

25  
SOLCTĐẠI HỌC BÁCH KHOA HÀ NỘI  
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

# Thị giác máy tính

## Bài 8: Phân loại ảnh

### Nội dung buổi học

- Tổng quan về thị giác ngữ nghĩa
- Bài toán phân lớp / nhận dạng ảnh
- Mô hình Bag-of-words
  - Mô hình
  - Vocabulary tree
- Đánh giá
- Phân lớp và lựa chọn siêu tham số?

# Thị giác ngữ nghĩa Semantic vision



3

Đây có phải đèn giao thông?  
(Phân lớp / nhận dạng – Classification / recognition)



4

## Con người xuất hiện ở đâu trên ảnh? (Phát hiện - Detection)



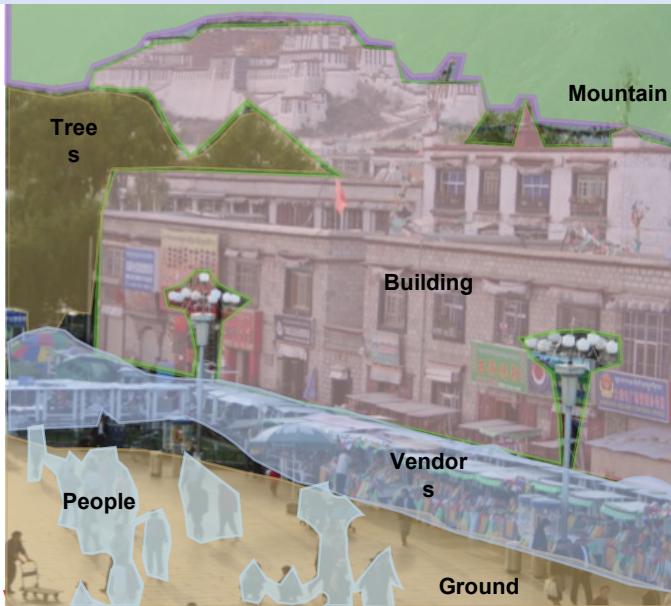
5

## Đây có đúng là địa điểm Potala? (Định danh - Identification)



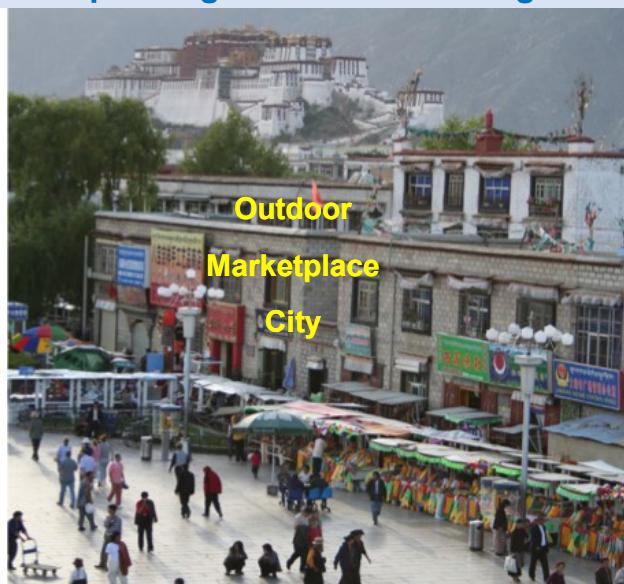
6

Có những gì trong bức ảnh này?  
(Phân vùng ngữ nghĩa - semantic segmentation)



7

Đây là loại khung cảnh gì?  
(Phân lớp khung cảnh - Scene categorization)



8

## Những người này đang làm gì?

(Nhận dạng hành động / sự kiện - Activity / Event Recognition)



## Nhận dạng đối tượng khó không?

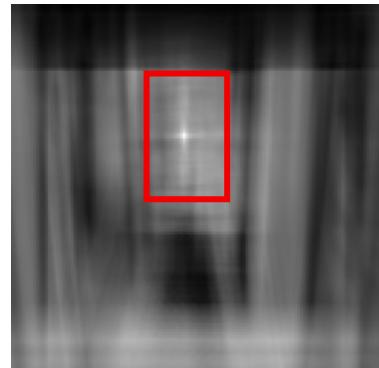
Đây là cái ghế



Hãy tìm cái ghế trong ảnh



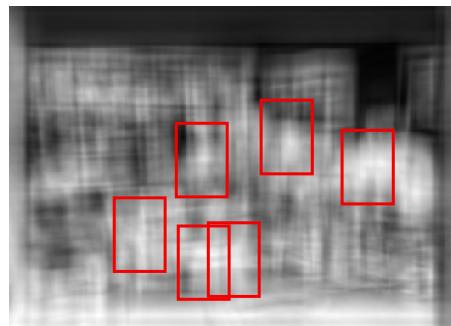
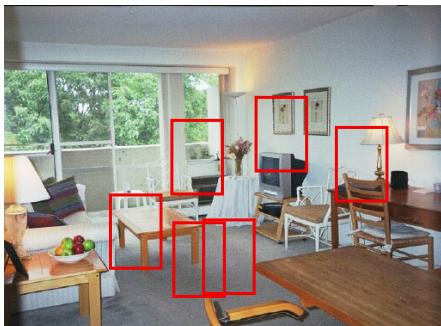
Độ tương quan sau khi chuẩn hóa



## Nhận dạng đối tượng khó khăn?



Tìm cái ghế trong bức ảnh này



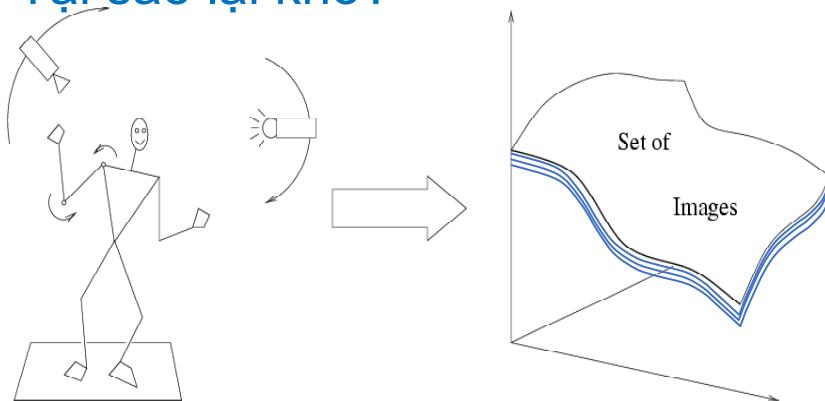
Quá nhiều nhiễu

Kỹ thuật template matching đơn giản không thể  
giải quyết được

## Ví dụ khác ...



## Tại sao lại khó?



Đa dạng:  
Vị trí camera  
Ánh sáng  
Hình dáng vật thay đổi



## Thách thức: các góc chụp khác nhau



Michelangelo 1475-1564



15

## Thách thức: ánh sáng thay đổi

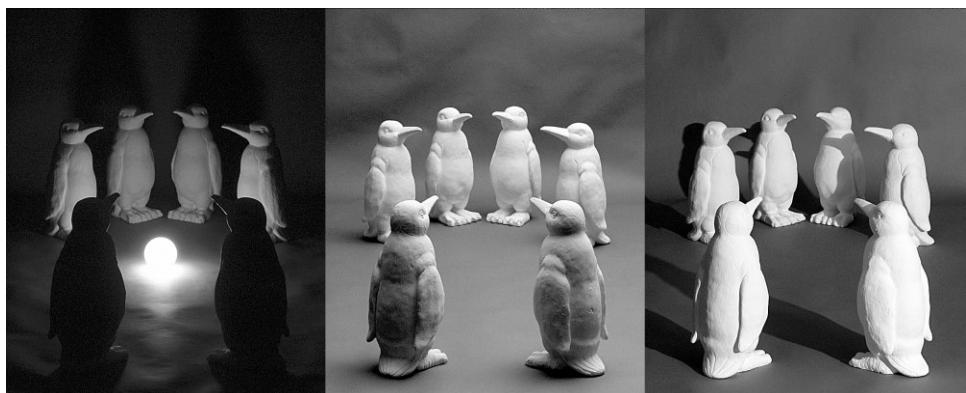


image credit: J. Koenderink

16

## Thách thức: kích thước đa dạng

and small things

from Apple.

