

# Machine Perception

## Lecture 6: Localization and SLAM (Part II)

Piotr Skrzypczyński

Institute of Robotics and Machine Intelligence  
Poznań University of Technology



---

POZNAN UNIVERSITY OF TECHNOLOGY

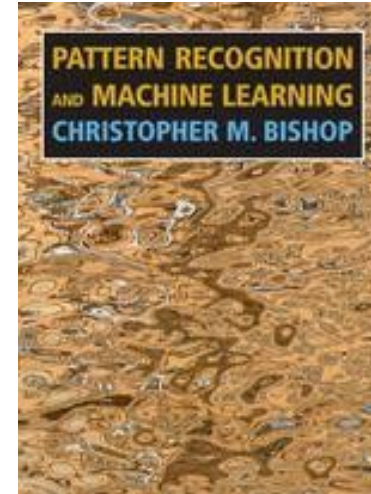
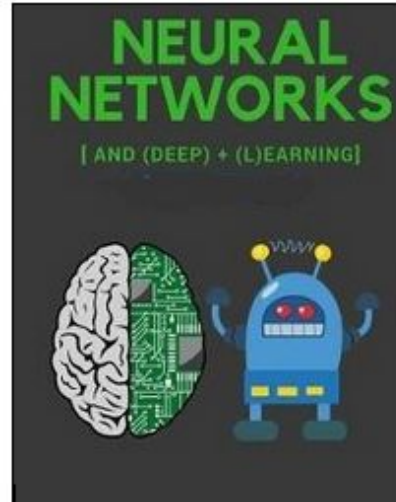
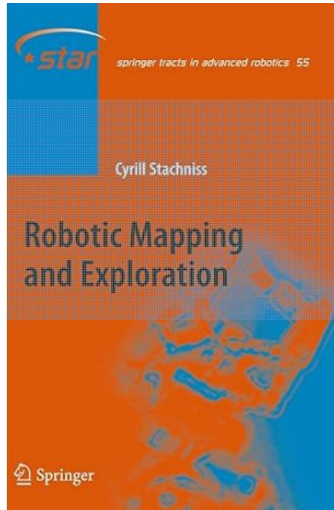
---

# Lecture outline

- Visual Place Recognition - definition, variants, role in SLAM.
- Approaches to Visual Place Recognition .
- Bag of Visual Words.
- Deep learning based approach - NetVLAD.
- Learning based place recognition in robotics.

# Literature

1. C. Stachniss, *Robotic Mapping and Exploration*, Springer Verlag, 2009.
2. M. Nielsen, *Neural Networks and Deep Learning* <http://neuralnetworksanddeeplearning.com/>
3. C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer-Verlag Berlin, Heidelberg, 2006.



# Place recognition

- Robotics : has the robot been to this place before ? Which images were taken around the same location ?
- Image retrieval : have I seen this image before? Which images in my database look similar to it ?

## Visual Place Recognition: A Survey

Stephanie Lowry, Niko Sünderhauf, Paul Newman, *Fellow, IEEE*, John J. Leonard, *Fellow, IEEE*, David Cox, Peter Corke, *Fellow, IEEE*, and Michael J. Milford, *Member, IEEE*

**Abstract**—Visual place recognition is a challenging problem due to the vast range of ways in which the appearance of real-world places can vary. In recent years, improvements in visual sensing capabilities, an ever-increasing focus on long-term mobile robot autonomy, and the ability to draw on state-of-the-art research in other disciplines—particularly recognition in computer vision and animal navigation in neuroscience—have all contributed to significant advances in visual place recognition systems. This paper presents a survey of the visual place recognition research landscape. We start by introducing the concepts behind place recognition—the role of place recognition in the animal kingdom, how a “place” is defined in a robotics context, and the major components of a place recognition system. Long-term robot operations have revealed that changing appearance can be a significant factor in visual place recognition failure; therefore, we discuss how place recognition solutions can implicitly or explicitly account for appearance change within the environment. Finally, we close with a discussion on the future of

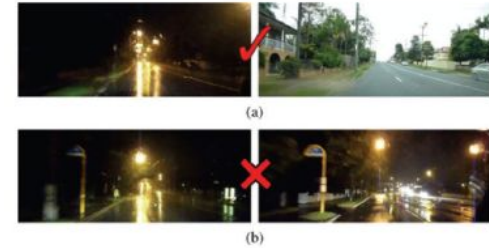


Fig. 1. Visual place recognition systems must be able to (a) successfully match very perceptually different images while (b) also rejecting incorrect matches between aliased image pairs of different places.

G.  
Ro  
IET  
lio

1. S. Lowry et al., "Visual Place Recognition: A Survey", IEEE Transactions on Robotics, vol. 32, no. 1, pp. 1-19, 2016

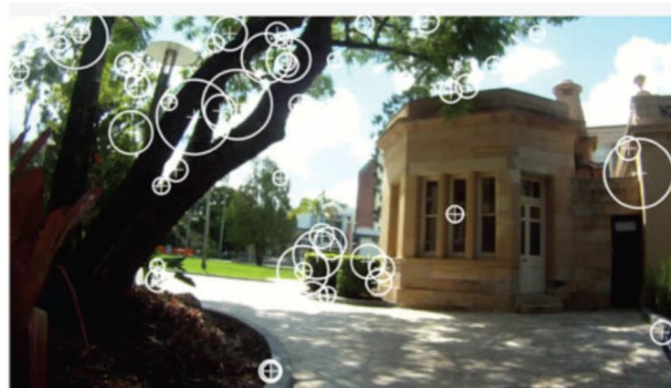
# Place recognition

- **Visual Place Recognition**
  - goal : query an image in a database of  $N$  images
  - complexity:  $NM^2$  feature comparisons (assumes each image has  $M$  features)
- **Appearance changes:** illumination, weather conditions, dynamic objects (people, cars,...), Viewpoint changes
- **Perceptual aliasing:** two different places may look similar (building, roads, ...)



