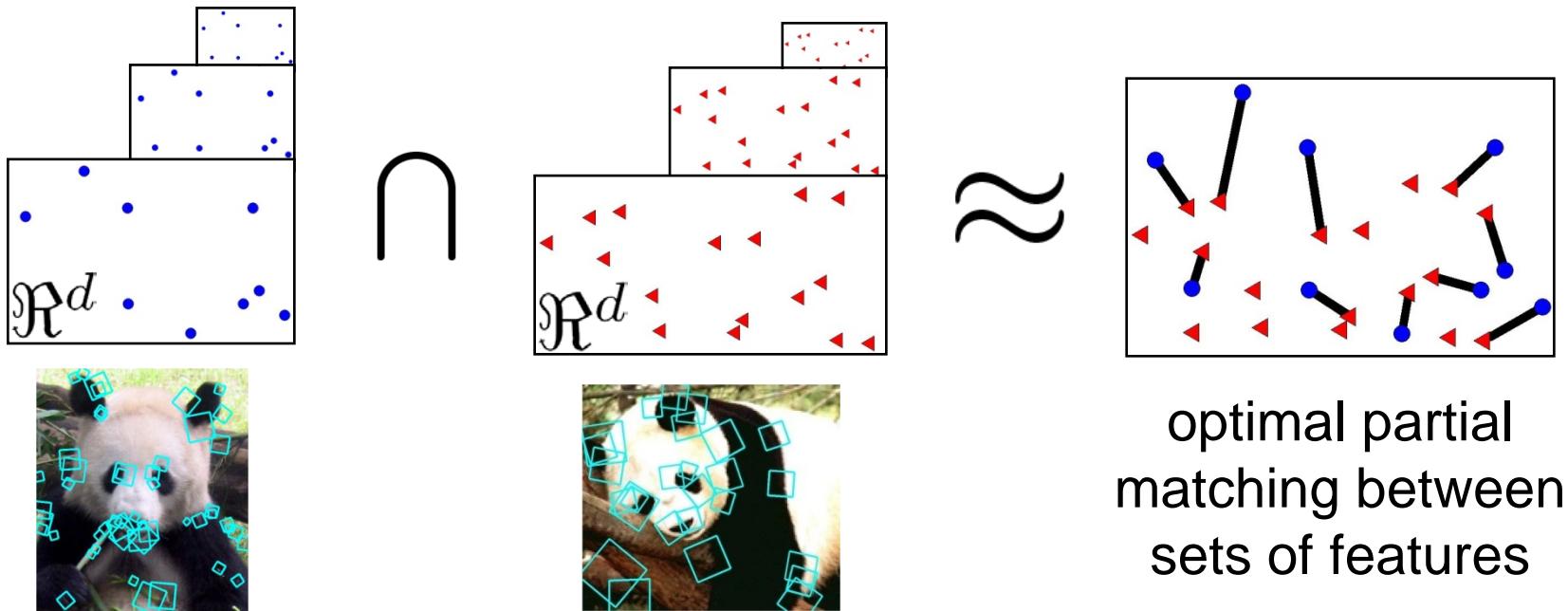


# Advanced Representation and Learning

Vinay P. Namboodiri

- Slide Credit to James Hays, Hakan Bilen, Kai Yu and Derek Hoiem

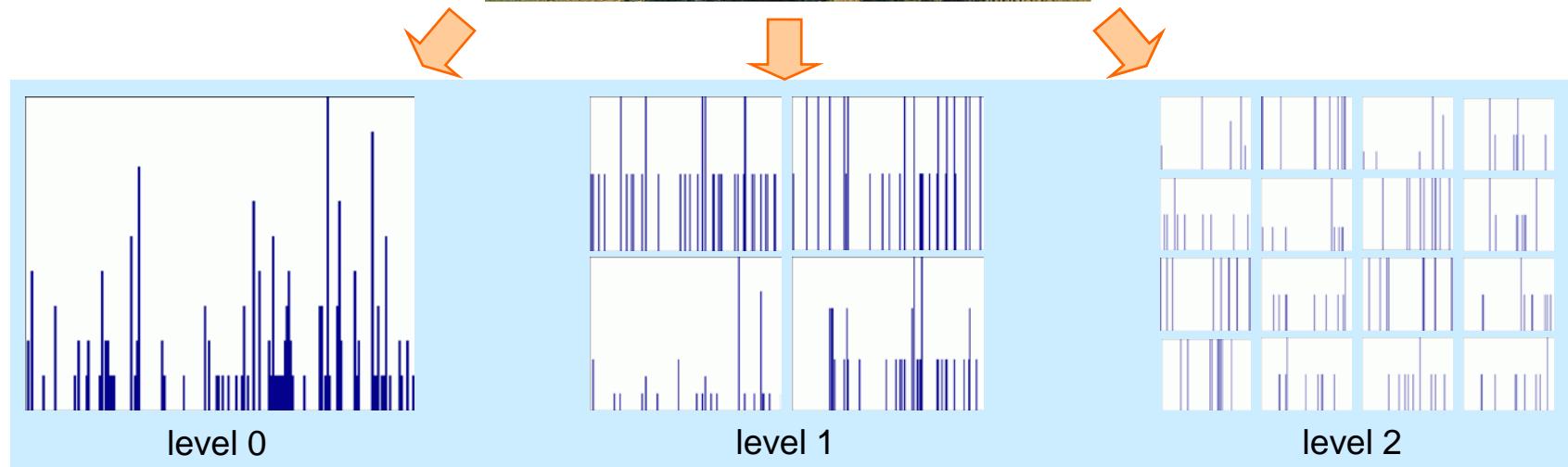
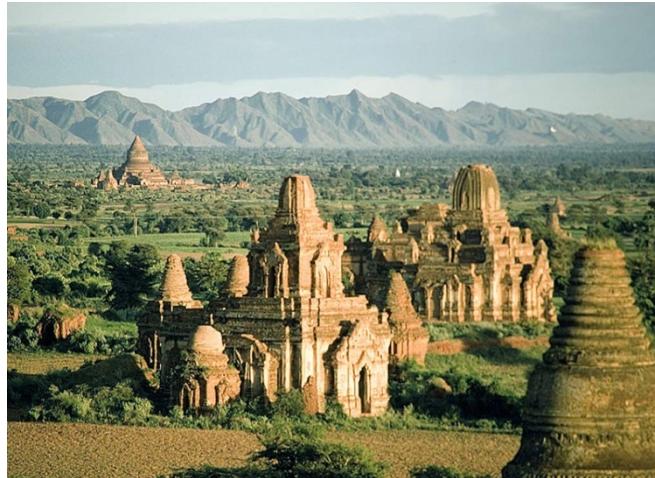
# Last class: Pyramid match kernel



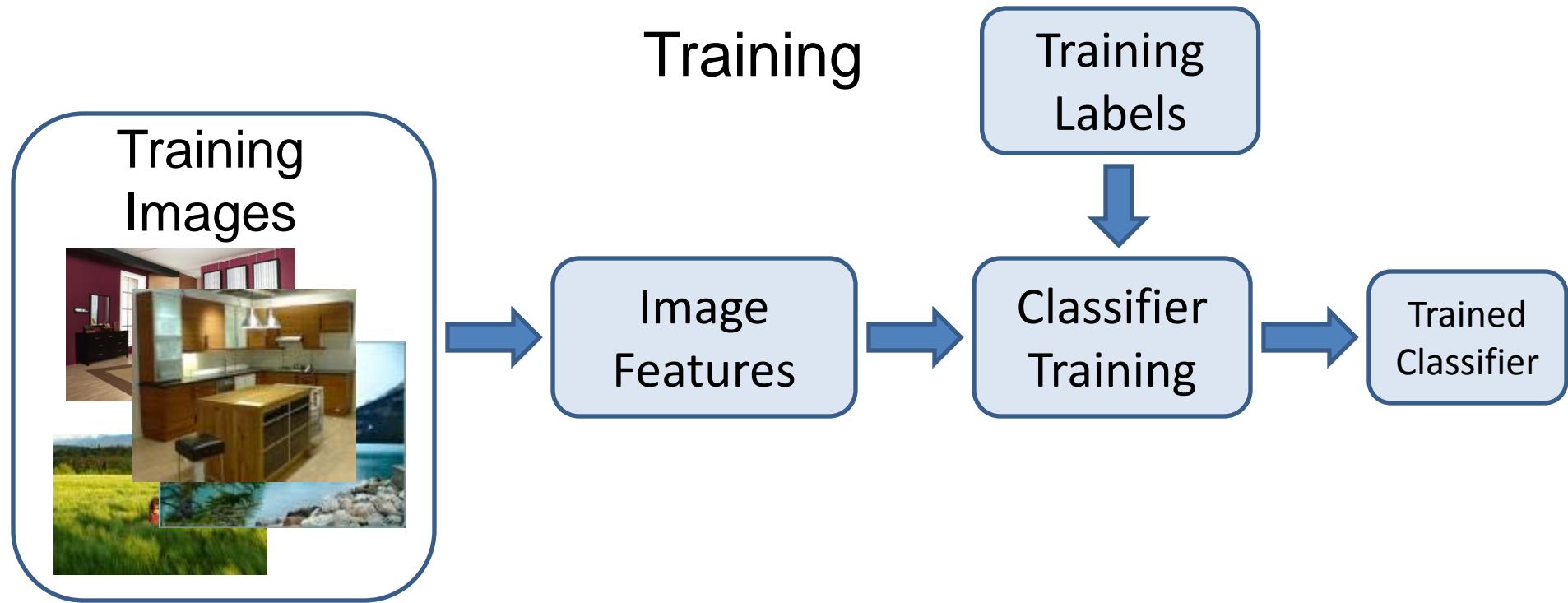
$$K_{\Delta} (\Psi(\mathbf{X}), \Psi(\mathbf{Y})) = \sum_{i=0}^L \underbrace{\frac{1}{2^i} \left( \mathcal{I}(H_i(\mathbf{X}), H_i(\mathbf{Y})) - \mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y})) \right)}_{\text{difficulty of a match at level } i}$$

