

Robot Vision

Visual Place Recognition

Dr. Chen Feng

cfeng@nyu.edu

ROB-UY 3203, Spring 2024



Overview

- Overview of Visual Place Recognition
- Bag of Visual Words (BoVW)
- VLAD

References

- Szeliski 2011
 - Section 14.4
- (BoVW)Jupyter notebook by Olga Vysotska:
https://github.com/ovysotska/in_simple_english/blob/master/bag_of_visual_words.ipynb
- (VLAD) Jégou, H., Douze, M., Schmid, C. and Pérez, P., 2010, June. Aggregating local descriptors into a compact image representation. In *2010 IEEE computer society conference on computer vision and pattern recognition* (pp. 3304-3311). IEEE.
- (NetVLAD) Arandjelovic, R., Gronat, P., Torii, A., Pajdla, T. and Sivic, J., 2016. NetVLAD: CNN architecture for weakly supervised place recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5297-5307).
- Kendall, A., Grimes, M. and Cipolla, R., 2015. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision* (pp. 2938-2946).

Overview of Visual Place Recognition

- We predict where the images are captured using database of

- Geo-tagged images
 - Visual place recognition



- 3D points and descriptors
 - Image based localization



Keywords Related to Visual Place Recognition

Structure from Motion

Visual SLAM

**Image-based
localization**

**Visual place
recognition**

Image retrieval

Image classification

Keywords Related to Visual Place Recognition

Structure from Motion

Bundle adjustment

Geometric
verification
(RANSAC)

ANN

Feature detection
& description

Image retrieval

Visual SLAM

Tracking

Loop detection
& closure

**Image-based
localization**

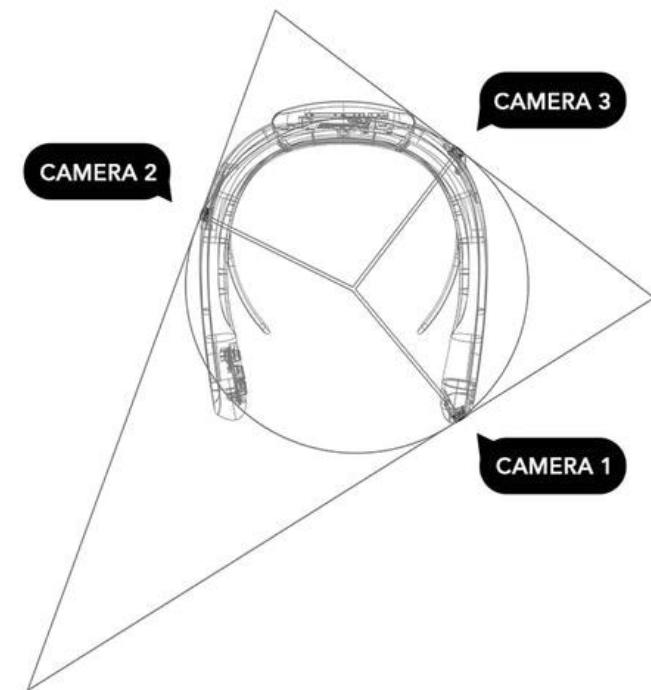
**Visual place
recognition**

Image description
& indexing

Image classification

Potential Application of Visual Place Recognition

- Accelerate incremental SfM and SLAM
 - Localization and navigation of robots and cars
- Integrate images to the real-world coordinates
- Acquiring images is becoming easier
 - e.g. Narrative Clip, FITT360



Street-level Visual Place Recognition

- Given a query image of a particular street or a building, we seek to find one, or at most a few, images in the geotagged database depicting the same place!