# Place recognition

- **Approaches**
  - Local descriptors
  - Global descriptors
  - Learning-based methods
- Use analogies from text retrieval:
- Visual Words
- Vocabulary of Visual Words
- “Bag of Words” (BoW) approach



(a)

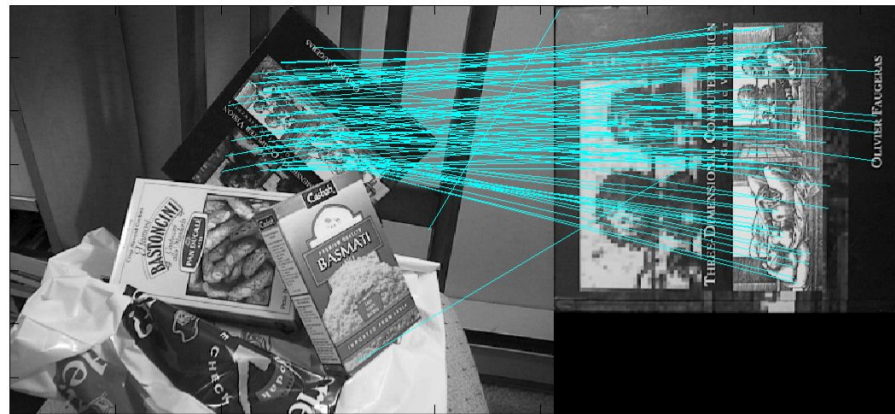


(b)



# Place recognition

- **Approaches**
  - Local descriptors
  - Global descriptors
  - Learning-based methods
- Use analogies from text retrieval:
- Visual Words
- Vocabulary of Visual Words
- “Bag of Words” (BoW) approach



# Global descriptors

## Early approaches:

- color histograms
- principal component analysis
- other statistics on edges, corners, and color patches

## GIST descriptor:

- image is filtered at different orientations and different frequencies to extract information from the image
- results are averaged to generate a compact vector that represents the “gist” of a scene



▶ Man-made open environment.



▶ Man-made closed urban environment.



▶ Perspective view of a man-made closed urban environment.  
Large space with small elements.



▶ Flat view of a man-made urban environment, vertically structured.



# Global vs. local descriptors

## Global descriptors:

- better at handling lighting conditions and seasonal variation
- more sensitive to viewpoint changes

## Local descriptors:

- allow estimating feature (and camera) geometry
- sensitive to lighting conditions and seasonal variations



► Man-made open environment.



► Man-made closed urban environment.



► Perspective view of a man-made closed urban environment.  
Large space with small elements.



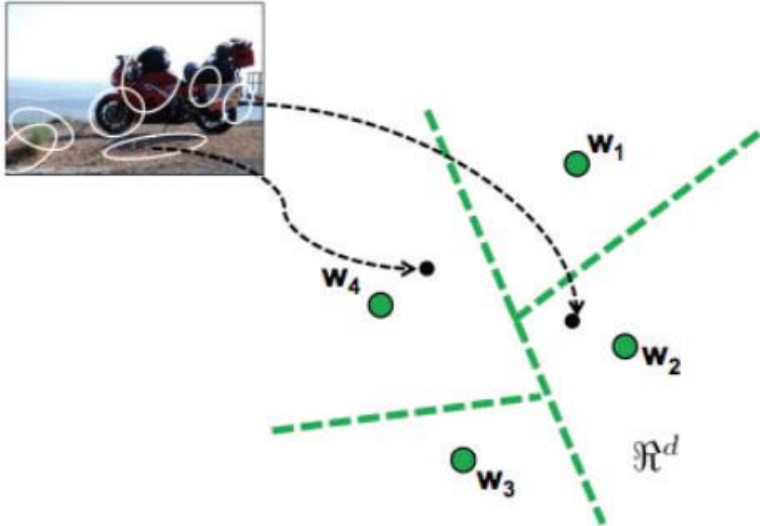
► Flat view of a man-made urban environment, vertically structured.

