



Computational Structures in Data Science



UC Berkeley EECS
Adj. Ass. Prof.
Dr. Gerald Friedland

Lecture #6: 'Guest Lecture'

September 30, 2016

<http://inst.eecs.berkeley.edu/~cs88>

Driving Scientific Discovery with Social Media Images and Videos



Groovin'. Snowball dances to the beat. Irena Schulz, Bird Lovers Only Rescue Inc

Dr. Gerald Friedland
fractor@icsi.berkeley.edu

Consumer-Produced Videos are Growing in the Internet

- YouTube claims 65k 100k video uploads per day or 48 72 hours every minute
- Youku (Chinese YouTube) claims 80k video uploads per day
- Facebook claims 415k video uploads per day!

3

Why do we care?

Consumer-Produced Multimedia allows empirical studies at never-before seen scale.



Spontaneous motor entrainment to music in multiple vocal mimicking species
A Schachner, TF Brady, IM Pepperberg, MD Hauser - Current Biology, 2009

4

Google "giving directions to a location"

Web Images Maps **Videos** News More Search tools

3 results (0.13 seconds)

Ad related to "giving directions to a location"

Maps & Driving Directions
driving-directions.easymaps.co/
 Enter Address Or Location. Free Maps & Directions w/Toolbar!

Key considerations for all maps from the Course Creating a Map with Illu...
www.lynda.com > ... > [Creating a Map with Illustrator](#)
 Are you giving directions to a location, or general information about the area. How much of the area should ...

Community Helpers - SlideServe
www.slideserve.com/kagami/community-helpers
 Aug 3, 2012
 This will help with giving directions to a location. Materials: Our maps A step by step direction route on chart ...

PPT – A Study on Wearable Computing PowerPoint presentation | free to download
www.powershow.com/.../A_Study_on_Wearable_Computing_powerpol...
 Technology which allows for the human and ... The concept of wearable computers attempts to bridge the 'interaction gap' ... Sprout. Spot. 17 /18. Conclusion .

Stay up to date on these results:
 • [Create an email alert for "giving directions to a location"](#)

Driving Directions & Maps
www.maps-directions.org/Directions
 Enter Starting Point & Destination. Get Directions. Quick & Easy.

YP.com Maps and Directions
www.yellowpages.com/
 Find & Discover Local Businesses on YP.COM

Directions To And From
www.getdrivingdirections.co/Directions
 Enter Address or Driving Location. Driving Directions & Maps w/Toolbar

Accept Online Donations
www.securegive.com/OnlineGiving
 Turn your giving into your website to grow your giving today!

Location Maps
www.myhomemsn.com/
 Get Access To Maps & Directions. Make MSN Your Homepage Today.

Map Quest Directions
shopping.yahoo.com/Books
 Great Deals on Map Quest Directions Shop Now and Save. Yahoo Shopping

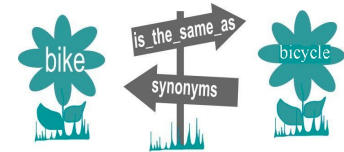
[See your ad here >](#)

Challenges I

User-provided tags are:

- sparse
- any language
- imply random context

Solution: Use the actual audio and video content for search.



Challenges II

Research to search the actual audio and video information is hindered by:

- YouTube videos not legally downloadable
- No reliable annotation
- Search in YouTube doesn't work (see Challenges I...)



The Multimedia Commons

100.2M Photos
800K Videos

Features for Machine Learning
(Visual, Audio, Motion, etc.)

User-Supplied Metadata
and New Annotations

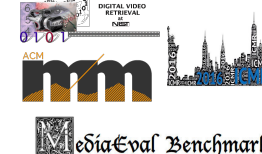
Tools for Searching,
Processing, and Visualizing

100M videos and images, and a growing pool of tools for research with easy access through Cloud Computing

Collaboration Between Academia and Industry:

