

# Discriminative Clustering for Image Co-segmentation

Armand Joulin

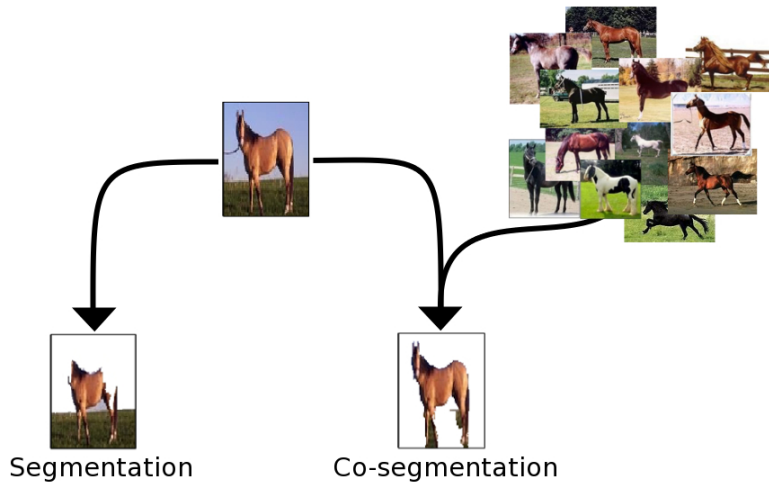
Francis Bach

Jean Ponce

INRIA Ecole Normale Supérieure, Paris

January 2010

# Introduction

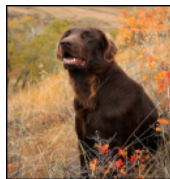
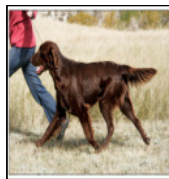
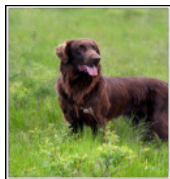


# Introduction

- ▶ **Task:** dividing simultaneously  $q$  images in  $k$  different segments
  - ▶ When  $k = 2$ , this reduces to dividing images into **foreground** and **background** regions.
- ▶ Our approach considers simultaneously the object recognition and the segmentation problems
  - ▶ **Semi-supervised discriminative clustering**
- ▶ Well-adapted to segmentation problems for 2 reasons :
  - ▶ Re-use existing features for supervised classification
  - ▶ Introduce spatial and local color-consistency constraints.

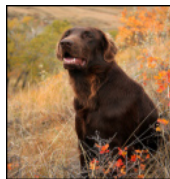
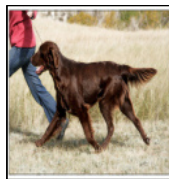
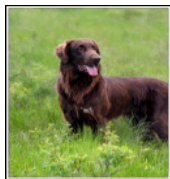


# Problem Notations



- ▶ Input:  $q$  images.
  - ▶ Each image  $i$  is reduced to a subsampled grid of  $n_i$  pixels
- ▶ For the  $j$ -th pixel (among the  $\sum_{i=1}^q n_i$  pixels), we denote by :
  - ▶  $c^j \in \mathbb{R}^3$  its color,
  - ▶  $p^j \in \mathbb{R}^2$  its position within the corresponding image,
  - ▶  $x^j$  an additional  $k$ -dimensional feature vector.

# Problem Notations

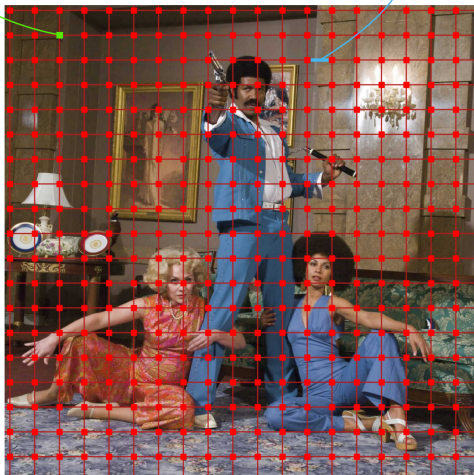


- ▶ Input:  $q$  images.
  - ▶ Each image  $i$  is reduced to a subsampled grid of  $n_i$  pixels
- ▶ For the  $j$ -th pixel (among the  $\sum_{i=1}^q n_i$  pixels), we denote by :
  - ▶  $c^j \in \mathbb{R}^3$  its color,
  - ▶  $p^j \in \mathbb{R}^2$  its position within the corresponding image,
  - ▶  $x^j$  an additional  $k$ -dimensional feature vector.
- ▶ **Goal:** find  $y =$  vector of size  $\sum_{i=1}^q n_i$  such that
  - ▶  $y_j = 1$  if the  $i$ -th pixel is in the foreground
  - ▶  $-1$  otherwise.

