

A Sketch-based Approach for Multimedia Retrieval

MS (by research) Thesis,

August, 2012 - November, 2016

Koustav Ghosal
Supervisor : Dr. Anoop Namboodiri

Centre for Visual Information Technology
IIIT Hyderabad

November 29, 2016

Outline

1 Motivation

2 Video Retrieval

3 Image Retrieval

4 Zero-Shot Learning

Outline

1 Motivation

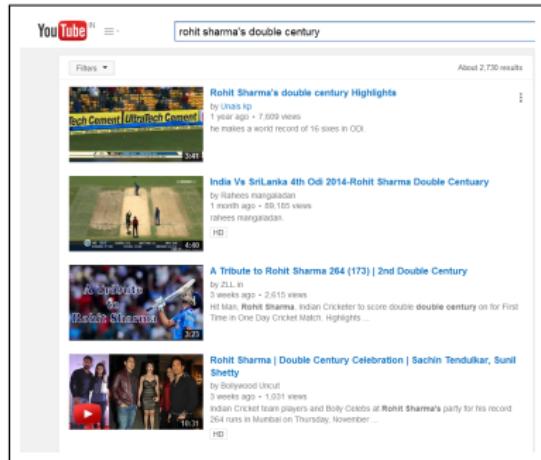
2 Video Retrieval

3 Image Retrieval

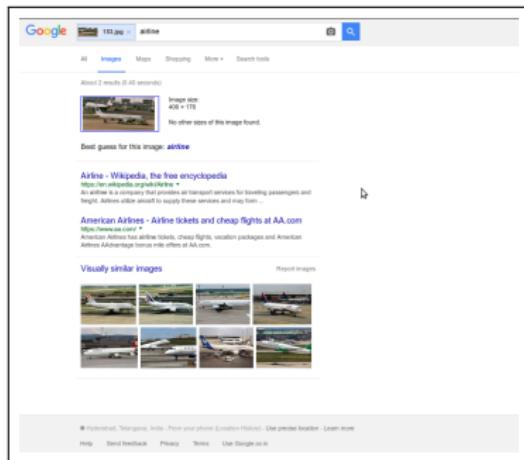
4 Zero-Shot Learning

Content Based Multimedia Retrieval

Textual and Example Based Queries



(a)



(b)

Figure: (a) Query by Text (b) Query by Example

Content Based Multimedia Retrieval

Limitations of Text based Query



Figure: Image Search

Metadata may not represent the original content.

Content Based Multimedia Retrieval

Limitations of Text based Query

"All those red coloured vehicles which came south and turned left"



Figure: Tracking

A more complicated event means a more complicated query.

Content Based Multimedia Retrieval

Limitations of Example based Query

Examples are not always available.
In fact, their absence being the reason for the search.

Why Sketch-based queries ?

Advantages

- Efficiently encodes information like shape, pose, colour, size etc. , all at once.
- A free-hand sketch is more convenient to draw than typing lengthy queries.
- A sketch is closer to the *content* of a video as compared to meta-data (tags, comments, captions).

Challenges in Sketch-based systems

Perceptual Variability

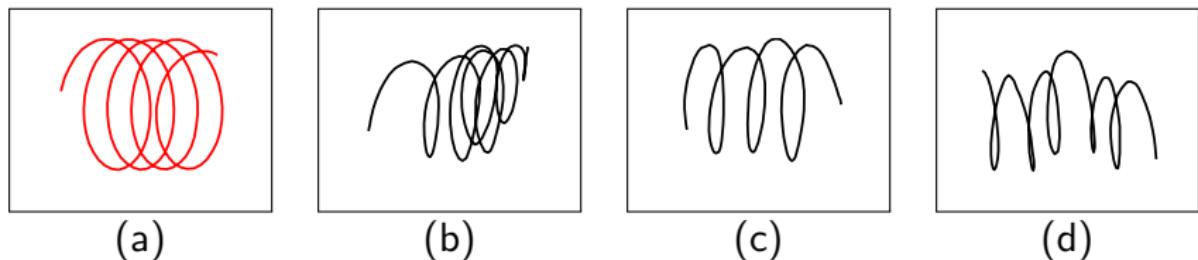


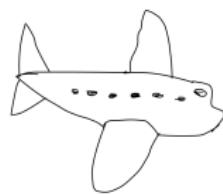
Figure: Different interpretations of the same trajectory. (a) Original (b),(c),(d)
User Inputs

Cognitive variation in motion perception in human beings .

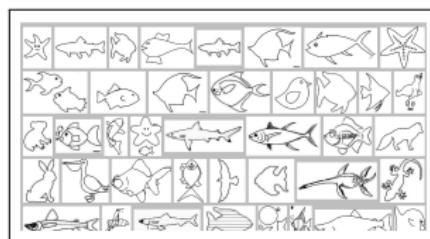
Challenges in Sketch-based systems

Multimodality

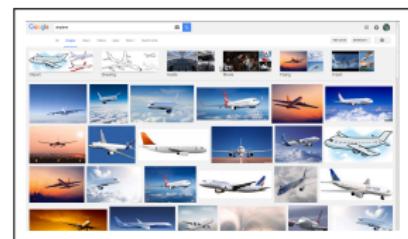
An image is a multi-channel dense representation of an object/scene.



(a)



(b)



(c)

Figure: (a) Query (b) Results (c) Desired Output

"A simple sketch is a high level sparse representation of the object/scene being searched for." [Li et al., 2015]

Challenges in Sketch-based systems

Scarcity of Data



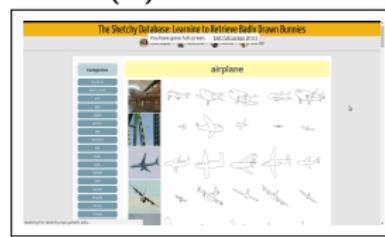
(a) ImageNet



(b) Caltech



(c) TU-Berlin



(d) Sketchy

Figure: Number of samples (a) 14,197,122 (b) 30,607 (c) 20,000 (d) 75,471

The system should be able to generalize to unknown classes.

Challenges in Sketch-based systems

Summary

- Perceptual Variability.
- Multimodality.
- Scarcity of Data.

