

Machine Learning, Pattern Recognition, Cross-modal Analysis and Fusion

Alex Hauptmann
School of Computer Science
Carnegie Mellon University,
Pittsburgh, PA, USA
alex@cs.cmu.edu
<http://www.cs.cmu.edu/~alex>

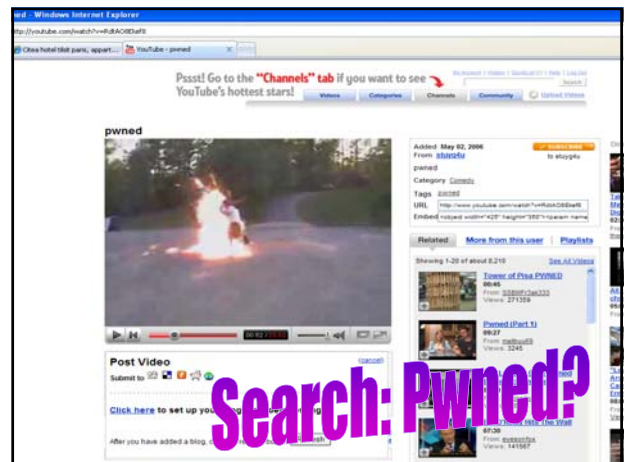
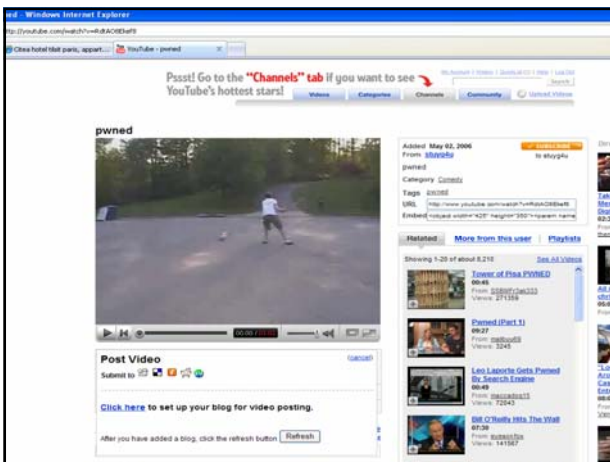
Carnegie Mellon

Outline

Part Two

- Bridging the Semantic Gap across Modalities
 - A Large Scale Ontology for Multimedia (LSCOM)
 - Retrieval Experiments with LSCOM
 - Active Learning of Semantic Concepts
- YouTube and Challenges for the Future

Carnegie Mellon



TRECVID Benchmark Evaluation

Goal is to promote progress in content-based retrieval from digital video via open, metrics-based evaluation

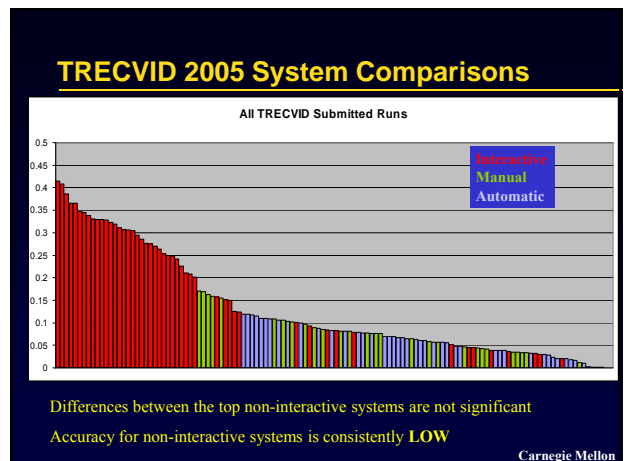
- Laboratory style evaluation, modeled on real world
- Emphasis on good science
- Evaluation is driven by participants (and feasibility)
 - More than 60 participating groups in 2007
- Collection is fixed, available in the spring
 - 50% data used for development, 50% for testing
- ~24 Test queries available in July, 1 month to submission