(Actual size)



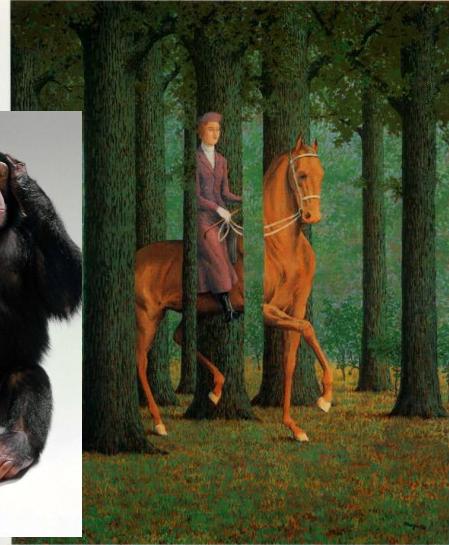
17

## Thách thức: đối tượng biến dạng



18

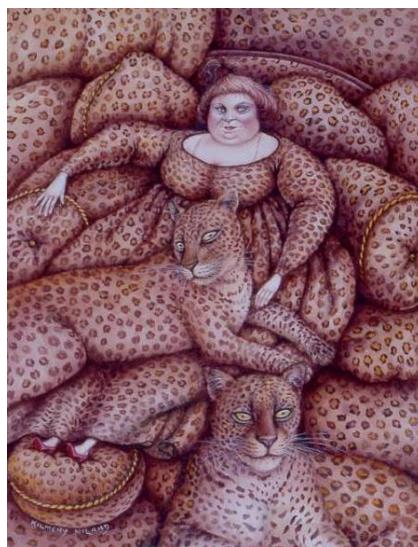
## Thách thức: che khuất



Magritte, 1957

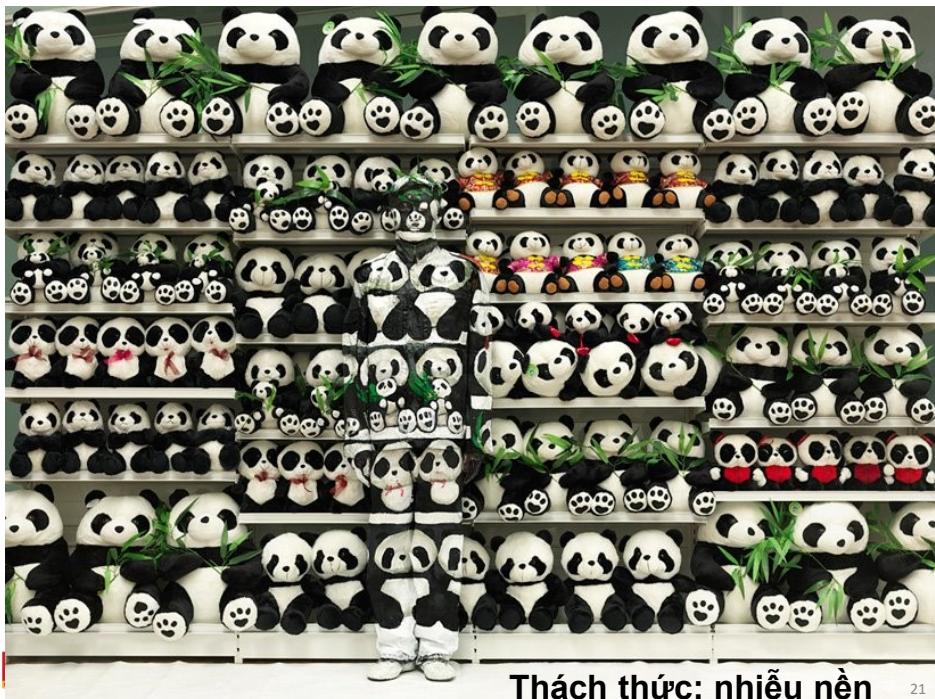
19

## Thách thức: nhiễu/ nền



Kilmeny Niland. 1995

20



## Thách thức: sự đa dạng về đối tượng trong từng lớp



Svetlana Lazebnik

## Phân lớp / nhận dạng ảnh

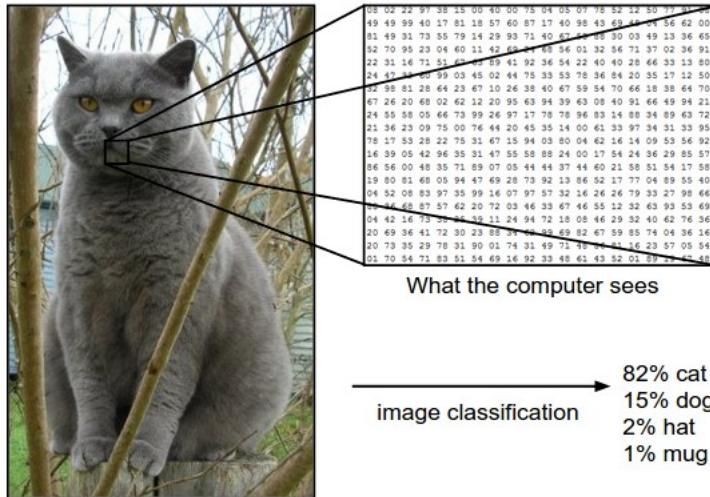
## Phân lớp / nhận dạng ảnh



(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}

→ cat

## Bài toán phân lớp ảnh



## Tiếp cận dựa trên dữ liệu

- Thu thập bộ dữ liệu ảnh kèm nhãn các lớp
- Sử dụng ML huấn luyện một bộ phân lớp
- Đánh giá bộ phân lớp trên ảnh test