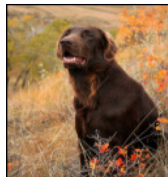
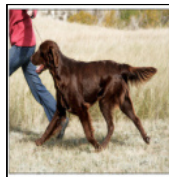
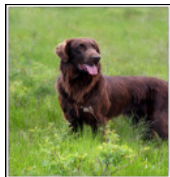
# Problem Notations

feature  $x_i$   
color  $c_i$   
position  $p_i$

spatial consistency  
based on  $\Delta c$  and  $\Delta p$



# Local consistency and discriminative clustering



- ▶ Co-segmenting images relies on two tasks :
  1. **Within an image**: maximize local spatial and appearance consistency (**normalized cuts**)
  2. **Over all images**: maximize the separability of two classes between different images (**semi-supervised SVMs**)



# Local consistency through Laplacian matrices

(Shi and Malik, 2000)

- ▶ Spatial consistency *within* an image  $i$  is enforced through a similarity matrix  $W^i$ 
  - ▶  $W^i$  is based on color features ( $c^j$ ) and spatial position ( $p^j$ )
  - ▶ Similarity between two pixels  $l$  and  $m$  within an image  $i$ :

$$W_{lm}^i = \exp(-\lambda_p \|p^m - p^l\|^2 - \lambda_c \|c^m - c^l\|^2), \quad (1)$$

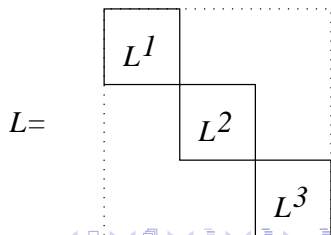
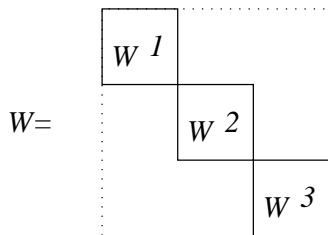
# Local consistency through Laplacian matrices

(Shi and Malik, 2000)

- ▶ Spatial consistency *within* an image  $i$  is enforced through a similarity matrix  $W^i$ 
  - ▶  $W^i$  is based on color features ( $c^j$ ) and spatial position ( $p^j$ )
  - ▶ Similarity between two pixels  $l$  and  $m$  within an image  $i$ :

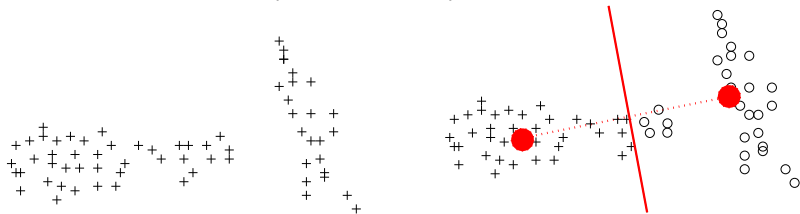
$$W_{lm}^i = \exp(-\lambda_p \|p^m - p^l\|^2 - \lambda_c \|c^m - c^l\|^2), \quad (1)$$

- ▶ Concatenate all similarity matrices into a block-diagonal matrix  $W$  (with  $W_i$  on its diagonal)
- ▶ Normalized **Laplacian** matrix  $L = I_n - D^{-1/2}WD^{-1/2}$



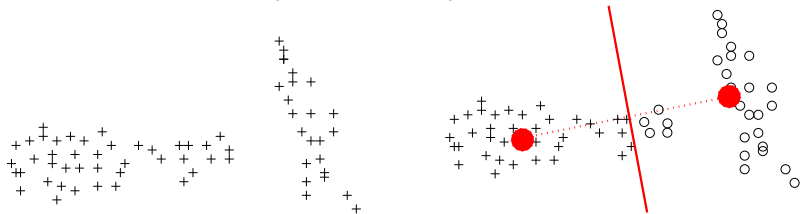
# Discriminative clustering

- ▶ **Generative** clustering (e.g., K-means)