Details at <http://www-nlpir.nist.gov/projects/trecvid>

Carnegie Mellon

TRECVID 2005 System Comparisons

All TRECVID Submitted Runs



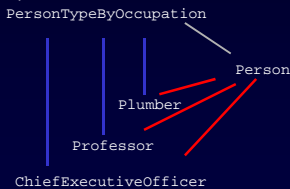
Differences between the top non-interactive systems are not significant
Accuracy for non-interactive systems is consistently **LOW**

Carnegie Mellon

Cyc can Suggest Extensions/Improvements

Use of Faceting Collections

- Higher-order collections that facet a class into sub-classes whose instances share a property
- Example:



Carnegie Mellon

Final Statistics of the Ontology

Size: 784 → 2638

Annotated: 450

Distribution through LSCOM.org

Carnegie Mellon

What are the Limits of Semantic Concepts

- How much can semantic concepts improve video search? What is the upper bound?
- How many do we need for good retrieval?
- Are there enough concepts out there?
- Can we find them?

Carnegie Mellon

Can Semantic Concepts provide significant boost to IR?

- Low level visual features are not sufficient to understand an image or video clip ("**Semantic Gap**")
 - Low-level: Texture, color, shape, interest points, motion, audio (SFFT, MelCep, Zero crossing, ...)
- Describe video through intermediate **semantic** concepts
 - Face, car, outdoors, boat, building, clouds, sky, water, ...
- Semantic concepts can be learned automatically
- Semantic concepts are useful for retrieval and summarization

Carnegie Mellon

Why are Semantic Concepts Important?

- What if we could detect a lot of concepts?
- Speech recognition analogy
 - ~~100 words~~ → ~~1000 words~~ → 20,000 words → 64,000 words
- LSCOM – A Large Scale Ontology for Multimedia
 - 2 year workshop to define and annotate 1000 concepts
 - Defined 850 concepts
 - Extended via ontology to ~2400 concepts,
 - Annotated 450 concepts on 70 hours of TV news
 - Available at www.LSCOM.org

Carnegie Mellon

Building a Semantic Concept Detector

- Get training data
 - Manual labels over a representative, large data set
- Extract feature vectors
 - Color, texture, shape, motion, feature points, text, audio
 - Very high dimensional, continuous features
- Run machine learning algorithm
 - SVM, AdaBoost, Decision Trees, Graphic Models, etc.

Presto!

But there is a lot of detail to be considered...

- Details matter: feature selection, parameter tuning, model combination

Carnegie Mellon

Video Data: TrecVid 2006 Development Corpus

- About 70 hours of multilingual broadcast news from October – November 2004
- English, Chinese, Arabic
- 61901 shots as the units of analysis
- Truth annotations for about 450 LSCOM concepts
- Evaluation by NIST on a separate test corpus for detection of a subset of the semantic concepts

Carnegie Mellon

39 Semantic Concepts (LSCOM-Lite)

1	Sports	20	Person
2	Entertainment	21	Government-Leader
3	Weather	22	Corporate-Leader
4	Court	23	Police-Security
5	Office	24	Military
6	Meeting	25	Prisoner
7	Studio	26	Animal
8	Outdoor	27	Computer-TV
9	Building	28	Flag-US
10	Desert	29	Airplane
11	Vegetation	30	Car
12	Mountain	31	Bus
13	Road	32	Truck
14	Sky	33	Boat-Ship
15	Snow	34	Walking-Running
16	Urban	35	People-Marching
17	Waterfront	36	Explosion-Fire
18	Crowd	37	Natural-Disaster
19	Face	38	Maps
		39	Charts

Carnegie Mellon

Annotated Concept Sets

- Trecvid 2006 development data
 - ~70hours English, Arabic, Chinese News
 - 62000 shots
- 3 Annotated Concept Sets:
 - LSCOM Lite
 - 39 concepts
 - Media Mill
 - 75 concepts that overlap with LSCOM
 - LSCOM
 - 300 concepts
 - Minimal frequency cutoff

Carnegie Mellon

A Speculative Scenario

What if we could detect a lot of concepts?

