Recent Advances and Open Issues of Digital Image/Video Search

Prof. Shih-Fu Chang

Digital Video and Multimedia Lab Columbia University

www.ee.columbia.edu/dvmm June 2007

Acknowledgement

Columbia University

 Winston Hsu, Wei Jiang, Lyndon Kennedy, Akira Yanagawa, Eric Zavesky, Dongqing Zhang

Kodak Research

- Alex Loui, Jiebo Luo
- IBM Research
 - Shahram Ebadollahi, Milind Naphade, Paul Natsev, John R. Smith, Lexing Xie
- CMU
 - Alex Hauptmann

Need of Image/Video Search

- Explosive growth of online image/video data, personal media, broadcast news videos, etc.
- 5 billion images on the Web, 31 million hours of TV programs each year
- Successful services like Youtube and Flickr
- Image/video search exciting opportunity





However... User Expectation in Practice

- "...type in a few words at most, then expect the engine to bring back the perfect results. More than 95 percent of us never use the advanced search features most engines include, ..." – *The Search*, J. Battelle, 2003
- Keyword search is the primary search method.
- User input is often limited, thus
 - Difficult to get detailed visual descriptions
 - Difficult to get user feedback for refined search



Google Zeitgeist publishes top keywords monthly

lewsmakers	Getaways	
 virginia tech knut yuri gagarin shaha riza kurt vonnegut 	 <u>hawaii</u> <u>dubai</u> <u>mexico</u> <u>chelsea</u> <u>london</u> 	
 reece - Top Gaining Q 1. <u>notis</u> 2. <u>holly valance</u> 3. <u>shrek</u> 	tueries: April 2007 6. <u>χαρτης</u> <u>θεσσαλονικης</u> 7. <u>ΜΕΤΑΦΡΑΣΗ</u> 8. <u>land rover</u>	 11. <u>chelsea fc</u> 12. <u>manchester united</u> 13. <u>Μύκονος</u> 14. <u>πλουταρχος</u>



Examples of Keyword Image Search



Example Search

• Text Query on Google: "Manhattan Cruise"



Image content analysis may help refine results





From amandajin





From nosilla q

Insufficient Precision of Social Tags

		precision
	Bronx-Whitestone Br.	1.00
York	Brooklyn Br.	0.38
	Chrysler Building	0.65
mark	Columbia University	0.30
	Empire State Building	0.18
S	Flatiron Building	0.70
	George Washington Br.	0.48
	Grand Central	0.37

Test

New

City

land

labe

Many tags from social networks are of low precision

(due to batch uploading?)

Times Square	0.56
Verrazano Narrows Br.	0.66
World Trade Center	0.13

An Interesting Paradigm: Image Tagging via Game Playing

(Von Ahn & Dabbish, CHI 04)

- Used in Goggle Image Labeler
- Use competitive games to motivate users
- Has attracted many participants for free!
 - Some users spent hours in a day
- Claim the potential of annotating the whole Web in just few months!
 - 5 Billion images





- Affected by personal background and proficiency
- Matched words tend to be obvious ones
 - Scoring system encourages guessing partner's thoughts, rather than indexing image content
 - Players sometimes just perform manual OCR
- An interesting paradigm, but perhaps not a complete s.-F. Chang, Communication
 12

Opportunity for Content Analysis: Large-Scale Auto. Image Tagging Framework

- Audio-visual features
- Surrounding text
- SVM or graph models
- Context fusion

 Rich semantic description based on content analysis



Statistical models

Large-Scale Concept Detectors
Publicly Available

- Columbia374
 - 374 baseline detectors for LSCOM multimedia ontology
- MediaMill
 - 491 concept detectors for LSCOM and MediaMill 101 Lexicons
- IBM MARVEL Search System
 - Trials with BBC, CNN
 - Real-time standalone detectors from IBM AlphaWorks
- Others ...

What Concept to Detect?