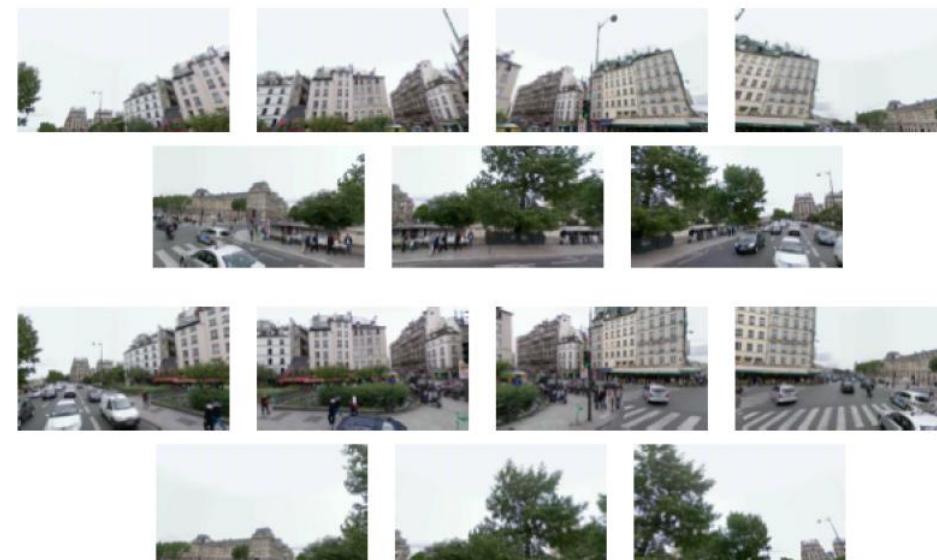


Database of Geo-Tagged Images

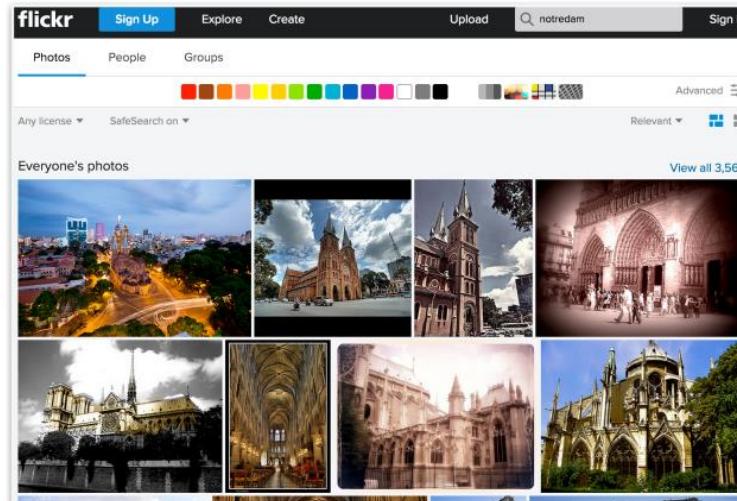
- Typical source of geo-tagged images:
 - Google Streetview (panorama)
 - Driving/walking with a ladybug camera



- We can use original panoramas or generate perspective cutouts



Database of Geo-Tagged Images



Landmarks

Internet Community Photos,
e.g. Flickr, Panoramio

- + Densely captured in different times!
(200M images on Flickr (2012))
- Not really depicting the place :)

See Part 2 for image based localization.

See also [Hays-CVPR'08, Hays-ICCV'09, Quack-CVIR'08].



Street-level

StreetView (panoramic) images

- + Almost all the streets on cities !
- A few images at the same locations
(different times available now!)



Visual place recognition by image retrieval

We are interested in “**a single query testing image**”
and “**a large image database**”.