Challenges in Sketch-based systems

Summary

- Perceptual Variability : Video Retrieval
- Multimodality : Image Retrieval
- Scarcity of Data : Zero Shot Learning

Outline

1 Motivation

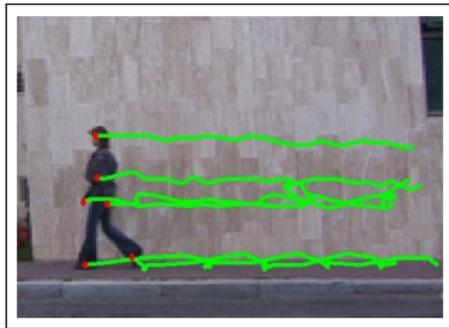
2 Video Retrieval

3 Image Retrieval

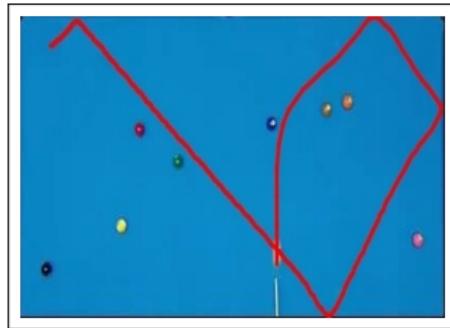
4 Zero-Shot Learning

Motion

Videos characterized by motion.



(a)



(b)

Figure: (a) Human Tracking (b) Sports Analysis

Features

Qualitative

Objective : Features minimize perceptual variability.

Qualitative Spatio Temporal Features

Features based on motion properties which tell us “how” rather than “how much”.

- Shape
- Direction
- Scale

Combines 3 different aspects of motion.

Aspects of Motion

Shape

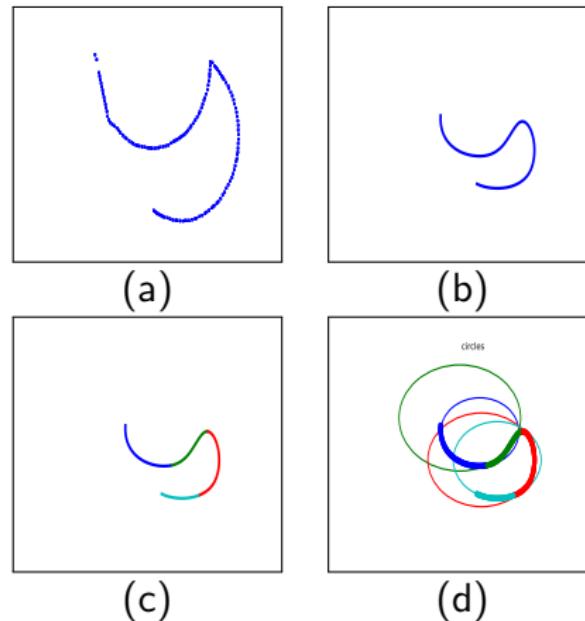
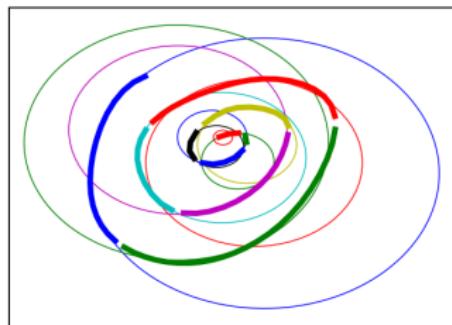


Figure: A sample motion with the corresponding m -segments: (a) Original (b) Smooth and Normalized (c) m -segments (d) Circle-Based Representation

Aspects of Motion

Shape

$$J = \min_{x_0, y_0, r} \sum_i^n x_i^2 + y_i^2 - 2x_0x_i - 2y_0y_i + x_0^2 + y_0^2 + r^2$$

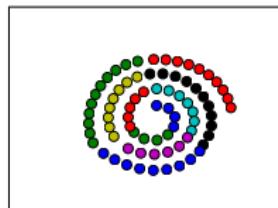


$S = (x_\mu, y_\mu, r, w, s)$
 $(x_\mu, y_\mu), r$ = center, radius of the circle.
 w = Slope of the segment.
 s = Normalized length of arc.

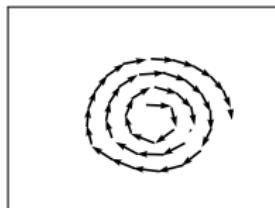
K-Means, Bag of Motion

Aspects of Motion

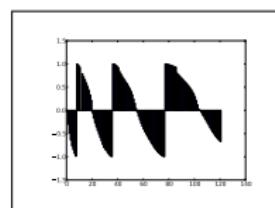
Scale, Direction



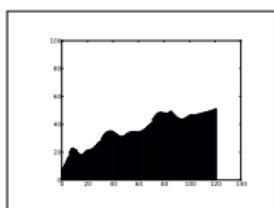
(a)



(b)



(c)



(d)

Figure: A Spiral Motion from our synthetic dataset
(a) Points sampled equidistantly in each segment
(b) Directions tracked for each equipoint segment
(c) Temporal Change of Direction
(d) Temporal Change of scale

$$\text{Trajectory Direction} = (\alpha_1, \alpha_2, \dots, \alpha_n)$$

$$\alpha_k = \sin \theta_k$$

$$\text{Trajectory Scale} = (d_1, d_2, \dots, d_n)$$

d_k = distance of the current segment from the mean.

Summary of Features

4 types

- **Bag of Motion** : Trajectory = **Histogram**.
- **Ordered Bag of Motion** : $Trajectory = (s_1, s_2, \dots, s_m)$, where
 $s_k = (x_\mu, y_\mu, r, w, s)_k$
- **Direction** : $= (\alpha_1, \alpha_2, \dots, \alpha_n)$
- **Scale** : $= (d_1, d_2, \dots, d_n)$

