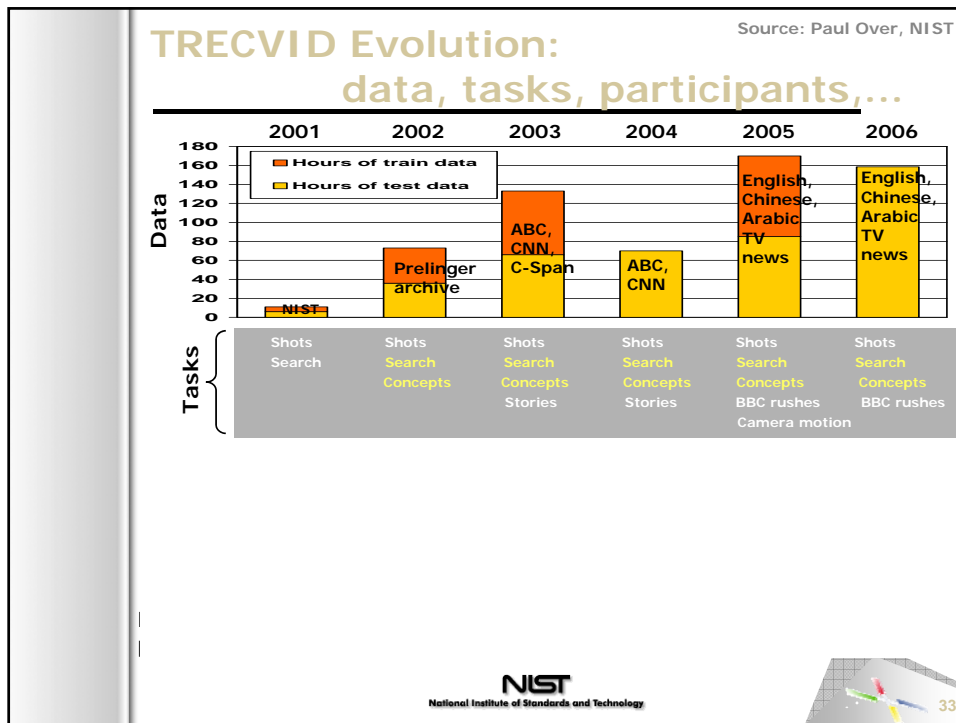
Source: Paul Over, NIST

Year	Hours of train data	Hours of test data
2001	~10	~10
2002	~70	~40
2003	~130	~50
2004	~70	~60
2005	~170	~100
2006	~170	~100

■ Hours of train data  
■ Hours of test data

2001: NIST  
 2002: Prelinger archive  
 2003: ABC, CNN, C-Span  
 2004: ABC, CNN  
 2005: English, Chinese, Arabic TV news  
 2006: English, Chinese, Arabic TV news

  
 National Institute of Standards and Technology









## Concept detection task

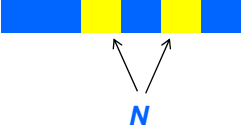
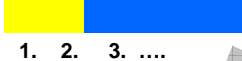
➤ Given:

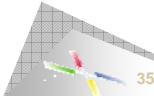
- ✓ a video dataset segmented into set of  $S$  unique shots
- ✓ set of  $N$  semantic concept definitions:

➤ Task:

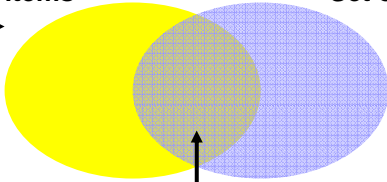
- ✓ How well can you detect the concepts?
- ✓ Rank  $S$  based on presence of concept from  $N$

$S$    $\Rightarrow$  


35




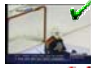

## Measuring uncertainty

Set of relevant items
Set of retrieved items

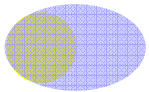


Set of relevant retrieved items

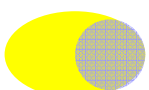
Results

1.  ✓
2.  ✗
3.  ✓
4.  ✓
5.  ✗

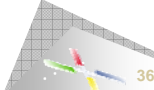
➤ Precision



➤ Recall



inverse relationship



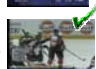



36

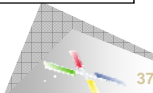
## TRECVID evaluation measures

- Classification procedure
  - ✓ Training: many hours of (partly) annotated video
  - ✓ Testing: many hours of **unseen** video
- Evaluation measure: **Average Precision**
  - ✓ Combines precision and recall
  - ✓ Averages precision after every relevant shot
  - ✓ Top of the ranked list most important