Best Case:

- Perfect concept detection (Oracle)
- Perfect concept combination (Oracle)

How well can you retrieve relevant shots (documents)

Carnegie Mellon

Computing the Optimal Weights

$$LLB = \max_{\lambda_i} AP \left(\sum_{i=1}^k \lambda_i f_i(D, Q) \right)$$

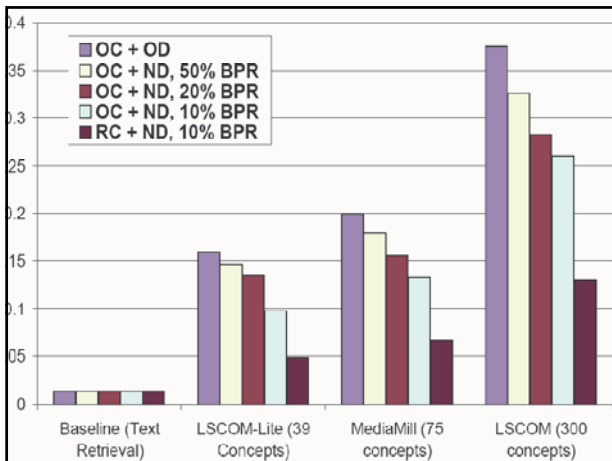
This is an upper bound under perfect conditions

Carnegie Mellon

Assumptions for Realistic Evaluation

- Noisy Detection - Assume MAP of 0.5, 0.2 or 0.1
 - TRECVID 2007 semantic concept detection results just under 0.2
- Approximate through precision/recall breakeven point at 50%, 20%, 10%
- Realistic combination (RC)
 - achieves only 50% of perfect (Oracle) combination
 - Based on prior TRECVID experiments

Carnegie Mellon



Extrapolation to thousands of concepts

- How many concepts do we need?
- Exponential function to fit curve
 - Only 3 points

Assumption:

- Things get harder as you add more concepts
 - Proportional to the difference between the current MAP and the upper limit of 1
i.e. the higher the current MAP, the less benefit a new concept offers

Carnegie Mellon

Extrapolation Assumptions

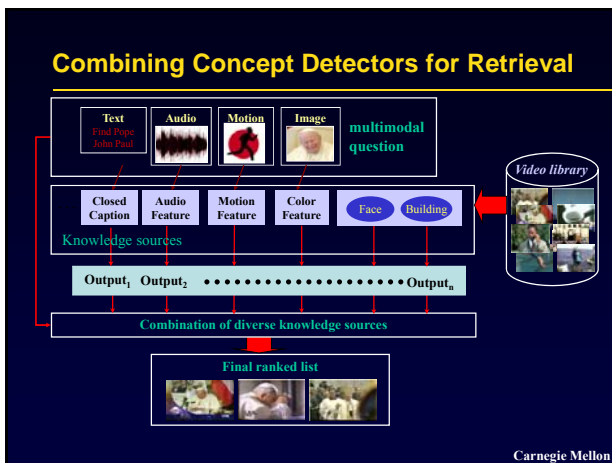
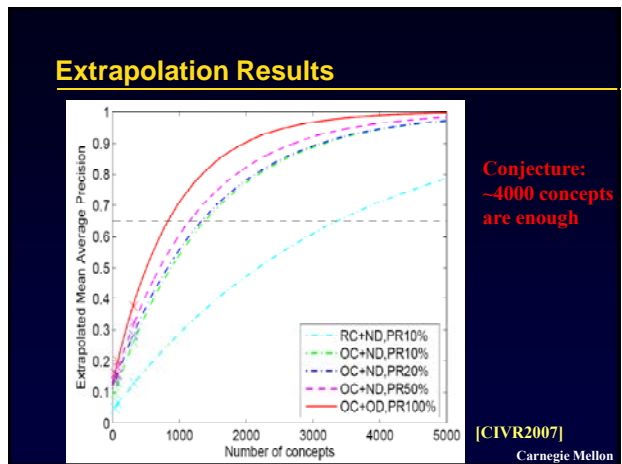
$$\frac{dm}{dx} \propto (1 - m)$$

where m is the value of MAP, x is the number of concepts, and the boundary condition is $m(\infty) = 1$.