Pipeline

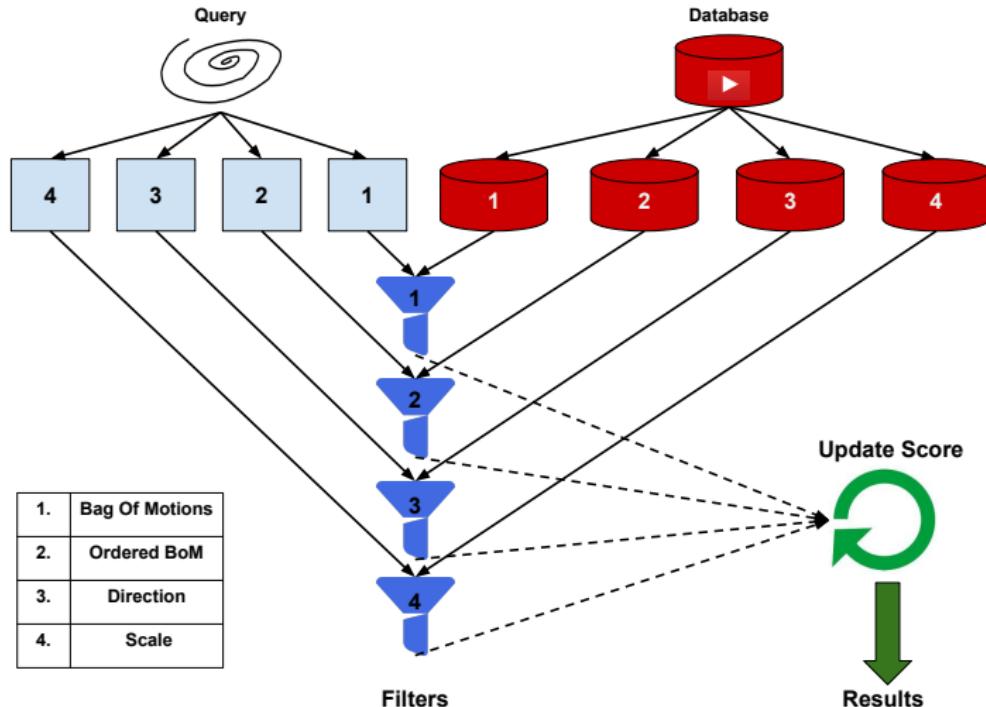
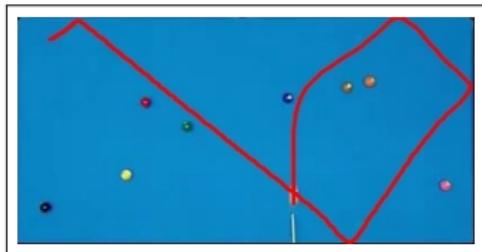
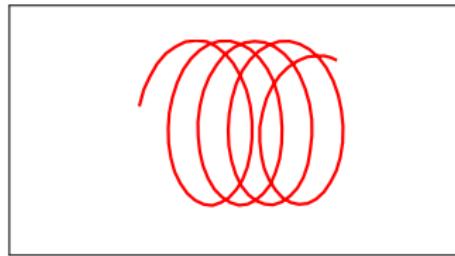


Figure: Cascaded Retrieval

Datasets



(a) Pool Dataset



(b) Synthetic Dataset.

Figure: (a) Five classes each containing 20 videos each. (b) Five classes containing 20 videos each. Thus the dataset had 200 videos

Results

A Sample Retrieval

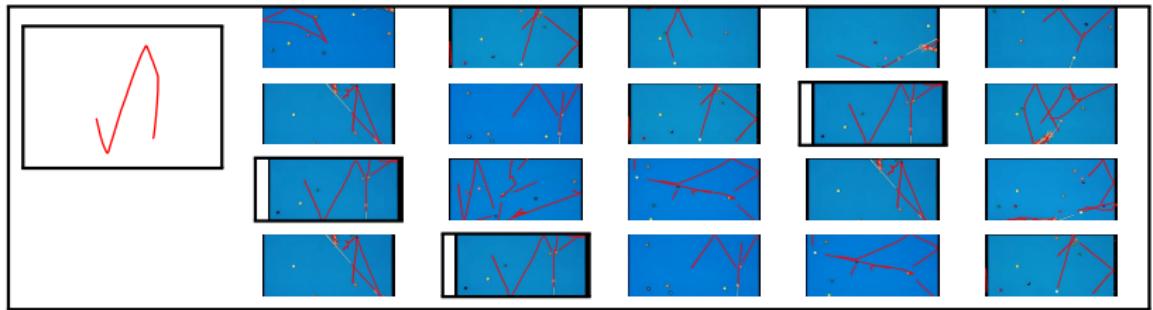
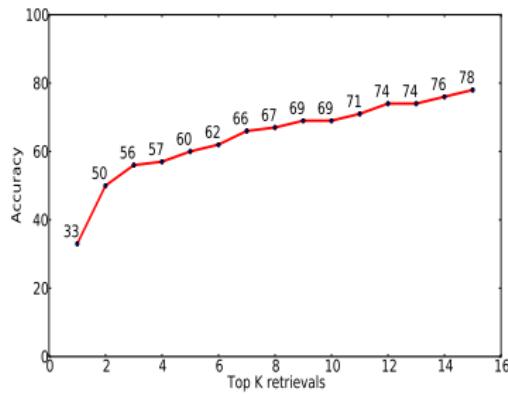


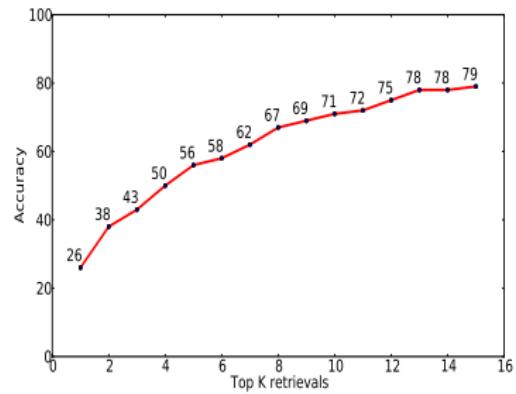
Figure: Qualitative Results

Results

Quantitative Results : Accuracy



(a) Pool Videos dataset

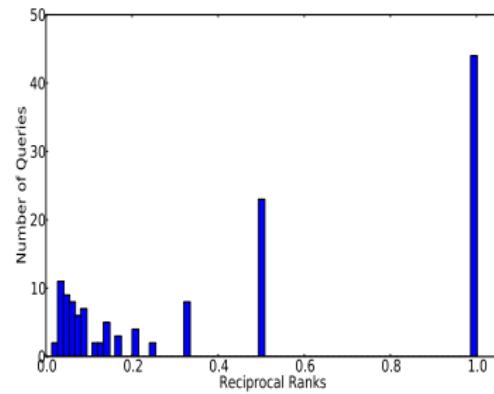


(b) Synthetic Motion dataset

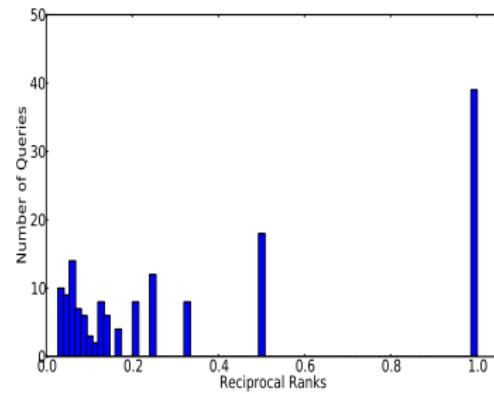
Figure: Accuracy at different top K retrievals.