number of new matches at level i

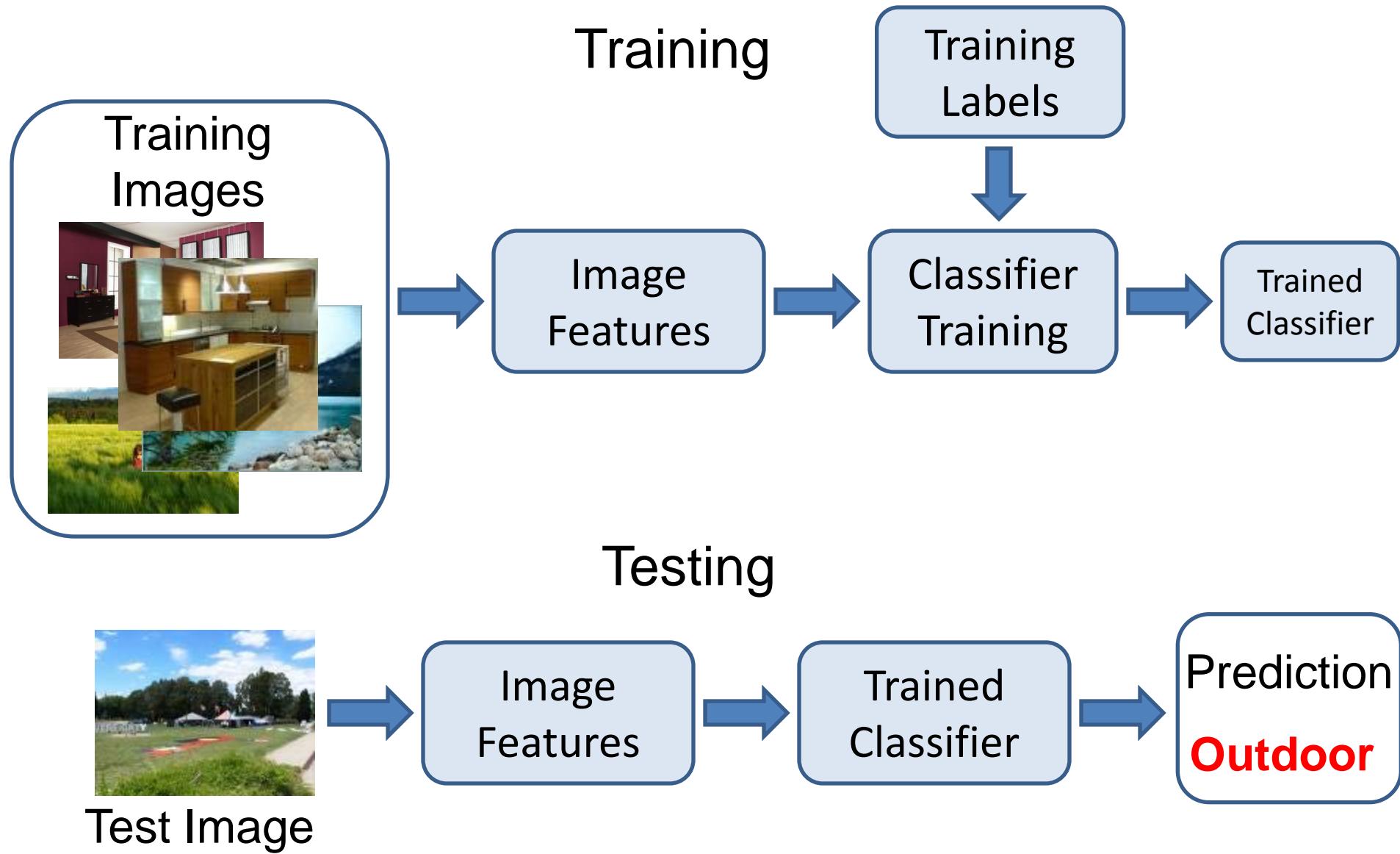
# Last Class: Spatial Pyramid



# Image Categorization



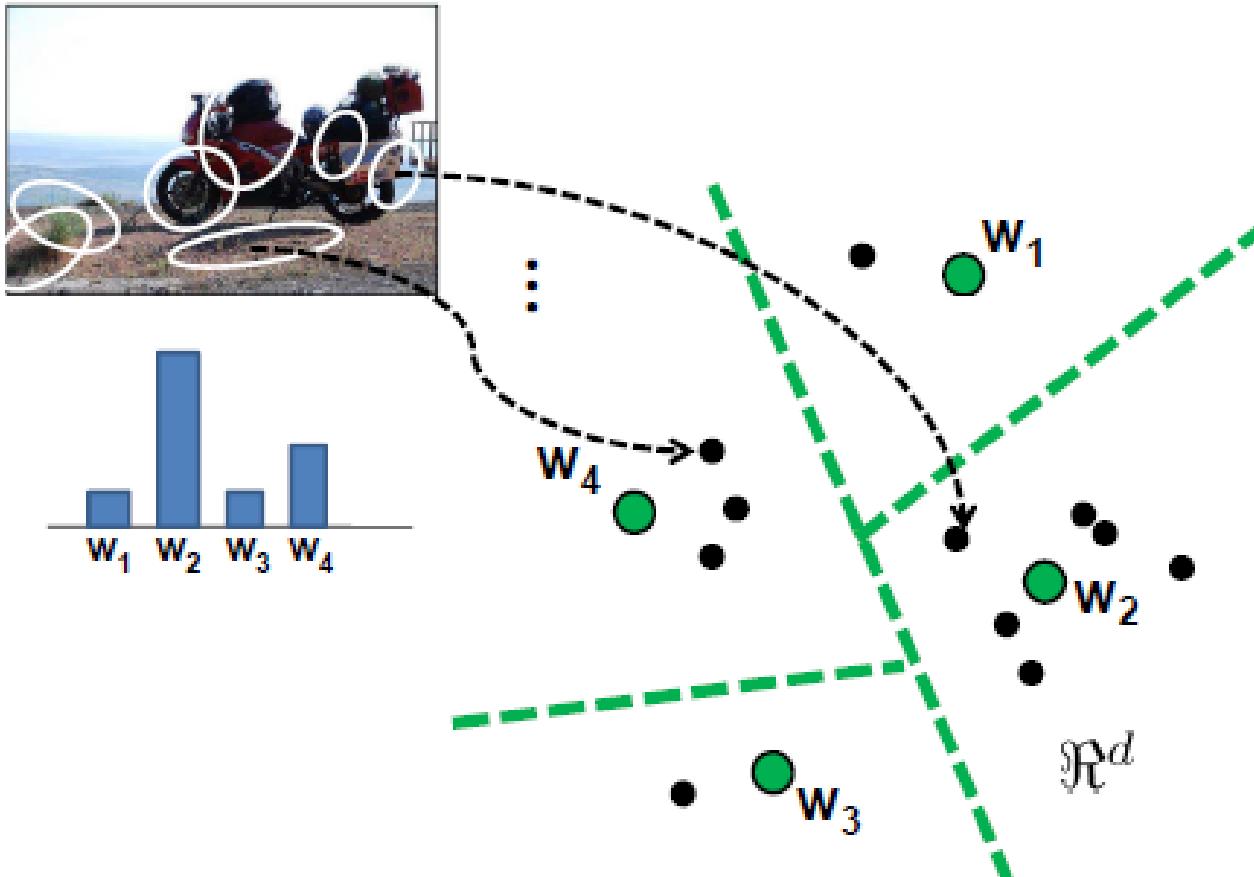
# Image Categorization



# Why do good recognition systems go bad?

- E.g. Why isn't our Bag of Words classifier at 90% instead of 70%?
- Training Data
  - Huge issue, but not necessarily a variable you can manipulate.
- Learning method
  - Probably not such a big issue, unless you're learning the representation (e.g. deep learning).
- Representation
  - Are the local features themselves lossy
  - What about feature quantization? That's VERY lossy.

# Standard Kmeans Bag of Words



[http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag\\_of\\_visual\\_words.pdf](http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag_of_visual_words.pdf)

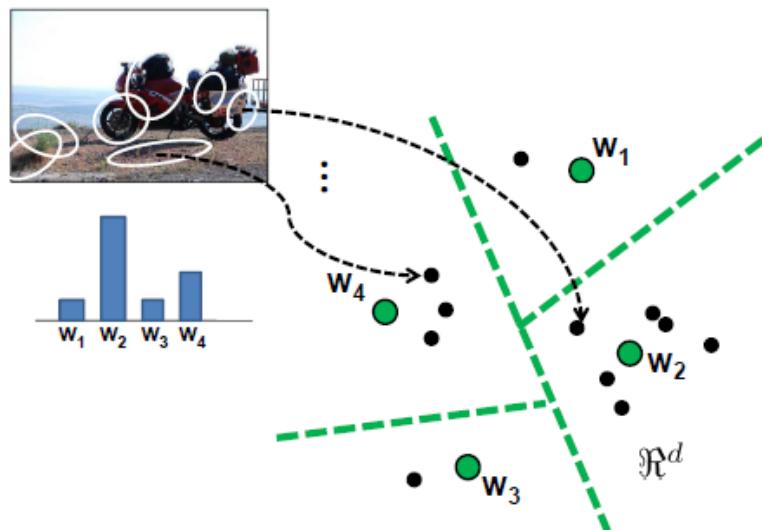
# Today's Class

- More advanced quantization / encoding methods that represent the state-of-the-art in image classification and image retrieval.
  - Soft assignment (a.k.a. Kernel Codebook)
  - VLAD
  - Locally Linear Coding
- Classification based on Latent Variables

# Motivation

*Bag of Visual Words* is only about **counting** the number of local descriptors assigned to each Voronoi region

Why not including **other statistics**?