# Local descriptors: Bag of Visual Words

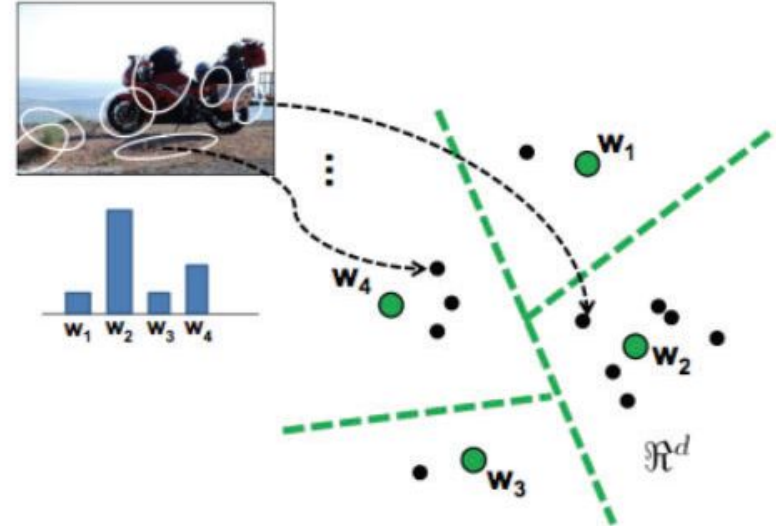
- Based on text retrieval and summarization methods

- 1) Extract features and descriptors in image
- 2) Discretize feature space (clustering)
- 3) Store the frequency of the features for each image

**Each cluster is a “visual word”**

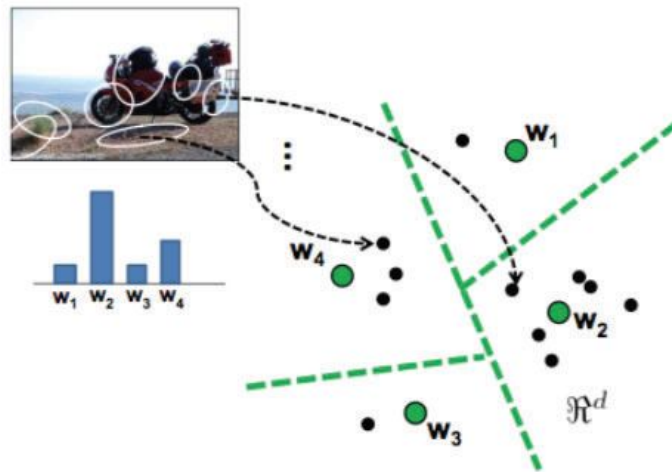
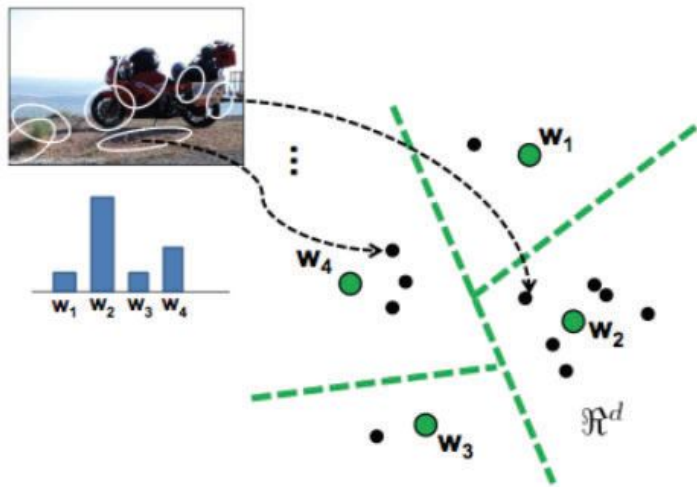


J. Sivic, A. Zisserman. Video Google: A text retrieval approach to object matching in videos. ICCV, 2003



# Local descriptors: Bag of Visual Words

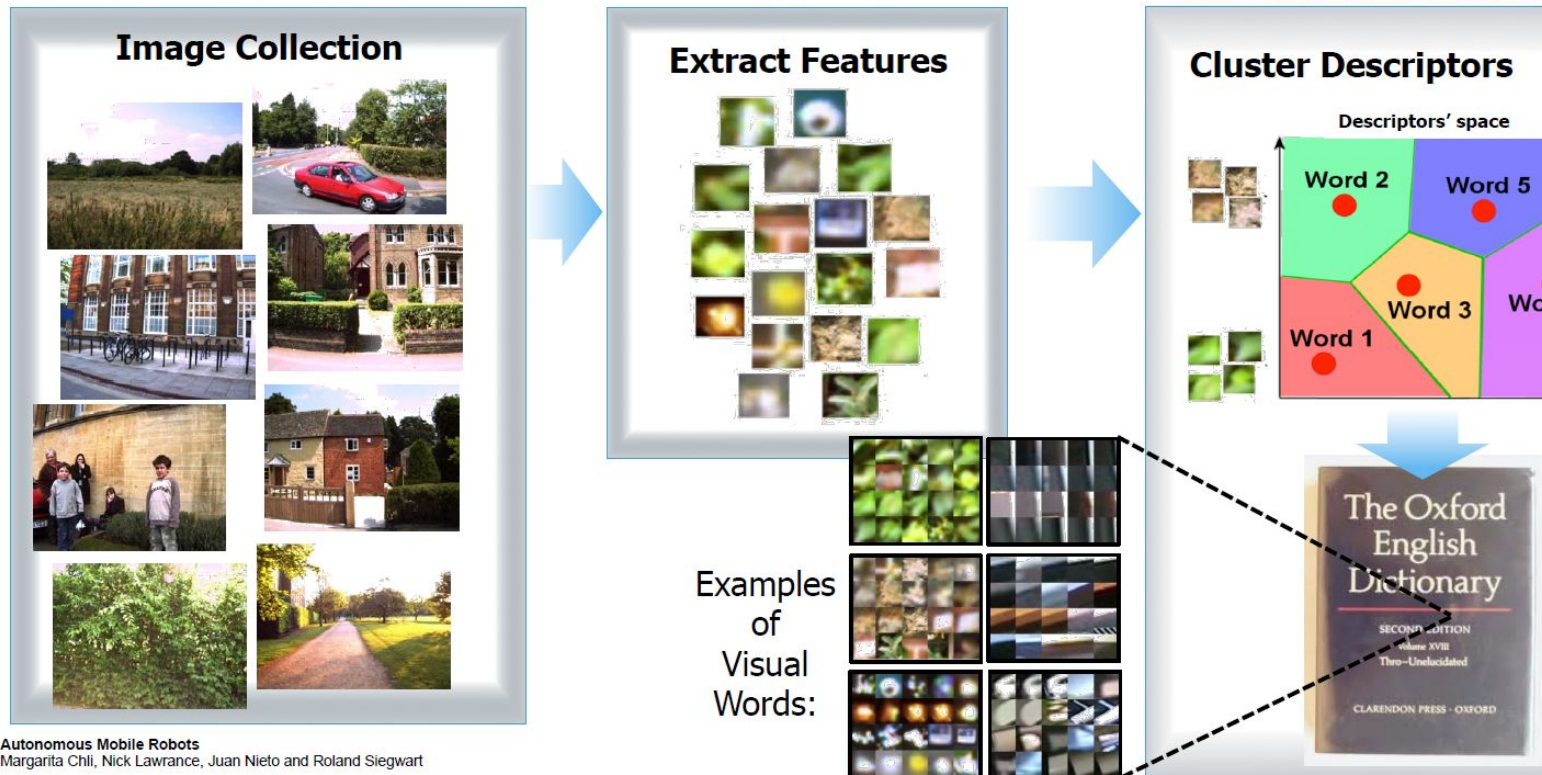
- Two images are compared based on the corresponding histogram (Hamming distances, other metrics, ...)
- Faster version: vocabulary tree
- Alternative: VLAD (Vector of Locally Aggregated Descriptors)