- One effort: Large Scale Concept Ontology for Multimedia (LSCOM)
 - Joint effort by news/intelligence analysts, librarians, researchers
 - Broadcast News Domain
 - Selection Criteria
 - useful, detectable, observable
 - 834 concepts defined, 449 concepts annotated
 - Labeled over 61,000 shots of TRECVID 2005 data set
 - **33** Million judgments collected, 100 person-month labor
 - Download by 170+ groups so far
 - http://www.ee.columbia.edu/dvmm/lscom/

LSCOM data set downloaded by 170+ groups

- Yahoo! Research
- Intel
- AT&T
- FXPAL
- University of Amsterdam
- Oxford University
- Nanyang Technological University, Singapore
- National Taiwan University
- Tsinghua University
- KDDI, Japan
- Dublin City University, Ireland
- University of Central Florida
- University of Texas, Austin
- UC Berkeley
- Others ...

Example LSCOM Concepts (449)

- Event/Activity (56 13%)
 - Airplane taking off, car crash, explosion, etc
- People (113 25%)
 - Person, male/female, firefighter, etc
- Location (89 20%)
 - Cityscape, hospital, airfield, etc
- Object (135 30%)
 - Vehicle, map, tank, power plant, etc
- Scene (49 10%)
 - Vegetation, urban, interview, etc
- Program (7 2%)
 - Entertainment, weather, finance, etc

Another Effort : Consumer Video Ontology (Kodak-Columbia, 2007)

	 Activity (6) Occasion (16) Scene (15) Object (25) People (11) Sound (14) 	Activity: Occasion : Scene: Object: People:
ł	Object Motion (3) Social (4)	Camera Motion: Object Motion:
Lexicon and annotation over 1300 consumer videos planned		

How to Obtain High-Quality Annotations?



- Label only one concept at a time
 - Instead of open text simultaneous labeling
 - Found to be more accurate
 - And faster!
- But hard to scale up
 - Throughput: 1-3 sec/image
 - A laborious process (100 person-month) for 33 M labels
 - No user incentive

Building Image Classifiers – starting from the basics



- General for all concepts, easy to implement
- Late fusion better than early fusion
- 374 baseline detectors (Columbia 374) released

Examples of Basic Image Features

grid layout + color moment



s.-F. Chang, Columbia ensions

Gabor texture





48 dimensions

edge direction histogram





73 dimensions

9-

http://www.ee.columbia.edu/ln/dvmm/columbia374/

Columbia374 Download (Version 1.0)

www.ee.columbia.edu/dvmm/columbia37

Google

✓ ++ ×

-

Columbia374

Columbia University's Baseline Detectors for 374 LSCOM Semantic Visual Concepts

Summary

Semantic concept detection represents a key requirement in accessing large collections of digital images/videos. Automatic detection of presence of a large number of semantic concepts, such as "person," or "waterfront," or "explosion", allows intuitive indexing and retrieval of visual content at the semantic level. Development of effective concept detectors and systematic evaluation methods has become an active research topic in recent years. For example, a major video retrieval benchmarking event, NIST TRECVID[1], has contributed to this emerging area through (1) the provision of large sets of common data and (2) the organization of common benchmark tasks to perform over this data.

However, due to limitations on resources, the evaluation of concept detection is usually much smaller in scope than is generally thought to be necessary for effectively leveraging concept detection for video search. In particular, the TRECVID benchmark has typically focused on evaluating, at most, 20 visual concepts, while providing annotation data for 39 concepts. Still, many researchers believe that a set of hundreds or thousands of concept detectors would be more appropriate for general video retrieval tasks. To bridge this gap, several efforts have developed and released annotation data for hundreds of concepts [2, 5, 6].