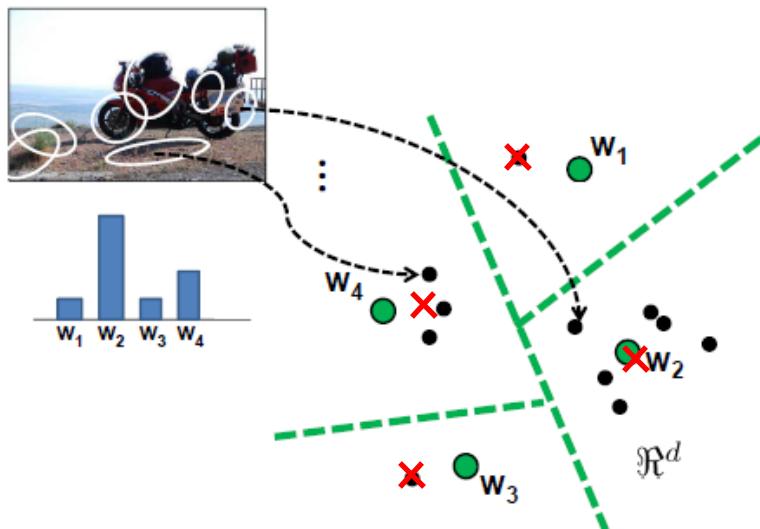
[http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag\\_of\\_visual\\_words.pdf](http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag_of_visual_words.pdf)

# Motivation

*Bag of Visual Words* is only about **counting** the number of local descriptors assigned to each Voronoi region

Why not include **other statistics**? For instance:

- mean of local descriptors X



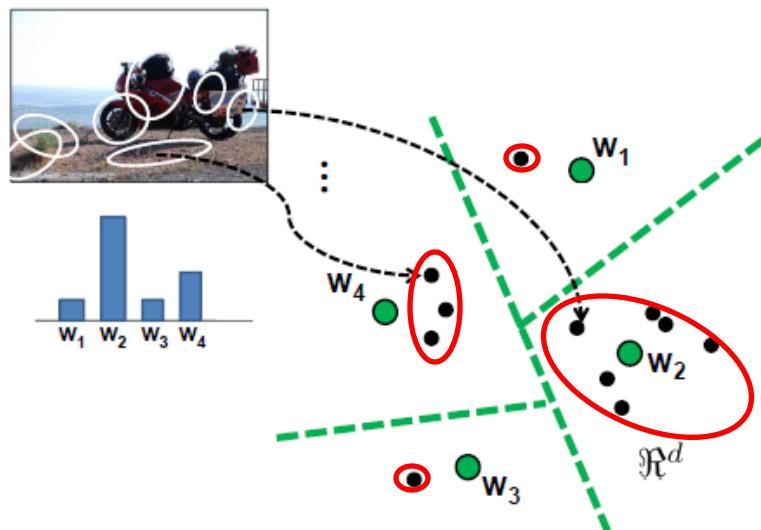
[http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag\\_of\\_visual\\_words.pdf](http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag_of_visual_words.pdf)

# Motivation

*Bag of Visual Words* is only about **counting** the number of local descriptors assigned to each Voronoi region

Why not including **other statistics**? For instance:

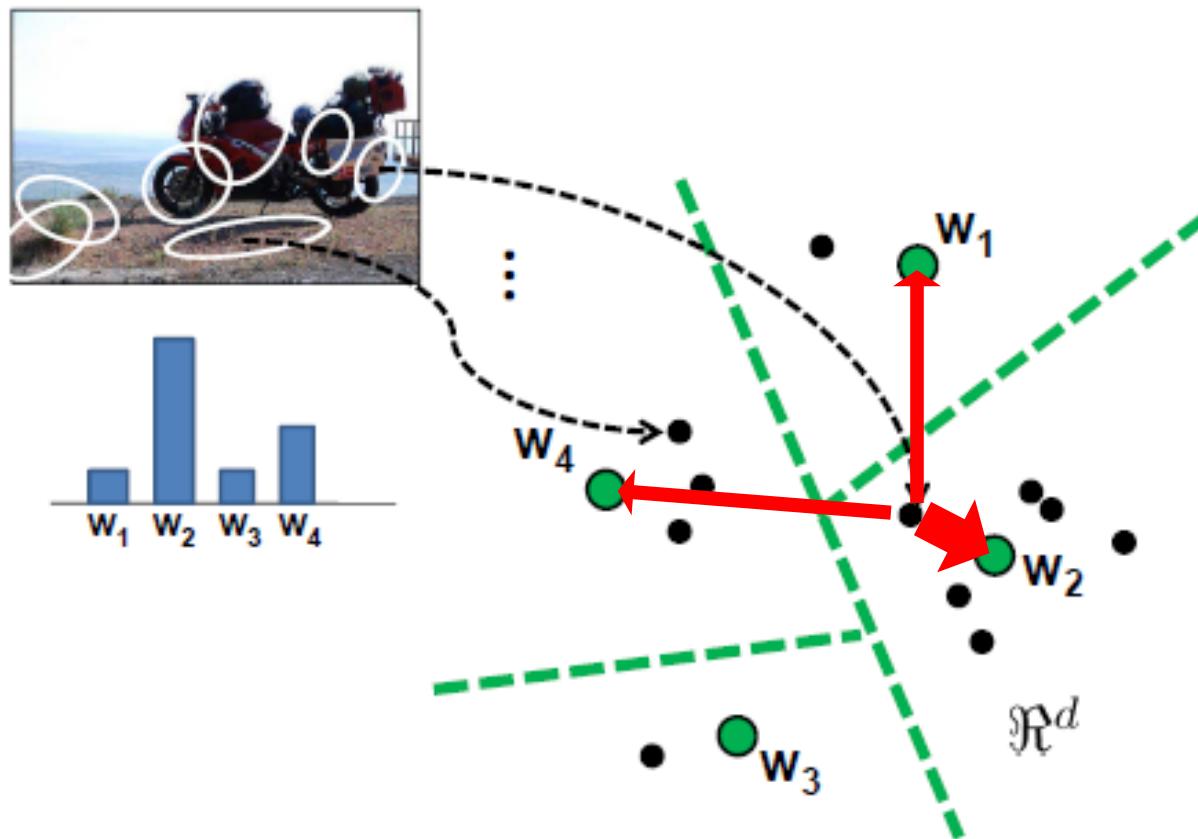
- mean of local descriptors
- (co)variance of local descriptors



[http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag\\_of\\_visual\\_words.pdf](http://www.cs.utexas.edu/~grauman/courses/fall2009/papers/bag_of_visual_words.pdf)

# Simple case: Soft Assignment

- Called “Kernel codebook encoding” by Chatfield et al. 2011. Cast a weighted vote into the most similar clusters.



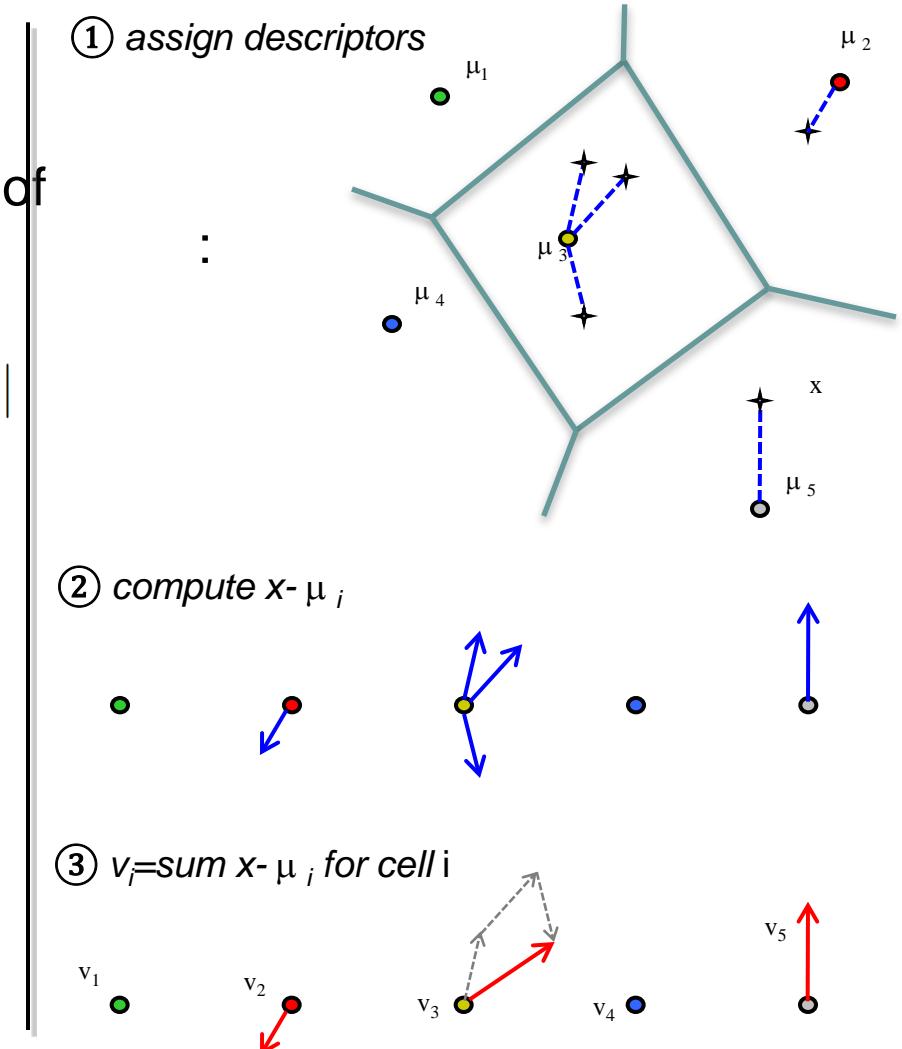
# The VLAD representation

Given a codebook  $\{\mu_i, i = 1 \dots N\}$ ,  
e.g. learned with K-means, and a set of  
local descriptors  $X = \{x_t, t = 1 \dots T\}$

- ① assign  $\text{NN}(x_t) = \arg \min_{\mu_i} \|x_t - \mu_i\|$

- ②③ compute:  $v_i = \sum_{x_t: \text{NN}(x_t) = \mu_i} x_t - \mu_i$

- concatenate  $v_i$ 's +  $\ell_2$  normalize



Jégou, Douze, Schmid and Pérez, "Aggregating local descriptors into a compact image representation", CVPR'10.