Results

Quantitative Results : Mean Reciprocal Rank



(a) Pool Videos dataset



(b) Synthetic Motion dataset

Figure: Mean Reciprocal Ranks

Outline

1 Motivation

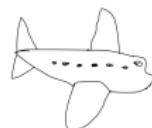
2 Video Retrieval

3 Image Retrieval

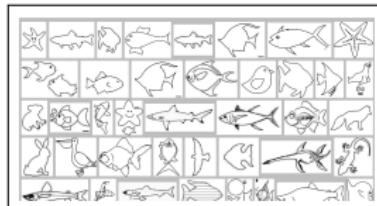
4 Zero-Shot Learning

Sparsity of Sketches

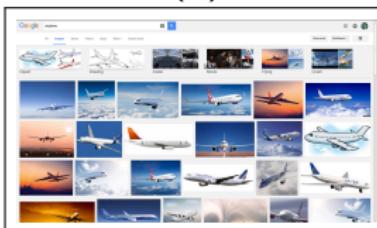
Two different modalities



(a)



(b)



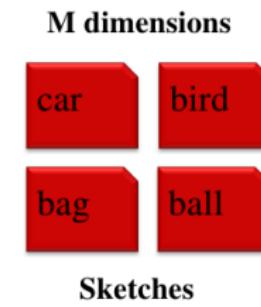
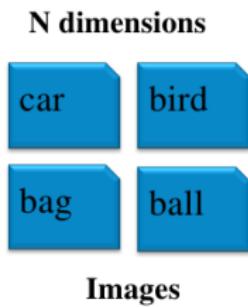
(c)

Figure: (a) Query (b) Results (c) Desired Output

Should not be compared directly.