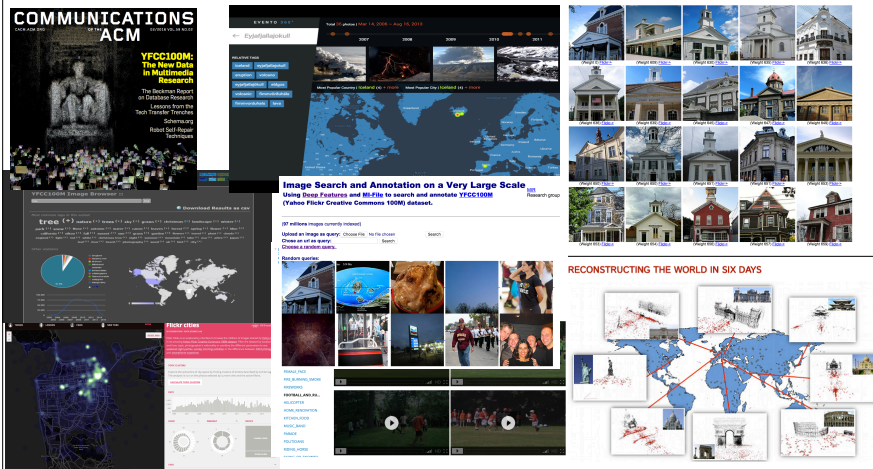


Benchmarks & Grand Challenges:



Supported in part by NSF Grant 1251276
 *BIGDATA: Small: DCM; DA: Collaborative Research;
 SMASH: Scalable Multimedia Content Analysis in a High-level Language"

The Multimedia Commons: An Open Infrastructure for Large-Scale Multimedia Research



<http://mmcommons.org>

Challenges II



Challenges I



10

Work on Multimedia Content Retrieval

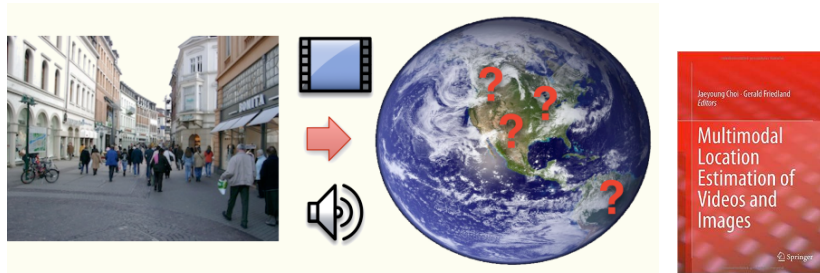
- Computer Vision: Focus on solving the AI problem, e.g. through object labeling
- Video Retrieval:
 - Computer Vision techniques
 - Motion
 - Audio
 - Metadata

Our Approaches to Content-based Video Search

- Focus on events (time and location)
- Combine text and image/video similarity searches and event search
- Try to 'translate' multimedia data into text

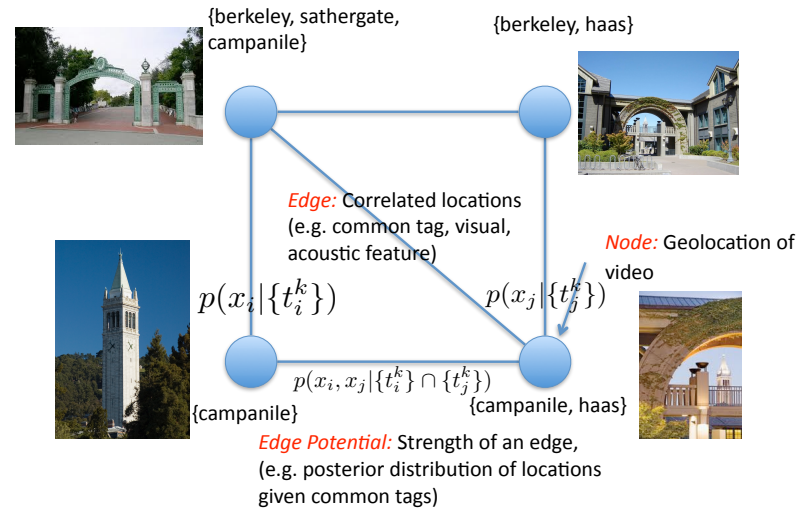
12

Events: Multimodal Location Estimation



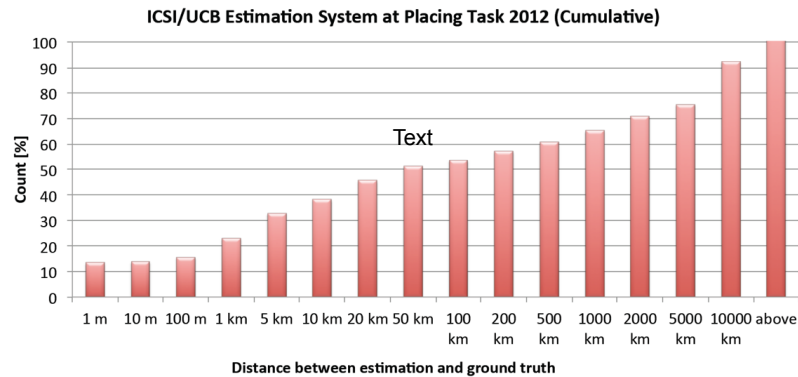
<http://mmle.icsi.berkeley.edu>

Intuition for the Approach



MediaEval Benchmark

MediaEval Benchmarking Initiative for Multimedia Evaluation
 The "multi" in multimedia: speech, audio, visual content, tags, users, context