$$AP = \frac{1/1 + 2/3 + 3/4 + \dots}{\text{Total Number of correct shots}}$$

**Results**


1.  ✓
2.  ✗
3.  ✓
4.  ✓
5.  ✗



With the MediaMill team


## Semantic Pathfinder @ TRECVID

**TRECVID 2004-2006 Benchmark Comparison**




Average Precision

The Good




The Bad

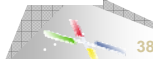


ill-defined / few examples

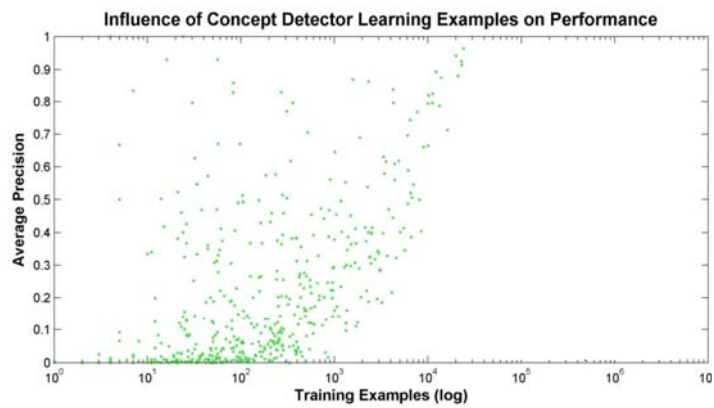
The Ugly



exploit TV repetition



## 491 detectors, a closer look



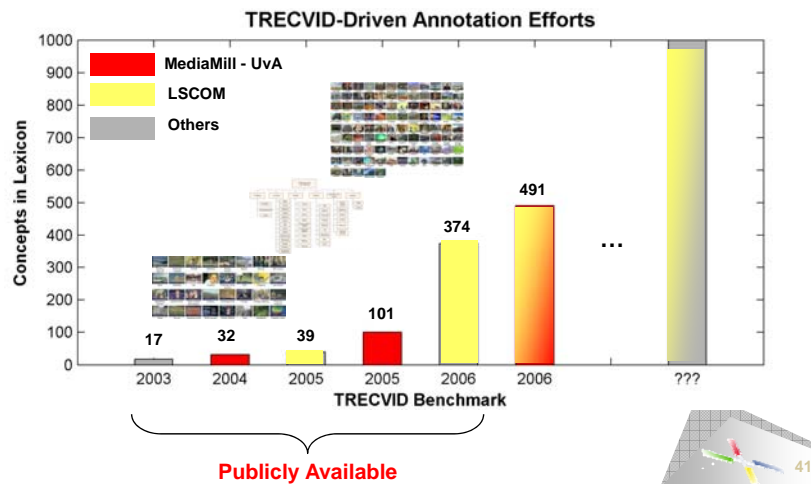
*The number of labeled image examples used at training time seems decisive in concept detector accuracy.*

## Demo time!



## Concept detector: requires examples

- TRECVID's collaborative research agenda has been pushing manual concept annotation efforts



## Concept definition

- **MM078-Police/Security Personnel**
  - ✓ Shots depicting law enforcement or private security agency personnel.




References:  
 Christel, Informedia, 2005  
 Volkmer et al, ACM MM 2005

## Collaborative annotation tool

TRECVID 2005

- Manual annotation by 100+ TRECVID participants
  - ✓ Incomplete, but reliable



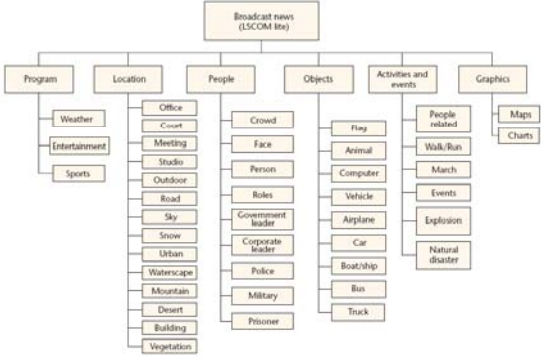
43

References:  
 Naphade et al, IEEE Multimedia 2006

## Manual annotations: LSCOM-lite

TRECVID 2005

- LSCOM:
  - ✓ Large Scale Annotation for Multimedia
  - ✓ Aims for ontology of 1,000 annotated concepts
- LSCOM-Lite: annotations for 39 semantic concepts
  - ✓ Used in TRECVID 2005 and 2006



```

    graph TD
      Root[Broadcast news (LSCOM lite)] --> Program
      Root --> Location
      Root --> People
      Root --> Objects
      Root --> Activities[Activities and events]
      Root --> Graphics

      Program --> Weather
      Program --> Entertainment
      Program --> Sports

      Location --> Office
      Location --> Court
      Location --> Meeting
      Location --> Studio
      Location --> Outdoor
      Location --> Road
      Location --> Sky
      Location --> Snow
      Location --> Urban
      Location --> Waterscape
      Location --> Mountain
      Location --> Desert
      Location --> Building
      Location --> Vegetation

      People --> Crowd
      People --> Face
      People --> Person
      People --> Roles
      People --> Government[Government leader]
      People --> Corporate[Corporate leader]
      People --> Police
      People --> Military
      People --> Prisoner

      Objects --> Flag
      Objects --> Animal
      Objects --> Computer
      Objects --> Vehicle
      Objects --> Airplane
      Objects --> Car
      Objects --> Boat[Boat/ship]
      Objects --> Bus
      Objects --> Truck

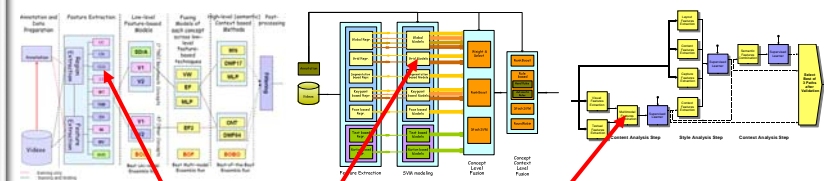
      Activities --> People[People related]
      Activities --> Walk[Walk/Run]
      Activities --> March
      Activities --> Events
      Activities --> Explosion
      Activities --> Natural[Natural disaster]

      Graphics --> Maps
      Graphics --> Charts
    
```