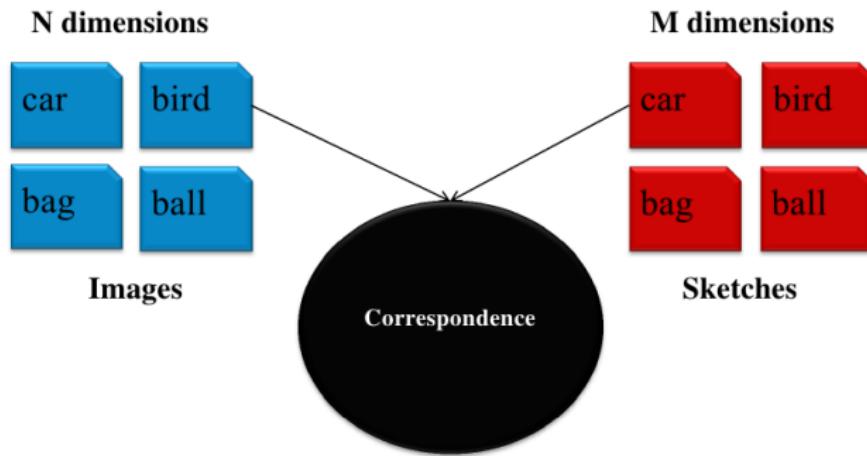
Our Model

Two Modalities

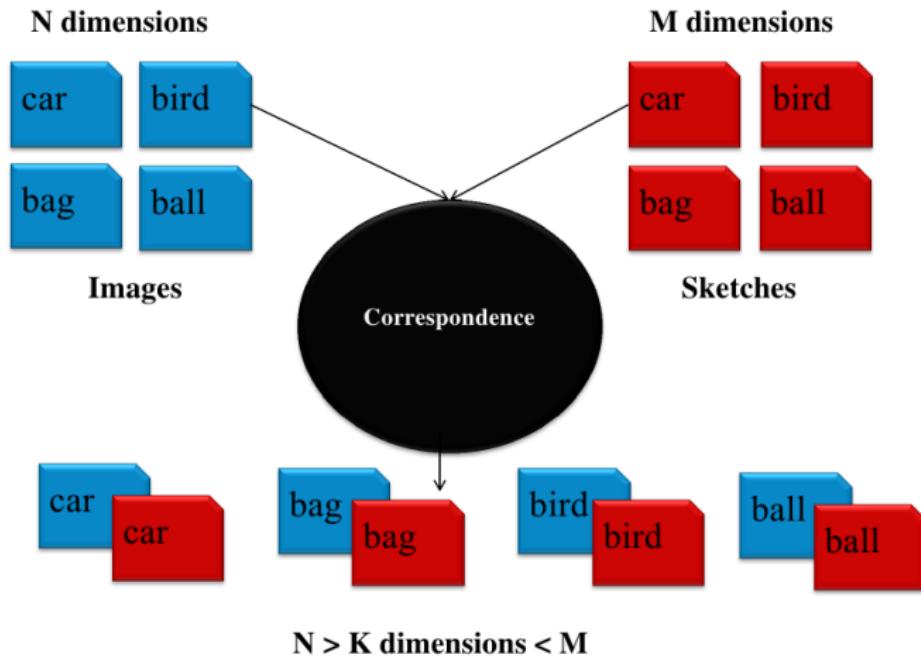


Our Model

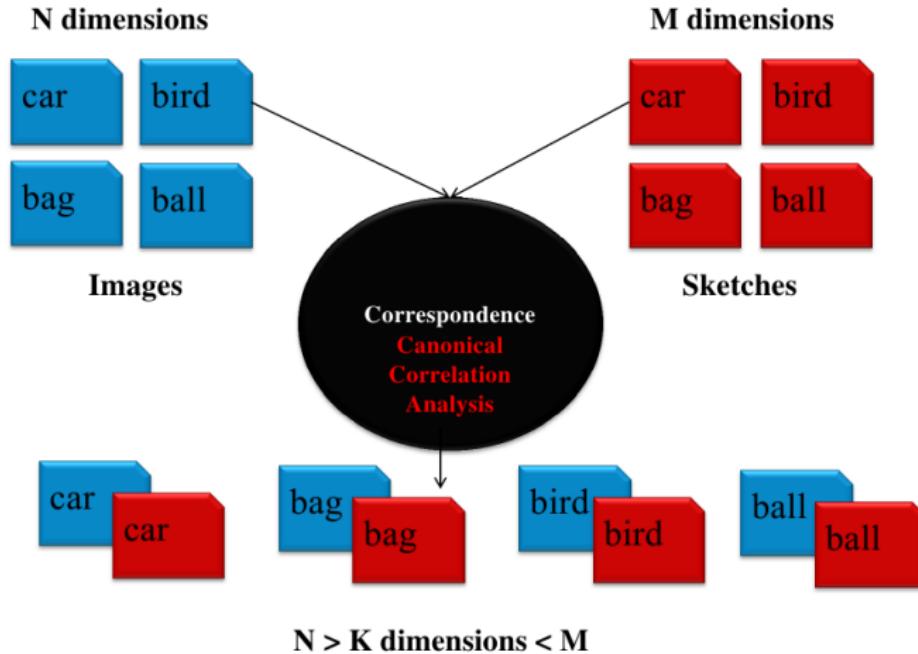
Correspondence



Our Model

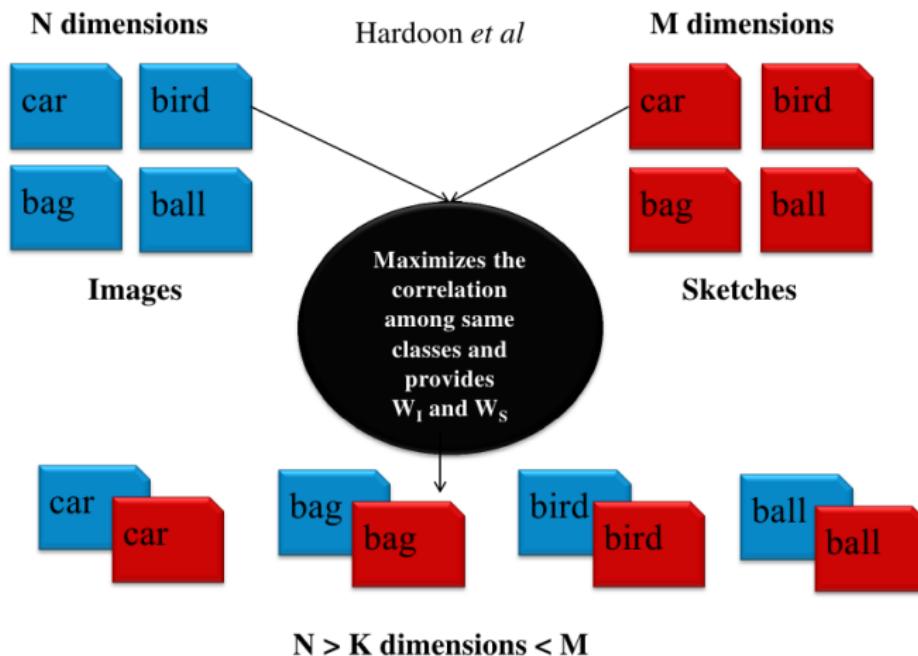


Our Model



Our Model

Two Modalities



Standard CCA

Given two sets I and S ,

$$I = (I_1, I_2, \dots, I_n)$$

$$S = (S_1, S_2, \dots, S_n)$$

we try to find two subspaces,

$$P^I = \langle W_I, I \rangle, P^S = \langle W_S, S \rangle$$

such that their correlation is maximized.

$$\begin{aligned}\rho &= \max_{W_I, W_S} \text{corr}(P^I, P^S) \\ &= \max_{W_I, W_S} \frac{\langle P^S, P^I \rangle}{\|P^S\| \|P^I\|}\end{aligned}\tag{1}$$

Standard CCA

continued...

As derived in [Hardoon et al., 2004], Equation 1 reduces to,

$$\rho = \max_{W_I, W_S} \frac{W_I' \text{Cov}_{IS} W_S}{\sqrt{W_I' \text{Cov}_{II} W_I} \sqrt{W_S' \text{Cov}_{SS} W_S}} \quad (2)$$

and the covariance matrix of (A_I, A_S) given by:

$$\text{Cov} = \mathbb{E} \left[\begin{pmatrix} A_I \\ A_S \end{pmatrix} \begin{pmatrix} A_I \\ A_S \end{pmatrix}' \right] = \begin{bmatrix} \text{Cov}_{II} & \text{Cov}_{IS} \\ \text{Cov}_{SI} & \text{Cov}_{SS} \end{bmatrix} \quad (3)$$

Equation 2 can be solved as an Eigen value problem for W_I and W_S .

Cluster CCA

As suggested in [Rasiwasia et al., 2014], we compute the covariance matrices as follows,

$$Cov_{IS} = \frac{1}{M} \sum_{c=1}^C \sum_{j=1}^{|I^c|} \sum_{k=1}^{|S^c|} I_j^c S_k^{c'} \quad (4)$$

$$Cov_{II} = \frac{1}{M} \sum_{c=1}^C \sum_{j=1}^{|I^c|} |S^c| I_j^c I_j^{c'} \quad (5)$$

$$Cov_{SS} = \frac{1}{M} \sum_{c=1}^C \sum_{k=1}^{|S^c|} |I^c| S_k^c S_k^{c'} \quad (6)$$

where $M = \sum_{c=1}^C |I^c||S^c|$, is the total number of pairwise correspondences across C classes.

Pipeline

Training and Testing

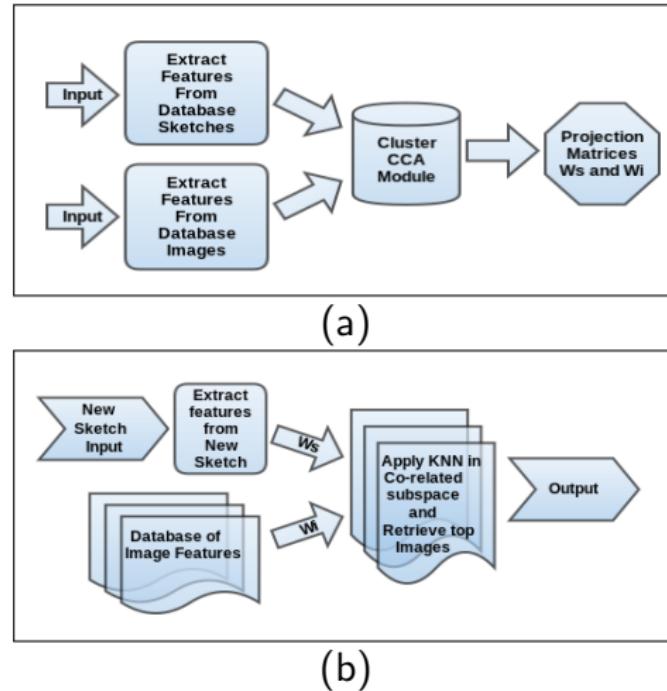
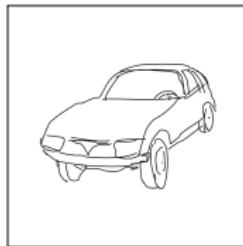
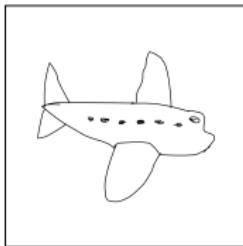