# A first example: the VLAD

A graphical representation of  $v_i = \sum_{x_t: \text{NN}(x_t) = \mu_i} x_t - \mu_i$



Jégou, Douze, Schmid and Pérez, "Aggregating local descriptors into a compact image representation", CVPR'10.

# Local Coordinate Coding (LCC)

Yu, Zhang & Gong, NIPS 09

Wang, Yang, Yu, Lv, Huang CVPR 10

- Dictionary Learning: k-means (or hierarchical k-means)

- Coding for  $x$ , to obtain its sparse representation  $a$

Step 1 – ensure locality: find the  $K$  nearest bases

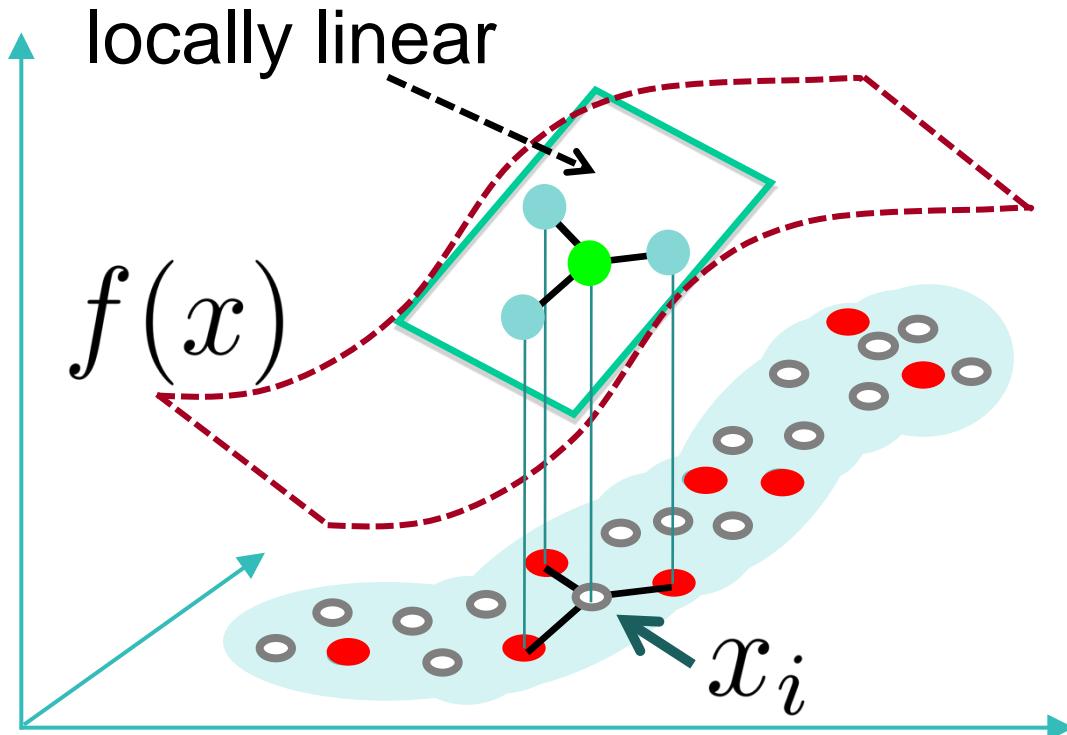
$$[\phi_j]_{j \in J(x)}$$

Step 2 – ensure low coding error:

$$\min_a \left\| x - \sum_{j \in J(x)} a_{i,j} \phi_j \right\|^2, \quad \text{s.t.} \quad \sum_{j \in J(x)} a_{i,j} = 1$$

# Function Approximation based on LCC

Yu, Zhang, Gong, NIPS 10



○ data points

● bases

# Object Classification with Latent Variables

IJCV 2013, BMVC 2011

H. Bilen, V.P. Namboodiri and L.V. Gool

VISICS, ESAT, KU Leuven

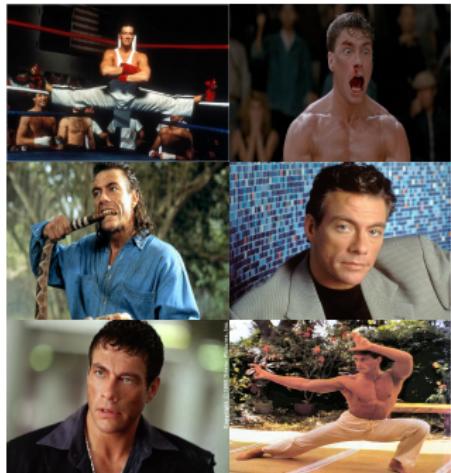
November 19, 2013



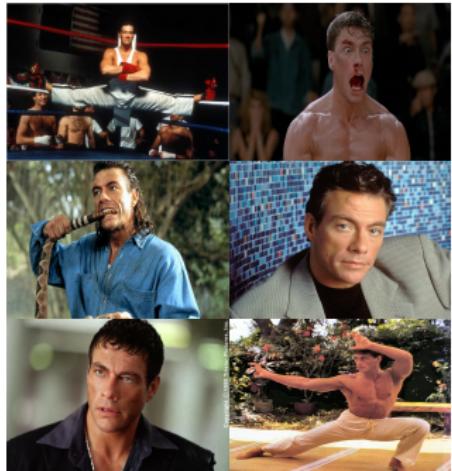
- 1 Recognition Problems
- 2 Motivation and Challenges
- 3 Contributions
- 4 Previous Work
- 5 Learning Spatial Pyramids



?

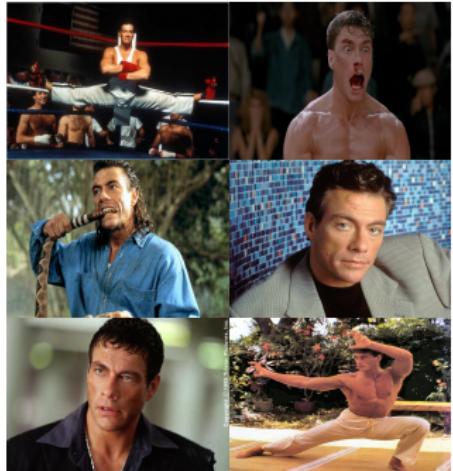


?



observed data



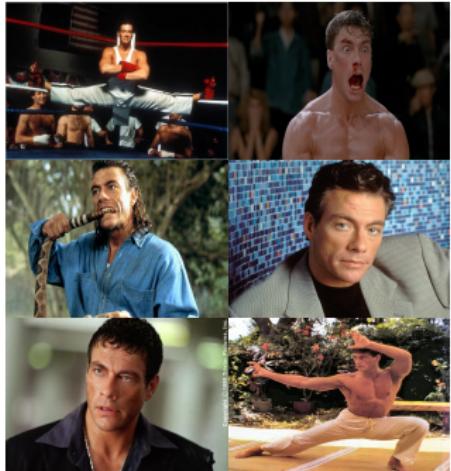


observed data

↔  
**LABEL**



Jean Claude Van Damme



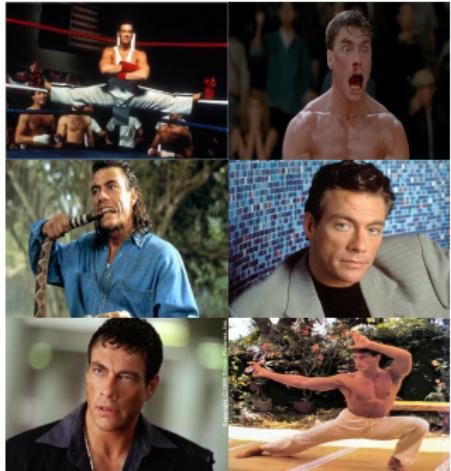
LABEL



Jean Claude Van Damme

observed data

- instance specific recognition



↔

LABEL



Jean Claude Van Damme

observed data

- instance specific recognition
- generic class recognition e.g. bicycle, person

# What is an object class?

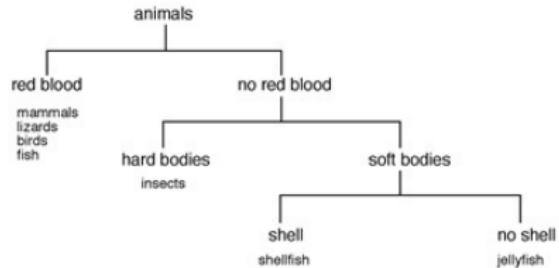
KU LEUVEN



# What is an object class?

## Classical view

- members share a set of properties



# What is an object class?

## Classical view

- members share a set of properties
- finding the common
- functional or visual?



image credits: [webdesignerdepot.com](http://webdesignerdepot.com)

# What is an object class?

## Classical view

- members share a set of properties
- finding the common
- functional or visual?



## Exemplar view

- shared properties  $\Rightarrow$  *similarity*



image credits: PASCAL VOC 2007 dataset

# What is an object class?

## Classical view

- members share a set of properties
- finding the common
- functional or visual?



## Exemplar view

- shared properties  $\Rightarrow$  *similarity*



image credits: PASCAL VOC 2007 dataset

# What is an object class?

## Classical view

- members share a set of properties
- finding the common
- functional or visual?