44

## TRECVID Criticism

- Focus is on the final result
  - ✓ TRECVID judges **relative** merit of indexing methods
  - ✓ Ignores repeatability of intermediate analysis steps
- Systems are becoming more complex
  - ✓ Typically combining several features and learning methods
- Component-based optimization and comparison impossible



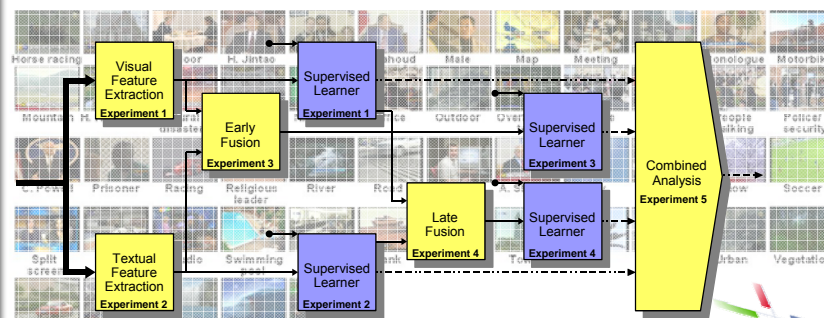
What is the contribution of these components?

45

## MediaMill Challenge

- **The Challenge provides**
  - ✓ Manually annotated lexicon of 101 semantic concepts
  - ✓ Pre-computed low-level multimedia features
  - ✓ Trained classifier models
  - ✓ Five experiments
  - ✓ Baseline implementation together with baseline results
- **The Challenge allows to**
  - ✓ Gain insight in intermediate video analysis steps
  - ✓ Foster repeatability of experiments
  - ✓ Optimize video analysis systems on a component level
  - ✓ Compare and improve upon baseline

• **The Challenge lowers threshold for novice multimedia researchers**



Online available: <http://www.mediamill.nl/challenge/>

## MediaMill Challenge

### ➤ Advantages

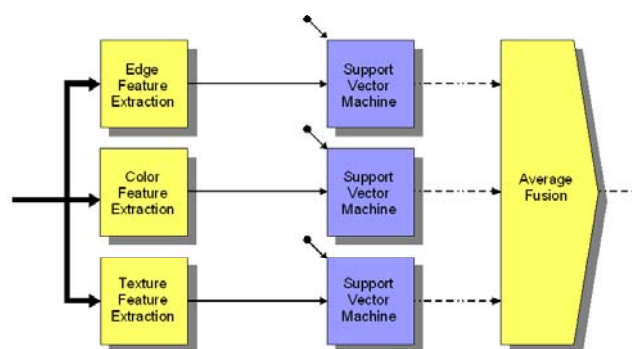
- ✓ For research
  - People can focus on the experiment for which they have the expertise without having to do all the processing
    - Pure computer vision
    - Pure natural language processing
    - Pure machine learning
    - .....
- ✓ For education
  - Students can do
    - large scale experiments
    - compare themselves to each other
    - ..... and to the state-of-the-art

47

## Columbia374

### ➤ Baseline for 374 concept detectors

- ✓ Focus is on visual analysis experiments



Online available:  
<http://www.ee.columbia.edu/ln/dvmm/columbia374/>

48

**Case study**

**Fabchannel.com**

**FABCHANNEL**  
LIVE CONCERT VIDEOS

- Fabchannel narrowcasts concerts from Amsterdam Paradiso and Melkweg venues
  - ✓ Currently +/- 700 concerts online
- Fabchannel request
  - ✓ What can you do with 45 hours of live concerts?
- Answer:
  - ✓ Let's try the semantic pathfinder to detect concert concepts



**Results for singer**

**FABCHANNEL**  
LIVE CONCERT VIDEOS



50



**FABCHANNEL**  
LIVE CONCERT VIDEOS

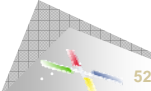
## Results for drummer



The screenshot shows a window titled 'MediaMill Semantic Video Search Engine'. It features a search bar at the top with buttons for 'Query', 'Grid', 'Cross', 'Sphere', 'Play', 'Info', and 'Sort'. Below the search bar is a small video player showing a drummer. The main area displays a grid of 18 video thumbnails, each with a unique ID (e.g., 'm002\_001\_001', 'm002\_002\_001', etc.). The thumbnails show various scenes of drummers performing. At the bottom right of the window, there are buttons for 'Current shot', 'Play', 'Bookmark', and 'Unbookmark'. A small graphic of a grid with colored lines and the number '51' is located in the bottom right corner of the slide.

## Conclusions

- An international community is building a bridge to narrow the semantic gap
  - ✓ Currently detects more than 500 concepts in broadcast video
  - ✓ Generalizes outside news domain
- Important lessons
  - ✓ No superior method for all concepts exists,
  - ✓ Best to **learn** optimal approach per concept
  - ✓ Best methods cover variation in features, classifiers, and concepts



## Concept detection challenges

- Show generality of approach over several domains
  - ✓ Show benefit of web-based image/video and annotations
- Show that concept classes work with less analysis
  - ✓ People, objects, setting
- Show benefit of using dynamic nature of video
  - ✓ Events
- Show that an ontology can help
  - ✓ How to connect logical relations to uncertain detectors?
- Show that 'iconological' concepts can be detected
  - ✓ E.g. funny, sarcastic, cozy, ...