Figure: Proposed pipeline (a) Training (b) Retrieval

Datasets



TU Berlin Dataset, 250 categories, 80 sample each category



Caltech 256, Pascal VOC 2007

Features

Table 1 : Summary of Features

Feature	Dimension	Source
CALTECH - SIFT	1000	VI-Feat. [Vedaldi and Fulkerson, 2008]
CALTECH - HOG	20000	VI-Feat [Vedaldi and Fulkerson, 2008]
CALTECH - CNN	4096	Krizhevsky <i>et al.</i> [Krizhevsky et al., 2012]
PASCAL - SIFT	1000	Guillaumin <i>et al.</i> [Guillaumin et al., 2010]
PASCAL - HOG	20000	VI-Feat [Vedaldi and Fulkerson, 2008]
PASCAL - CNN	4096	Krizhevsky <i>et al.</i> [Krizhevsky et al., 2012]
TU-BERLIN - SIFT-Like	501	Eitz <i>et al.</i> [Eitz et al., 2012]
TU-BERLIN - HOG	20000	VI-Feat [Vedaldi and Fulkerson, 2008]
TU-BERLIN - Fisher	250000	Rosalia <i>et al.</i> [Schneider and Tuytelaars, 2014]
TU-BERLIN - CNN	4096	Yang <i>et al.</i> [Yang and Hospedales, 2015]

Results

Qualitative

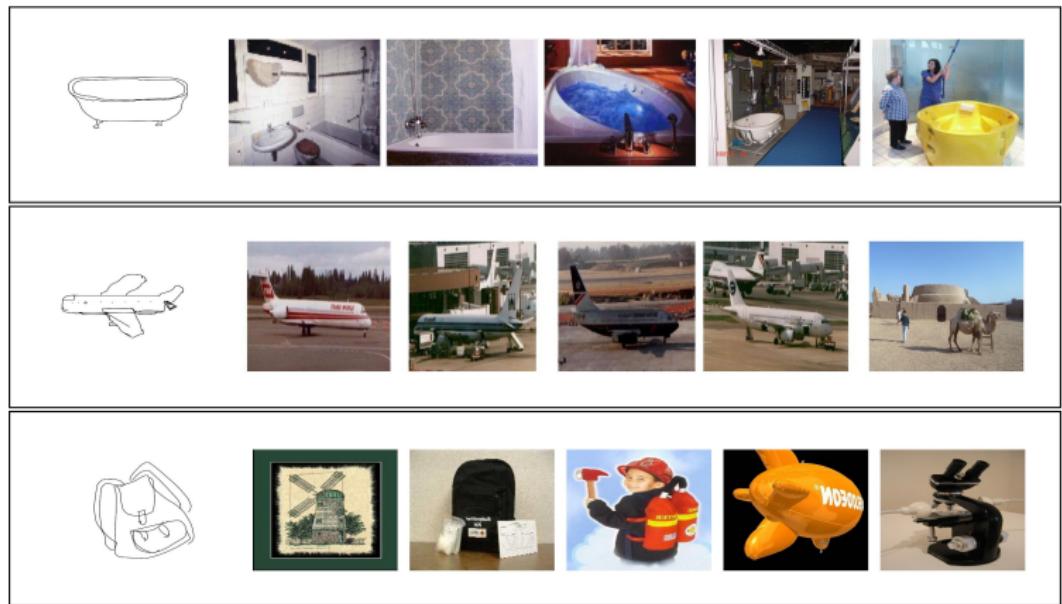


Figure: Example Retrievals From Our System

Results

Quantitative

Table : Mean Average Precision (MAP) for Image-Sketch feature combinations

Dataset	SIFT-SIFT	SIFT-HOG	SIFT-Fisher	HOG-SIFT	HOG-HOG	HOG-Fisher	CNN-CNN
Caltech	0.06	0.03	0.20	0.14	0.02	0.01	0.20
Pascal	0.13	0.12	0.05	0.18	0.09	0.06	0.06

Table : Performance improvement in mAP values

Dataset	Features	Before CCA	After CCA
Caltech	SIFT-Fisher	0.01	0.20
Caltech	CNN-CNN	0.01	0.20
Pascal	HOG-SIFT	0.01	0.18
Pascal	SIFT-SIFT	0.06	0.13