Example training set



## Các bước huấn luyện

Ảnh huấn  
luyện



Đặc trưng  
ảnh

## Các bước huấn luyện

Ảnh huấn  
luyện

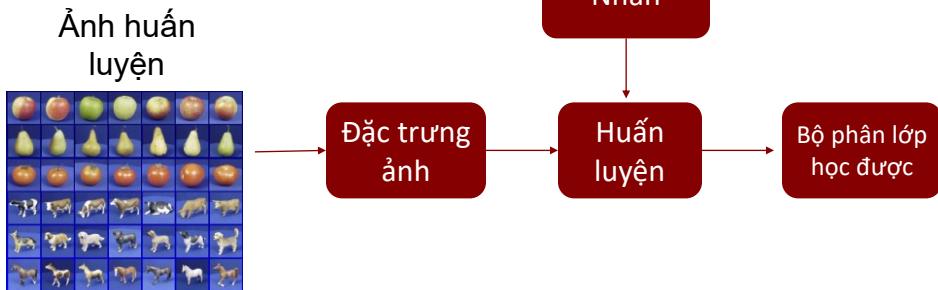


Đặc trưng  
ảnh

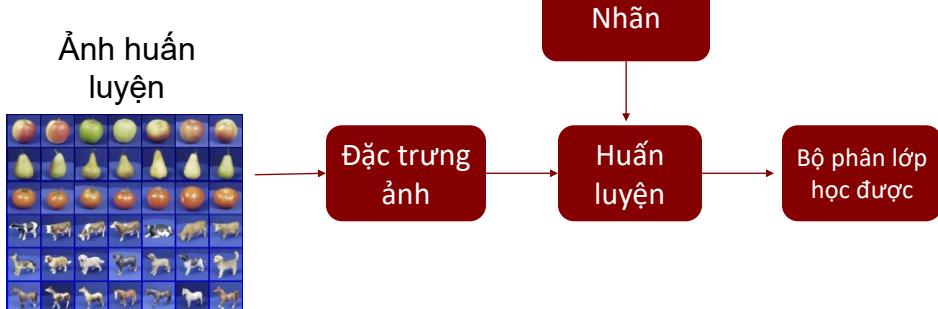
Nhãn

Huấn  
luyện

## Các bước huấn luyện



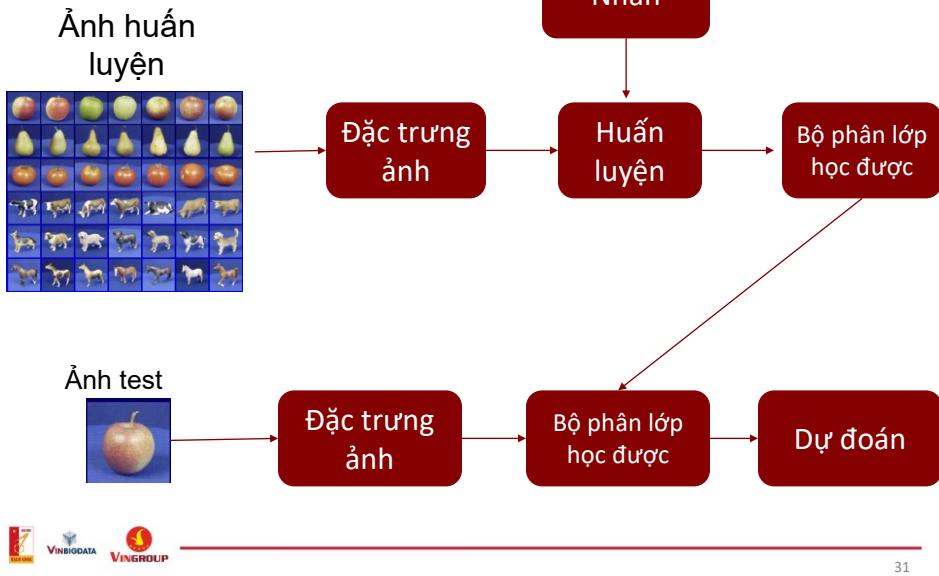
## Các bước huấn luyện



### Ảnh test



## Các bước huấn luyện



## Trích chọn đặc trưng (nhắc lại)

- Hai loại đặc trưng được trích chọn từ ảnh:
  - Đặc trưng cục bộ và toàn cục
- **Đặc trưng toàn cục:**
  - Mô tả toàn bộ ảnh như 1 đối tượng
  - Đặc trưng đường biên, đặc trưng hình dạng, đặc trưng kết cấu
    - Ví dụ: Invariant Moments (Hu, Zernike), Histogram Oriented Gradients (HOG), PHOG, and Co-HOG, ...
- **Đặc trưng cục bộ:**
  - Mô tả đặc trưng cục bộ mô tả từng vùng nhỏ trong ảnh, từng vùng cục bộ của đối tượng (điểm đặc trưng trong ảnh).
  - Biểu diễn đặc trưng kết cấu/màu sắc trong mỗi vùng cục bộ ảnh
    - Ví dụ: SIFT, SURF, LBP, BRISK, MSER và FREAK, ...

## Mô hình Bag of words

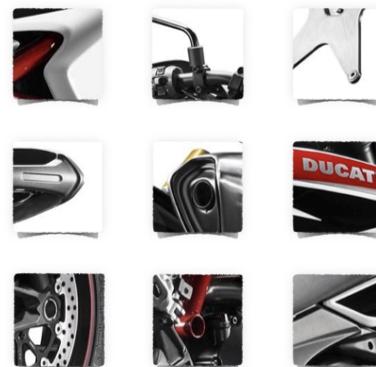


33

Một số đặc trưng cục bộ giàu thông tin



Một đối tượng



Tập hợp các đặc trưng cục bộ  
(bag-of-features)



- Giải quyết tốt vấn đề che khuất
- Bất biến với kích thước
- Bất biến với phép xoay

34

## CalTech6 dataset



class	<b>bag of features</b>	<b>bag of features</b>	<b>Parts-and-shape model</b>
	Zhang et al. (2005)	Willamowski et al. (2004)	Fergus et al. (2003)
airplanes	<b>98.8</b>	97.1	90.2
cars (rear)	98.3	<b>98.6</b>	90.3
cars (side)	<b>95.0</b>	87.3	88.5
faces	<b>100</b>	99.3	96.4
motorbikes	<b>98.5</b>	98.0	92.5
spotted cats	<b>97.0</b>	—	90.0

BoW có độ chính xác khá tốt cho bài toán phân lớp ảnh

Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)



35

## Bag-of-features

Biểu diễn một dữ liệu (văn bản, vân ảnh, hình ảnh) dưới dạng lược đồ (histogram) của các đặc trưng

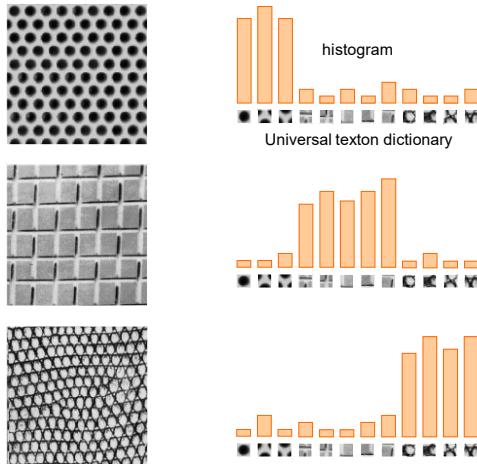
Một ý tưởng cũ

(e.g., nhận dạng vân ảnh và truy vấn thông tin)

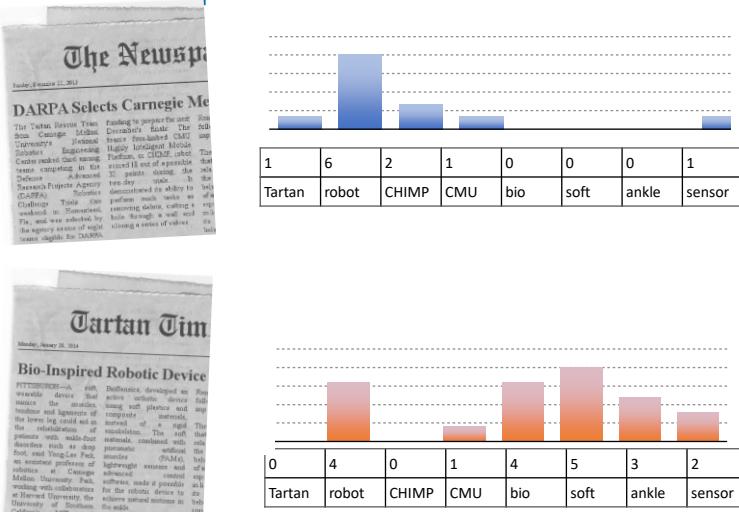


36

## Texture recognition



## Vector Space Model



Một văn bản (datapoint) biểu diễn dưới dạng véc-tơ tần suất xuất hiện các từ (feature)

$$\mathbf{v}_d = [n(w_{1,d}) \ n(w_{2,d}) \ \cdots \ n(w_{T,d})]$$

$n(\cdot)$  Tần suất xuất hiện

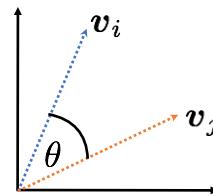
Lược đồ theo các từ

Độ tương đồng giữa hai văn bản?

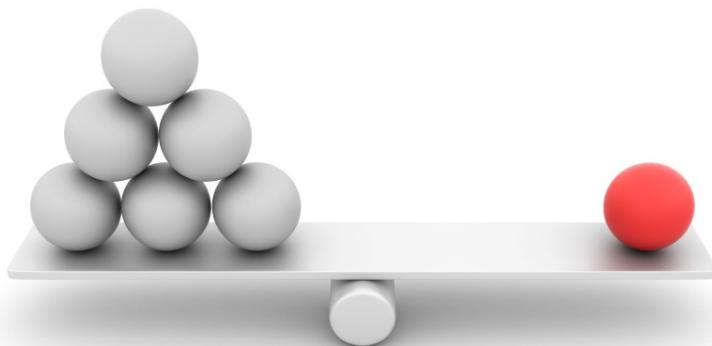
Có thể dùng độ đo cosine

$$d(\mathbf{v}_i, \mathbf{v}_j) = \cos \theta$$

$$= \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|}$$



Nhưng không phải các từ đều quan trọng như nhau



## TF-IDF

Term Frequency Inverse Document Frequency

$$\mathbf{v}_d = [n(w_{1,d}) \ n(w_{2,d}) \ \cdots \ n(w_{T,d})]$$

Đánh trọng số của các từ theo quy tắc sau:

$$\mathbf{v}_d = [n(w_{1,d})\alpha_1 \ n(w_{2,d})\alpha_2 \ \cdots \ n(w_{T,d})\alpha_T]$$

$$n(w_{i,d})\alpha_i = n(w_{i,d}) \log \left\{ \frac{\text{term frequency}}{\text{inverse document frequency}} \right\}$$

$$= n(w_{i,d}) \log \left\{ \frac{D}{\sum_{d'} \mathbf{1}[w_i \in d']} \right\}$$

