

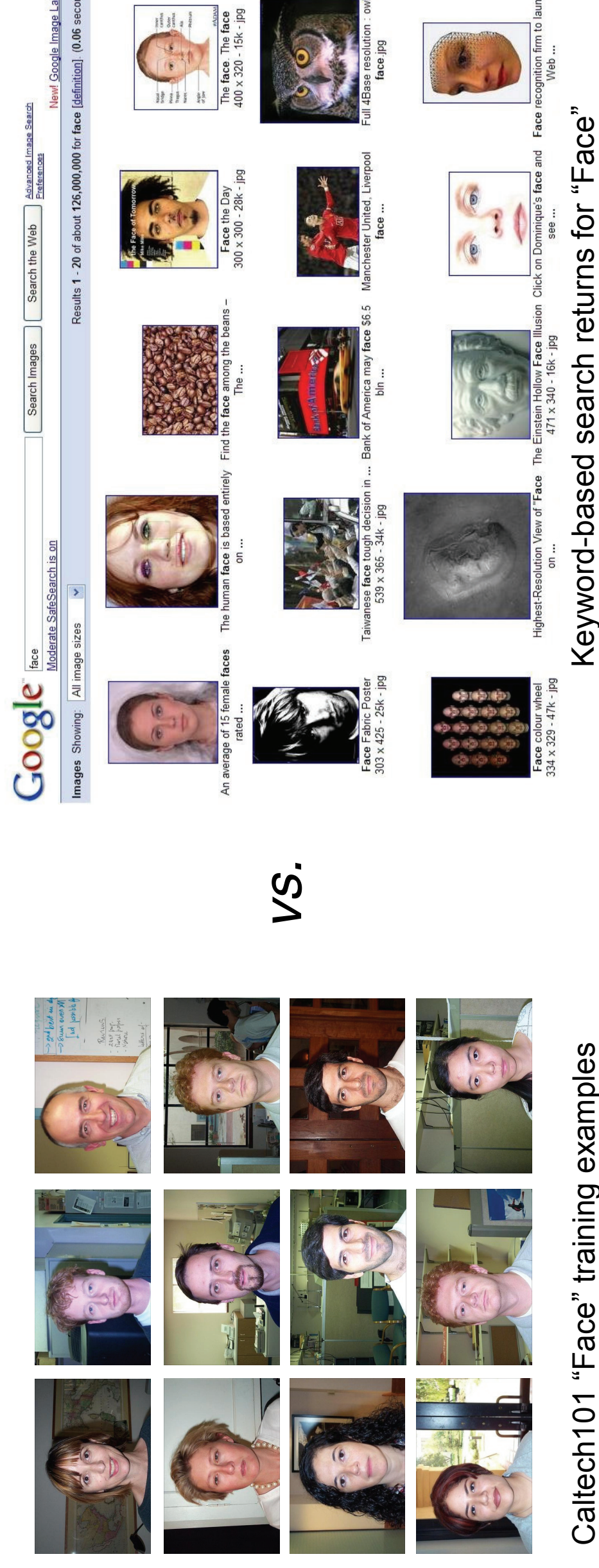
Multiple-Instance Learning for Weakly Supervised Object Categorization

Keywords to Visual Categories:
Sudheendra Vijayanarasimhan and Kristen Grauman
Department of Computer Sciences, University of Texas at Austin

Problem

Carefully prepared labeled datasets are useful for learning visual category models, but they are expensive to obtain.

Images on the Web indexed using keywords are available in large quantities, but do not necessarily portray the query term and are much more variable than typical benchmark training sets.



How can we learn reliable category models from images returned by keyword search?

