

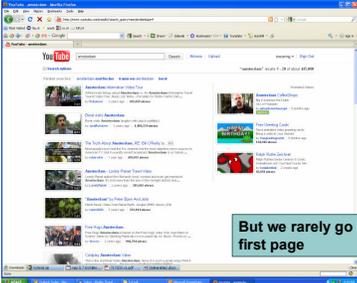
Multimedia Analytics: Exploration of Large Multimedia Collections

Marcel Worring

With contributions from
Ork de Rooij, Daan Odijk, Cees Snoek, Koen van de Sande

Informatics Institute
University of Amsterdam

The Internet



YouTube every minute >50 hours of video uploaded

But we rarely go beyond the first page

Broadcasting archives

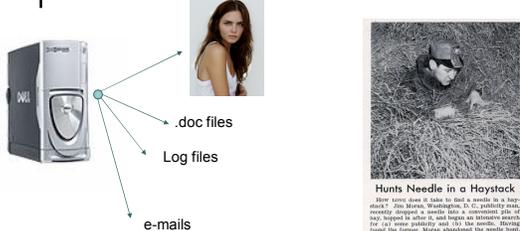


BEELD EN GELUID

Sound and Vision Archive

Yearly 15.000 hours video added

Forensics



10 Tb is a large set of potential locations for child porn, but how to find it?

We need to find all needles

Intelligence

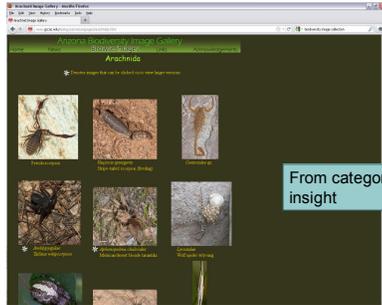



Hsdl.org

Not obvious what you are looking for, any clue would help

publicintelligence.net

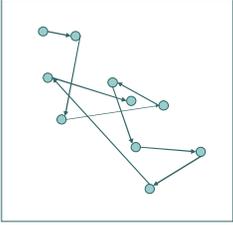
Science



From categorization to insight

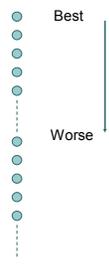
The basis: ranking of data

Some query defines starting point and order



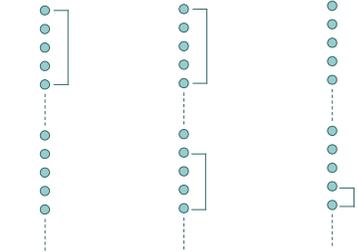
An image/video collection

Result



Tasks

Most of the techniques are based on providing a ranking, different applications different criteria of success.



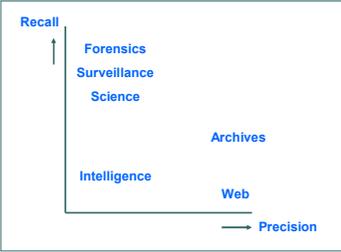
High precision

High recall - many items

High recall - few items

Precision versus recall

Search requirements in retrieval vary



Recall

Forensics
Surveillance
Science

Archives

Intelligence

Web

Precision

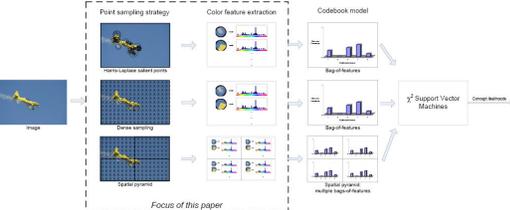
"Insight"

What is a category?

- Definition
 - A class or group of things, people, etc. possessing some quality or qualities in common; a division in a system of classification
- So for images
 - Could be in the content or in the metadata
- Current focus
 - The visual content

Concept detection

References: van de Sande, CIVR 2008, Snoek Trecvid 2008



Image

Flatt sampling strategy

Color feature extraction

Codebook model

χ^2 Support Vector Machines

Concept names

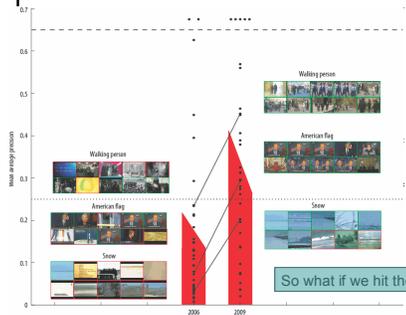
Focus of this paper

One of the best detection pipelines

50 - 500 concepts
Wildlife, Face,
People marching,

Quality is steadily improving

Ref: Snoek IEEE Computer 2010



mean average precision

0.7

0.6

0.5

0.4

0.3

0.2

0.1

0

2006

2009

Walking person

American flag

Sewer

So what if we hit the target line?