53

## Using concept detectors

- “We are now seeing researchers starting to use the confidence values from concept detectors, within the shot retrieval process and this appears to be the roadmap for future work in this area.”
  - ✓ Alan Smeaton, Information Systems, 32(4):545-559, 2007

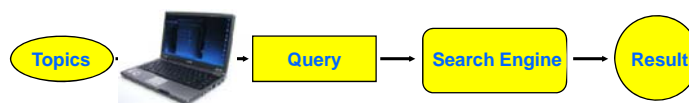
54

## Measure concept detector influence

- Hypothesis 1:
  - ✓ *Increasing the number of concept detectors in a lexicon improves video retrieval accuracy.*
  
- Hypothesis 2:
  - ✓ *Combining concept detectors from a lexicon improves video retrieval accuracy.*

55

## TRECVID automatic search task



- Automatically solve search topic
- Return 1,000 ranked shot-based results
- Evaluate using Average Precision
  
- TRECVID 2005
  - ✓ 85 hrs test set – Chinese, Arabic, English TV News
  - ✓ 24 search topics

56

## Topic examples



Find shots of one or more helicopters in flight.



Find shots of a hockey rink with at least one of the nets fully visible from some point of view.



Find shots of an office setting, i.e., one or more desks/tables and one or more computers and one or more people



Find shots of a group including at least four people dressed in suits, seated, and with at least one flag.

57

## Influence of lexicon size

- Lexicon = 363 machine learned concept detectors
- Procedure
  1. Set bag size  $B = 10$ ;
  2. Select random bag of  $B$  detectors from lexicon
  3. Determine maximum performance for each search topic
  4.  $B += 10$ ;
  5. Go back to step 2.
- Repeat the process 100 times
  - ✓ Reduces random positive and negative effects

58

### Influence of lexicon size

TRECVID 2005

Mean Average Precision

Concept Detectors in Lexicon

Linear increase for first 60 detectors

- Size matters
  - ✓ Lexicon of 150 detectors comes close to maximum performance
- Some detectors perform well for specific topics
  - ✓ Tennis game detector for “find two visible tennis players”
- Substantial number of detectors not accurate enough yet
  - ✓ Only small increase when more than 70 detectors are used

### Influence of detector combination

- How to combine multiple detectors?
  - ✓ Experiment: pair-wise oracle fusion

Office  $\lambda$

Computer  $1 - \lambda$

- Improvement for 20 out of 24 topics
- Increase per topic as high as 89%
- Overall increase 10%

### Typical results

The diagram illustrates search results for two queries. The first query, 'Find shots of a graphic map of Iraq, location of Baghdad marked - not a weather map.', yields two results: 'Best' (a map with Baghdad marked) and '2nd Best' (a map with overlaid text). The second query, 'Find shots of George Bush entering or leaving a vehicle (e.g., car, van, airplane, helicopter, etc) (he and vehicle both visible at the same time)', yields two results: 'Best' (a shot of rocket propelled grenades) and '2nd Best' (a shot of Iyad Allawi). A red text box asks 'How to select relevant detectors automatically?'.

Find shots of a graphic map of Iraq, location of Baghdad marked - not a weather map.

Best + 2nd Best

Maps + Overlaid Text

How to select relevant detectors automatically?

Find shots of George Bush entering or leaving a vehicle (e.g., car, van, airplane, helicopter, etc) (he and vehicle both visible at the same time)

rocket propelled grenades + Iyad Allawi ?

61

### Problem statement

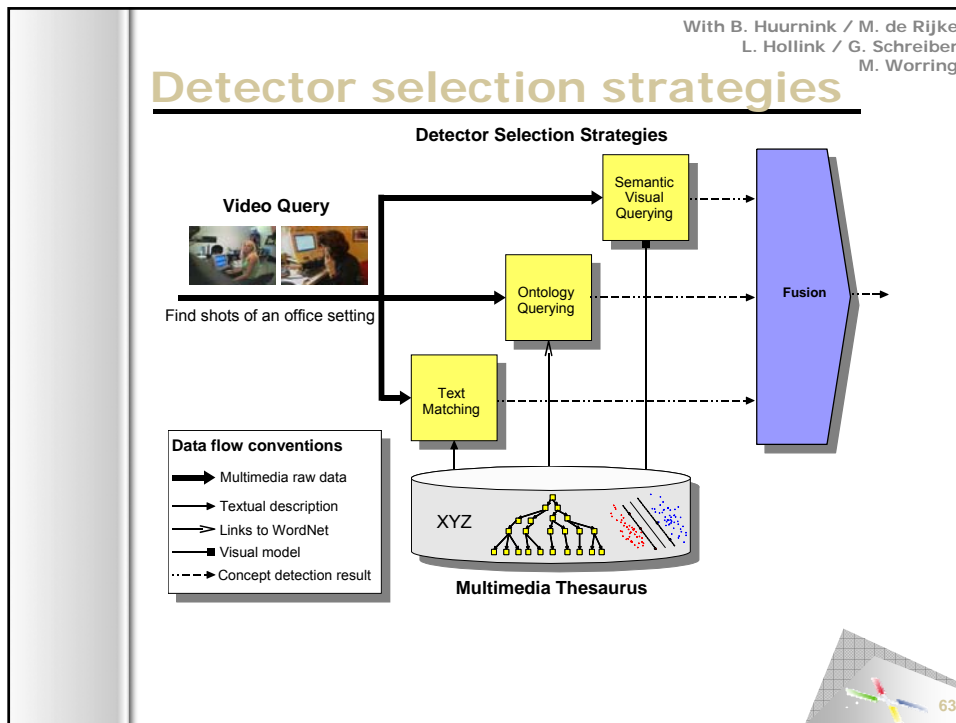
The flowchart shows the process of searching for office settings. It starts with 'Topics' (represented by an oval) which leads to a 'Query' (represented by a box). The 'Query' is processed by a 'Search Engine' (represented by a box), resulting in a 'Result' (represented by a circle). A feedback loop arrow points from the 'Result' back to the 'Query'.