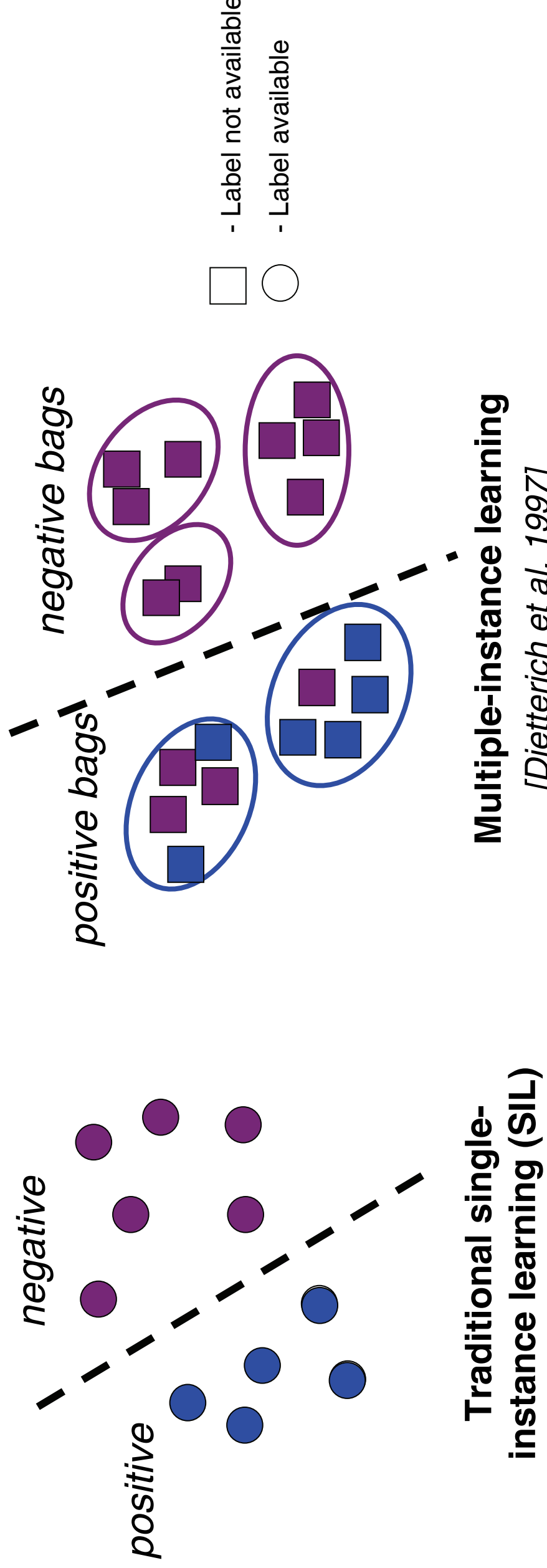
Our approach

We consider a *multiple-instance visual category learning* scenario to directly obtain discriminative models for specified categories.

- Large-margin solution with constraints to accommodate expected sparsity of good examples
- Iteratively improve multiple-instance classifier by automatically refining the representation

Previous approaches cluster to find visual themes [Sivic et al. 2005, Fergus et al. 2005, Li et al. 2007], but it is impossible to guarantee any clusters will correspond to desired topic. Others apply models known to work well with correctly labeled data [Fergus et al. 2004, Schroff et al. 2007], but extremely noisy data is problematic.

Multiple-Instance Learning (MIL)



In MIL, training examples are *bags* of instances, and we only know that each positive bag has at least one true positive instance. The classifier must predict labels for novel *instances*.

Algorithm Overview

Collect set of positive bags \mathcal{X}_p from multiple search engines using the category name translated into multiple languages.

Collect set of negative bags \mathcal{X}_n from images with known labels, or unrelated searches.

Let X denote a bag, x denote an instance, and let $\mathcal{X}_n = \{x|x \in X \in \mathcal{X}_n\}$ be all negative instances.

Sparse MIL (sMIL)

To begin, we solve a large-margin decision problem with constraints as suggested in [Bunescu & Mooney, 2007] to reflect the fact that positive bags may contain as few as one true positive:

minimize: $\frac{1}{2} \|w\|^2 + \frac{C}{|\mathcal{X}_n|} \sum_{x \in \tilde{\mathcal{X}}_n} \xi_x + \frac{C}{|\mathcal{X}_p|} \sum_{X \in \mathcal{X}_p} \xi_X$

subject to: $w \phi(x) + b \leq -1 + \xi_x, \quad \forall x \in \tilde{\mathcal{X}}_n$

$$w \frac{\phi(X)}{|X|} + b \geq \frac{2 - |X|}{|X|} - \xi_X, \quad \forall X \in \mathcal{X}_p$$
$$\xi_x \geq 0, \xi_X \geq 0.$$

where $\phi(X) = \sum_{x \in X} \phi(x)$.

The sparse MIL problem is convex, and reduces to single-instance learning when bags are of size 1.