We start posing more difficult questions

Multimedia Analysis

Sometimes it works

Rover Mini Cooper



Concepts with score higher than 0.2

Car 0.715

Ground_Vehicles 0.555

Vehicle 0.387

- Tags
- Views
- Model
- Date and Time
- Location
- Concepts

Sometimes it doesn't

il decollo



Concepts with score higher than 0.2

Adult 0.203

Building 0.358

Indoor 0.311

Male_Person 0.223

- Tags
- Views
- Model
- Date and Time
- Location
- Concepts

Automatic versus interactive

(Semi-) interactive

Recall ↑

High recall retrieval

Forensics
Surveillance
Science

Intelligence

Archives

Web

Explorative search → Precision

Automatic

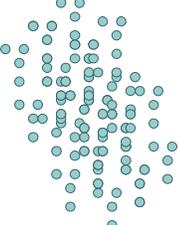
"Insight"

Man versus machine

What humans are good at

-  Recognize
-  Interpret context
-  Associate

What machines are good at

-  Bulk processing

Multimedia Analytics

- o Aims
 - To combine the best of both worlds
- o For
 - Datasets that are too large to be handled by humans alone and too complex to be handled by machines only
- o With the constraint
 - That we need to see an image before we can make a judgment on it

Chinchor, Thomas, Wong, Christel, Ribarsky CG&A 2010

Definition

Multimedia Analytics

=

Multimedia Analysis

+

Visual Analytics

Visual Analytics

Visual Analytics is the science of analytical reasoning facilitated by **interactive visual interfaces**



The field started as a result of 9/11

An overwhelming amount of traces and information.

But the scope of visual analytics is much broader and applies to any scenario where large data collections have to be used.

Automatic versus interactive

(Semi-) interactive

Recall

High recall retrieval

Forensics
Surveillance
Science

Archives

Intelligence

Web

Explorative search

Precision

Automatic

"Insight"

Threads

- Definition
 - A thread is a linked sequence of shots in a specified order, based upon an aspect of their content.
- Static threads
 - Are pre-computed
- Dynamic threads
 - Are created on the fly

Navigation support: Threads

One possible thread

And another one
Basically different queries

Currently viewed shot

Snoek, Worring, Koelma, Smeulders
IEEE Trans MM 2007

Beyond the query result

- When I pose a (semantic) query
 - The result is a one dimensional ranked list
 - So one dimension of the display can be used for another purpose, the CrossBrowser

or the sphere based variation

Rank

Before Time After

de Rooij, Worring IEEE Trans MM 2010

The RotorBrowser

Even more directions to go through

Threads + time

Current shot

Very good for explorative browsing

Automatic versus interactive

(Semi-) interactive

Recall

High recall retrieval

Forensics
Surveillance
Science

Archives

Intelligence

Web

Explorative search

Precision

Automatic

"Insight"

The guidelines we learned

- o Guideline 1
 - The user has to be able to efficiently inspect the setting, scene, and all objects and their motions within a shot for relevance and the user should be able to adapt the level of detail given according to the complexity.
- o Guideline 2
 - The user has to be able to view the temporal context of video shots in order to determine the relevance of the current shot.
- o Guideline 3
 - The user must be able to label multiple shots at the same time, when this is more efficient.
- o Guideline 4
 - The system should show unexplored areas within the collection.

Guidelines

- o Guideline 5
 - The system should always yield a default navigation path.
- o Guideline 6
 - The user has to be able to efficiently inspect the potential relevance of the navigation options in the navigation context.
- o Guideline 7
 - The system should aid the user to return to earlier navigation options by visual or spatial means.
- o Guideline 8
 - The interface should not switch between different panes.
- o Guideline 9
 - The interface must use a clear mapping between navigation and visualization.