This yields:

$$m(x) = 1 - \exp(ax + b)$$

Carnegie Mellon



Probabilistic Model for Video Retrieval

- What is the probability *this* shot (document) is an answer for *this* query?

Estimate $P(Y=1|D,Q)$

- Basic model: **logistic regression** for multimedia IR
 - For every document D and query Q ,

$$P(Y = 1 | D, Q) = \sigma \left(\sum_{i=1}^N \lambda_i P(S_i | D, Q) \right)$$

where $\sigma(x) = (1 + e^{-x})^{-1}$

- Learn the combination weights λ_i from past retrieval results

Carnegie Mellon

“Ranking” Logistic Regression

- Problem: logistic regression is designed for **classification**
 - Simplify the QA problem to binary classification
 - Several drawbacks, e.g. bias to irrelevant documents

- Our model: “ranking” logistic regression

$$\max_{\lambda} \sum_{q \in Q} \sum_{d \in D} \sum_{d_2 \in D} \log \sigma \left(\sum_{i=0}^N \lambda_i [P(S_i | d_1, q) - P(S_i | d_2, q)] \right)$$

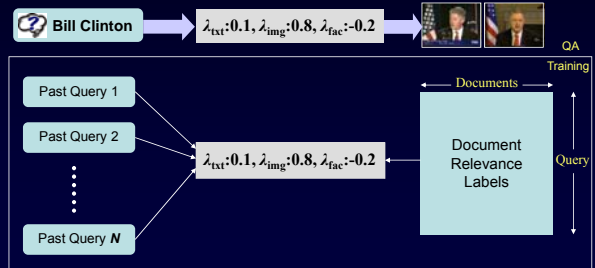
- Maximize the **prediction gaps** between positive/negative examples
- This is a lower bound of the **average precision**

- Approximate inference
 - Reduce from $O(n^2)$ to the same complexity as logistic regression

Carnegie Mellon

Query Independent Combination

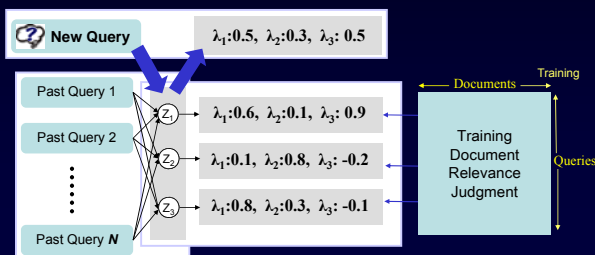
- Learn a single combination function from training queries
- Uniform** combination function for all other queries



Carnegie Mellon

Probabilistic Latent Query Analysis (pLQA) [Yan06]

- Each query can be described with a mixture of latent query classes
- Each query class is associated with a set of combination weights
- Query description is used to predict to which query class it belongs



Carnegie Mellon

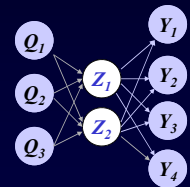
Basic pLQA (BpLQA)

- BpLQA: learn the query classes directly from training data
 - Introduce new latent class variables Z in the logistic regression model
 - Query classification and combination optimization in a unified framework

$$P(y | D, Q) = \sum_z P(z | Q; \mu) \cdot \sigma \left(y \sum_{i=0}^N \lambda_i f_i(D, Q) \right)$$

- Advantages over query-class combination

- No **manual definition** on query classes
- Allow **mixture** of query classes for one query
- Able to find **optimal number** of query classes using model selection criteria, e.g. AIC, BIC



Carnegie Mellon

Adaptive pLQA (ApLQA) Extension to Unseen Queries

- Problem:** $P(z|Q;\mu)$ is associated with each training query
- Goal:** extend basic pLQA to handle unseen queries outside the training collection
- ApLQA: represent $P(z|Q;\mu)$ as a soft-max function over a vector of predefined query features q_i
 - e.g., whether the text query contain a specific person name
 - e.g., how many objects are mentioned in the query
 - e.g., the mean of text retrieval scores

$$P(y | D, Q) \propto \sum_z \exp \left(\sum_{i=0}^L \mu_i q_i \right) \cdot \sigma \left(y \sum_{i=0}^N \lambda_i f_i(D, Q) \right)$$