J. Choi, G. Friedland, V. Ekambaram, K. Ramchandran: "Multimodal Location Estimation of Consumer Media: Dealing with Sparse Training Data," in Proceedings of IEEE ICME 2012, Melbourne, Australia, July 2012.

An Experiment

Listen!

- Which city was this recorded in?
 Pick one of: Amsterdam, Bangkok, Barcelona, Beijing, Berlin, Cairo, CapeTown, Chicago, Dallas, Denver, Duesseldorf, Fukuoka, Houston, London, Los Angeles, Lower Hutt, Melbourne, Moscow, New Delhi, New York, Orlando, Paris, Phoenix, Prague, Puerto Rico, Rio de Janeiro, Rome, San Francisco, Seattle, Seoul, Siem Reap, Sydney, Taipei, Tel Aviv, Tokyo, Washington DC, Zuerich
- Solution: Tokyo, highest confidence score!

Evento360: Search with Combined Textual, Visual, and Acoustic Features



'Translate Multimedia': Scenario

Empirical Study: How do Children learn to catch a ball?



Example Video



<https://www.youtube.com/watch?v=o6QXcP3Xvus>

Properties of Consumer-Produced Videos of Multimedia Commons

- Visuals: No constraints in angle, number of cameras, cutting, editing
- Audio: 70% heavy noise, 50% speech, any language, 40% dubbed, 3% professional content
- Metadata: geotags correlated with technology adaptation, tags in high part of Zipf distribution

Analyzing the Audio Track



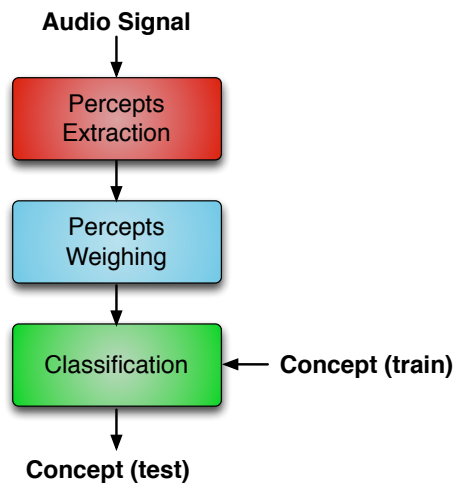
21

Approach

- Extract “audible units” aka percepts.
- Determine which percepts are common across a set of videos we are looking for but uncommon to others.
- Similar to text document search.

22

Conceptual System Overview



23

Percepts Extraction

- High number of initial segments
- Features: MFCC19+D+DD+MSG
- Minimum segment length: 30ms
- Train Model(A,B) from Segments A,B belonging to Model(A) and Model(B) and compare using BIC:

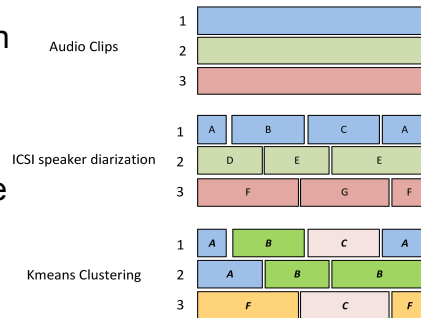
$$\log p(X|\Theta) - \frac{1}{2}\lambda K \log N$$

- Derived from Speaker Diarization

24

Percepts Dictionary

- Percepts extraction works on a per-video basis
- Use k-means to unify percepts across videos in the same set and build „prototype percepts“
- Represent video sets by supervectors of prototype percepts = “words”



25

Questions...

- How many unique “words“ define a particular concept?
- What’s the occurrence frequency of the „words“ per set of video?
- What’s the cross-class ambiguity of the „words“?
- How indicative are the highest frequent „words“ of a set of videos?