ForkBrowser

Navigation possibilities

Direct mapping of keyboard to visualization

Initial query

similarity thread 1

similarity thread 2

Micons to show content, can be zoomed

Timeline as context

History

GO TO CURRENT SHOT

BOOKMARK

Active Zooming

Quick inspection of a large part of a thread

1. Inspecting navigation possibilities
2. Label many at the same time

ForkBrowser depicting scene from the Klokhuis program

active zooming based on similarity

The only screen switch in the system, but different task

ForkBrowser

Navigation possibilities obtained through active learning

Initial query

similarity thread 1

similarity thread 2

Timeline

History

GO TO CURRENT SHOT

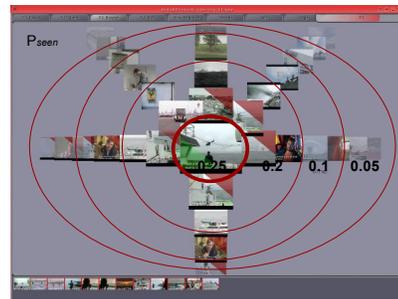
BOOKMARK

Browsing through results

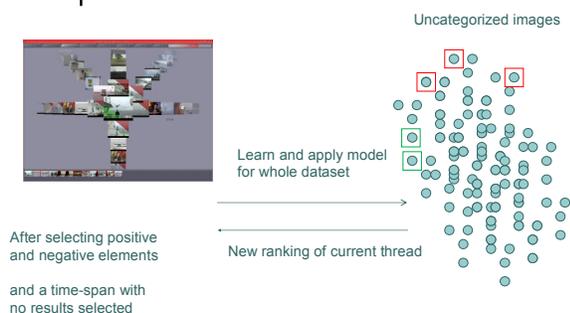
User marks relevant results



Negative results marked automatically



Update navigation path



Refined default navigation path



Interactive Benchmarking

o TRECVID

- Search topics defined by NIST
- Participants have fixed amount of time to search through the collection



Visual evaluation of the result by NIST

Experimental setup

- o Dataset: TRECVID 2008
 - 200 hours of video
 - 35766 individual shots
- o Three experiments
 - Study of potential benefit of threads
 - Determining the potential of relevance feedback in this dataset
 - Trecvid participation

Example tasks in TRECVID

Task 1: Find pictures of "Job Cohen"



Only 3 clusters: full frontal face, 2 people, and table of people.
Type B - specific visual example.

Task 2: Find boats or ships



Much variety: top-down, closeup, inside boats, with horizon.
Type C - generic visually diverse w/ concept "boat-ship".

Task 3: Find closeups of hands



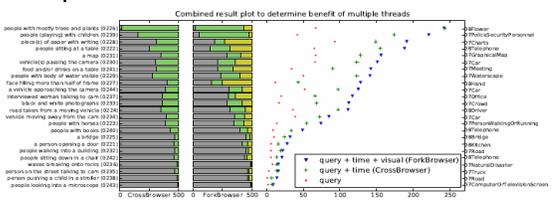
Extremely high variety of results, both closeup and further away.
Type C - generic visually diverse w/ concept "hands".

Task 4: RW pictures of airplanes



In sky and on ground, mostly from war documentaries, bad quality, many angles.
Type D - generic visually diverse, no direct concept match.

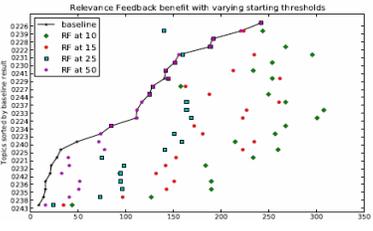
Evaluation



Grey: Query
Time: Green
Visual similar: Yellow

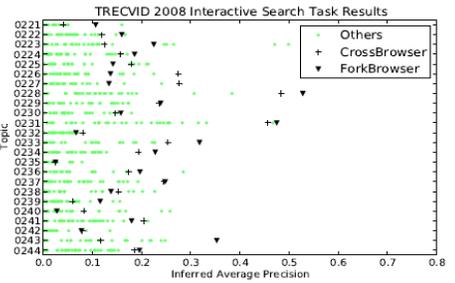
But first using simulated users
-Judge shot based on ground truth
-Follow thread with most results if possible
-Otherwise follow query thread

Evaluation



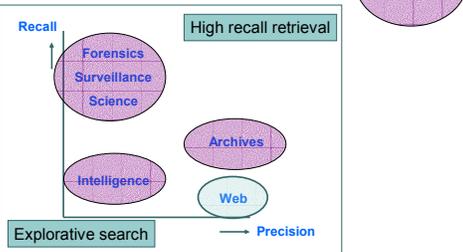
Activating RF at X indicating that X times a non-relevant item was found (500 User Interaction Steps).