## Exemplar view

- shared properties  $\Rightarrow$  *similarity*
- difficult to evaluate
- requires a huge collection



image credits: PASCAL VOC 2007 dataset

## Classical view

- members share a set of properties
- finding the common
- functional or visual?



motorbike

not motorbike

## Exemplar view

- shared properties  $\Rightarrow$  *similarity*
- difficult to evaluate
- requires a huge collection

## In practice

- closed, well defined sets
- train and test sets

# What is an object class?

## Classical view

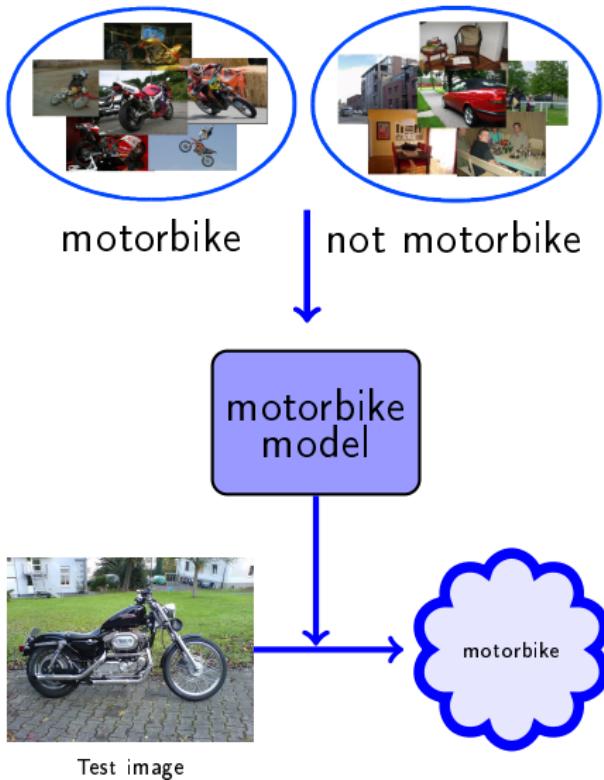
- members share a set of properties
- finding the common
- functional or visual?

## Exemplar view

- shared properties  $\Rightarrow$  *similarity*
- difficult to evaluate
- requires a huge collection

## In practice

- closed, well defined sets
- train and test sets



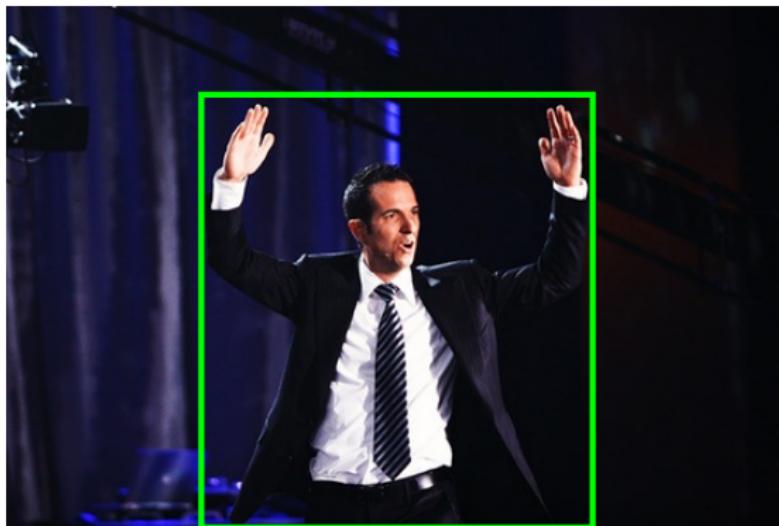
---

→ localization

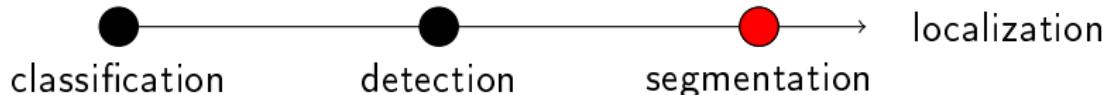
classification → localization



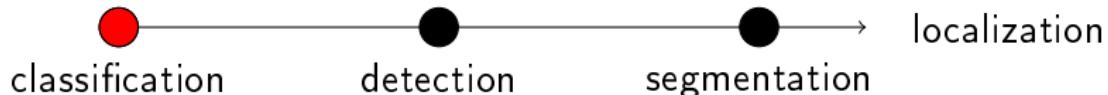
*label*



*bounding box and label*



*label* at each pixel



*label*

# Why does it matter?

KU LEUVEN



## Intellectual curiosity

- Can we extract semantic information from a signal?
- Can we make machines see?



## Intellectual curiosity

- Can we extract semantic information from a signal?
- Can we make machines see?



## Applications

- Assistive technologies
- Semantic image and video search
- Surveillance
- Medical imaging
- Human-computer interaction
- ...

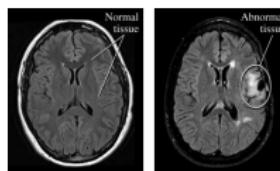


Figure 1

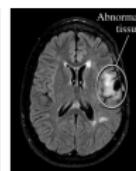


Figure 2



image credits: blog.timesunion.com, renttycoons.com, Intermountain Medical Imaging, digitaltrends.com, iwatchsystems.com

## Sources of variation

- illumination
- viewpoint
- background clutter
- occlusion
- intra-class variability
- intra-class vs. inter-class variability

## Sources of variation

- illumination
- viewpoint
- background clutter
- occlusion
- intra-class variability
- intra-class vs. inter-class variability



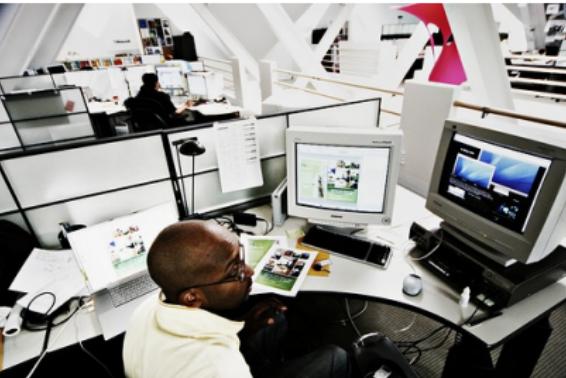
## Sources of variation

- illumination
- viewpoint
- background clutter
- occlusion
- intra-class variability
- intra-class vs. inter-class variability



## Sources of variation

- illumination
- viewpoint
- **background clutter**
- occlusion
- intra-class variability
- intra-class vs. inter-class variability



## Sources of variation

- illumination
- viewpoint
- background clutter
- **occlusion**
- intra-class variability
- intra-class vs. inter-class variability



## Sources of variation

- illumination
- viewpoint
- background clutter
- occlusion
- intra-class variability
- intra-class vs. inter-class variability



## Sources of variation

- illumination
- viewpoint
- background clutter
- occlusion
- intra-class variability
- intra-class vs. inter-class variability



Shih-Tzu

## Maltese Dog



develop invariant representations

## develop invariant representations

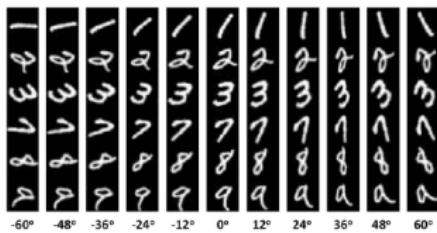
- pros: no extra parameters
- cons: less discriminative

7	2	1	0	4	1	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	5	1	2	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	4	6	4	3
7	0	2	9	1	7	3	2	9	7
7	6	2	7	8	4	4	7	3	6
3	6	9	3	1	4	1	7	6	9

## develop invariant representations

- pros: no extra parameters
- cons: less discriminative

## model variances explicitly

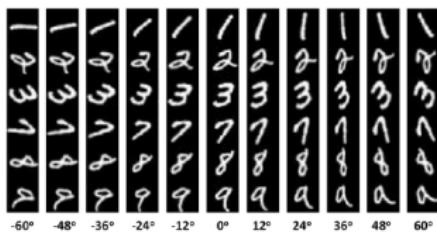


## develop invariant representations

- pros: no extra parameters
- cons: less discriminative

## model variances explicitly

- pros: more discriminative
- cons:
  - extra parameters

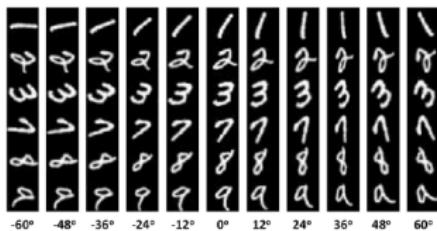


## develop invariant representations

- pros: no extra parameters
- cons: less discriminative

## model variances explicitly

- pros: more discriminative
- cons:
  - extra parameters
  - annotation



- time-consuming, expensive, not much fun!

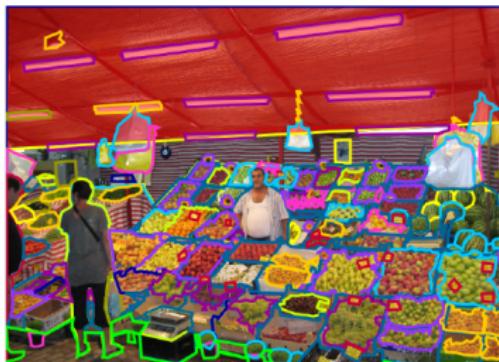


image credits: Barriosa and Torralba

- time-consuming, expensive, not much fun!



image credits: Barriosa and Torralba

improve classification task by

improve classification task by

- modeling various sources of variation,

improve classification task by

- modeling various sources of variation,
- without requiring any annotation but only class labels,

improve classification task by

- modeling various sources of variation,
- without requiring any annotation but only class labels,
- incorporating them into image representation as **latent variables**,

improve classification task by

- modeling various sources of variation,
- without requiring any annotation but only class labels,
- incorporating them into image representation as **latent variables**,
- jointly learning both classifier parameters and latent parameters in a discriminative setting

## Object Classification

- state of the art classification methods

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations
- various quantization and encoding methods are compatible with our work

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations
- various quantization and encoding methods are compatible with our work

## Latent Variable Models

- generative methods (GMMs, HMMs)

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations
- various quantization and encoding methods are compatible with our work