Query image



Database of geo-tagged images



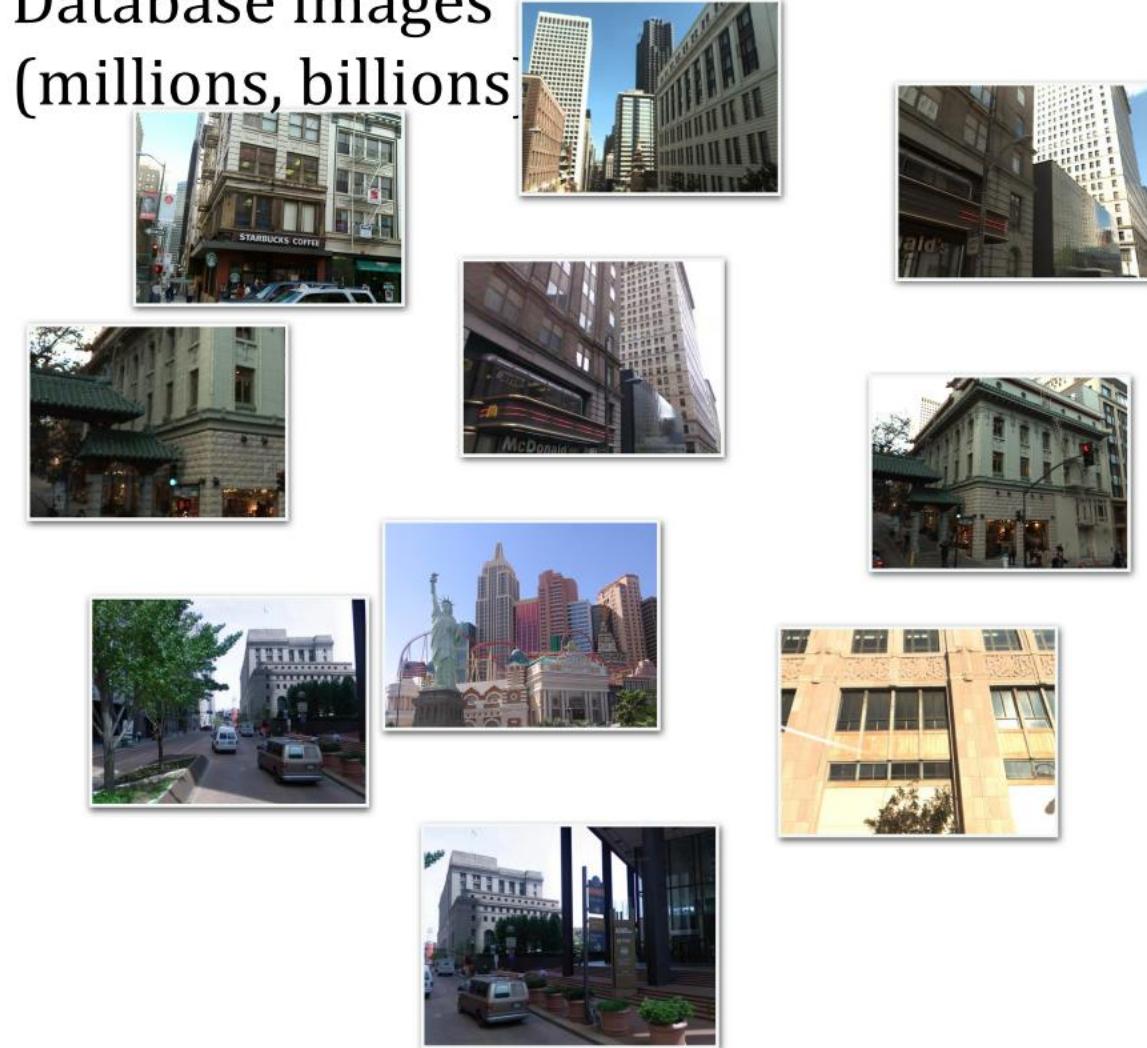


A standard image retrieval [Philbin-CVPR07]

Query image



Database images
(millions, billions)