Find shots of an office setting

Topics → Query → Search Engine → Result

➤ How to translate query topic to concept detectors?

62



## Influence of detector selection combi

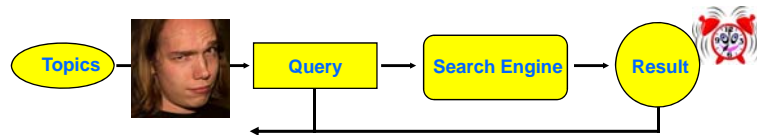
TRECVID 2005

➤ Individual selection strategies seem comparable  
 ✓ But, **oracle** combination of selection strategies pays off!

	Best	Text Matching	Ontology Querying	Visual Querying
Find shots of a tall building (with more than 5 floors above the ground)	Office Building	Tower	Building	Grass
Find shots of an office setting, i.e., one or more desks/tables and one or more computers and one or more people	Computer	Computer	Office	Emile Lahoud
Find shots of one or more palm trees.	Weapon	Tropical Setting	Tree	Fire Weapon

## TRECVID interactive search task

- So many choices for retrieval...
  - ✓ Why not let user decide interactively?



<http://trecvid.nist.gov/>

65

References:  
Carnegie Mellon University

## 'Classic' Infromedia system

- First multimodal video search engine



66



References:  
FxPal

## MediaMagic

➤ Focus on the story level

The screenshot shows the MediaMagic interface. On the left is a large grid of video thumbnails labeled 'A'. Below it is a search bar containing the text 'condoleezza rice' labeled 'B'. To the right of the search bar is a search results pane labeled 'D' which includes a video player showing a news segment. Below the search bar are checkboxes for 'Show related words' and 'Show segment text', and radio buttons for 'Text matching', 'Fuzzy', 'Blind', and 'Exact'. Below the search bar is a row of thumbnails labeled 'C'. To the right of the search bar is a list of search results labeled 'E', with a detailed view of a video segment labeled 'F'. The interface also includes a 'Clear' button and a 'Find shots of Condoleezza Rice' section at the bottom.

67

References:  
IBM

## IBM MARVeI

➤ A web based system

The screenshot shows the IBM MARVeI web interface. The page title is 'MARVeI Multimedia Analysis and Retrieval System (IBM T. J. Watson Research Center) - Microsoft Windows'. The interface includes a search bar, a 'Random Retrieval' section, and a grid of search results. Each result includes a video thumbnail and a text snippet. The search results are labeled 'A' through 'E'. The interface also includes a 'Search' button and a 'Apply' button. The URL 'http://mp7.watson.ibm.com/marvel/' is displayed at the bottom.

http://mp7.watson.ibm.com/marvel/

68

References:  
Oulu University

## Cluster-temporal browsing

➤ Using that result are typically similar/close in time

The screenshot shows a web-based video retrieval interface. At the top, there's a search bar and navigation options. Below it, a grid of video thumbnails is displayed, each with a small text snippet. The interface is titled 'Cluster-temporal browsing' and includes a 'References' section for 'Oulu University'. The video thumbnails are arranged in a grid, and the text snippets provide context for each video frame.

References:  
Dublin City University

## Físchlár

➤ Optimized for use by “real” users

The screenshot displays the 'Físchlár' video retrieval interface. It features a 'QUERY PANEL' on the left with search filters and a 'SEARCH RESULT' section in the center showing a grid of video thumbnails. A 'SAVED POINT' section is visible on the right. The interface is designed for user interaction and includes various search and navigation tools.

References:  
NUS & ICT-CAS

## VisionGo

➤ Extremely fast and efficient

	current	processed
positive	39	6
negative	33	84
submitted	0	0
labeled	72	90
unlabeled	928	910

71

References:  
Carnegie Mellon University

## Extreme video retrieval

➤ Observation

- ✓ Correct results are retrieved, but not optimally ranked
- ✓ If user has time to scan results exhaustively, retrieval is a matter of watching, selecting, and sorting **quickly**