Iterative refinement of positive bags

A limitation of the above objective is that the summed constraints mean each bag is mapped to the mean of its components; we would like instead for a bag to be represented mainly by the true positives within it.

We propose a new iterative refinement scheme for MIL that biases the bag representation towards instances that are likely to be positive:

A positive bag is represented by a weighted sum of its instances.

Initially all instances receive equal weight.

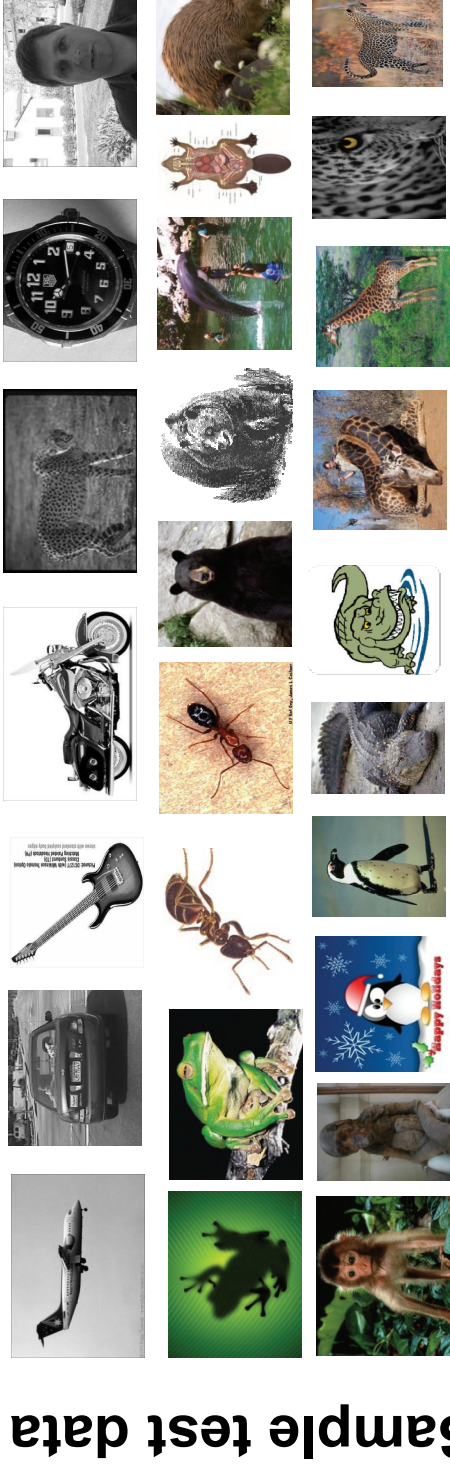
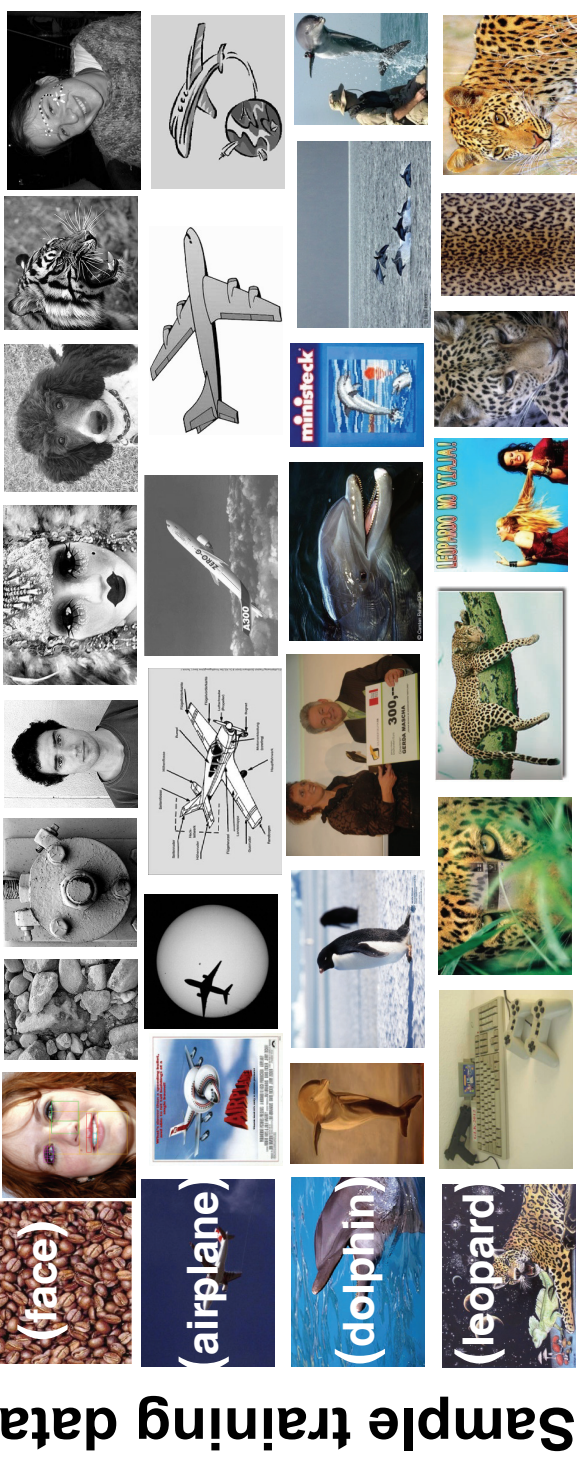
Weights are updated based on the distance of each instance from the current hyperplane.