And @ TRECVID



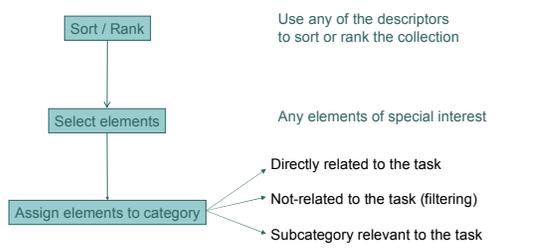
Automatic versus interactive

(Semi-) interactive



Automatic

High recall search processes

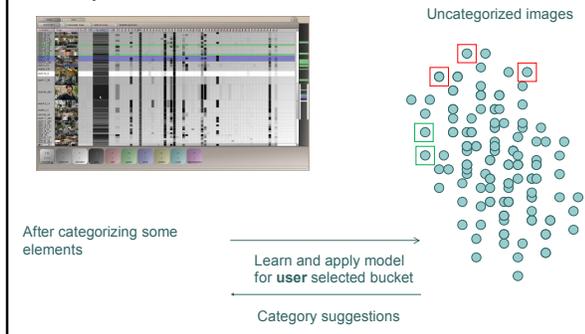


Use any of the descriptors to sort or rank the collection

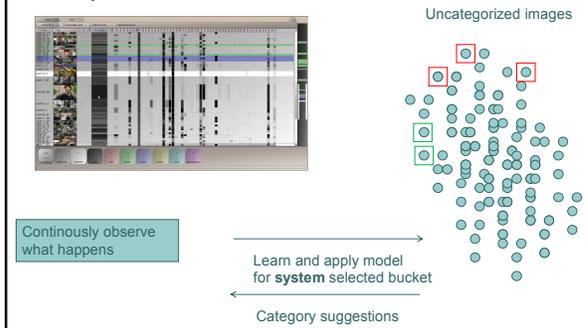
Any elements of special interest

- Directly related to the task
- Not-related to the task (filtering)
- Subcategory relevant to the task

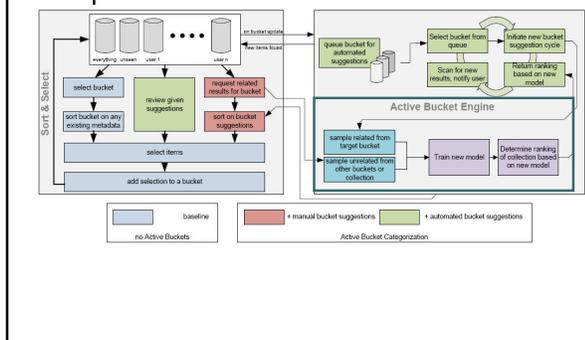
On demand suggestions



Unobtrusive assistance



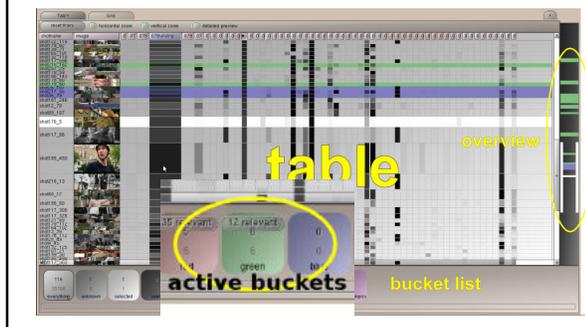
Active Buckets



Sampling and learning

- Two strategies
 - Nearest Mean: fast and simple
 - Compute mean of elements in bucket
 - Rank collection with respect to distance to the mean
 - Support Vector Machines: more expensive, more accurate
 - Dynamically learn new model from examples in the bucket and randomly selected (presumed) negatives