Outline

1 Motivation

2 Video Retrieval

3 Image Retrieval

4 Zero-Shot Learning

Zero-Shot Learning

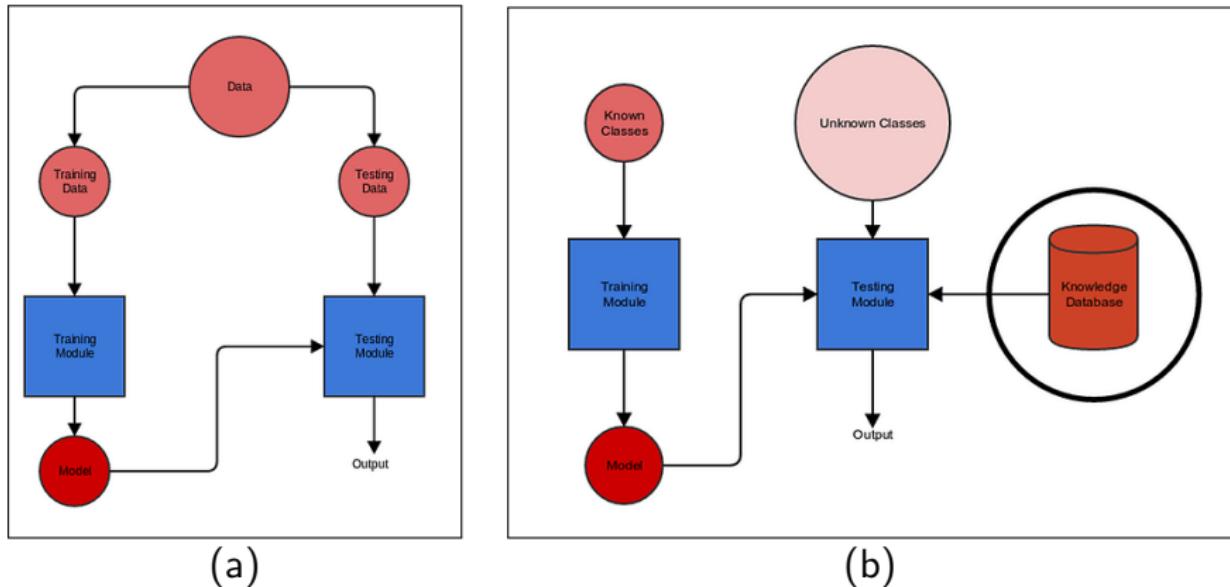


Figure: (a) Standard Classifier (b) Zero-Shot Classifier

Knowledge Database acts as an Oracle, who knows everything

Word 2 Vec

- Word 2 Vec by Mikolov *et al.* [Mikolov et al., 2013] is a vector space, where semantically similar words are mapped together.
- Apples* and *Oranges* are closer than *Apples* and *Mumbai*.
- The distance between two classes in Word 2 Vec space represents their semantic similarity.

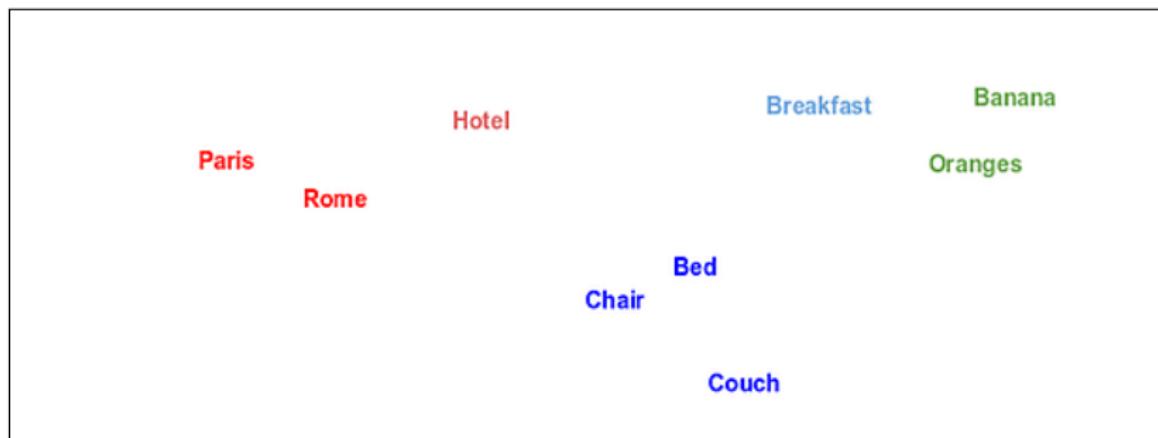


Figure: An artistic expression of Word2Vec vector space.

Training

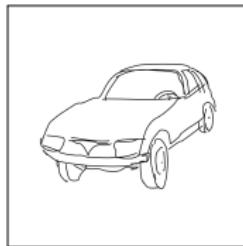
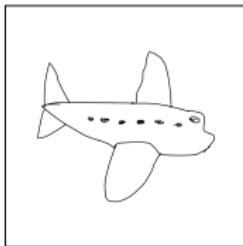
Objective Function

We train a CNN as in [Yu et al., 2015] for TU Berlin and as in [Chatfield et al., 2014] for Caltech256 dataset and replace the soft-max layer as follows.

$$J(\theta) = \sum_{y \in \mathcal{Y}} \sum_{x^i \in \mathcal{X}_y} \|w_y - g(x^i)\|^2$$

A close miss is penalized less than a distant miss.

Datasets



TU Berlin Dataset, 250 categories, 80 sample each category



Caltech 256, Pascal VOC 2007

Results

Qualitative

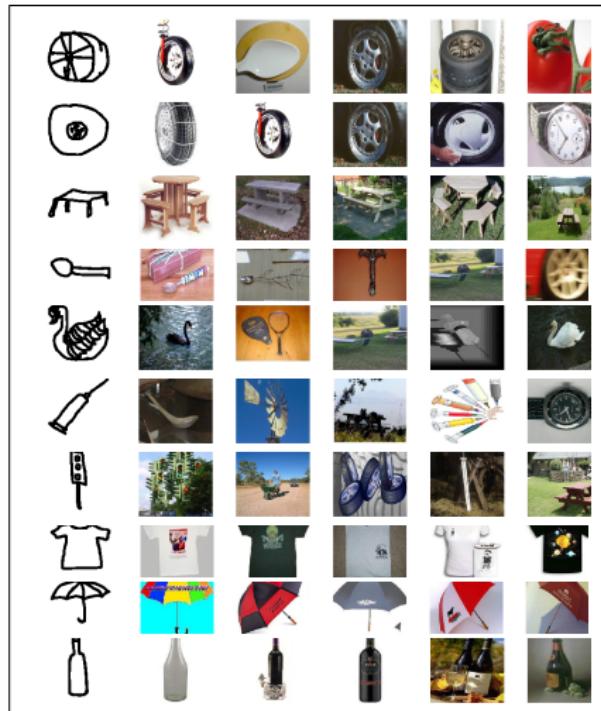


Figure: Example Queries

Experiments

- Sketch Classification.
- Uni-Modal Sketch Retrieval.
- Cross-Modal Retrieval.
- Zero Shot Retrieval.

Results

Classification

Algorithm	Features	Classifier	Accuracy
TU-Berlin Dataset			
Yang <i>et al.</i> [Yu et al., 2015]	CNN-Ensemble	SOFT-MAX	74.9%
Yang <i>et al.</i> [Yu et al., 2015]	CNN-Single	SOFT-MAX	72.6%
Schneider <i>et al.</i> [Schneider and Tuytelaars, 2014]	FISHER	SVM	63.1%
Eitz <i>et al.</i> [Eitz et al., 2012]	BOW	SVM	56%
Proposed	DFSR	RANDOM FOREST	70.22%

Table: Classification Results show that our features perform reasonably well almost at par with the state of the art methods.