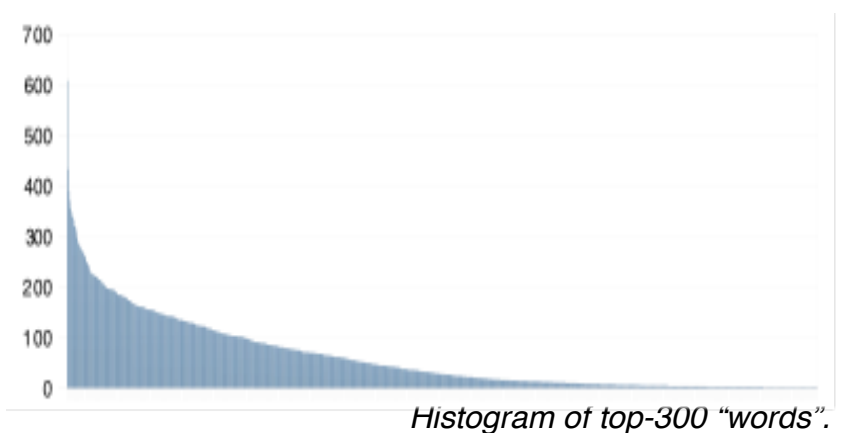
26

Properties of “Words”

- Sometimes same “word” describes more percepts (homonym)
 - Sometimes same percepts are described by the different “words” (synonym)
 - Sometimes multiply “words” needed to describe one percepts
- => Problem?

27

Distribution of “Words”



Long-Tailed Distribution (~ Zipf)

28

Zipf?



29

TF/IDF

$$TF(c_i, D_k) = \frac{\sum_j n_j P(c_i = c_j | c_j \in D_k)}{\sum_j} \quad IDF(c_i) = \log \frac{|D|}{\sum_k P(c_i \in D_k)}$$

- $TF(c_i, D_k)$ is the frequency of “word” c_i in concept D_k .
- $P(c_i = c_j | c_j \in D_k)$ is the probability that “word” c_i equals c_j in concept D_k
- $|D|$ is the total number of concepts
- $P(c_i \in D_k)$ is the probability of “word” c_i in concept D_k

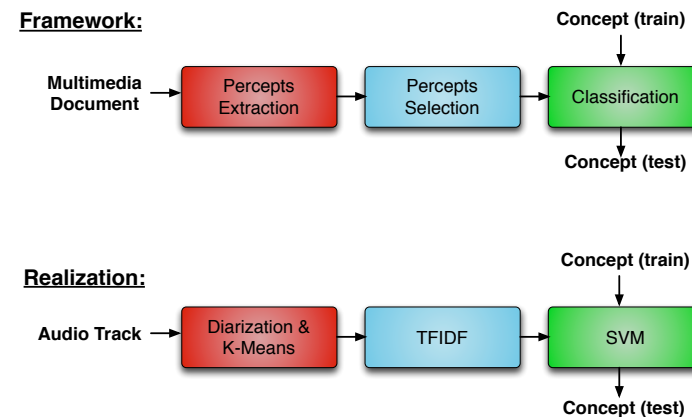
30

Classify the Words

- Have: New input video and set of representative videos
- Need: Does this belong to the same set
- Classifier takes 300 tuples of (“words“, TF-IDF values) as input
- Use SVM with Intersection Kernel (IKSVM) / Deep Learning

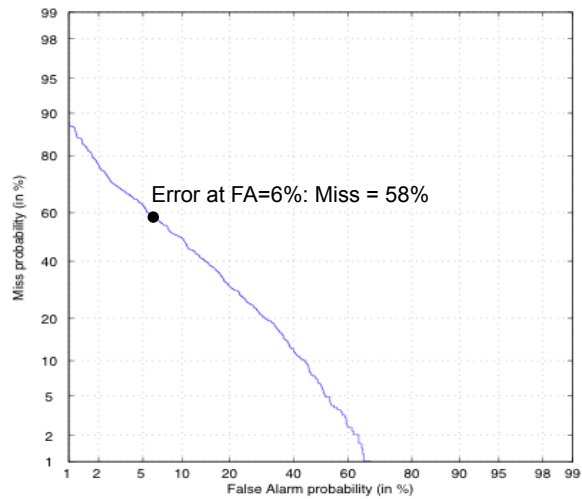
31

System Overview



32

Audio-Only Detection in TRECVID MED 2011



33

Visualization of Zipfian Percepts

- Top-1 percepts very representative of concept.



34

Thank You! Questions?

35