Active Buckets interface



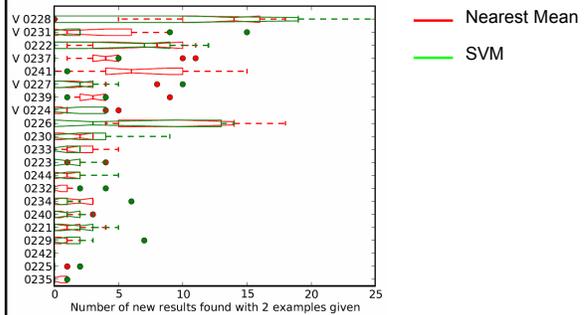
Experimental setup

- Dataset: TRECVID 2008
 - 200 hours of video
 - 35766 individual shots
- Two experiments
 - Determining the potential of relevance feedback in this dataset
 - User experiment comparing the three strategies

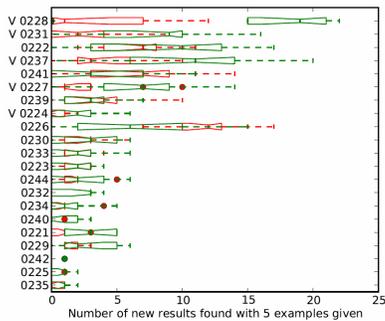
Experiment 1: Potential

- Protocol
 - 10 times random selection of N positive elements
 - Learning using the two strategies
 - Nearest Mean Classifier
 - approx 0.5 seconds per bucket
 - Support Vector Machines
 - approx 2-6 seconds per bucket
- Measure
 - Number of new correct examples in the top 50 ranked results

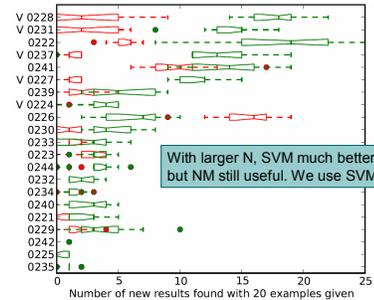
Learning experiment: N=2



Learning experiment N=5



Learning experiment N=20



User experiment

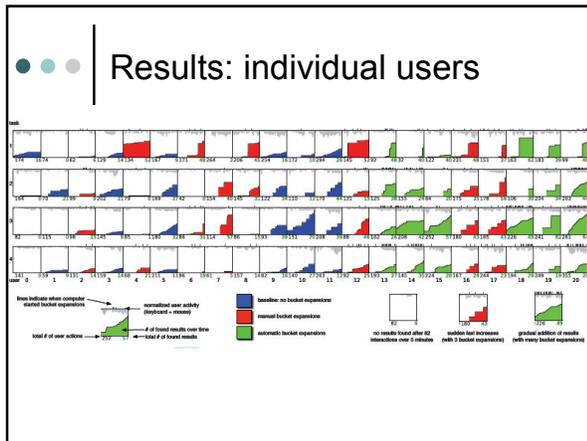
- 21 student users in three groups
 - Baseline
 - Interactive retrieval only
 - Passive buckets
 - Bucket expansion on demand only
 - Active buckets
 - Unobtrusive bucket expansion
- 5 minutes per topic

Results: elements found

Task	Baseline	Manual	Automatic
1	13 ± 7	36 ± 19 *	46 ± 8 * °
2	24 ± 14	23 ± 19	32 ± 8 * °
3	27 ± 24	31 ± 18	53 ± 18 * °
4	13 ± 8	19 ± 7 *	28 ± 5 * °
Average	20 ± 16	27 ± 18 *	40 ± 15 * °
Participants	8	7	6

- * significant at the p=0.01 level compared to baseline
- ° significant at the p=0.01 level compared to manual

Task 1: specific, high visual similarity
 Task 2: generic visually diverse, concept available
 Task 3: generic visually diverse, concept available
 Task 4: generic visually diverse, no concept available



Discussion: task 1

- Specific, visually similar
 - Active bucket users able to select available results in one go
 - When a few results are found active buckets provide the rest of the cluster

Task 1: Find pictures of "Job Cohen"

Only 3 clusters: full frontal face, 2 people, and table of people.
Type B - specific visual example.