Results

Uni-Modal Retrieval

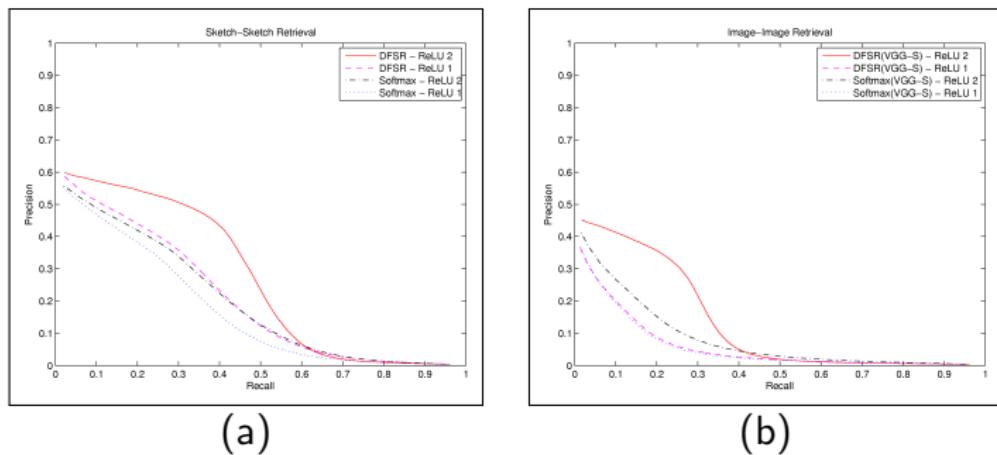


Figure: Uni-modal retrieval : (a) Sketch Modality (b) Image Modality

Results

Cross-Modal Retrieval

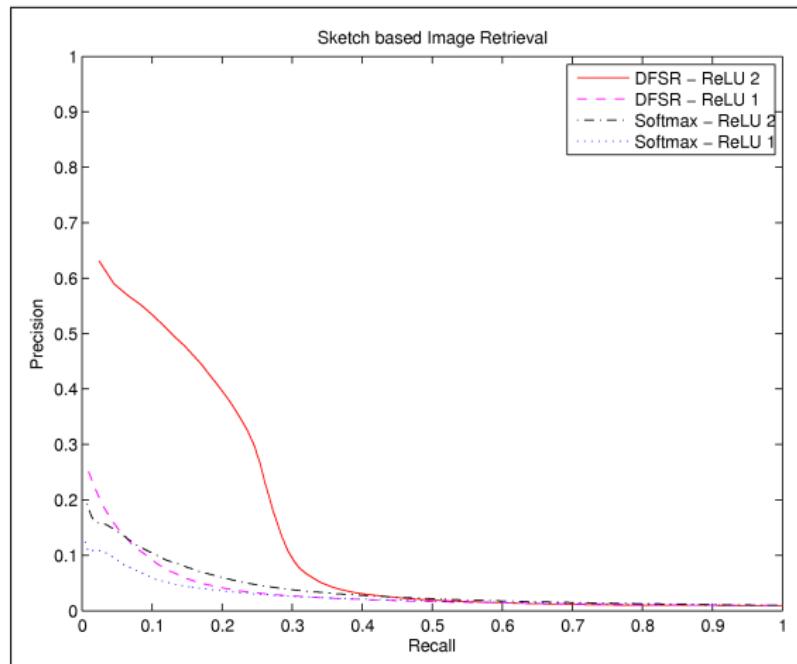


Figure: Cross Modal Retrieval

Results

Zero Shot Retrieval

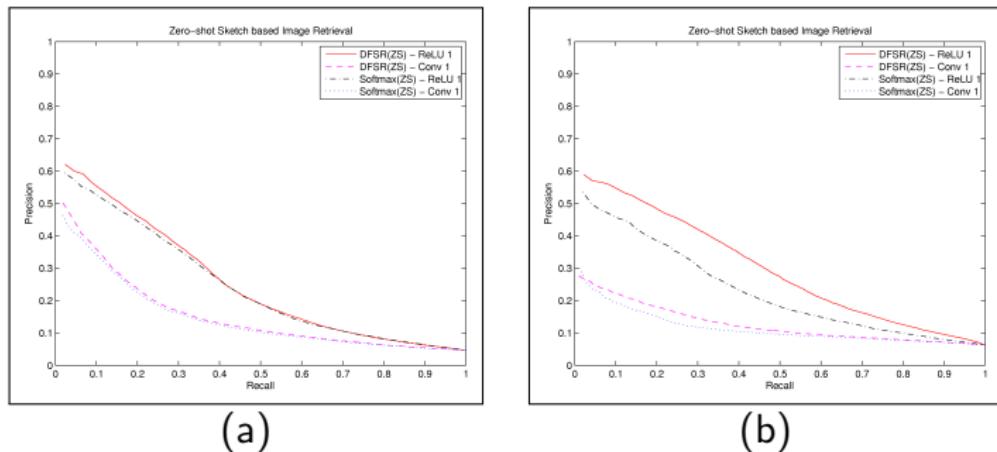


Figure: Zero-Shot Cross-Modal Retrieval : (a) PR curve for the worst performing partition. It can be observed that DFSR outperform the features proposed by Yang *et al.* [Yang and Hospedales, 2015]. (b) PR curve for the best performing partition.

Summary



Figure: Related Domains

- Qualitative features for video retrieval.
- Cluster CCA for Multi Modal Image Retrieval.
- Deep Features for Semantic Retrieval.

References I

- Chatfield, K., Simonyan, K., Vedaldi, A., and Zisserman, A. (2014). Return of the devil in the details: Delving deep into convolutional nets. In *British Machine Vision Conference*.
- Eitz, M., Hays, J., and Alexa, M. (2012). How do humans sketch objects? *ACM Trans. Graph.*
- Guillaumin, M., Verbeek, J., and Schmid, C. (2010). Multimodal semi-supervised learning for image classification. In *CVPR*.
- Hardoon, D., Szedmak, S., and Shawe-Taylor, J. (2004). Canonical correlation analysis: An overview with application to learning methods. *Neural computation*.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *NIPS*.
- Li, Y., Hospedales, T. M., Song, Y.-Z., and Gong, S. (2015). Free-hand sketch recognition by multi-kernel feature learning. *CVIU*.

References II

- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Rasiwasia, N., Mahajan, D., Mahadevan, V., and Aggarwal, G. (2014). Cluster canonical correlation analysis. In *AI Statistics*.
- Schneider, R. G. and Tuytelaars, T. (2014). Sketch classification and classification-driven analysis using fisher vectors. *TOG*.
- Vedaldi, A. and Fulkerson, B. (2008). VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/>.
- Yang, Y. and Hospedales, T. M. (2015). Deep neural networks for sketch recognition. *arXiv preprint arXiv:1501.07873*.
- Yu, Q., Yang, Y., Song, Y.-Z., Xiang, T., and Hospedales, T. M. (2015). Sketch-a-net that beats humans. In *BMVC*.

Thank you.
koustav.ghosal@research.iiit.ac.in