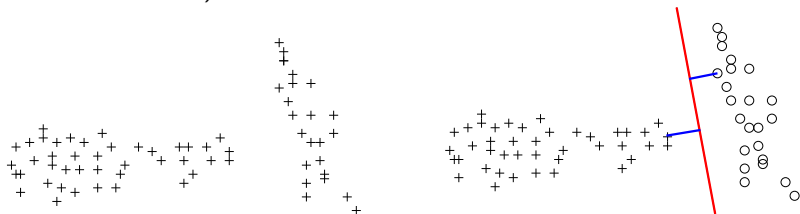


# Discriminative clustering

- ▶ **Generative** clustering (e.g., K-means)



- ▶ **Discriminative** clustering (Xu et al., 2002, Bach and Harchaoui, 2007)



# Discriminative clustering

- ▶ Discriminative clustering framework based on **positive definite kernels**
- ▶ Histograms of features  $\Rightarrow$  kernel matrix  $K$  based on the  $\chi^2$ -distance:

$$K_{lm} = \exp \left( - \lambda_h \sum_{d=1}^k \frac{(x_d^l - x_d^m)^2}{x_d^l + x_d^m} \right), \quad (2)$$

- ▶ Equivalent to mapping each of our  $n$   $k$ -dimensional vectors  $x^j$ ,  $j = 1, \dots, n$  into a high-dimensional Hilbert space  $\mathcal{F}$  through a **feature map**  $\Phi$ , so that  $K_{ml} = \Phi(x^m)^T \Phi(x^l)$

## Discriminative clustering

- ▶ Minimize with respect to both the predictor  $f$  and the labels  $y$  (Xu et al., 2002):

$$\frac{1}{n} \sum_{j=1}^n \ell(y_j, f^T \Phi(x^j)) + \lambda_k \|f\|^2, \quad (3)$$

where  $\ell$  is a loss function.

## Discriminative clustering

- ▶ Minimize with respect to both the predictor  $f$  and the labels  $y$  (Xu et al., 2002):

$$\frac{1}{n} \sum_{j=1}^n \ell(y_j, f^T \Phi(x^j)) + \lambda_k \|f\|^2, \quad (3)$$

where  $\ell$  is a loss function.

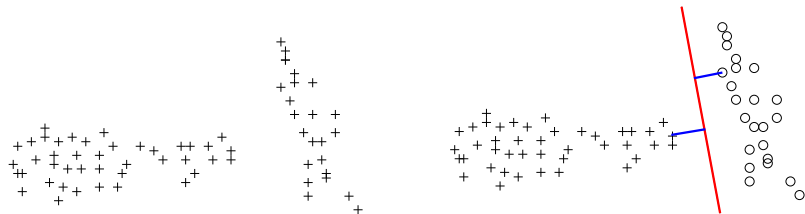
- ▶ Square loss function ( $\ell(a, b) = (a - b)^2$ ), solution in closed form (Bach and Harchaoui, 2007)

$$g(y) = \min_f \frac{1}{n} \sum_{j=1}^n \ell(y_j, f^T \Phi(x^j)) + \lambda_k \|f\|^2 = \text{tr}(Ayy^T)$$

where  $A = \lambda_k (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) (n \lambda_k I_n + K)^{-1} (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T)$ .

- ▶ **Linear in**  $Y = yy^T \in \mathbb{R}^{n \times n}$

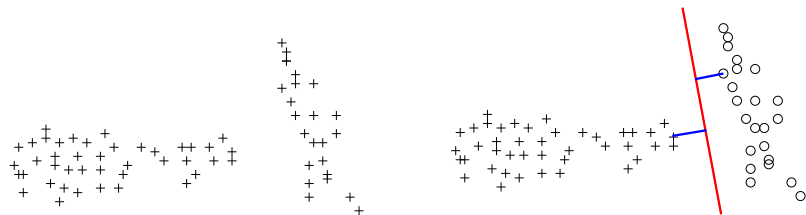
## Cluster size constraints



- ▶ Putting all pixels into a single class leads to perfect separation
  - ▶ Constrain the number of elements in each class (Xu et al., 2002)



# Cluster size constraints



- ▶ Putting all pixels into a single class leads to perfect separation
  - ▶ Constrain the number of elements in each class (Xu et al., 2002)
- ▶ Multiple images:
  - ▶ constrain the number of elements of each class in each image to be upper bounded by  $\lambda_1$  and lower bounded by  $\lambda_0$ .
  - ▶ Denote  $\delta_i \in \mathbb{R}^n$  the indicator vector of the  $i$ -th image