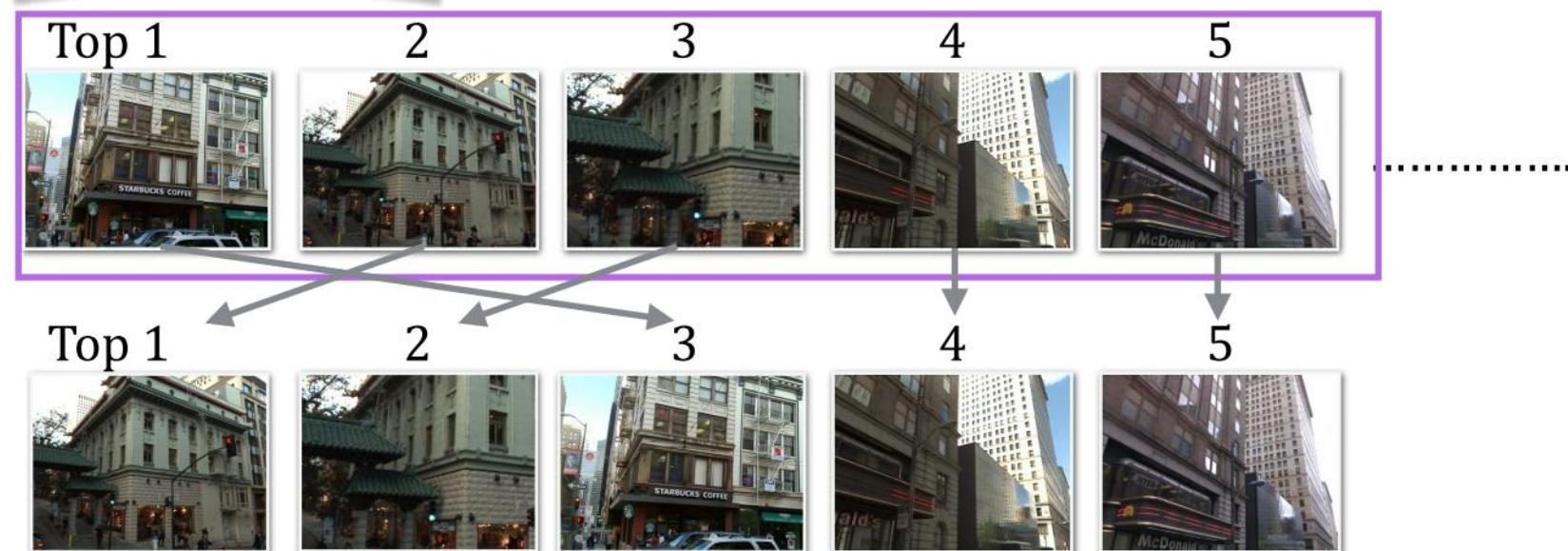


A standard image retrieval [Philbin-CVPR07]

Query image



Step1: Initial ranking/shortlisting
Step2: Re-ranking to improve the list



Step1 should be fast -> BoVW, VLAD, Fisher vectors
Step2 can be a bit more costly -> geometric verification

Visual Place Recognition Pipeline



Database images



Offline

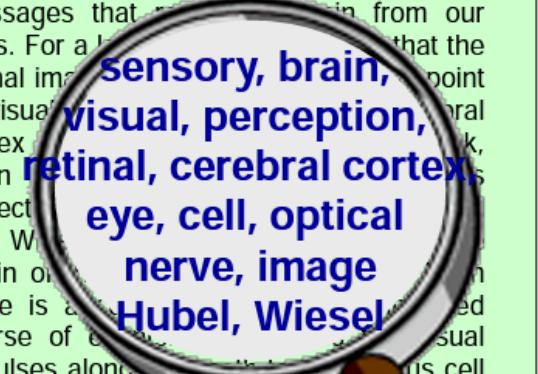
1. Feature detection & description
2. Training visual vocabulary
3. Image description

4. Feature detection & description
5. Image description
6. Initial ranking
7. Re-ranking with geometric verification



Bag-of-Visual-Words(BoVW)

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach the brain from our eyes. For a long time it was believed that the retinal image was processed directly in the visual cortex. However, in 1960, Hubel and Wiesel demonstrated that the visual system is much more complex than previously thought. They found that the visual system consists of several layers of nerve cells. In the retina, there is a layer of ganglion cells which project their axons to the optic nerve. The optic nerve then enters the brain and terminates in the lateral geniculate nucleus. From here, the optic tract projects to the optic radiation, which then connects to the optic cortex in the occipital lobe. Hubel and Wiesel also discovered that the visual system is organized into columns of cells, each column being responsible for a specific type of visual information. For example, one column might be responsible for detecting vertical edges, while another might be responsible for detecting horizontal edges. This organization allows the visual system to process visual information in parallel, allowing for rapid and efficient processing.



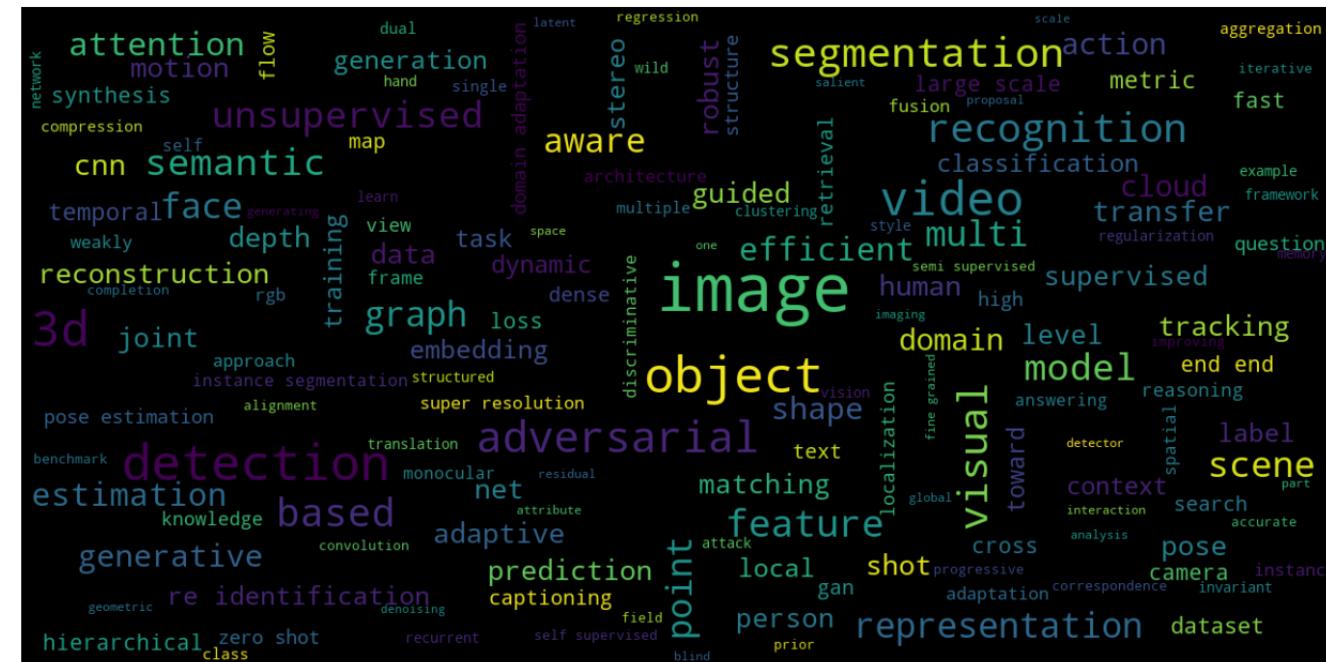
China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a 15% jump in exports. China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value