➤ Push the **user** to the max = **very demanding!**

- ✓ ~~Rapid~~-serial visual presentation
- ✓ Adjust browser to depth of results

72

## Futuristic video retrieval

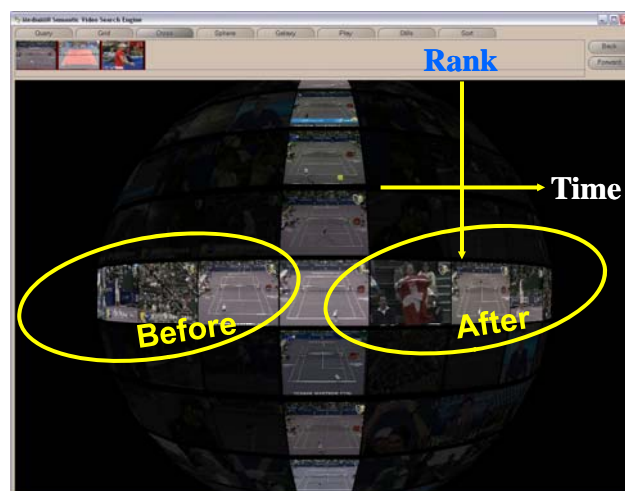


Jonathan Wang, Carnegie Mellon University

73

With the MediaMill team

## CrossBrowsing through results



74

References:  
de Rooij, CIVR 2008

## ForkBrowser

75

## Demo time!

76


RotorBrowser

CrossBrowser

## NIST TRECVID benchmark

anno 2001

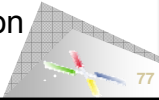
- Benchmark objectives
  - ✓ Promote progress in video retrieval research
  - ✓ Provide common dataset (shots, recognized speech, key frames)
  - ✓ Use open, metrics-based evaluation
- Large international field of participants



- ✓ and the 70 others...

- Currently the **de facto standard** for evaluation

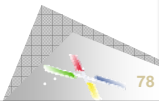
<http://trecvid.nist.gov/>



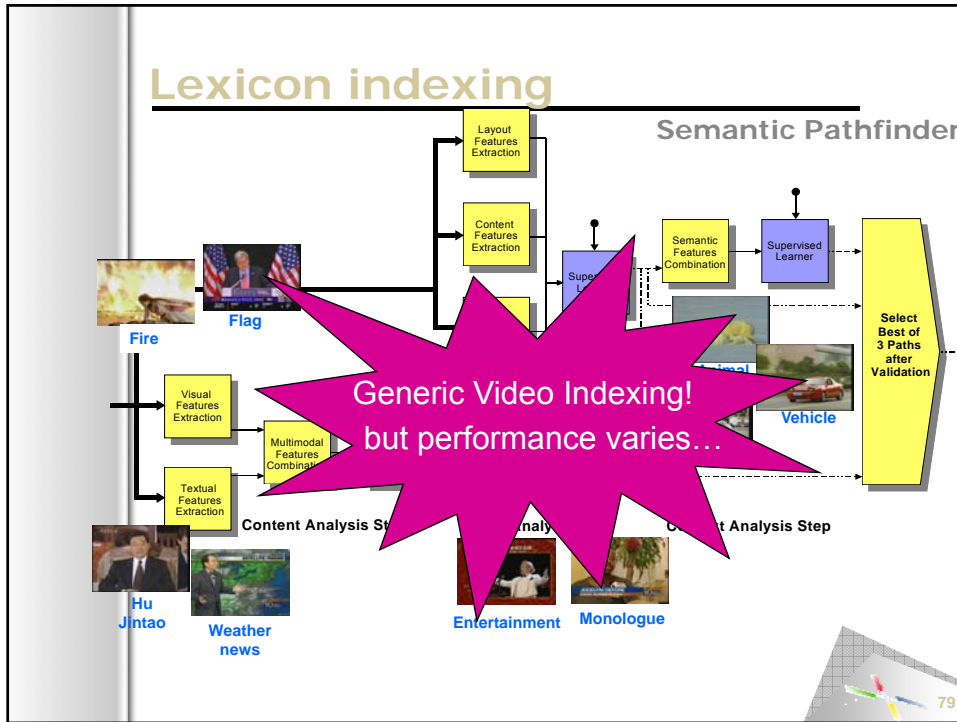
77

## Experimental Setup

- Experiment 1
  - ✓ TRECVID 2004 (64 hrs test set – English TV News)
  - ✓ **Lexicon with 32 learned concepts** (where others use max. 10)
  - ✓ All other components 'standard'
- Experiment 2
  - ✓ TRECVID 2005 (85 hrs test set – Chinese, Arabic, English TV News)
  - ✓ **Lexicon with 101 learned concepts** (where others use max. 39)
  - ✓ Added advanced display (CrossBrowser)