Carnegie Mellon

Parameter Estimation for ApLQA

- Bayesian information criterion (BIC) to select the optimal number of query classes
- Maximum likelihood estimation using the EM algorithm
 - Initialize combination weights from automatically selected queries
 - E-step:** compute the expectation of latent query type Z

$$P(z_i | Q, D) = \frac{\exp \left(\sum_{i=0}^L \mu_i q_i \right) \cdot \sigma \left(y \sum_{i=0}^N \lambda_i f_i(Q, D) \right)}{\sum_{i=0}^L \exp \left(\sum_{i=0}^L \mu_i q_i \right) \cdot \sigma \left(y \sum_{i=0}^N \lambda_i f_i(Q, D) \right)}$$

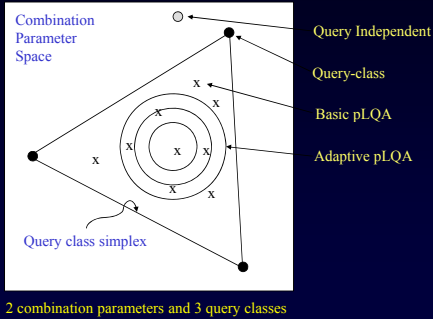
- M-step:** estimate the parameters using posterior probabilities of Z

$$\lambda_i = \arg \max_{\lambda} \sum_{i=0}^N \sum_{j=0}^L P(z_j | Q, D) \log \left[\sigma \left(y \sum_{i=0}^N \lambda_i f_i(Q, D) \right) \right]$$

$$\mu_j = \arg \max_{\mu} \sum_{i=0}^L P(z_i | Q, D) \log \left[\frac{1}{Z} \exp \left(\sum_{i=0}^L \mu_i q_i \right) \right]$$

Carnegie Mellon

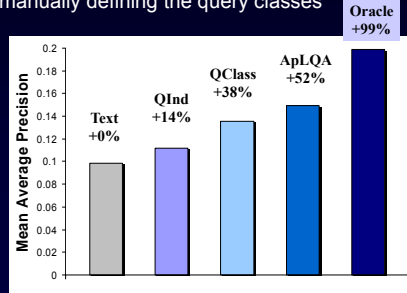
Geometric Interpretation



Carnegie Mellon

Video Retrieval Performance

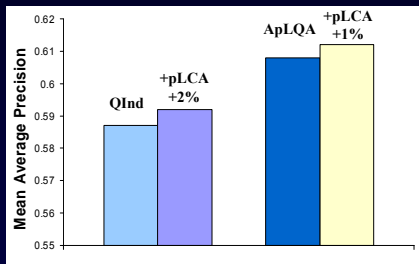
- Average retrieval performance on TREC'02-'05 test sets
- ApLQA can achieve a higher performance without manually defining the query classes



Carnegie Mellon

Meta-Search Performance

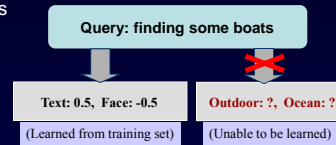
- Performance of pLQA on TREC-8 testing sets
 - pLQA on 2 retrieval algorithms w. 5 additional ranking features



Carnegie Mellon

Limitations of pLQA

- Due to a **limited amount** of training data, the weights of many ranking features are simply set to zero
- More than 80% of the ranking features are ignored in TREC'05
- However, these features are useful for some specific queries

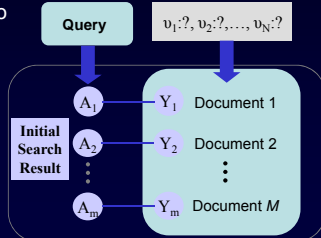


Let's try relevance feedback (local content analysis)

Carnegie Mellon

Probabilistic Local Context Analysis (pLCA) [Yan06]

- Goal:** automatically leverage useful features for **current query**
- Method:** assume combination parameters u of "unlearned" ranking features to be **latent variables** instead of zero