Quick Guide to Columbia374

Columbia374 Citation:

Akira Yanagawa, Shih-Fu Chang, Lyndon Kennedy and Winston Hsu, "Columbia University's Baseline Detectors for 374 LSCOM Semantic Visual Concepts", Columbia University ADVENT Technical Report #222-2006-8, March 12, 2007. [pdf]

1. Visual Faetures & Lists

Download the Visual Factures and the lists in Columbia374. (219 MB file. Expands to 695 MB on disk.)

Models for LIBSVM (Ver. 2.81)

Download the Models trained by LIBSVM (Ver. 2.81) in Columbia374. (3.5 GB file. Expands to 4 GB on disk.)

3. Scores

Performance of baseline models in TRECVID 2006



Performance Metric – Average Precision (AP)

#	hit	R	Р	
I ₁	+	0.10	1.00	
	_	_		
I ₃	+	0.20	0.67	
I ₄	+	0.30	0.75	
۱ ₅	- <u>-</u>	<u> </u>	<u>-</u>	
I ₆	-	-	-	
۱ ₇	+	0.40	0.57	
۱ ₈	-	-	-	
I ₉	-	-	-	
I ₁₀	+	0.50	0.5	
J ₁₁	-	-	-	

- Rank detection results
- Compute precision/recall values at each hit point
- Average precision at different recall values
- Mean average precision (MAP)
 - Mean of APs across concepts
 - Top TRECVID system MAP ~0.3



Comparison with Text-based Classifiers

- Extract text features from closed caption, ASR or translated ASR (for foreign videos)
- Apply classifiers : Naïve Bayes, SVM, Max. Entropy ...



Text classifiers usually lower than visual models by 3-5 times.

Comparing Semantic Classification (text vs. visual) (1)



Visual concept search – "boat"

(images from TRECVID)

Comparing Semantic Classification (text vs. visual) (2)



Concept search results – "car"

Example: good detectors for LSCOM concept

bridge waterfront crowd explosion fire US flag Military personnel Search Results Search Results: Search Deputter Search Results: WHAT WETDOW LODE NUMBER OF THE PARTY MHAT VETOMY LOOKS IS NAME PERMIT MOST PALLILY, NY INVIE DECOME THE



Complexity

Slow training, fast detection

- Training time over 61K shots per SVM
 - 20 mins (color), 3.5 mins (texture), 1.2 mins (edge)
 - About 1 month over 20 PCs to train Columbia374
- Feature extraction
 - 0.4 sec (color), 0.4 sec (texture), 2.8 sec (edge) per image
- Testing
 - < 30 ms per image</p>



Advanced Features and Models



Representation and Learning





(Zhang & Chang CVPR 06)

Learning Graph Object Model



BUT

- Finding the correspondence of parts and computing matching probability are NP-complete
- Use advanced machine learning methods: Loopy Belief Propagation, and Gibbs Sampling plus Belief Optimization



Boosted Conditional Random Fields



S.-F. Chang, Columbia U.

Predict When will Context Help

Concept Fusion (CF) does not always help:

- Complex inter-conceptual relationships difficult to estimate from limited training samples
- Strong classifiers may suffer by fusion with inaccurate context

Use CF for concept C_i , only if C_i is weak or with strong context

$$E(C_i) > \lambda$$
 or

weak self

$$\frac{\sum_{C_j, j \neq i} I(C_j; C_i) E(C_j)}{\sum_{C_j, j \neq i} I(C_j; C_i)} < \beta$$

Strong context

 $I(C_i; C_j)$ -- mutual information between C_i and C_j $E(C_i)$ -- error rate of independent detector for C_i

Fuse Context Only When it is Predicted to Help

- Apply context fusion smartly
- 16 concepts are predicted to improve
- 14 actually improve, 4 missed





-36- >>digital video | multimedia lab>>>

Example: Road Detector

Road



digital video | multimedia laboratary



Combine Context-Fusion with User Input

automatic content recognition
 +
 context fusion
 +
 user input

Instead of Automatic Tagging 100% of Concepts



Tags:Person: 79%Clinton: 75%US Flag: 80%Podium: 70%Give speech: 75%Press conference: 65%

What if we can ask users to help 1-2 labels?