# Place recognition (Bag of Visual Words)

- Collect a large enough dataset that is representative of all possible images that are relevant to your application (e.g., for automotive place recognition, you may want to collect million of street images sampled around the world)
- Extract features and descriptors from each image and map them into the descriptor space (e.g., for SIFT, 128 dimensional descriptor space)
- Cluster the descriptor space into  $K$  clusters
- The centroid of each cluster is a visual word.
- This is computed by taking the arithmetic average of all the descriptors within the same cluster, e.g., for SIFT, each cluster contains SIFT features that are very similar to each other;
- The visual word then is the average all the SIFT descriptors in that cluster

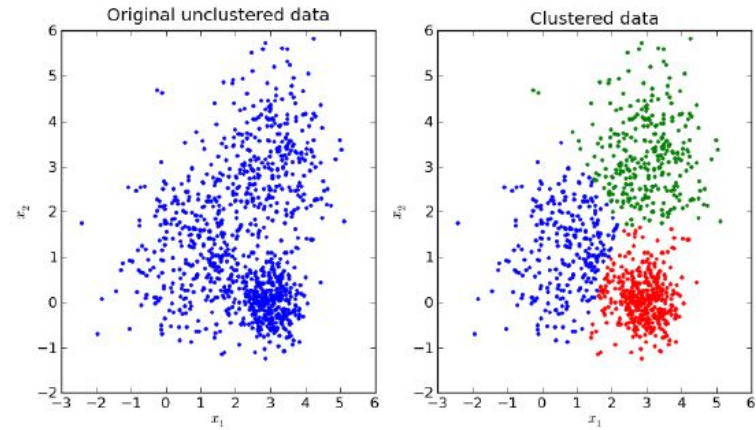
# Place recognition (Bag of Visual Words)



# Place recognition (Bag of Visual Words)

- The **k-means** clustering algorithm partitions  $n$  data point into  $k$  clusters in which each data point  $x$  belongs to the cluster  $S_i$  with center  $m_i$
- It minimizes the sum of squared Euclidean distances between points  $x$  and their nearest cluster centers  $m_i$
- **Algorithm:**
  1. Randomly initialize  $k$  cluster centers
  2. Iterate until convergence:
  3. Assign each data point  $x_j$  to the nearest center  $m_i$
  4. Recompute each cluster center as the mean of all points assigned to it

$$D(X, M) = \sum_{i=1}^k \sum_{x \in S_i} (x - m_i)^2$$





# Place recognition (Bag of Visual Words)

- The **Image Vocabulary** is a data structure that lists all extracted visual words
- Each visual word is assigned a unique identifier (an integer number)
- Each word in the image vocabulary points to a list of images (from the entire image database) in which that word appears
- If the database grows, the vocabulary is updated accordingly

Alphabetical Index	
<b>A</b>	Aperture, 103
<b>B</b>	Band, Lower, 31
<b>C</b>	Cap Screw, Hexagon Socket Head, 125
<b>D</b>	Die, Butt Plate, 33
<b>E</b>	Ejector (Cartridge), 88
<b>F</b>	Follower, 74
<b>G</b>	Gas Cylinder, 123
<b>H</b>	Hammer, 14-15
<b>I</b>	Impeller, 103
<b>J</b>	Joint, 103
<b>K</b>	Key, 103
<b>L</b>	Latch, 97
<b>M</b>	Master Pin, 98
<b>N</b>	Nut, 103
<b>O</b>	Oil, 103
<b>P</b>	Pin, 103
<b>Q</b>	Quench, 103
<b>R</b>	Rod, 103
<b>S</b>	Screw, 103
<b>T</b>	Trigger, 103
<b>U</b>	Unit, 103
<b>V</b>	Valve, 103
<b>W</b>	Washer, 103
<b>X</b>	X-ray, 103
<b>Y</b>	Yoke, 103
<b>Z</b>	Zinc, 103