Carnegie Mellon

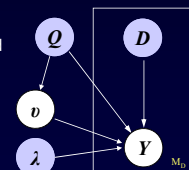
Model and Parameter Estimation

- Jointly optimize **weights of query-specific sources u** and **answer relevance Y** in initial retrieved shots

$$\max_{\gamma, \nu} P(\nu | Q) \cdot \prod_{j=1}^{M_0} \sigma \left(\sum_{i=0}^N y_i \lambda_i P(S_i | D_j, Q) + \sum_{i=1}^M y_i \nu_i P(S_i^* | D_j, Q) \right)$$

(Prior) (Query Independent) (Query Specific)

- Iteratively maximize **Y** and **u**
 - Y :** top-ranked shots as positive and others as negative
 - u :** feature selection + regularized logistic regression



Carnegie Mellon

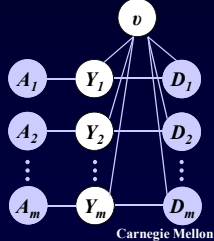
pLCA: Undirected Model and Parameter Estimation

- Compute the posterior probability of document relevance Y given initial results A based on an undirected graphical model

$$P(\bar{y} | \bar{a}; \bar{D}, Q) = \frac{1}{Z} \int \prod_l P(v_l | Q; v_l^0) \cdot \prod_{j=1}^{M_y} \exp\left(y_j a_j + y_j \sum_l v_l f_l(D_j, Q)\right) dv$$

- Variational inference

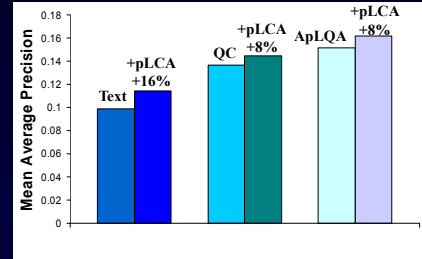
- Approximate the posterior distribution of $p(y, v | a)$ by a family of variational distributions $q(y, v)$ where y and v are independent
- Iteratively maximize a variational lower bound of the log-likelihood function until it converges



Carnegie Mellon

Video Retrieval Performance

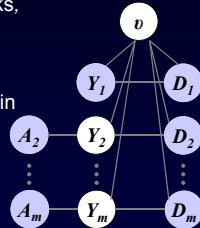
- Retrieval performance of pLCA on TREC'03-05
 - pLCA on 3 retrieval algorithms w. ~50 additional ranking features



Carnegie Mellon

pLCA with Relevance Feedback

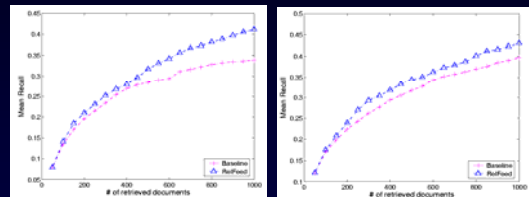
- Users could provide their own relevance judgments in an interactive retrieval interface.
- Given these relevance feedbacks, we can use a similar variational inference technique to update parameters except that some relevance variables Y_i are fixed in this case.



Carnegie Mellon

Results of pLCA Relevance Feedback

- Performance in terms of average recall
- Update combination parameters every 50 shots



Baseline: text retrieval

Baseline: ApLQA

Carnegie Mellon

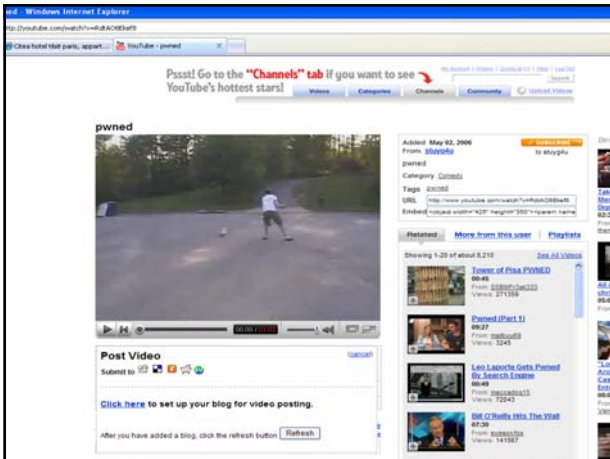
Extreme Video Retrieval (Interactive)