where $y_i = w \phi(x_i) + b$

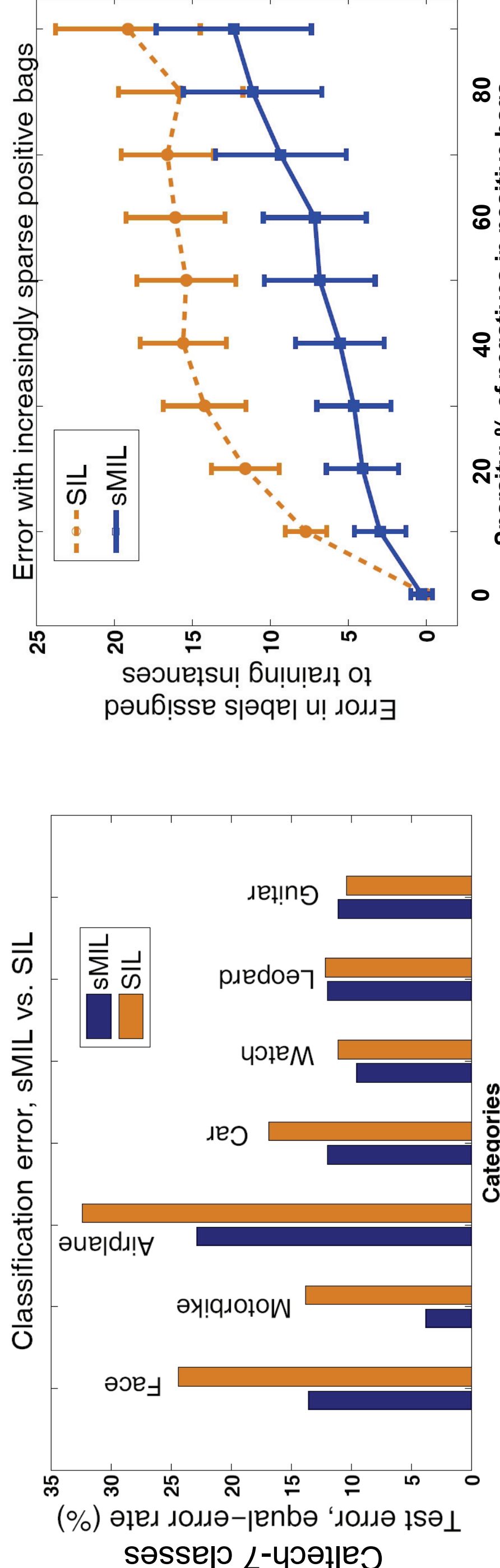
$y_m = \operatorname{argmax}_{x_i \in X} y_i$

Results

We evaluate our method for two kinds of tasks: categorizing novel examples, and re-ranking keyword search returns according to predicted relevance.