## Image vocabulary

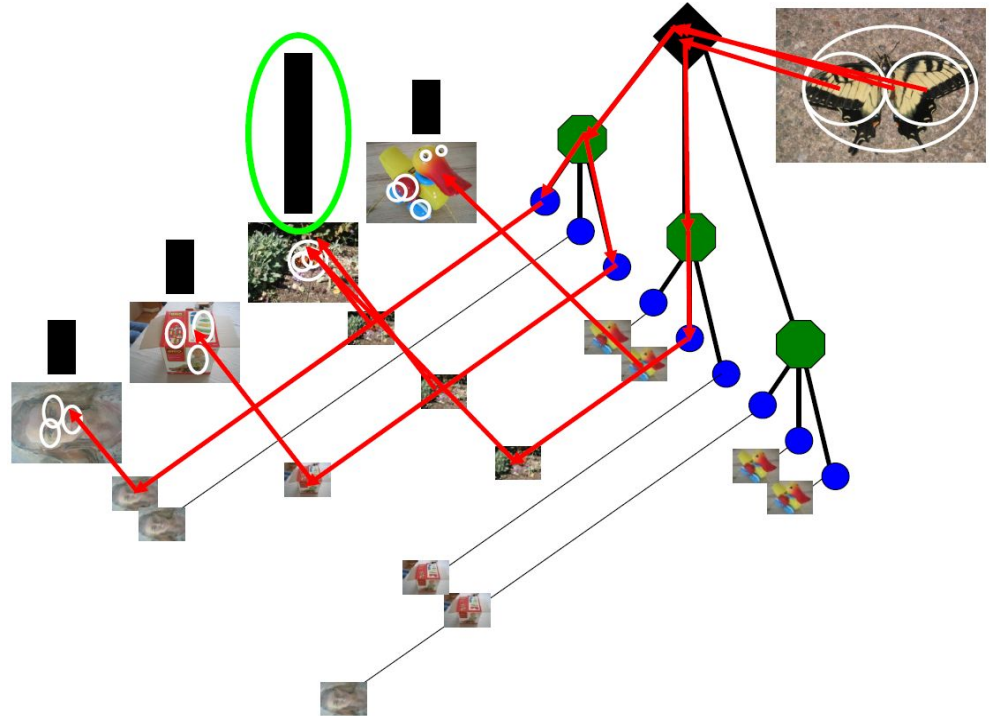
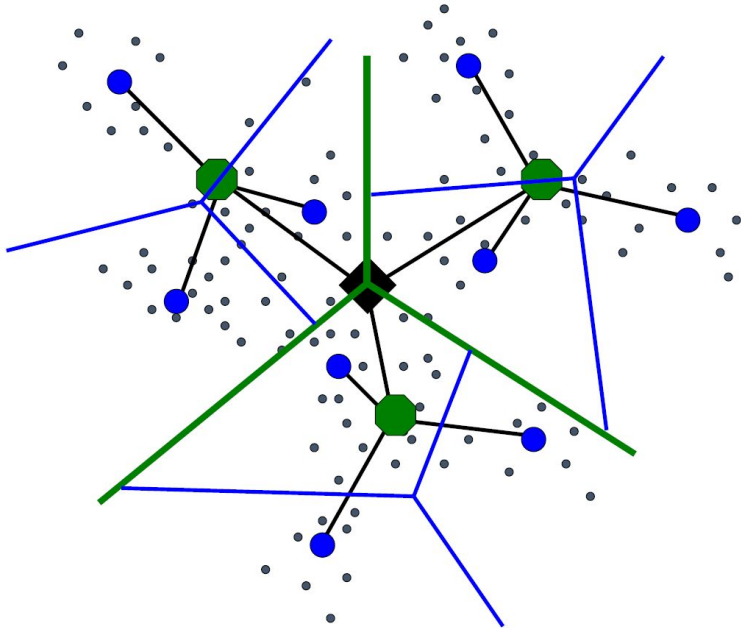
Visual words

List of images in which this word appears

0	
...	
101	
102	
103	
104	
105	
...	