- Maximizing the Synergy between Man and Machine
- Maximizes information transfer: System \rightarrow Human
 - Up to 10 keyframes/second
 - 1 - 9 images per page (dynamically adjustable)
 - Presentation intervals
 - Automatic: System presents next page after X msecs (adjustable)
 - Manual: Press key for next page
 - Click where relevant shot is seen
- Relevance feedback based on user clicks
- Automatic retrieval baseline for ranking order
- Requires full concentration: 3000-5000 keyframes in 15 minutes
- Second best interactive system in TrecVid 2005
- Fourth best in TrecVid 2006
- Best system in Video Olympics 2007 (NUS/CAS)

Carnegie Mellon

Challenges for the Future

Carnegie Mellon



YouTube

- Annotations as user tags and descriptions
- Not very accurate, but abundantly available
- Exploit other information
 - Author
 - User
 - Comment
 - Video Responses
- Google scale

Carnegie Mellon

Near Term Research (5 -10 years)

- Leverage the web
 - 'Public' data
 - Social tagging
 - Games
- Bridging the Semantic Gap (perhaps using Concepts)
 - Better retrieval
- Human computer interaction
 - better AV search interfaces
- Combining world knowledge with AV analysis for search and summaries
 - Ontologies?
- "Event Search" (R. Jain)
- Privacy, digital rights
- Data Integrity
- Data scale increased to 'Google-scale'

Carnegie Mellon

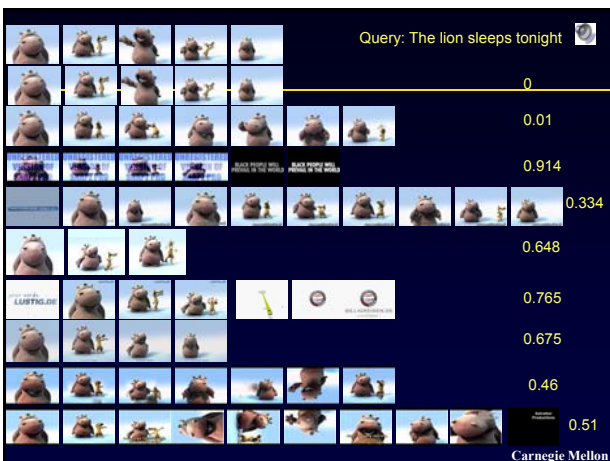
Data Integrity of Web Content

Near-duplicate and tamper detection for security, copyright, ease of browsing

Types of Near-Duplicates and Tampering:

- Reformatting (frames rate, frame size,
- Mixing old video with new audio (or vice versa)
- Image variations: color, lighting, cropping of frame
- Adding/Deleting frames, mostly at the head and tail
- Insertion of new pieces
- Re-editing for different versions: 1'08", 2'50"
- Overlay text and logos
- Recomposition of image components

Carnegie Mellon

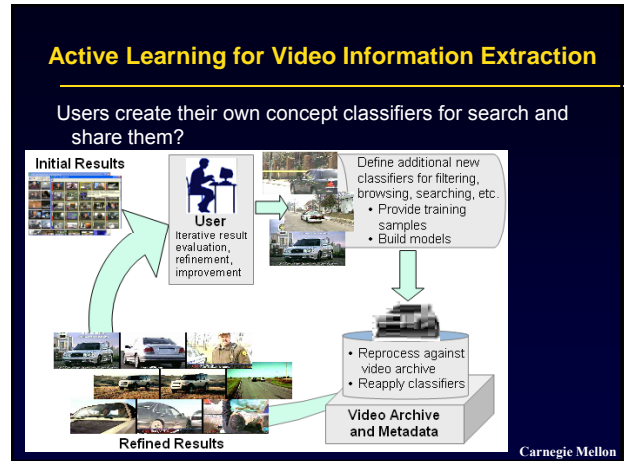


Carnegie Mellon

Human Computer Interaction: Better AV Search Interfaces

- Principles for dealing with errorful analysis
- Relevance Feedback, active learning, incremental learning
- Collective learning (ESP Game)
- User profiles
 - Never happens beyond a trivial way
 - No long term users in research, systems change too fast