## Latent Variable Models

- generative methods (GMMs, HMMs)
- discriminative methods (Hidden CRFs [Wang@CVPR06], Latent SVM [Felzenswalb@PAMI10, Yu@ICML09])

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations
- various quantization and encoding methods are compatible with our work

## Latent Variable Models

- generative methods (GMMs, HMMs)
- discriminative methods (Hidden CRFs [Wang@CVPR06], Latent SVM [Felzenswalb@PAMI10, Yu@ICML09])
- in both cases parameter learning is non-convex

## Object Classification

- state of the art classification methods
  - better feature quantization and encoding, LLC [Wang@CVPR09]
  - more powerful statistics, I. Fisher Encoding [Perronnin@ECCV10]
  - multiple features [Gehler@ICCV09], classifiers [Torresani@ECCV10]
- static representations
- various quantization and encoding methods are compatible with our work

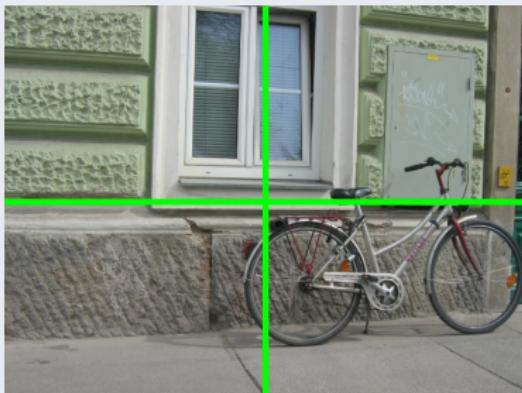
## Latent Variable Models

- generative methods (GMMs, HMMs)
- discriminative methods (Hidden CRFs [Wang@CVPR06], Latent SVM [Felzenswalb@PAMI10, Yu@ICML09])
- in both cases parameter learning is non-convex
  - EM (generative)
  - CCCP (discriminative)

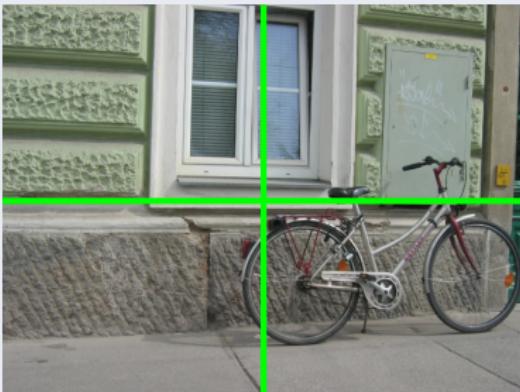
## Spatial Pyramids [Lazebnik@CVPR06]



## Spatial Pyramids [Lazebnik@CVPR06]

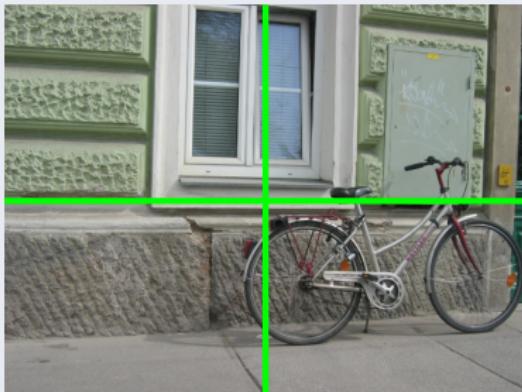


## Spatial Pyramids [Lazebnik@CVPR06]



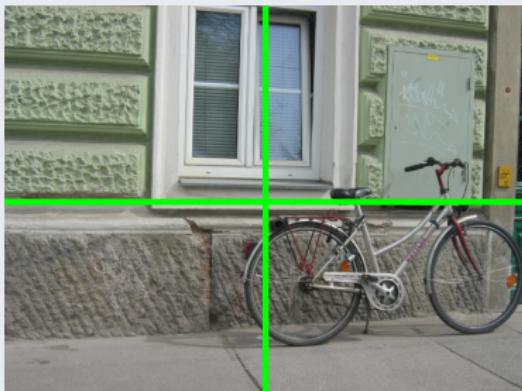
- pros: adds spatial information to the orderless BoW

## Spatial Pyramids [Lazebnik@CVPR06]



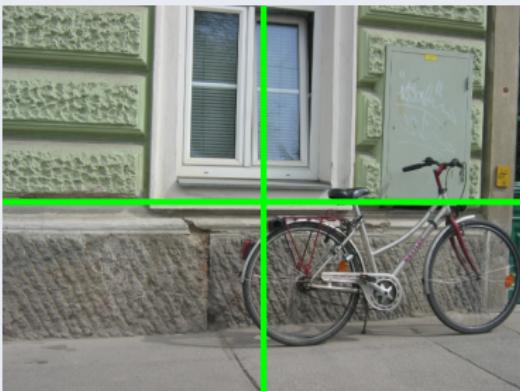
- pros: adds spatial information to the orderless BoW
- cons:

## Spatial Pyramids [Lazebnik@CVPR06]



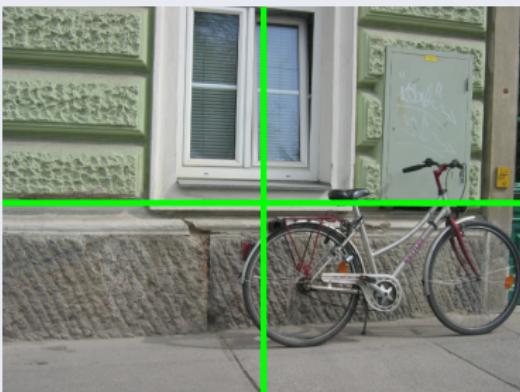
- pros: adds spatial information to the orderless BoW
- cons:
  - static, not optimal

## Spatial Pyramids [Lazebnik@CVPR06]



- pros: adds spatial information to the orderless BoW
- cons:
  - static, not optimal
  - features in cells are not aligned between images

## Spatial Pyramids [Lazebnik@CVPR06]

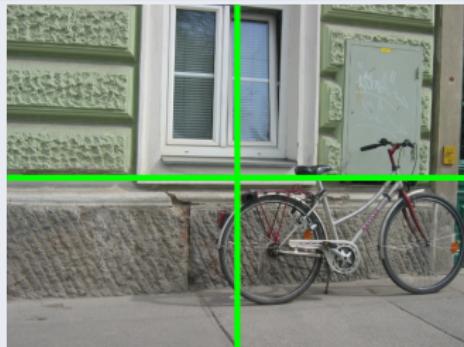


- pros: adds spatial information to the orderless BoW
- cons:
  - static, not optimal
  - features in cells are not aligned between images
  - includes discriminative and non-representative features