(down-weights common terms)

## BOW cho bài toán phân lớp ảnh

## Mô hình túi từ

- Đặc trưng cục bộ ~ 1 từ
- 1 ảnh ~ 1 văn bản
- Sử dụng mô hình vector: ảnh = tần suất xuất hiện của các từ trực quan



**Dictionary Learning:**  
Học Visual Words bằng phân cụm

**Encode:**  
Xây dựng Bags-of-Words (BOW) vectors  
cho từng ảnh

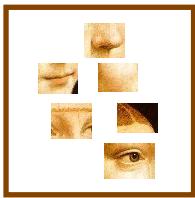
**Classify:**  
Huấn luyện và test dữ liệu mới sử dụng BOWs

---

## Dictionary Learning:

### Học Visual Words bằng phân cụm

1. Trích xuất đặc trưng (e.g., SIFT) từ các ảnh



---

## Dictionary Learning:

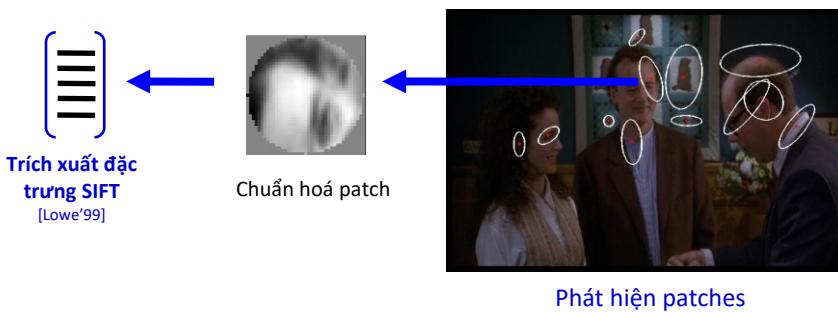
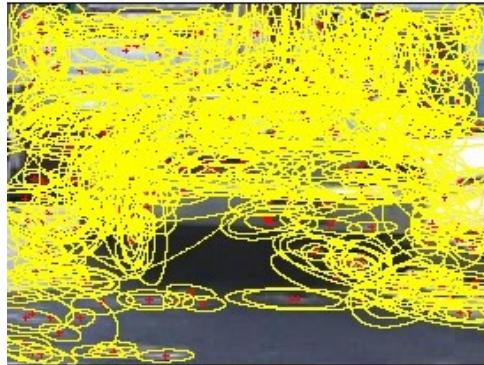
### Học Visual Words bằng phân cụm

2. Học visual dictionary (e.g., phân cụm K-means)

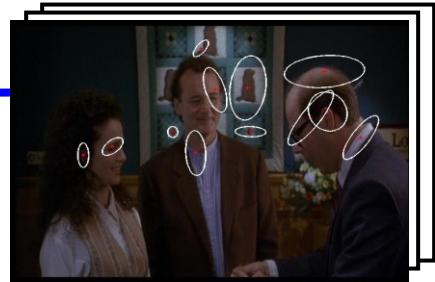


## What kinds of features can we extract?

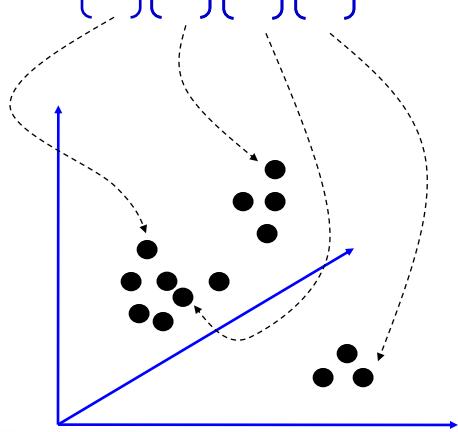
- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005
- Interest point detector
  - Csurka et al. 2004
  - Fei-Fei & Perona, 2005
  - Sivic et al. 2005
- Other methods
  - Random sampling (Vidal-Naquet & Ullman, 2002)
  - Segmentation-based patches (Barnard et al. 2003)