Carnegie Mellon



- ### Combining World Knowledge with AV Analysis
- Better understanding for better search
 - What is the program (content)
 - What is new, what is redundant
 - What does the user want
 - What summaries are useful, effective and can be automatically generated
 - Keyframes, storyboard, skims, text summaries, ...
 - In text analysis, ontologies (i.e. world knowledge) have NOT been a big success story!
 - This doesn't work at any reasonable scale
- Carnegie Mellon

- ### Ontologies Won't Solve These Problems
- Hand-constructed
 - As good as person building it
 - Brittle with missing and extraneous information
 - Don't match the data in the archive
 - Can't deal with the high error rates from AV analysis
-
- Carnegie Mellon

- ### Event Search (R. Jain)
- Richer analysis of importance and temporal aspects
- Calendar, blog, webpage, both "push" and "pull" technology
- Not clear exactly how it should work
- Carnegie Mellon

- ### Exploiting the Next Level of Scale
- GB connectivity everywhere
 - Terabyte personal storage
 - Device size, mobility/location, ubiquity
 - Partnering with Industry
 - to get large data sets,
 - long term studies,
 - fielded applications
 - with appropriate hardware
- Carnegie Mellon

Immediate challenges

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

Carnegie Mellon

Immediate challenges

- How to leverage concept detectors for search?
 - How to present detectors to users?
 - How to select the correct detectors?
 - How to combine concept detectors?
 - How to combine concept selection methods?
- What if there is no text to search
 - Consider home video domain for example
- How to balance semantic coverage and anticipated performance of detectors for a specific query?

Carnegie Mellon

Economic Challenges

- Money (e.g. from advertising like Google)
- Security surveillance (industrial, public spaces)
- DVR, TV, cell-phones
 - large market saturated with high resolution, large screen devices, high bandwidth connections and nothing worthwhile to display

Automatic analysis for search concepts
– At least for “computer assisted” search

Carnegie Mellon

Social Challenges

- Personal recordings
 - Multimedia phones (with camera, video recorder, GPS)
 - Seagate: Terabyte on the person in 5+ years
 - Digital cameras, camcorders
 - Add GPS, time, biometrics, temperature, other sensor data
- Web2.0: YouTube, Google Video, ...

How to archive, categorize, search and leverage this flood of data

Carnegie Mellon

Political and Cultural Challenges

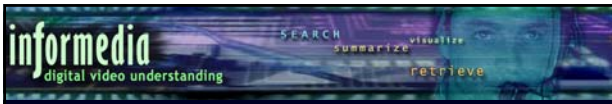
- Cultural archives
 - Preserve cultural identity in homogenizing, global world
- Access to news and information from other countries, languages
 - Broadening Perspectives
 - Communication and understanding across cultures and societies
 - Better informed citizens → more stable democracies
- Intelligence/Military Analysis
 - Foreign news imagery analysis
 - Satellite imagery, unmanned aerial vehicles
 - Web video
- Healthcare
 - Analysis of imagery
 - Combine with other sensor data
 - Longitudinal video (not just AV, but other sensors also)
 - Reduce costs
 - Unified, verifiable records

Carnegie Mellon

The Fundamental (BIG) Challenge

- Understanding video in general way
 - How do we see? Understand?
 - Given any video, describe it in words that are useful for search and summary
- Can computers be better than people
 - Calculators, chess, spelling (?), ... AV analysis and search
- Knowing when we don't know an answer
- Mind reading
 - understand a search question in the user's mental context

Carnegie Mellon

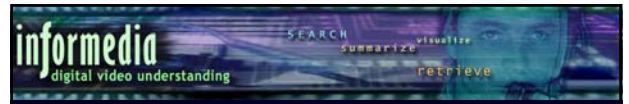


informedia
digital video understanding

SEARCH summarize visualize retrieve

Questions?

Carnegie Mellon



informedia
digital video understanding

SEARCH summarize visualize retrieve

Questions?

Carnegie Mellon