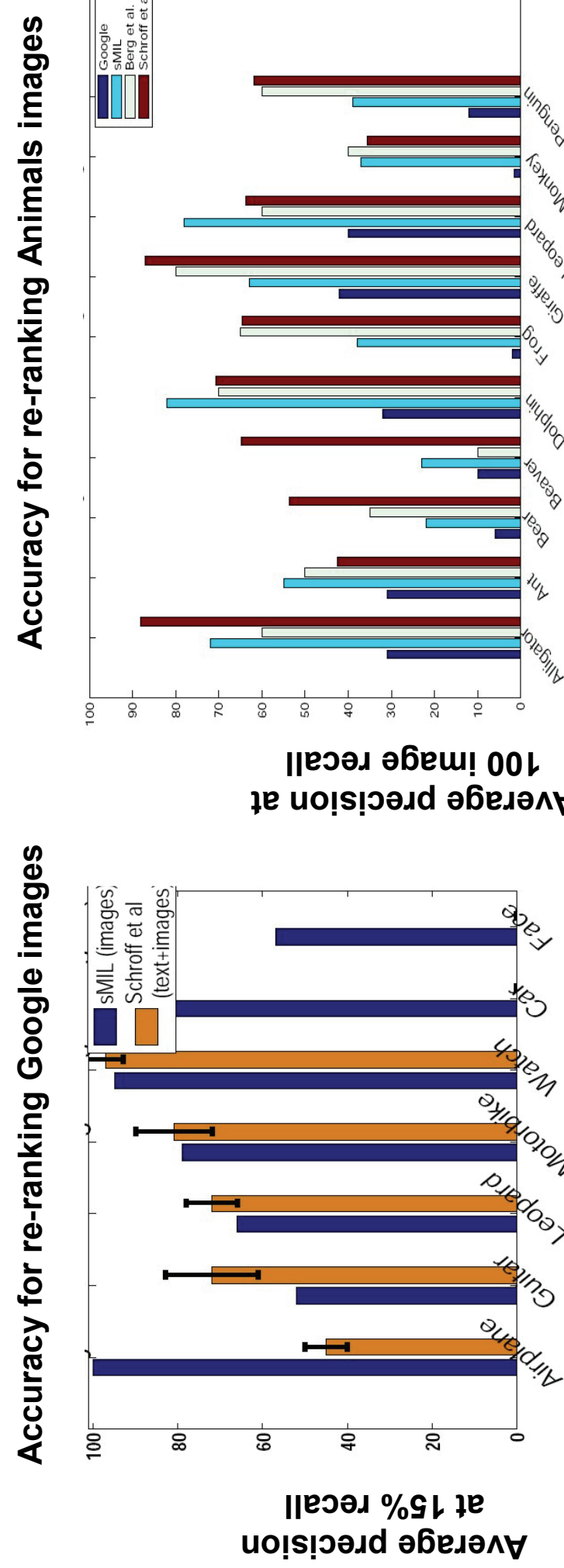
We represent images with bags of local SIFT features from four interest operators (Harris-Affine, DoG, edges, Kadir & Brady), and use an RBF kernel. Features were chosen to directly compare to previous work.



MIL is better suited to the sparse, noisy training data than the SIL baseline, and degrades much more gradually when given fewer true positives.

Amnt. of manual supervision:	ing labels	ing labels +segment.	true ing labels	none	none	none	none
Source of training data	Caltech	Caltech	Google	Caltech	Google	Google	Caltech
Method	Fergus et al. 2005	Ogilby et al. 2004	SIL-SVM	SVM'	sMIL	sMIL	sMIL
Category	Airplane	7.0	11.1	-	4.9	5.0	15.5
	Car (rear)	9.7	8.9	6.1	10.7	10.7	22.9
	Face	3.6	6.5	-	21.8	5.5	16.0
	Leopard	10.0	-	11.1	-	20.7	23.1
	Motorbike	6.7	7.8	6.0	15.4	3.8	13.0
	Guitar	-	-	-	-	6.2	12.4
	Wrist watch	-	-	-	-	31.8	3.8
Average error	-	-	-	-	7.3	19.9	8.2
	-	-	-	-	9.5	17.59	11.27
	-	-	-	-	-	-	12.14

Equal error rates on Caltech test data. Our results improve the state-of-the-art in unsupervised recognition, and are even competitive with fully supervised techniques.



Using image content alone, our approach provides accurate re-ranking results, and for some classes improves precision more than methods employing both text and image features.