Active Tagging





(Jiang, Chang, Loui, ICIP 06)



ICIP 2006

Examples of Best Questions



Best Questions to Ask User?

Active Tagging Consistently Help

- Performance gain increases with more user input
- Magic number here is 4



Power of Concept-based Representation





-43- >>digital video | multimedia lab>>>

Power of Semantic Indexing: Text-to-Concept Search



Mapping search topics to concepts

TRECVID search topics



Concept Search Demo

- Concept search case 1 (<u>link</u>)
- Concept search case 2 (<u>link</u>)
- Multimodal search (<u>link</u>)

Other Applications: Semantic Mining

Top results from text query: "Manhattan Cruise"



Decipher dominant visual concepts behindeach search topic



Related concepts have high mutual information with target labels.

Mutual information:

$$I(T;C) = \sum_{T} \sum_{C} P(T,C) \log \frac{P(T,C)}{P(T)P(C)}$$

- T: initial text search score
- C: concept detection score (quantized)
- Include both positive and negative correlations

Examples: discovered concepts per query

Query Topic	Positive Concepts	Negative Concepts
(158) Find shots of a helicopter in flight.	Weapons, Explosion Fire, Airplane, Smoke, Sky	Person, Civilian Person, Face, Talking, Male Person, Sitting
(164) Find shots of boats or ships.	Waterscape, Outdoor, Sky, Fire, Exploding Ordnance	Person, Civilian Person, Face, Adult, Suits, Ties
(179) Find shots of Saddam Hussein.	Court, Politics, Government Leader, Suits, Lawyer, Judge	Desert, Protest, Crowd, Mountain, Outdoor, <mark>Animal</mark>

Three main types of topic-concept relationships:

Generally Present Obvious, expected relationships. Concepts that are always together.

News Story Relationships unique to news story. Reranking uniquely powerful.

Mistaken Relationships due to detection errors. Rare, but still helpful.

Another Application: Event Detection

 Typical approach to event detection: detect object of interest, track over time, model spatio-temporal dynamics



Difficult for events without explicit objects and motions, e.g., "riot", "street battle", "flight combat", "parade"





New Approach: Concept-based Event Representation (1)





- Map image frames to confidence scores in the semantic concept space
- A video is associated with a bag of concept features



bag of concept features

Concept-based Event Representation (2)



Video Events as Distinct Traces in the Concept Space







Red: video 1, Green: video 2, Yellow: other videos

DVMM Lab, Columbia University

From Representation to Event Detection



Earth-Mover's Distance (EMD)

[Xu & Chang, CVPR 07]

Matching? Video 2 Video 1 concepts concepts **EMD** finds optimal alignment flows between - - -. . . features in the bags



-56- >>digital video | multimedia lab>>>

HMM+SVM to Model Event Traces in the Concept Space

[Xie et al, ICME 06]





-57- >>digital video | multimedia lab>>>

Performance Testing Using TRECVID Data





Significant performance gain by applying temporal models (EMD and HMM)

Landscape of Semantic Image/Video Indexing



Open Issues and Opportunities

- Systematic Way of Extending Ontology?
- Cross-Domain Model Development:
 - General detectors for News, Consumer, Web?
 - Sharing of data and models?
- Complexity and speed
 - Feature selection, fast training
- Map user queries to visual concepts
 - Wordnet for visual search?

More Information

Columbia DVMM Lab

- http://www.ee.columbia.edu/dvmm
- 374 SVM-based concept detectors available soon
- Online video search demos
- LSCOM lexicon and annotation
 - http://www.ee.columbia.edu/lscom
- Columbia374 concept detectors
 - http://www.ee.columbia.edu/columbia374

Columbia Video Search System

http://www.ee.columbia.edu/cuvidsearch



Prototype includes 160 hours, 3 languages (English, Arabic, Chinese), 6 channels