# Problem formulation

- ▶ Combining:
  - ▶ **spatial consistency** through Laplacian matrix  $L$
  - ▶ **discriminative cost** through matrix  $A$  and cluster size constraints

$$\min_{y \in \{-1,1\}^n} \operatorname{tr}\left(\left(A + \frac{\mu}{n}L\right)yy^T\right),$$

subject to  $\forall i, \lambda_0 \mathbf{1}_n \leq (yy^T + \mathbf{1}_n \mathbf{1}_n^T)\delta_i \leq \lambda_1 \mathbf{1}_n.$

- ▶ Combinatorial optimization problem
  - ▶ Convex relaxation with semi-definite programming (Goemans and Williamson, 1995)

## Optimization - Convex Relaxation

$$\min_{y \in \{-1,1\}^n} \text{tr}\left(\left(A + \frac{\mu}{n}L\right)yy^T\right),$$

$$\text{subject to } \forall i, \lambda_0 \mathbf{1}_n \leq (yy^T + \mathbf{1}_n \mathbf{1}_n^T) \delta_i \leq \lambda_1 \mathbf{1}_n.$$

- ▶ Reparameterize problem with  $Y = yy^T$
- ▶  $Y$  referred to as the **equivalence matrix**
  - ▶  $Y_{ij} = 1$  if points  $i$  and  $j$  belong to the same cluster
  - ▶  $Y_{ij} = -1$  if points  $i$  and  $j$  do not belong to the same cluster
- ▶  $Y$  is symmetric, positive semidefinite, with diagonal equal to one, and unit rank.

# Optimization - Convex Relaxation

- ▶ Denote by  $\mathcal{E}$  the *elliptope*, i.e., the convex set defined by:

$$\mathcal{E} = \{Y \in \mathbb{R}^{n \times n}, Y = Y^T, \text{diag}(Y) = \mathbf{1}_n, Y \succeq 0\},$$

- ▶ Reformulated optimization problem :

$$\begin{aligned} & \min_{Y \in \mathcal{E}} \text{tr}\left(Y\left(A + \frac{\mu}{n}L\right)\right), \\ & \text{subject to } \forall i, \lambda_0 \mathbf{1}_n \leq (Y + \mathbf{1}_n \mathbf{1}_n^T) \delta_i \leq \lambda_1 \mathbf{1}_n \\ & \text{rank}(Y) = 1 \end{aligned}$$

- ▶ Rank constraint is not convex
- ▶ **Convex relaxation by removing the rank constraint**

# Optimization

$$\begin{aligned} & \min_{Y \in \mathcal{E}} \operatorname{tr}(Y(A + \frac{\mu}{n}L)), \\ & \text{subject to } \forall i, \lambda_0 \mathbf{1}_n \leq (Y + \mathbf{1}_n \mathbf{1}_n^T) \delta_i \leq \lambda_1 \mathbf{1}_n \end{aligned}$$

- ▶ SDP: semidefinite program (Boyd and Vandenberghe, 2002)
- ▶ General purpose toolboxes would solve this problem in  $O(n^7)$
- ▶ Bach and Harchaoui (2007) considers a partial dualization technique that scales up to thousands of data points.
- ▶ To gain another order of magnitude: optimization through low-rank matrices (Journée et al, 2008)

## Efficient low-rank optimization (Journée et al, 2008)

- ▶ Replace constraints by penalization  $\Rightarrow$  optimization of a convex function  $f(Y)$  on the elliptope  $\mathcal{E}$ .
- ▶ Empirically: global solution has low rank  $r$
- ▶ Property: a local minimum of  $f(Y)$  over the rank constrained elliptope

$$\mathcal{E}_d = \{Y \in \mathcal{E}, \text{rank}(Y) = d\}$$

is a global minimum of  $f(Y)$  over  $\mathcal{E}$ , if  $d > r$ .

## Efficient low-rank optimization (Journée et al, 2008)

- ▶ Replace constraints by penalization  $\Rightarrow$  optimization of a convex function  $f(Y)$  on the elliptope  $\mathcal{E}$ .
- ▶ Empirically: global solution has low rank  $r$
- ▶ Property: a local minimum of  $f(Y)$  over the rank constrained elliptope

$$\mathcal{E}_d = \{Y \in \mathcal{E}, \text{rank}(Y) = d\}$$

is a global minimum of  $f(Y)$  over  $\mathcal{E}$ , if  $d > r$ .