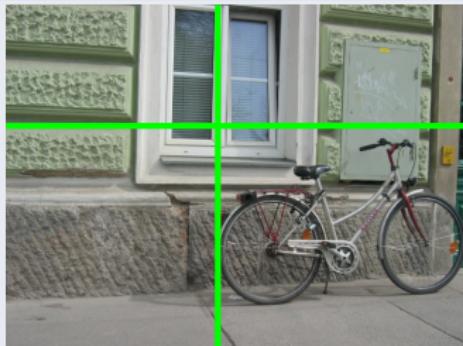
Split



Split

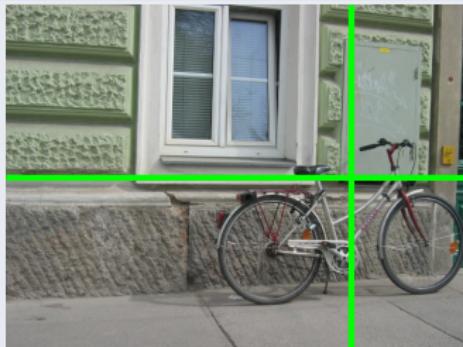


Split



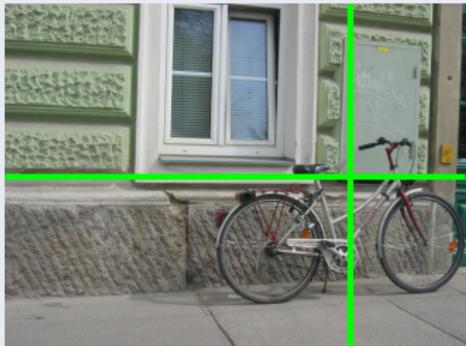
# More Flexible Image Representation

Split

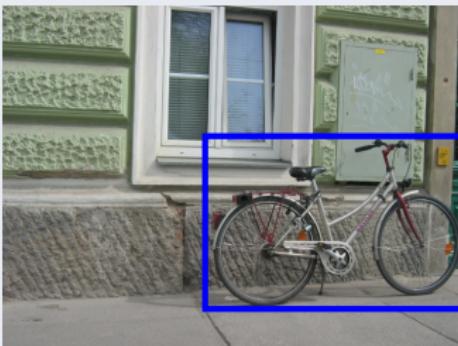


# More Flexible Image Representation

Split



Crop

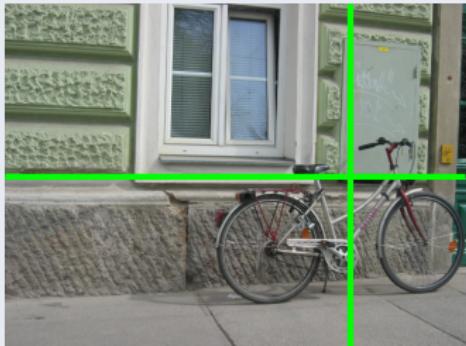


# More Flexible Image Representation

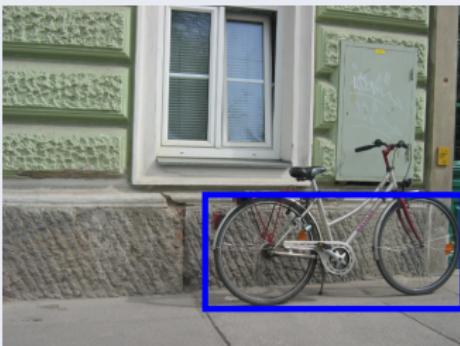
KU LEUVEN



Split



Crop

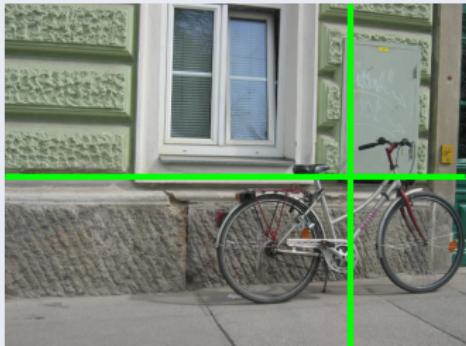


# More Flexible Image Representation

KU LEUVEN



Split

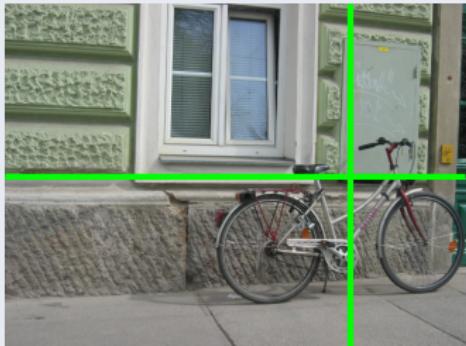


Crop

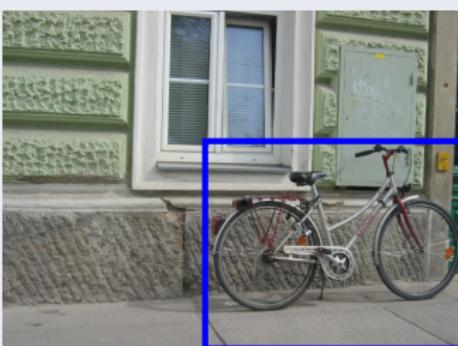


# More Flexible Image Representation

Split



Crop

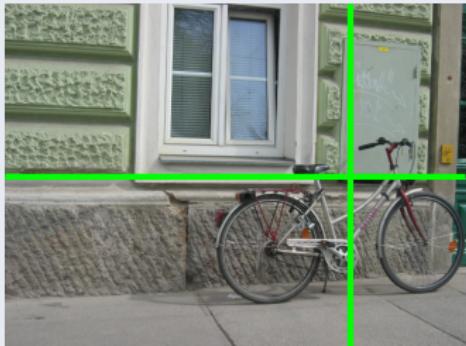


# More Flexible Image Representation

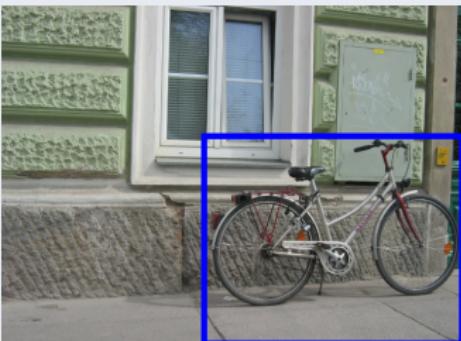
KU LEUVEN



Split



Crop



Crop - Uniform Split

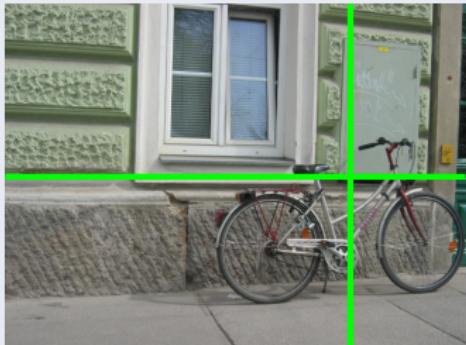


# More Flexible Image Representation

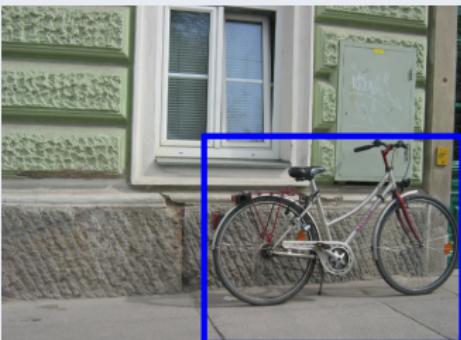
KU LEUVEN



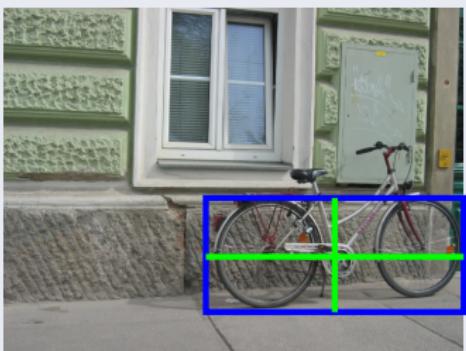
Split



Crop



Crop - Uniform Split



# More Flexible Image Representation

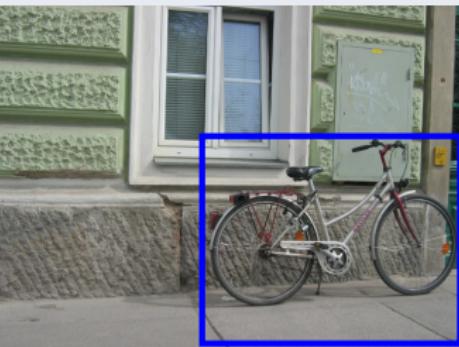
KU LEUVEN



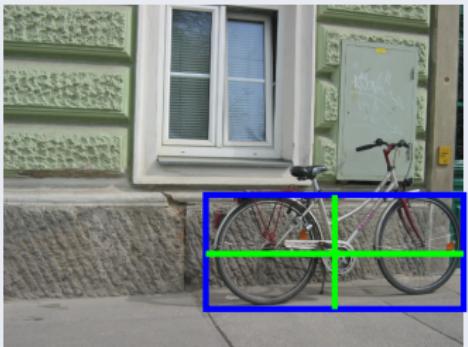
Split



Crop



Crop - Uniform Split

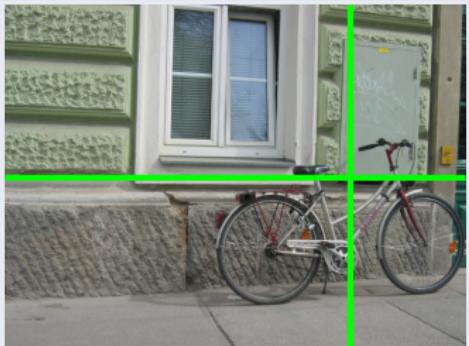


Crop - Split

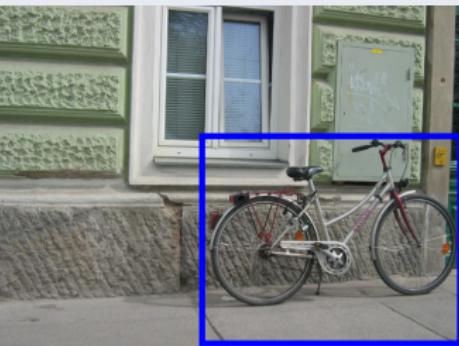


# More Flexible Image Representation

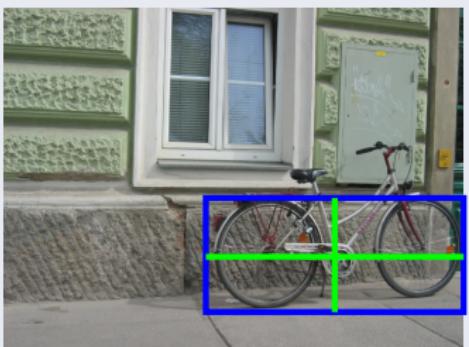
Split



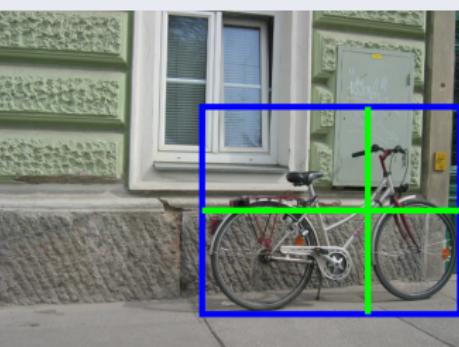
Crop



Crop - Uniform Split



Crop - Split



- learning classifier parameters and image windows/non-uniform grid are dependent tasks

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- image windows/non-uniform grid (crop/split) as latent parameters

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- image windows/non-uniform grid (crop/split) as latent parameters
- alternating optimization between classification parameters and latent parameters

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- image windows/non-uniform grid (crop/split) as latent parameters
- alternating optimization between classification parameters and latent parameters
- is non-convex and requires initialization of classification or latent parameters

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- image windows/non-uniform grid (crop/split) as latent parameters
- alternating optimization between classification parameters and latent parameters
- is non-convex and requires initialization of classification or latent parameters
- is guaranteed to decrease the objective function and to converge to a local minimum

- learning classifier parameters and image windows/non-uniform grid are dependent tasks
- disjoint learning is not optimal in terms of classification performance

## Latent (Structural) SVM [Yu@ICML09]

- image windows/non-uniform grid (crop/split) as latent parameters
- alternating optimization between classification parameters and latent parameters
- is non-convex and requires initialization of classification or latent parameters
- is guaranteed to decrease the objective function and to converge to a local minimum
- We use Latent SVM for learning

Training pairs  $\{(x_i, y_i, h_i)\} \in \mathcal{X} \times \mathcal{Y} \times \mathcal{H}$

- $x_i$  are images/videos,  $y_i$  are class labels,
- $h_i$  are coordinates of image windows/non-uniform grids.
- only  $x_i$  and  $y_i$  are given,  $h_i$  are not given

Training pairs  $\{(x_i, y_i, h_i)\} \in \mathcal{X} \times \mathcal{Y} \times \mathcal{H}$

- $x_i$  are images/videos,  $y_i$  are class labels,
- $h_i$  are coordinates of image windows/non-uniform grids.
- only  $x_i$  and  $y_i$  are given,  $h_i$  are not given

Discriminant function that measures the matching quality of the triplet:

$$f : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathcal{R}$$

$$f_w(x, y, h) = w^T \psi(x, y, h)$$

Training pairs  $\{(x_i, y_i, h_i)\} \in \mathcal{X} \times \mathcal{Y} \times \mathcal{H}$

- $x_i$  are images/videos,  $y_i$  are class labels,
- $h_i$  are coordinates of image windows/non-uniform grids.
- only  $x_i$  and  $y_i$  are given,  $h_i$  are not given

Discriminant function that measures the matching quality of the triplet:

$$f : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathcal{R}$$

$$f_w(x, y, h) = w^T \psi(x, y, h)$$

- $\psi(x, y, h)$  is a histogram of quantized SIFT features which is obtained through the image window or non-uniform grid  $h$ .