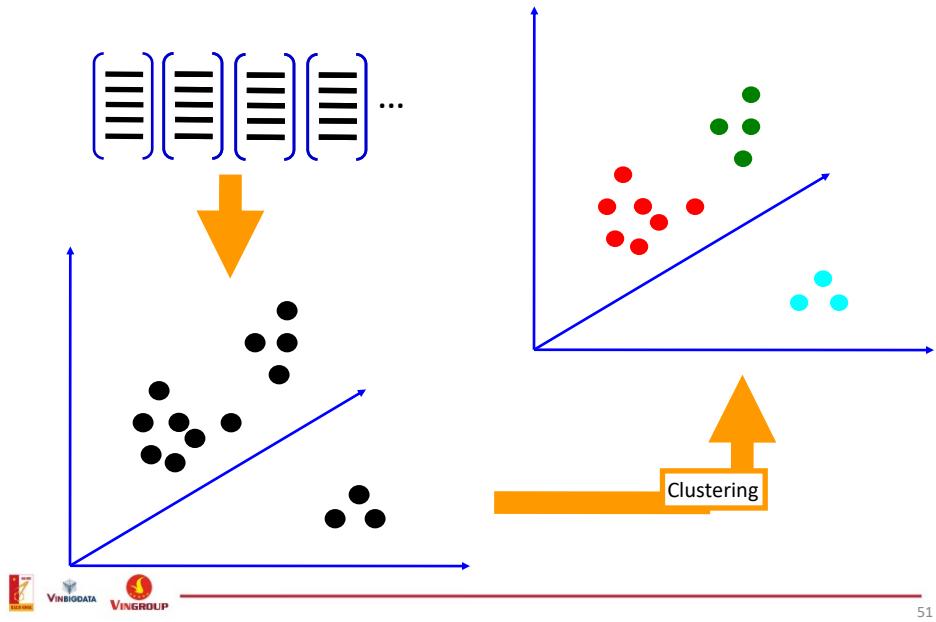


$$\left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \dots$$

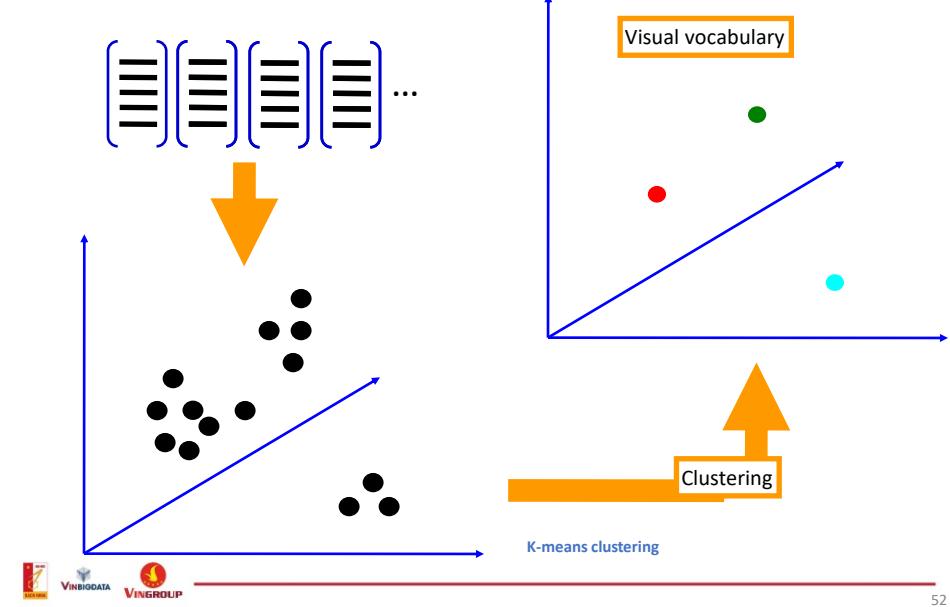


$$\left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \dots$$



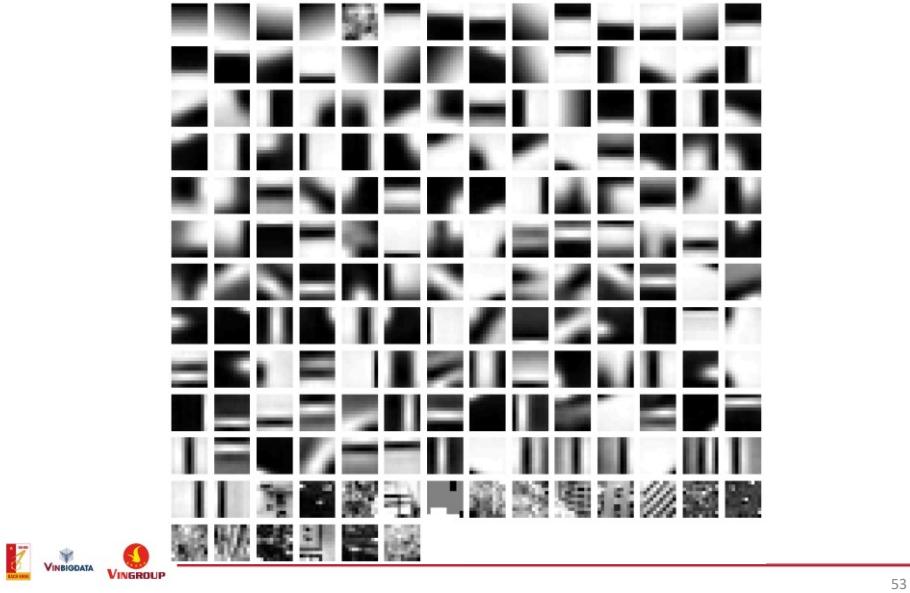


51



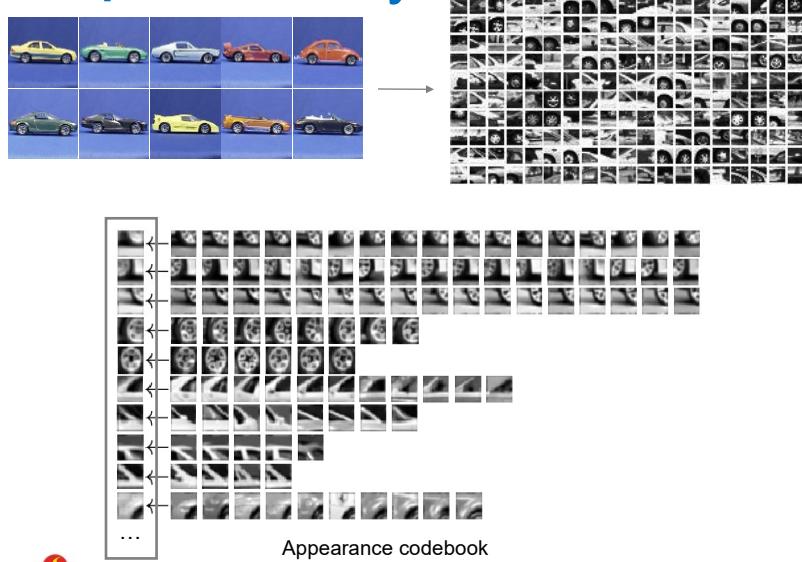
52

## Ví dụ visual dictionary



53

## Ví dụ dictionary



Source: B. Leibe



# Một ví dụ khác dictionary



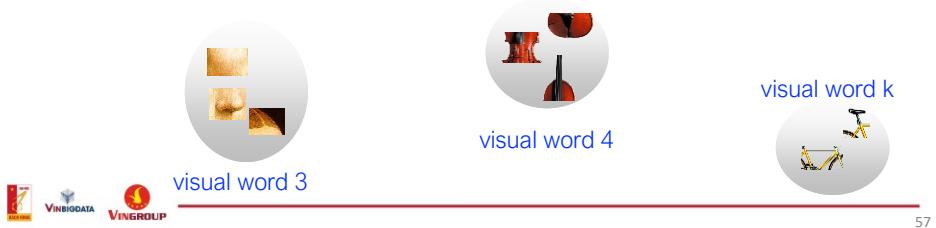
**Dictionary Learning:**  
Học Visual Words bằng phân cụm

**Encode:**  
Xây dựng Bags-of-Words (BOW) vectors  
cho từng ảnh

**Classify:**  
Huấn luyện và test dữ liệu mới sử dụng BOWs

**Encode:**

Xây dựng Bags-of-Words (BOW) vectors cho từng ảnh

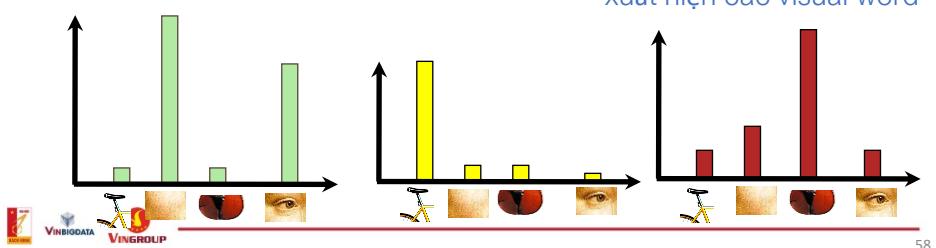


57

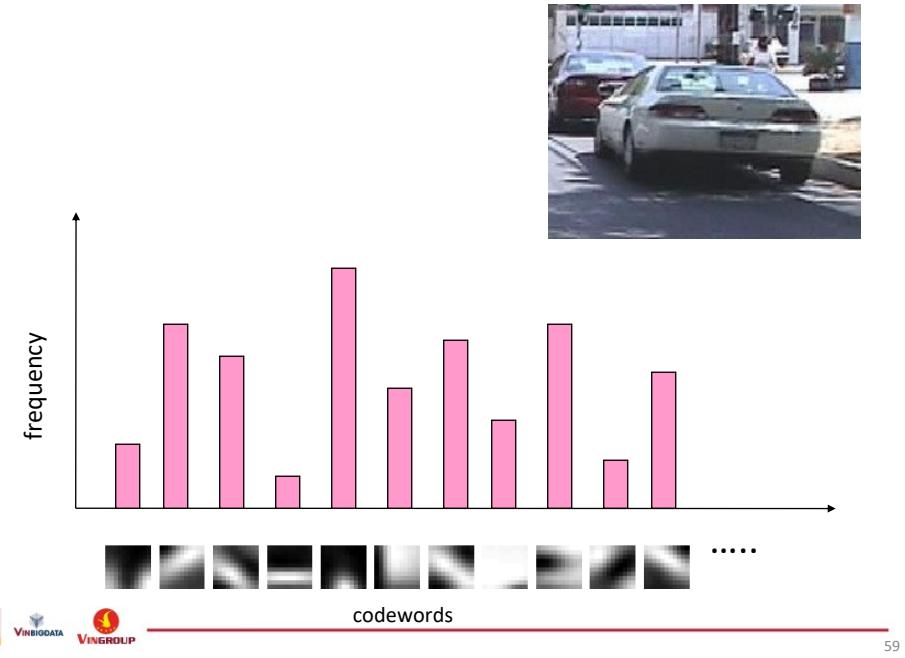
**Encode:**

Xây dựng Bags-of-Words (BOW) vectors cho từng ảnh

2. Lược đồ: đếm tần suất xuất hiện các visual word



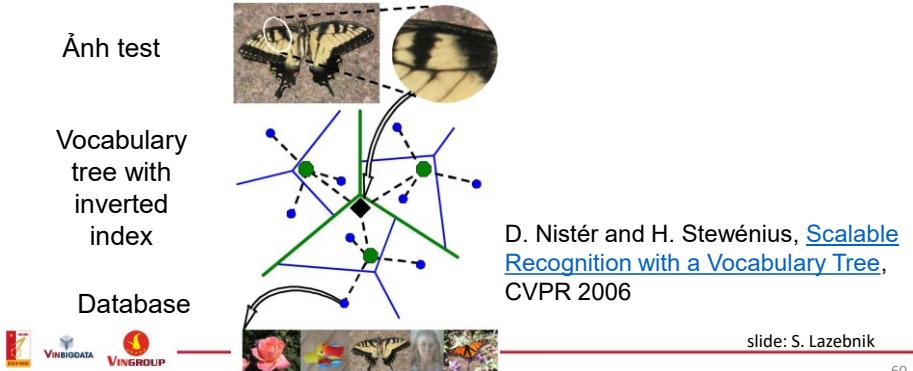
58



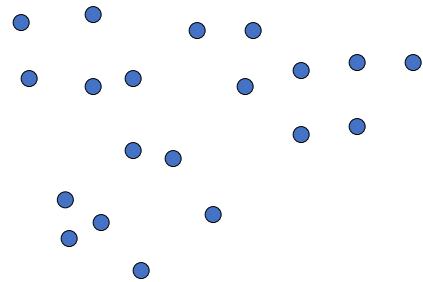
## Khả năng mở rộng cho CSDL lớn

- Làm sao để truy vấn một ảnh trong CSDL hàng triệu ảnh?