78

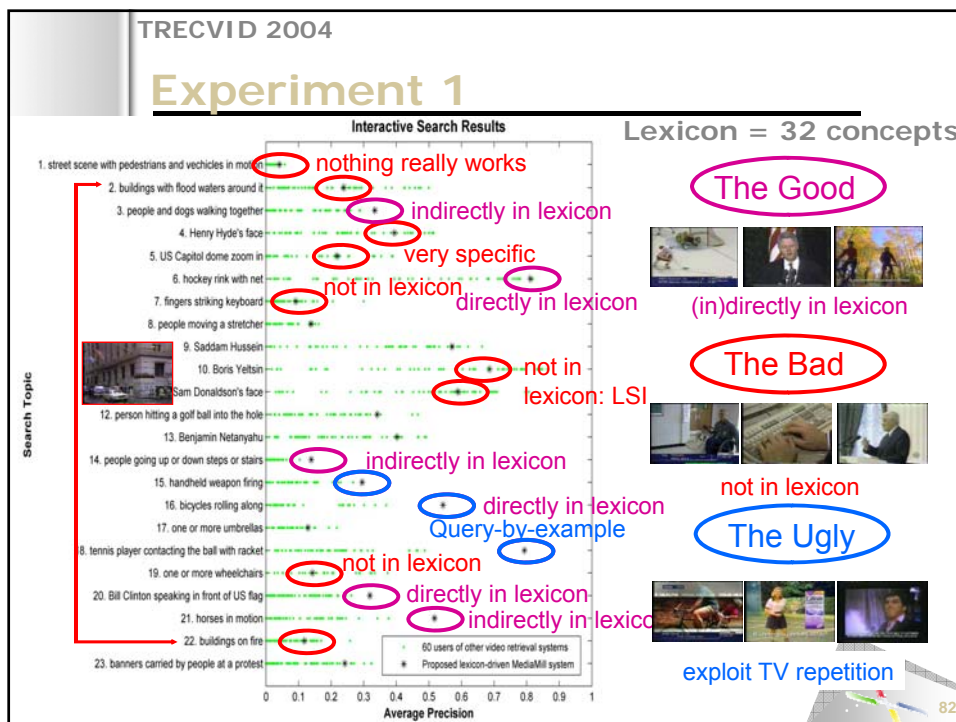


## Query selection

The screenshot shows a search interface with three main query methods highlighted by pink ovals:

- Query-by-keyword:** Located at the top, showing a search bar and language options.
- Query-by-concept:** A central grid of video thumbnails representing different concepts.
- Query-by-example:** A section at the bottom showing a video frame and a 'Visually similar to' search option.

... yields a ranking of the data





TRECVID 2005

## Learned lexicon of 101 concepts

The grid contains 101 concepts, each with a representative image and a label: Aircraft, I. Allawi, Anchor, Animal, Y. Arafat, Baseball, Basketball, Beach, Bicycle, Bird, T. Blair, Boat, Building, Bus, G. Bush jr., G. Bush sr., Candle, Car, Cartoon, Chair, Charts, B. Clinton, Cloud, Corp. leader, Court, Crowd, Cycling, Desert, Dog, Drawing, Drawing & Cartoon, Duo-anchor, Entertainment, Explosion, Face, Female, Fire weapon, Fish, Flag, Flag USA, Food, Football, Golf, Government building, Government leader, Graphics, Grass, Horse, Horse racing, House, Indoor, H. Jintao, J. Kerry, E. Lahoud, Male, Map, Meeting, Military, Monologue, Motorbike, Mountain, H. Nasrallah, Natural disaster, News paper, Night fire, Office, Outdoor, Overlaid text, People, People marching, People walking, Police security, C. Powell, Prisoner, Racing, Religious leader, River, Road, Screen, A. Sharon, Sky, Smoke, Snow, Soccer, Split screen, Sports, Studio, Swimming pool, Table, Tank, Tennis, Tower, Tree, Truck, Urban, Vegetation, Vehicle, Violence, Waterfall, Waterscape, Weather.

83

TRECVID 2005

## Experiment 2

Lexicon = 101 concepts

Interactive Search Results

Search Topic

Average Precision

48 users of other video retrieval systems  
 Proposed lexicon-driven MediaMill system

1. Condoleezza Rice (not in lexicon)  
 2. Iyad Allawi (poor concept detection)  
 3. Omar Karami  
 4. Hu Jintao  
 5. Tony Blair  
 6. Mahmoud Abbas  
 7. graphic map of Iraq, Baghdad marked (concept specification fails)  
 8. two visible tennis players on the court  
 9. people shaking hands  
 10. helicopter in flight  
 11. George W. Bush entering or leaving a vehicle  
 12. something on fire with flames and smoke  
 13. people with banners or signs  
 14. people entering or leaving a building (concept combination fails)  
 15. a meeting with a large table and people  
 16. a ship or boat  
 17. basketball players on the court  
 18. one or more palm trees  
 19. an airplane taking off  
 20. a road with one or more cars  
 21. one or more military vehicles  
 22. a tall building  
 23. a goal being made in a soccer match  
 24. office setting