Analogy to Text Documents

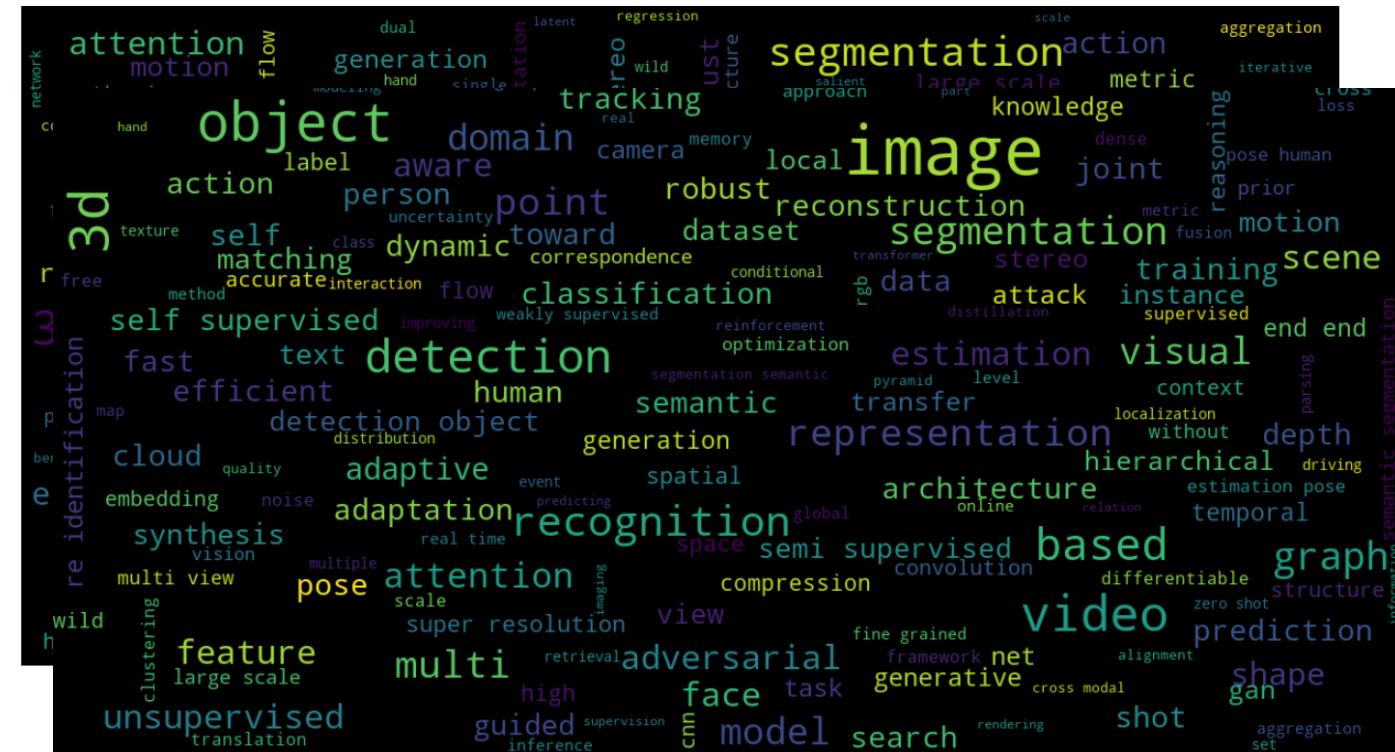
BoVW Origin: Bag-of-Words Model

- Orderless document representation: frequencies of words from a dictionary (Salton&McGill(1983))