- Efficient putative match generation
  - Fast nearest neighbor search, inverted indexes



## Vocabulary Tree là gì?

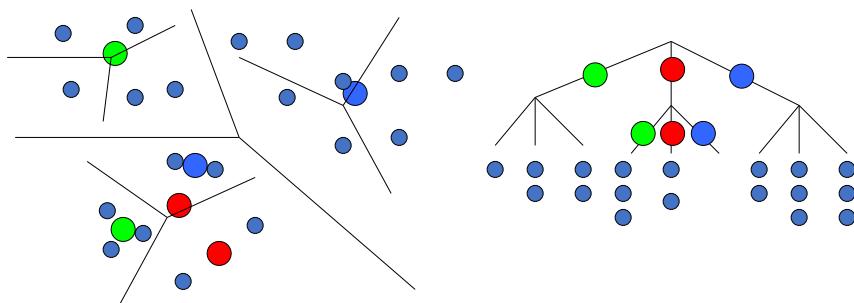


Nister and Stewenius CVPR 2006



61

## Vocabulary Tree là gì?

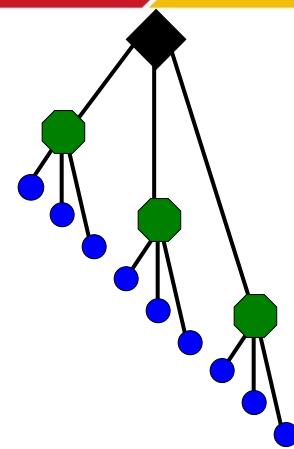


- Nhiều bước K-Means để xây dựng decision tree (offline)
- Truy vấn cây online

Nister and Stewenius CVPR 2006



62

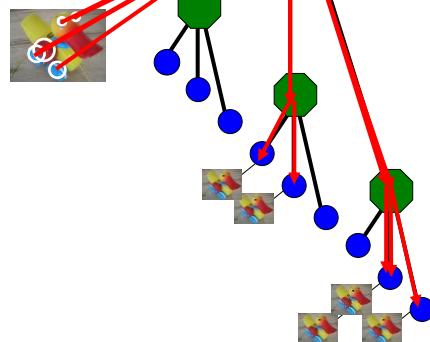


Slide credit: D. Nister

63



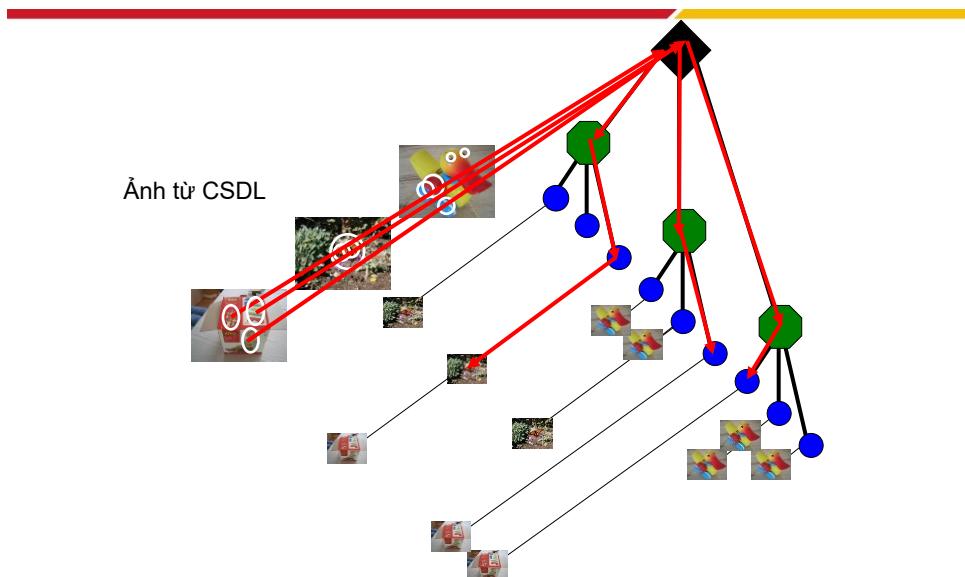
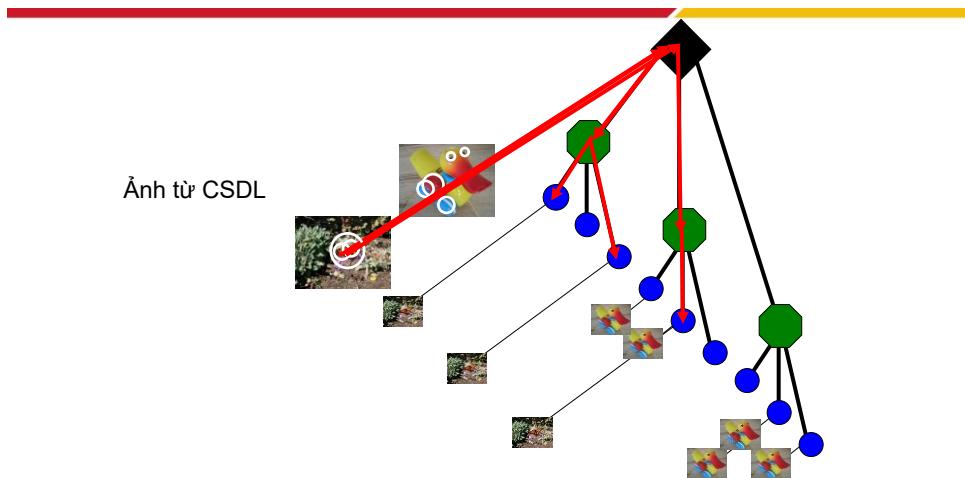
Ảnh từ CSDL