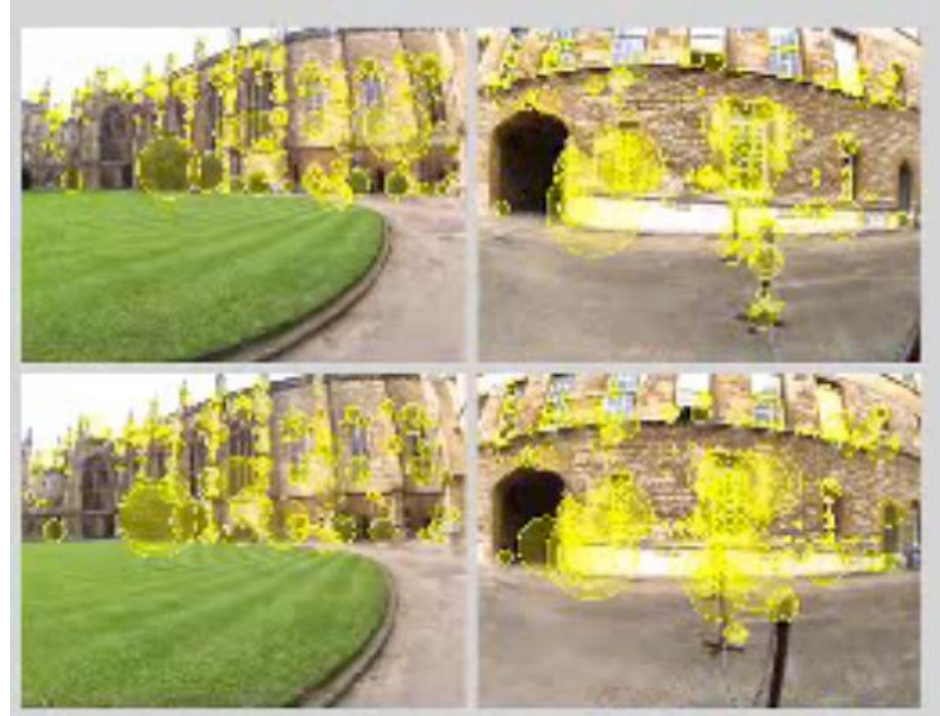
# Place recognition (Bag of Visual Words)

- Within the vocabulary, each visual word points to a list of images where that word occurs.
- During retrieval, each feature contributes to update the voting array.
- The image with most votes is returned.



# Place recognition (FAB-MAP)

- Place recognition for robot localization using stereo images
- Build the visual vocabulary using SURF features
- Probabilistic model of the world:
- World = a set of discrete places
- Place = a set of consecutive images
- At a new frame, compute:  $P(\text{being at a known place})$ ,  $P(\text{being at a new place})$
- Captures the dependencies of visual words to distinguish the most characteristic structure of each scene (using the Chow-Liu tree)



# Place recognition (FAB-MAP)

[robots.ox.ac.uk/~mjc/appearance\\_based\\_results.htm](http://robots.ox.ac.uk/~mjc/appearance_based_results.htm)



(a)  $p=0.9989$



(b)  $p=0.999$



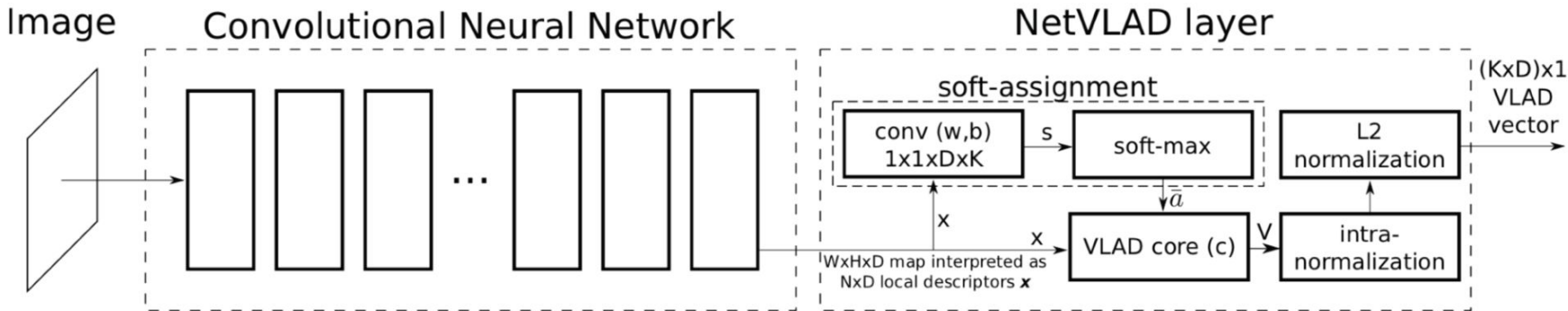
(c)  $p=0.999$

# Place recognition (NetVLAD)

- Earlier approaches using AlexNet or similar and use layers activations as descriptors
- **NetVLAD:**
  - CNN-based approach
  - Trained on the task of place recognition

## How to get labeled data ?

- a large dataset of panoramic images from the Google StreetView
- positions based on their (noisy) GPS
- seasonal variations
- illumination changes