Training pairs  $\{(x_i, y_i, h_i)\} \in \mathcal{X} \times \mathcal{Y} \times \mathcal{H}$

- $x_i$  are images/videos,  $y_i$  are class labels,
- $h_i$  are coordinates of image windows/non-uniform grids.
- only  $x_i$  and  $y_i$  are given,  $h_i$  are not given

Discriminant function that measures the matching quality of the triplet:

$$f : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathcal{R}$$

$$f_w(x, y, h) = w^T \psi(x, y, h)$$

- $\psi(x, y, h)$  is a histogram of quantized SIFT features which is obtained through the image window or non-uniform grid  $h$ .
- $w$  is the parameter vector of the latent SVM.

$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$

$$\min_w \left\{ \underbrace{\Omega(w)}_{\text{regularization}} + C \sum_{i=1}^n \left[ \underbrace{\max_{\hat{y}_i, \hat{h}_i} [f_w(x_i, \hat{y}_i, \hat{h}_i) + \Delta(y_i, \hat{y}_i)]}_{\text{most violated constraint}} - \underbrace{\max_{h_i^*} [f_w(x_i, y_i, h_i^*)]}_{\text{best latent variable}} \right] \right\}$$

- Loss function  $\Delta(y_i, \hat{y}_i)$



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$

$$\min_w \left\{ \underbrace{\Omega(w)}_{\text{regularization}} + C \sum_{i=1}^n \left[ \underbrace{\max_{\hat{y}_i, \hat{h}_i} [f_w(x_i, \hat{y}_i, \hat{h}_i) + \Delta(y_i, \hat{y}_i)]}_{\text{most violated constraint}} - \underbrace{\max_{h_i^*} [f_w(x_i, y_i, h_i^*)]}_{\text{best latent variable}} \right] \right\}$$

- Loss function  $\Delta(y_i, \hat{y}_i)$
- Iterate

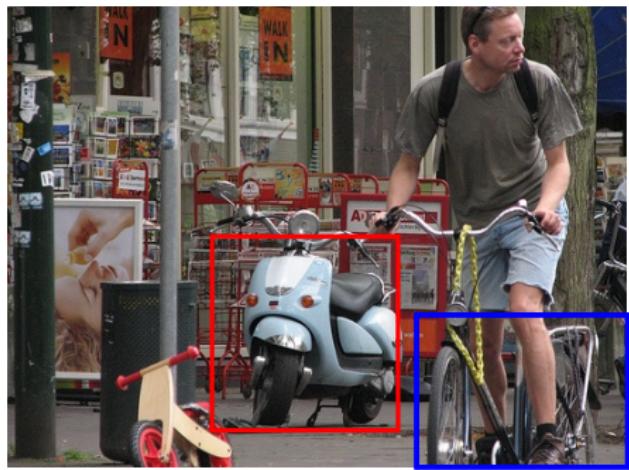


$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$

$$\min_w \left\{ \underbrace{\Omega(w)}_{\text{regularization}} + C \sum_{i=1}^n \left[ \underbrace{\max_{\hat{y}_i, \hat{h}_i} [f_w(x_i, \hat{y}_i, \hat{h}_i) + \Delta(y_i, \hat{y}_i)]}_{\text{most violated constraint}} - \underbrace{\max_{h_i^*} [f_w(x_i, y_i, h_i^*)]}_{\text{best latent variable}} \right] \right\}$$

- Loss function  $\Delta(y_i, \hat{y}_i)$
- Iterate

1 ...



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$

$$\min_w \left\{ \underbrace{\Omega(w)}_{\text{regularization}} + C \sum_{i=1}^n \left[ \underbrace{\max_{\hat{y}_i, \hat{h}_i} [f_w(x_i, \hat{y}_i, \hat{h}_i) + \Delta(y_i, \hat{y}_i)]}_{\text{most violated constraint}} - \underbrace{\max_{h_i^*} [f_w(x_i, y_i, h_i^*)]}_{\text{best latent variable}} \right] \right\}$$

- Loss function  $\Delta(y_i, \hat{y}_i)$
- Iterate

1 ...  
2 ...



$$g_w(x) = \underbrace{\arg \max_{\hat{y} \in \mathcal{Y}, \hat{h} \in \mathcal{H}} f_w(x, \hat{y}, \hat{h})}_{\text{prediction}}$$

$$\min_w \left\{ \underbrace{\Omega(w)}_{\text{regularization}} + C \sum_{i=1}^n \left[ \underbrace{\max_{\hat{y}_i, \hat{h}_i} [f_w(x_i, \hat{y}_i, \hat{h}_i) + \Delta(y_i, \hat{y}_i)]}_{\text{most violated constraint}} - \underbrace{\max_{h_i^*} [f_w(x_i, y_i, h_i^*)]}_{\text{best latent variable}} \right] \right\}$$

- Loss function  $\Delta(y_i, \hat{y}_i)$
- Iterate
  - 1 ...
  - 2 ...
- until convergence



## Object Classification

- Graz-02, PASCAL VOC 2007,  
Caltech 101
- Densely extracted SIFT local  
features

## Object Classification

- Graz-02, PASCAL VOC 2007, Caltech 101
- Densely extracted SIFT local features

## Action Classification

- Activities of Daily Living
- Harris3d detector + HOF descriptor [Laptev@CVPR08]

## Object Classification

- Graz-02, PASCAL VOC 2007, Caltech 101
- Densely extracted SIFT local features

## Action Classification

- Activities of Daily Living
- Harris3d detector + HOF descriptor [Laptev@CVPR08]

		Graz-02	VOC-07	Caltech-101	Activities
Baseline	BoW	$86.95 \pm 1.4$	49.86	$61.25 \pm 0.9$	79.33
	SP	$88.05 \pm 1.4$	54.74	$72.68 \pm 1.2$	88.00
Ours	crop	$88.40 \pm 1.1$	51.82	$62.16 \pm 1.0$	72.00
	split	$88.58 \pm 1.3$	55.32	$73.33 \pm 1.0$	88.00
	crop-uni-split	$90.38 \pm 1.9$	56.26	$75.31 \pm 0.7$	<b>90.67</b>
	crop-split	<b>90.62 ± 1.8</b>	<b>57.05</b>	74.93 ± 0.9	88.67

## Object Classification

- Graz-02, PASCAL VOC 2007, Caltech 101
- Densely extracted SIFT local features

## Action Classification

- Activities of Daily Living
- Harris3d detector + HOF descriptor [Laptev@CVPR08]

		Graz-02	VOC-07	Caltech-101	Activities
Baseline	BoW	$86.95 \pm 1.4$	49.86	$61.25 \pm 0.9$	79.33
	SP	$88.05 \pm 1.4$	54.74	$72.68 \pm 1.2$	88.00
Ours	crop	$88.40 \pm 1.1$	51.82	$62.16 \pm 1.0$	72.00
	split	$88.58 \pm 1.3$	55.32	$73.33 \pm 1.0$	88.00
	crop-uni-split	$90.38 \pm 1.9$	56.26	$75.31 \pm 0.7$	<b>90.67</b>
	crop-split	<b>90.62 ± 1.8</b>	<b>57.05</b>	74.93 ± 0.9	88.67

- The best results are obtained with the crop-uni-split or crop-split.

## Object Classification

- Graz-02, PASCAL VOC 2007, Caltech 101
- Densely extracted SIFT local features

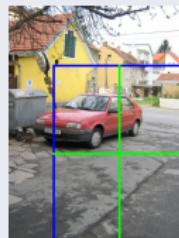
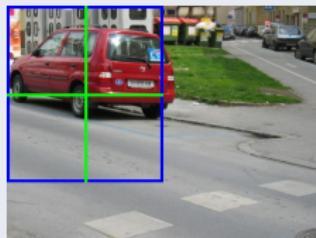
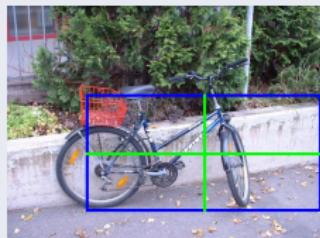
## Action Classification

- Activities of Daily Living
- Harris3d detector + HOF descriptor [Laptev@CVPR08]

		Graz-02	VOC-07	Caltech-101	Activities
Baseline	BoW	$86.95 \pm 1.4$	49.86	$61.25 \pm 0.9$	79.33
	SP	$88.05 \pm 1.4$	54.74	$72.68 \pm 1.2$	88.00
Ours	crop	$88.40 \pm 1.1$	51.82	$62.16 \pm 1.0$	72.00
	split	$88.58 \pm 1.3$	55.32	$73.33 \pm 1.0$	88.00
	crop-uni-split	$90.38 \pm 1.9$	56.26	$75.31 \pm 0.7$	<b>90.67</b>
	crop-split	<b><math>90.62 \pm 1.8</math></b>	<b>57.05</b>	$74.93 \pm 0.9$	88.67

- The best results are obtained with the crop-uni-split or crop-split.
- Both crop-uni-split and crop-split are always better than SP.

## crop-uniform split examples



## crop-split examples