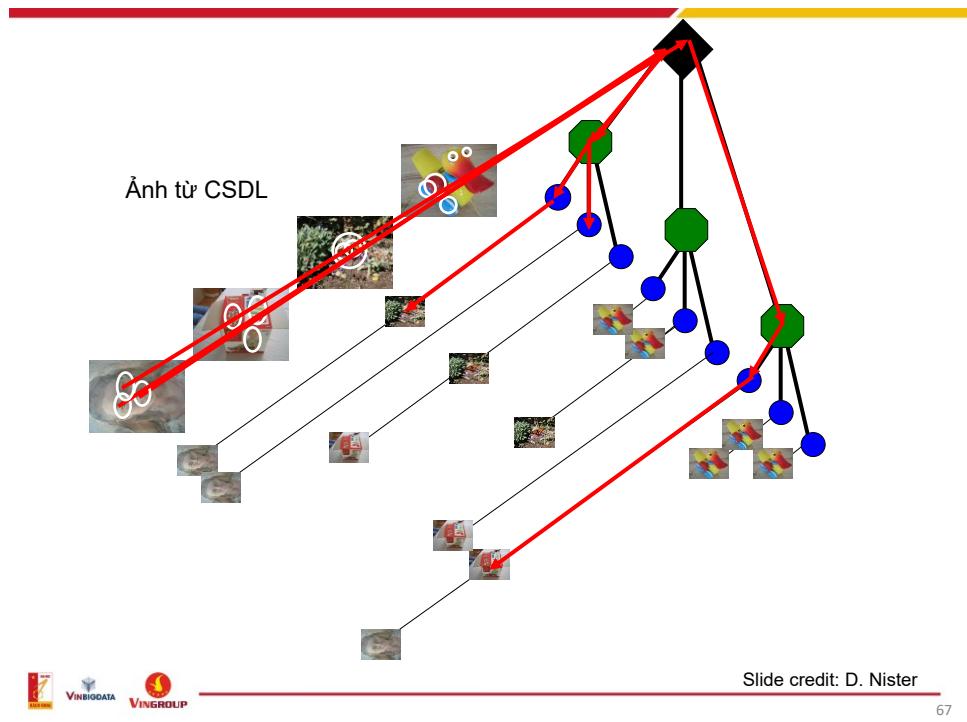


Slide credit: D. Nister

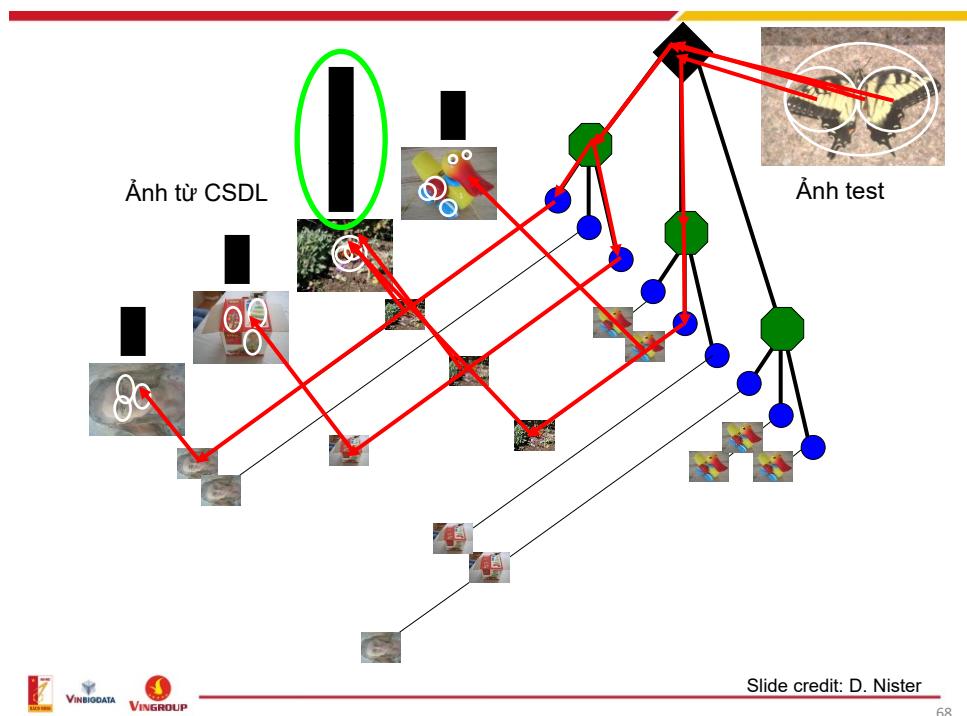
64







67



68

## Dictionary Learning:

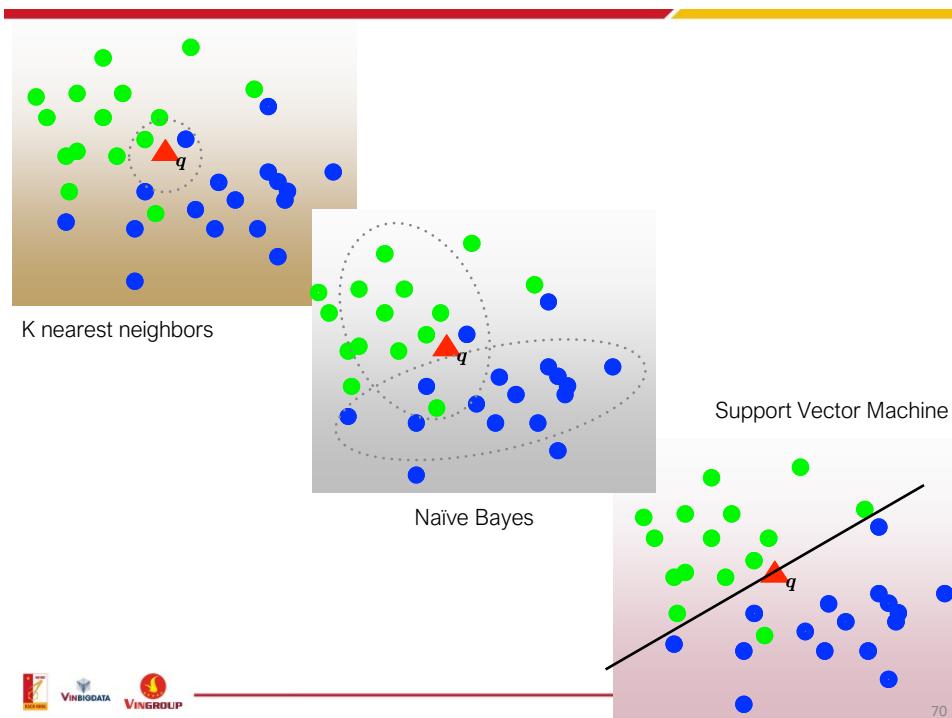
Học Visual Words bằng phân cụm

**Encode:**  
Xây dựng Bags-of-Words (BOW) vectors  
cho từng ảnh

**Classify:**  
Huấn luyện và test dữ liệu mới sử dụng BOWs



69



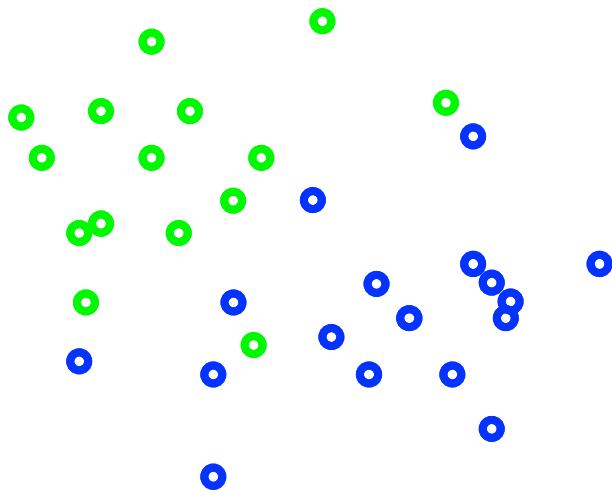
70

## K nearest neighbors



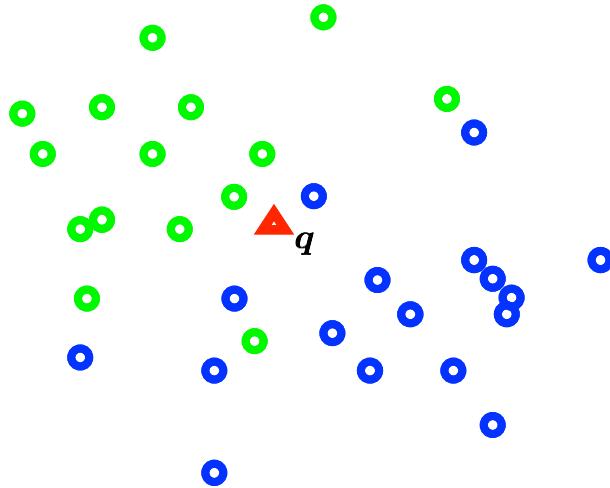
71

Distribution of data from two classes



---

Distribution of data from two classes

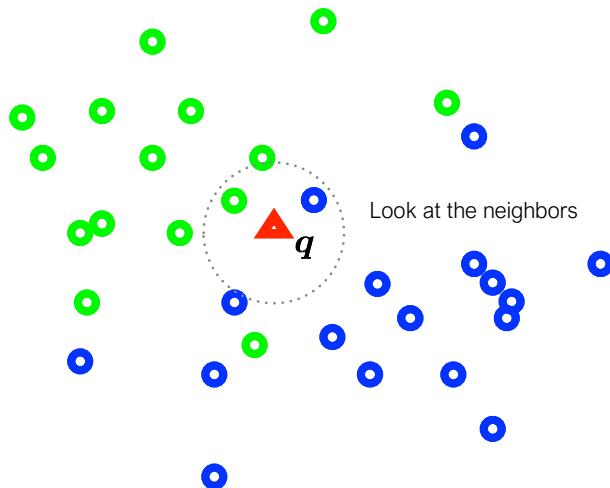


*Which class does  $q$  belong too?*

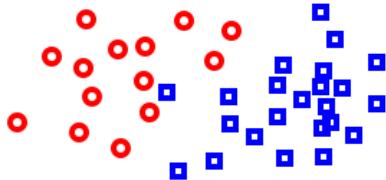
---

---

Distribution of data from two classes



# K-Nearest Neighbor (KNN) Classifier

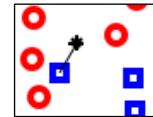


Non-parametric pattern classification approach

Consider a two class problem where each sample consists of two measurements ( $x, y$ ).