# Place recognition (NetVLAD)

- Mimic the classical pipeline with deep learning

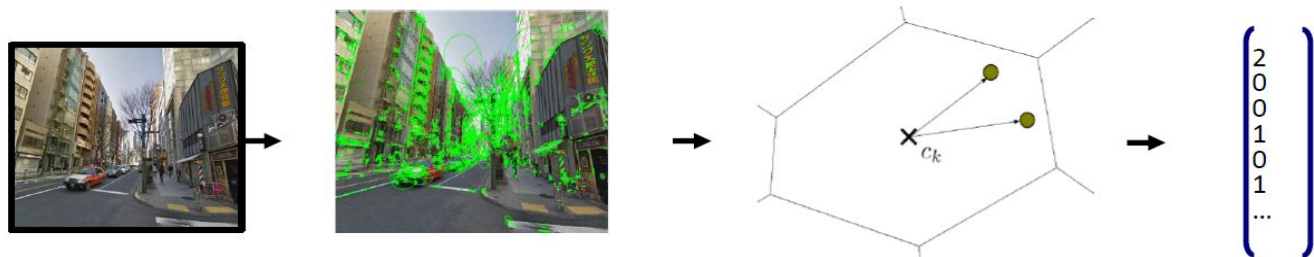
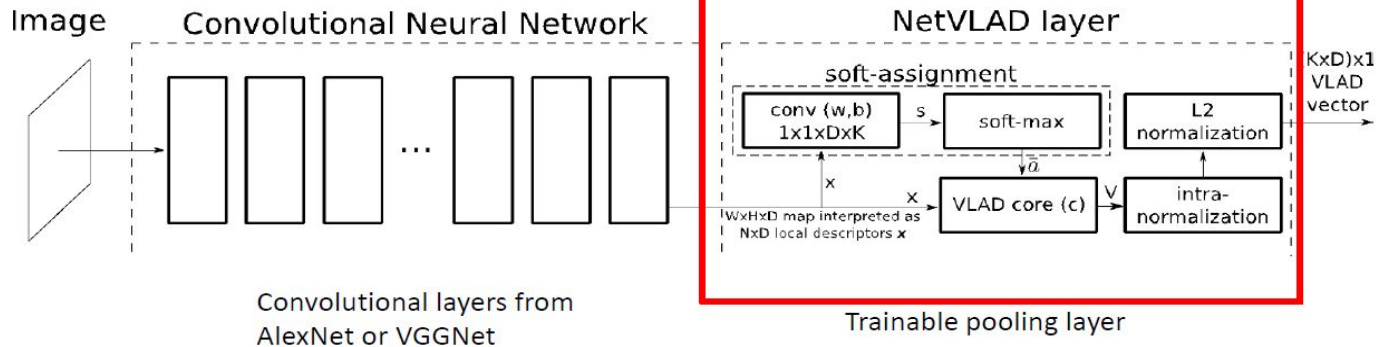


Image I

Extract local  
features (SIFT)

Aggregate  
(BoW, VLAD, FV)

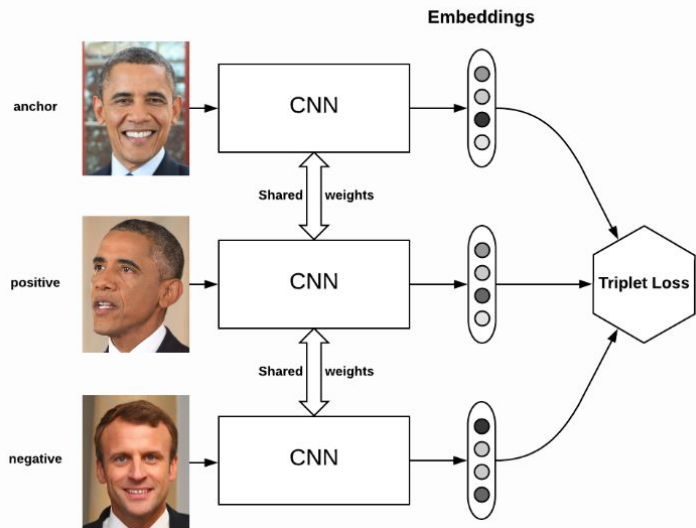
$F(I)$





# Place recognition (NetVLAD)

- NetVlad loss - triplet loss formulation



$$D_p = ||F_{\theta}(\text{city street}) - F_{\theta}(\text{city street})||^2$$

Matching samples

$$D_n = ||F_{\theta}(\text{city street}) - F_{\theta}(\text{forest})||^2$$

Non matching samples

$$L_{\theta} = \sum_{\text{samples}} \max(D_{p(\theta)} + \underbrace{m}_{\text{margin}} - D_{n(\theta)}, 0)$$

# Place recognition (NetVLAD)

- Code, dataset and trained network available online:  
<http://www.di.ens.fr/willow/research/netvlad/>