Discussion: task 2, 3

- Generic, concept available
 - Presence of appropriate concept makes fully interactive retrieval easier
- Task 2
 - Active buckets more results at once
 - In total same amount of results
- Task 3
 - Manually triggered more at once, but active more results in total

Task 2: Find boats or ships

Much variety: top-down, closeup, inside boats, with horizon.
Type C - generic visually diverse w/ concept "ship".

Task 3: Find closeups of hands

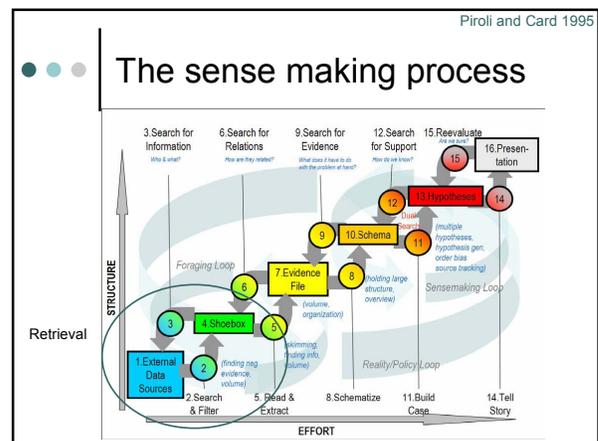
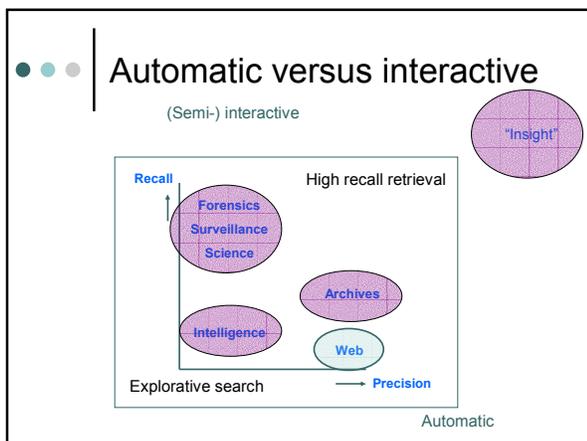
Extremely high variety of results, both closeup and further away.
Type C - generic visually diverse w/ concept "hands".

Discussion: task 4

- Generic, no concept available
 - Having to select a combination of concepts is more difficult for users
 - User used passive buckets when not appropriate
 - Active buckets of clear benefit here

Task 4: B/W pictures of airplanes

In sky and on ground, mostly from war documentaries, bad quality, many angles.
Type D - generic visually diverse, no direct concept match.



Characteristics of insight

- o *Complex*
 - Insight is complex, involving all or large amounts of the given data in a synergistic way, not simply individual data values.
- o *Deep*
 - Insight builds up over time, accumulating and building on itself to create depth.
 - Insight often generates further questions and, hence, further insight.
- o *Qualitative*
 - Insight is not exact, can be uncertain and subjective, and can have multiple levels of resolution.

Characteristics of insight

- o *Unexpected*
 - Insight is often unpredictable, serendipitous, and creative.
- o *Relevant.*
 - Insight is deeply embedded in the data domain, connecting the data to existing domain knowledge and giving it relevant meaning.
 - It goes beyond dry data analysis, to relevant domain impact.

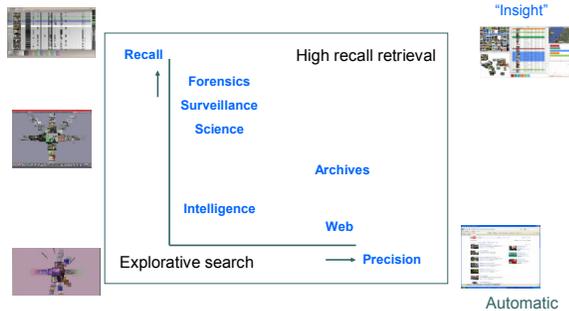
What is needed?

- o Different ways
 - to explore the data
 - to view the data
- o Categories
 - denoting current understanding
- o Coordinated views
 - When a change is made in one view it should be reflected in the whole visualization

Towards insight



Conclusion



m.worring@uva.nl