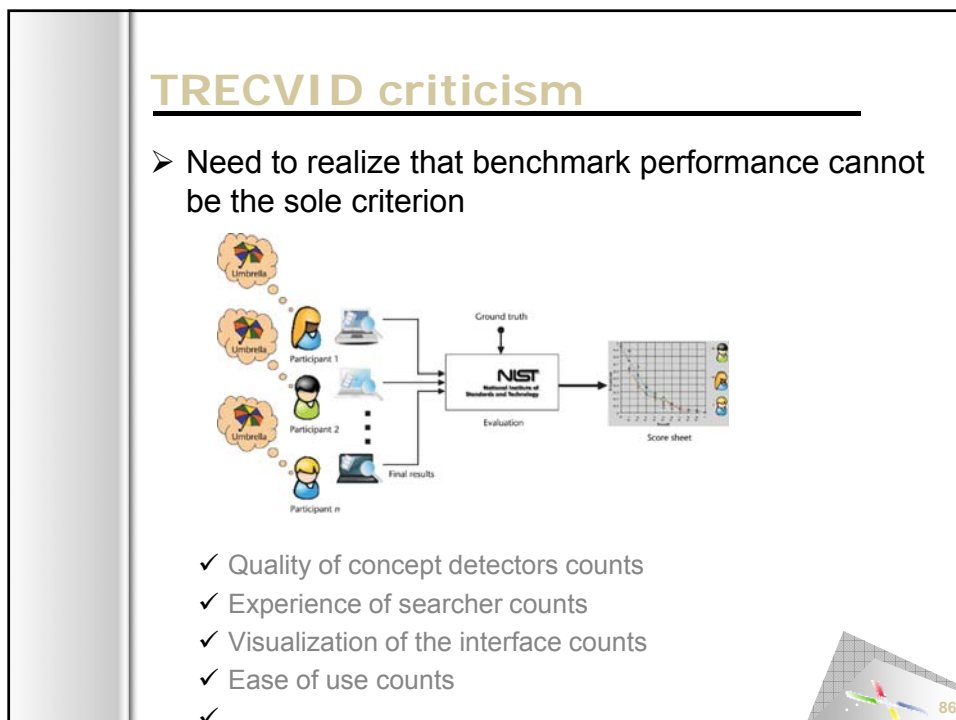
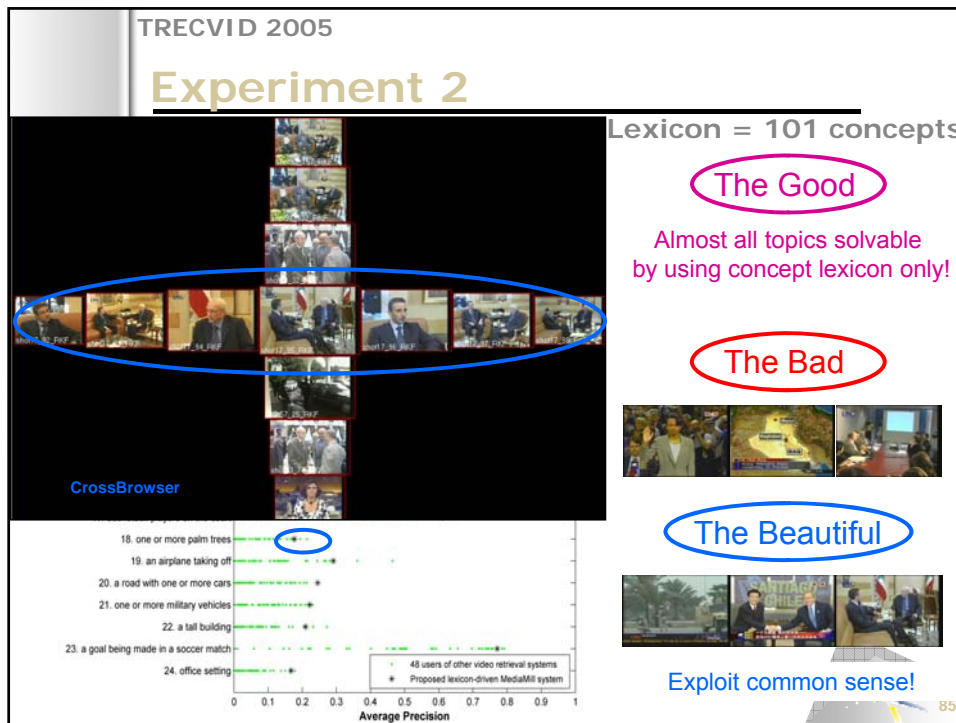
The Good  
 Almost all topics solvable by using concept lexicon only!


The Bad

The Beautiful

Exploit common sense!

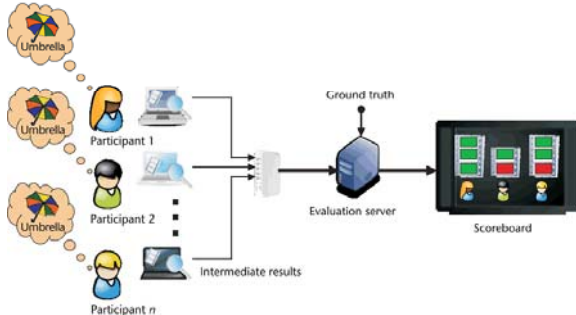
84





## VideOlympics

- Promote multiple facets of interactive video search
  - ✓ Real-time interactive video search ‘competition’
  - ✓ Simultaneous exposure of multiple video retrieval systems
  - ✓ Audience obtains complete overview of possibilities and limitations of state-of-the-art



87



## Participants

Showcase Event





Carnegie Mellon













88

## Setup

VIDEOLYMPICS WORLD IN BELIEF *MSE*

Showcase Event

Time: 89.296  
Positive: 63  
Negative: 23

Contest running...

One display

TRECVID like queries  
A result is submitted as soon as it is found

## Video trailer

VIDEOLYMPICS

<http://www.VideOlympics.org>

## Conclusions

---

- We should use
  - ✓ Generic semantic indexing
  - ✓ Large lexicons
- For browsing
  - ✓ The large lexicons are the best entrance
  - ✓ Ranking and time are most important dimensions
- TRECVID = priceless
  - ✓ Fosters a common research agenda and international collaboration
  - ✓ Many valuable resources available online
  - ✓ Offers you an opportunity to do it yourself easily

## Concept retrieval challenges

---

- How to leverage concept detectors for search?
  - ✓ How to select the best detectors automatically?
  - ✓ How to combine concept selection methods?
  - ✓ How to balance semantic coverage and anticipated performance of detectors for a specific query?
- How to help the user in browsing
  - ✓ How to present detectors and their uncertainty to users?
  - ✓ How to select the best detector on the fly?
  - ✓ How to present thread space in an intuitive manner?
- Generalization to unstructured domains
  - ✓ Consider YouTube domain for example
  - ✓ How to use contextual metadata?

Further information 

[www.MediaMill.nl](http://www.MediaMill.nl)



 93