Query

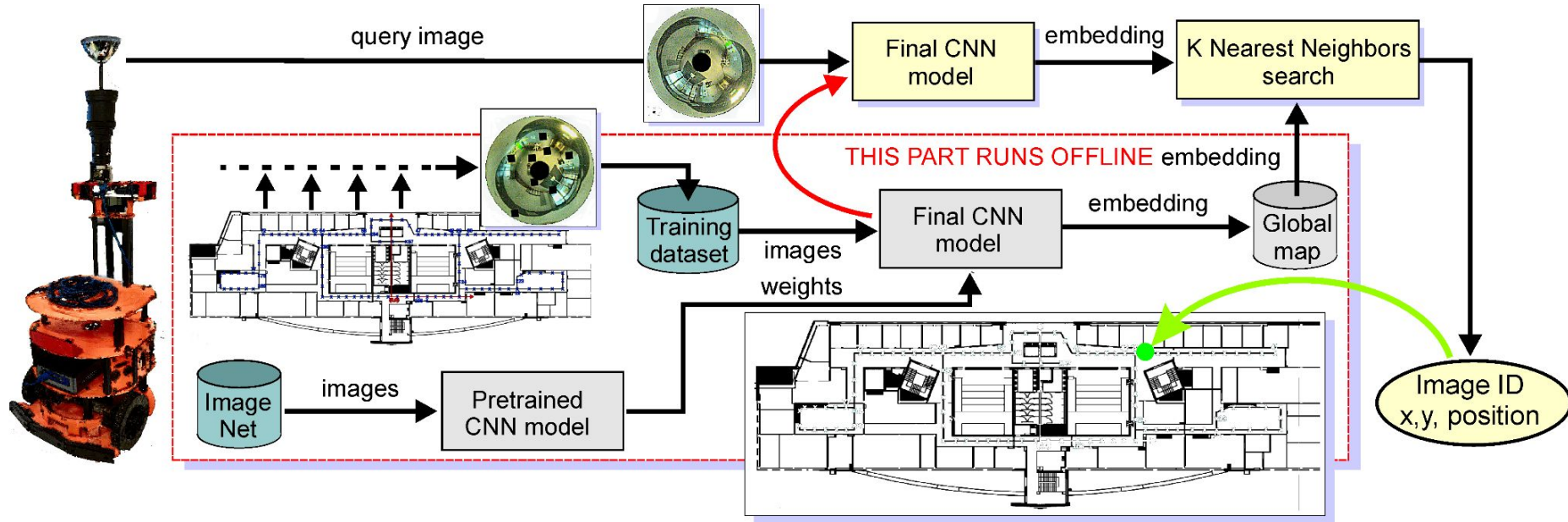


Top result



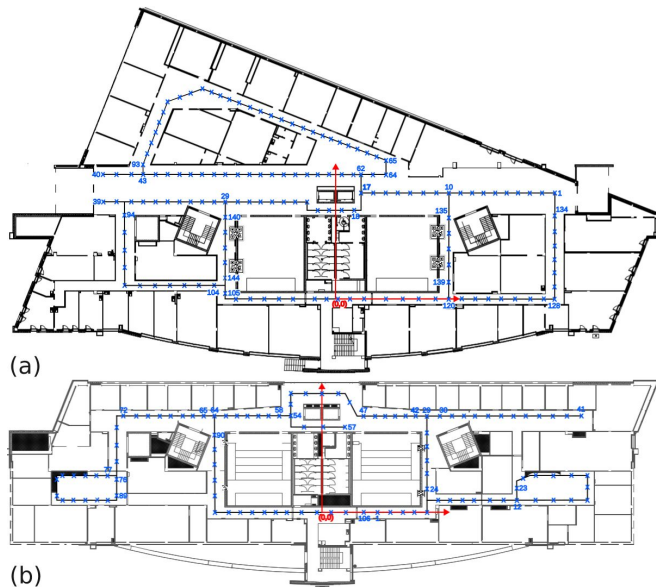
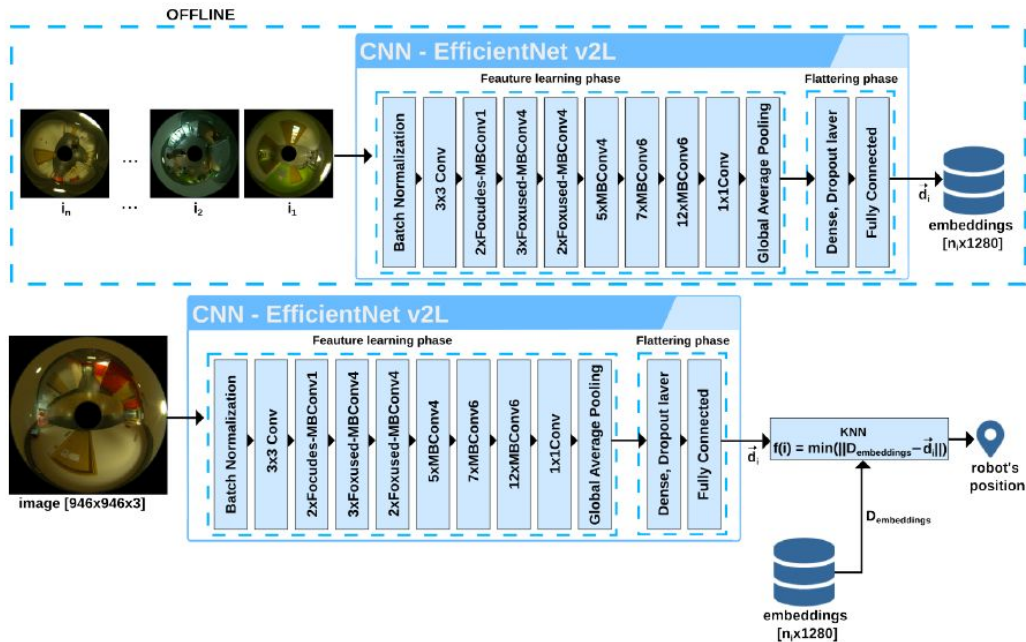
# Place recognition with omnidirectional images

- Appearance-based localization: visual place recognition with omnidirectional images obtained from a catadioptric camera



# Place recognition with omnidirectional images

- Diagram of the CNN-based image description blocks that produce embeddings used as global descriptors in the localization system. The global map is built from  $n_i$  images converted to embedding vectors  $d_i$  that are stored in the map (global descriptors).



# Outcome of the lecture

- A brief review of the Visual Place recognition approaches.
- More detailed presentation of the Bag of Visual Words idea and FAB-MAP as a localisation system that uses the BoW concept.
- NetVLAD as an example of end-to-end trainable place recognition system..
- Example of a simple place recognition system that uses embeddings. .



---

**POZNAN UNIVERSITY OF TECHNOLOGY**

---