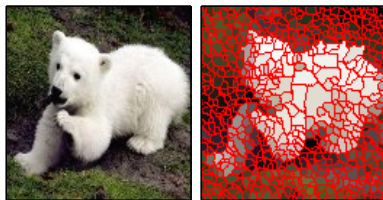
- ▶ Adaptive procedure to automatically find  $r$
- ▶ Manifold-based trust-region method for a given  $d$  (Absil et al., 2008)

## Low-rank optimization (Journée et al., 2008)

- ▶ **Final (combinatorial) goal:** minimize  $f(Y)$  over the rank-one constrained ellipsope  $\mathcal{E}_1 = \{Y \in \mathcal{E}, \text{rank}(Y) = 1\}$
- ▶ **Convex relaxation:** minimize  $f(Y)$  over the unconstrained ellipsope  $\mathcal{E}$
- ▶ **Subproblems:** minimize  $f(Y)$  over the rank- $d$  constrained ellipsope  $\mathcal{E}_d = \{Y \in \mathcal{E}, \text{rank}(Y) = d\}$  for  $d \geq 2$ 
  - ▶ It is a Riemannian manifold for  $d \geq 2$
  - ▶ If  $d$  is large enough, there is no local minima
  - ▶ Find a local minimum with trust-region method
- ▶ **Adaptive procedure:**
  - ▶ Start with  $d = 2$
  - ▶ Find local minimum over  $\mathcal{E}_d = \{Y \in \mathcal{E}, \text{rank}(Y) = d\}$
  - ▶ Check global optimality condition
  - ▶ Stop or augment  $d$

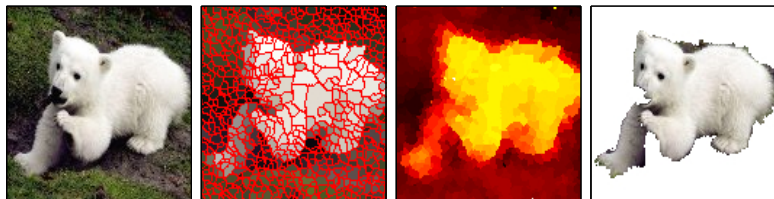


# Preclustering



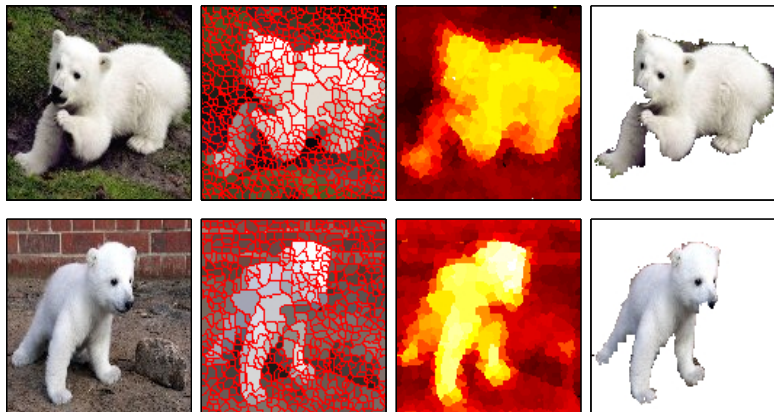
- ▶ Cost function  $f$  uses a full  $n \times n$  matrix  $A + (\mu/n)L$   
⇒ memory issues
- ▶ To reduce the total number of pixels
  - ▶ superpixels obtained from an oversegmentation of our images  
(watershed, Meyer, 2001)

# Rounding



- ▶ In order to retrieve  $y \in \{-1, 1\}$  from our relaxed solution  $Y$ , we compute the largest eigenvector  $e \in \mathbb{R}^n$  of  $Y$ .
- ▶ Final clustering is  $y = \text{sign}(e)$ .
- ▶ Other techniques could be used (e.g., rounding)
- ▶ Additional post-processing to remove some artefacts

## Method overview (co-segmentation on two bear images)



- ▶ From left to right: input images, over-segmentations, scores obtained by our algorithm and co-segmentations.