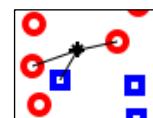
For a given query point  $q$ , assign the class of **the nearest neighbor**

$k = 1$



Compute the  **$k$  nearest neighbors** and assign the class by majority vote.

$k = 3$



## Nearest Neighbor is competitive

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

Test Error Rate (%)	
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
K-NN, Tangent Distance, 16x16	1.1
K-NN, shape context matching	0.67
1000 RBF + linear classifier	3.6
SVM deg 4 polynomial	1.1
2-layer NN, 300 hidden units	4.7
2-layer NN, 300 HU, [deskewing]	1.6
LeNet-5, [distortions]	0.8
Boosted LeNet-4, [distortions]	0.7

### MNIST Digit Recognition

- Handwritten digits
- 28x28 pixel images:  $d = 784$
- 60,000 training samples
- 10,000 test samples

Yann LeCunn

## What is the best distance metric between data points?

- Typically Euclidean distance
- Locality sensitive distance metrics
- Important to normalize.  
Dimensions have different scales

## How many K?

- Typically k=1 is good
- Cross-validation (try different k!)



77

## Distance metrics

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + \cdots + (x_N - y_N)^2} \quad \text{Euclidean}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{x_1 y_1 + \cdots + x_N y_N}{\sqrt{\sum_n x_n^2} \sqrt{\sum_n y_n^2}} \quad \text{Cosine}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_n \frac{(x_n - y_n)^2}{(x_n + y_n)} \quad \text{Chi-squared}$$



78

## Distance metrics

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

- Two most commonly used special cases of p-norm

$$\|x\|_p = \left( |x_1|^p + \cdots + |x_n|^p \right)^{\frac{1}{p}} \quad p \geq 1, x \in \mathbb{R}^n$$



79

## CIFAR-10 and NN results

Example dataset: **CIFAR-10**

**10 labels**

**50,000 training images**

**10,000 test images.**



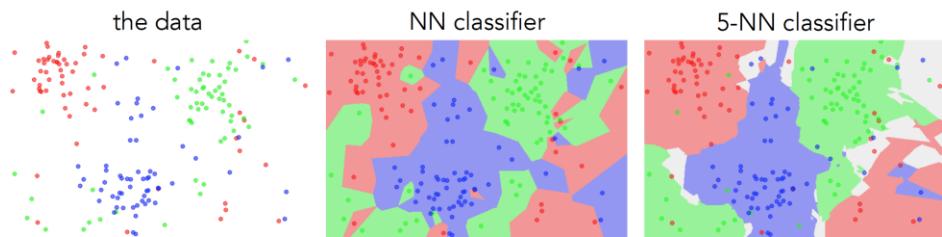
For every test image (first column),  
examples of nearest neighbors in rows



81

## k-nearest neighbor

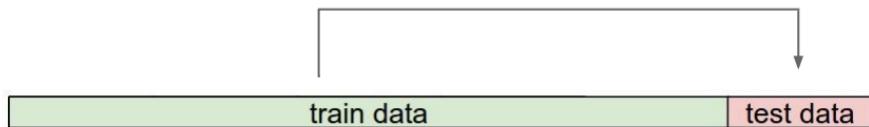
- Find the k closest points from training data
- Labels of the **k points** “vote” to classify



## Hyperparameters

- What is the best distance to use?
- What is the best value of k to use?
- i.e., how do we set the hyperparameters?
- Very problem-dependent
- Must try them all and see what works best

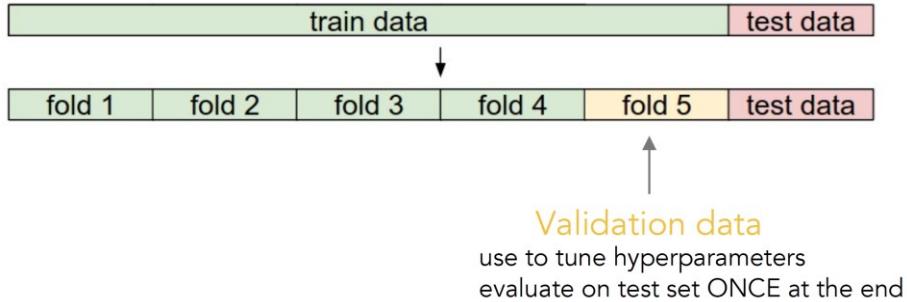
Try out what hyperparameters work best on test set.



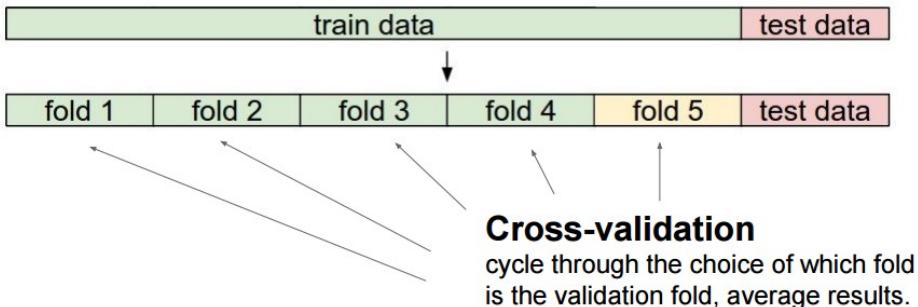
Trying out what hyperparameters work best on test set:  
Very bad idea. The test set is a proxy for the generalization performance!  
Use only **VERY SPARINGLY**, at the end.

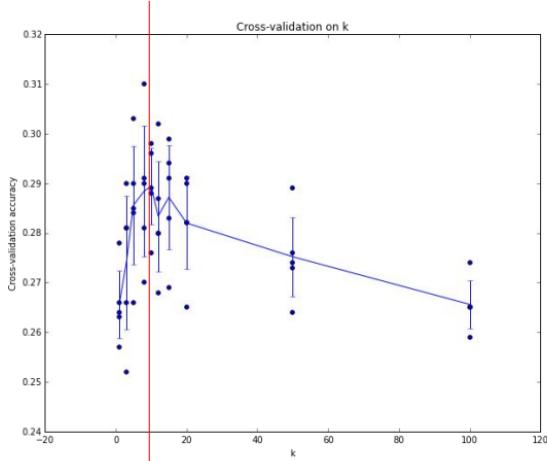


# Validation



# Cross-validation





Example of  
5-fold cross-validation  
for the value of  $k$ .

Each point: single  
outcome.

The line goes  
through the mean, bars  
indicated standard  
deviation

(Seems that  $k \approx 7$  works best  
for this data)



## How to pick hyperparameters?

- Methodology
  - Train and test
  - Train, validate, test
- Train for original model
- Validate to find hyperparameters
- Test to understand generalizability



# kNN

## Pros

- simple yet effective

## Cons

- search is expensive (can be sped-up)
- storage requirements
- difficulties with high-dimensional data



90

# kNN -- Complexity and Storage

- N training images, M test images
- Training:  $O(1)$
- Testing:  $O(MN)$
- Hmm...
  - Normally need the opposite
  - Slow training (ok), fast testing (necessary)



91

# Độ đo đánh giá kết quả

## Confusion matrix

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	TP	FN
<u>True 0</u>	FP	TN

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	hits	misses
<u>True 0</u>	false alarms	correct rejections

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

## Precision versus recall

- Accuracy
  - correct predictions / total predictions
- Precision:
  - how many of the object detections are correct?
- Recall:
  - how many of the ground truth objects can the model detect?
  - True Positive Rate (TPR)
- F1 score:
 
$$2 * precision * recall / (precision + recall)$$

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

## Tài liệu tham khảo

- Kristen Grauman (CS 376: Computer Vision, Spring 2018, The University of Texas at Austin)
- Ioannis Yannis, Gkioulekas (16-385 Computer Vision, Spring 2020, CMU)