# Results

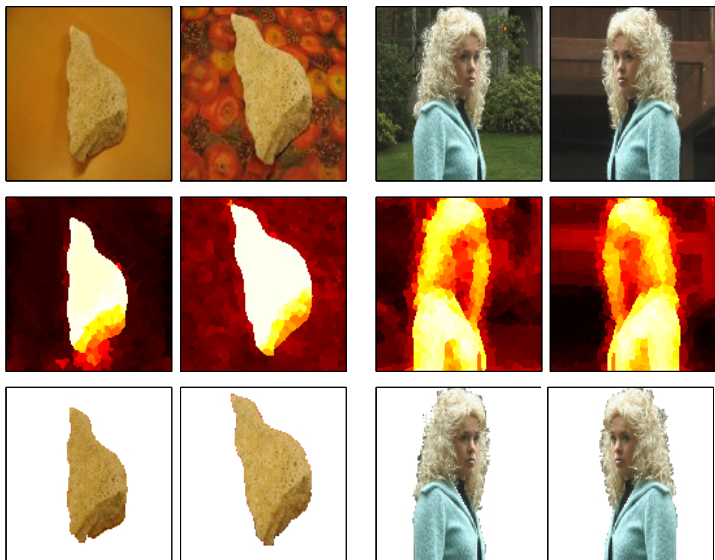
Results on two different problems :

- ▶ **Simple problems:** images with foreground objects which are identical or very similar in appearance and with few images to co-segment
- ▶ **Hard problems:** images whose foreground objects exhibit higher appearance variations and with more images to co-segment (up to 30).

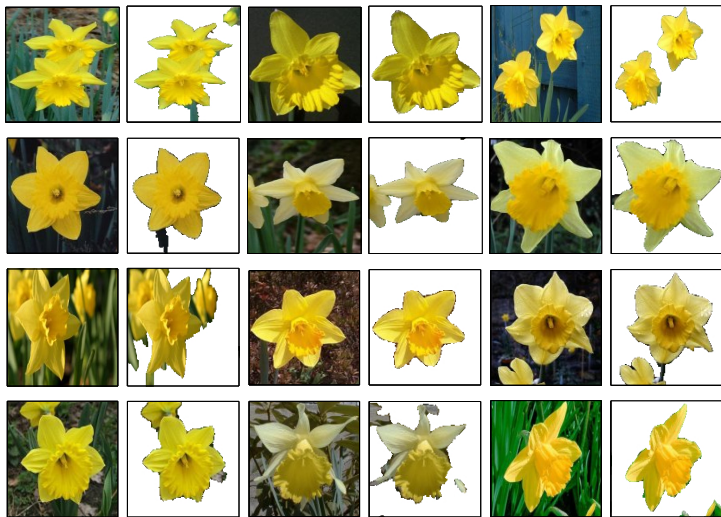
## Results - similar objects



## Results - similar objects



## Results - similar objects

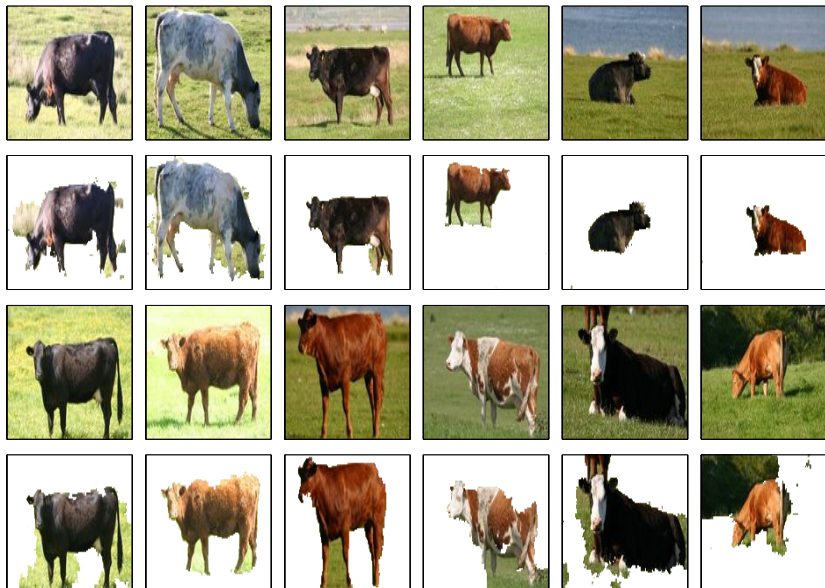


## Results - similar classes - Faces





## Results - similar classes - Cows



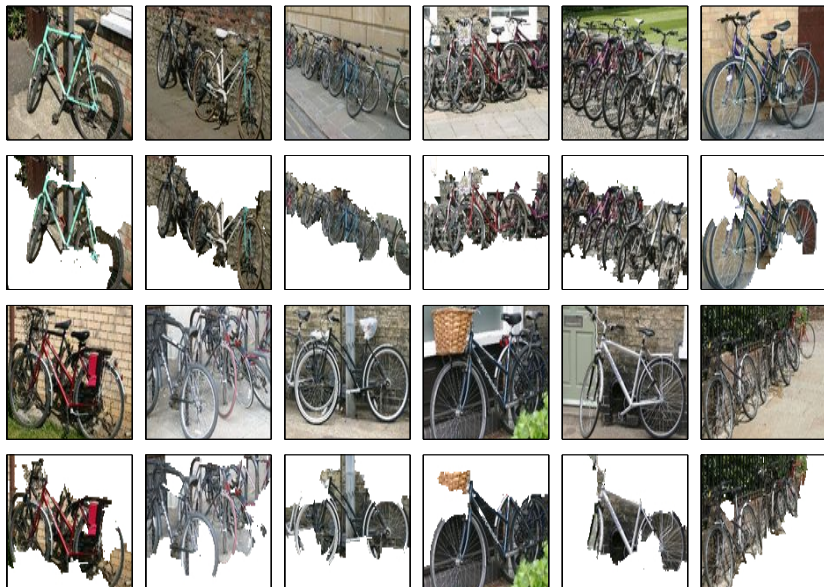
# Results - similar classes - Horses



## Results - similar classes - Cats



# Results - similar classes - Bikes



## Results - similar classes - Planes



## Comparison with MN-cut (Cour, Benezit, and Shi, 2005)

- Segmentation accuracies on the Weizman horses and MSRC databases.

class	#	cosegm.	independent	Ncut	uniform
Cars (front)	6	87.65 $\pm$ 0.1	<b>89.6 <math>\pm</math>0.1</b>	51.4 $\pm$ 1.8	64.0 $\pm$ 0.1
Cars (back)	6	<b>85.1 <math>\pm</math>0.2</b>	83.7 $\pm$ 0.5	54.1 $\pm$ 0.8	71.3 $\pm$ 0.2
Face	30	<b>84.3 <math>\pm</math>0.7</b>	72.4 $\pm$ 1.3	67.7 $\pm$ 1.2	60.4 $\pm$ 0.7
Cow	30	<b>81.6 <math>\pm</math>1.4</b>	78.5 $\pm$ 1.8	60.1 $\pm$ 2.6	66.3 $\pm$ 1.7
Horse	30	<b>80.1 <math>\pm</math>0.7</b>	77.5 $\pm$ 1.9	50.1 $\pm$ 0.9	68.6 $\pm$ 1.9
Cat	24	<b>74.4 <math>\pm</math>2.8</b>	71.3 $\pm$ 1.3	59.8 $\pm$ 2.0	59.2 $\pm$ 2.0
Plane	30	73.8 $\pm$ 0.9	62.5 $\pm$ 1.9	51.9 $\pm$ 0.5	<b>75.9 <math>\pm</math>2.0</b>
Bike	30	<b>63.3 <math>\pm</math>0.5</b>	61.1 $\pm$ 0.4	60.7 $\pm$ 2.6	59.0 $\pm$ 0.6

## Comparing co-segmentation with independent segmentations



- ▶ From left to right: original image, multiscale normalized cut, our algorithm on a single image, our algorithm on 30 images.

# Conclusion

- ▶ Co-segmentation through **semi-supervised discriminative clustering**
  1. **Within an image**: maximize local spatial and appearance consistency (**normalized cuts**)
  2. **Over all images**: maximize the separability of two classes between different images (**semi-supervised SVMs**)



# Conclusion

- ▶ Co-segmentation through **semi-supervised discriminative clustering**
  1. **Within an image**: maximize local spatial and appearance consistency (**normalized cuts**)
  2. **Over all images**: maximize the separability of two classes between different images (**semi-supervised SVMs**)
- ▶ Future work
  - ▶ Add negative images
  - ▶ More than 2 classes
  - ▶ Feature selection
  - ▶ Scale up to hundred of thousands
  - ▶ Change the loss function