BoVW Origin: Bag-of-Words Model

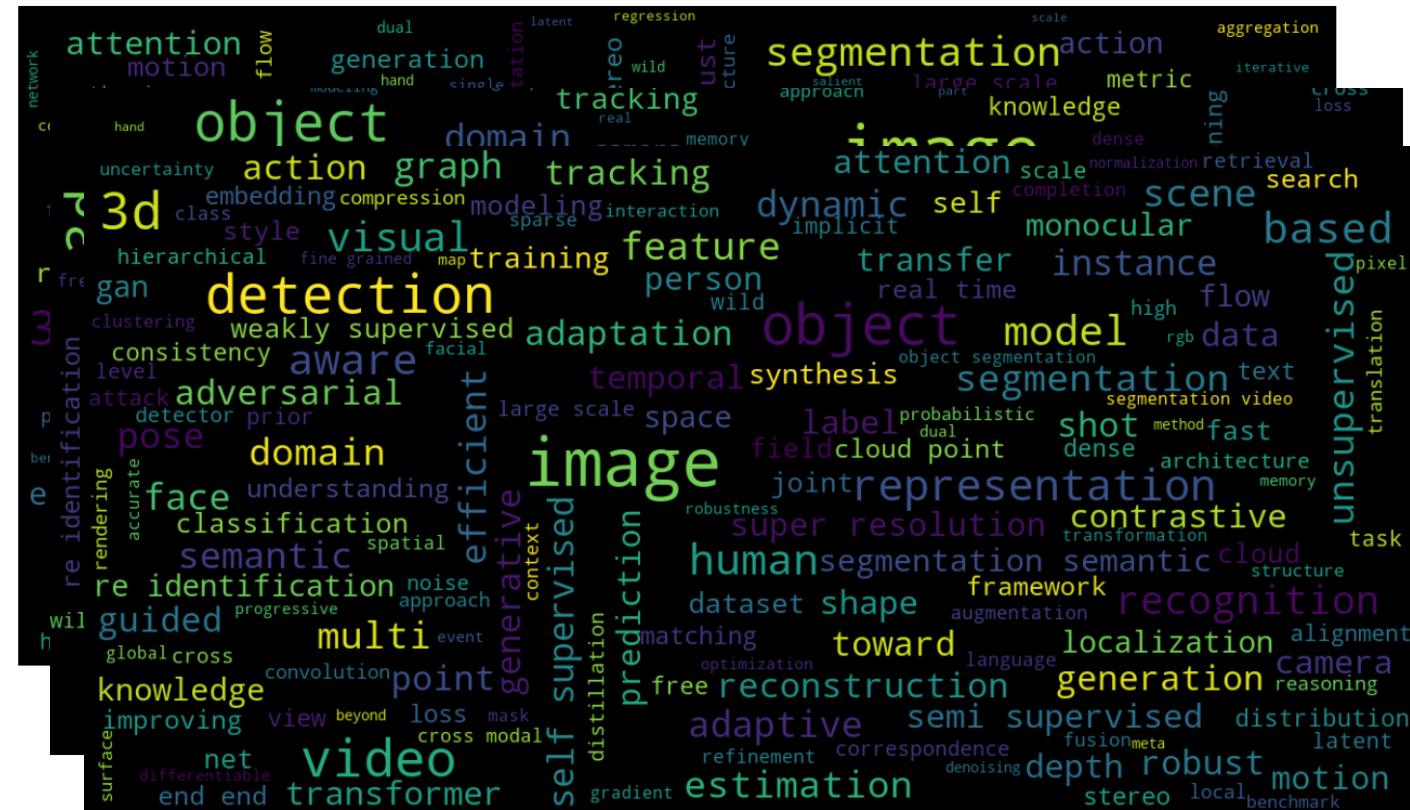
- Orderless document representation: frequencies of words from a dictionary (Salton&McGill(1983))



CVPR - 2020

BoVW Origin: Bag-of-Words Model

- Orderless document representation: frequencies of words from a dictionary
(Salton&McGill(1983))



CVPR - 2021



BoVW Pipeline

1. Extract features (e.g., SIFT)
2. Learn “visual dictionary” (e.g., K-Means Clustering)
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”



Image



Bag of "Visual Words"



BoVW Pipeline: Visual Overview

1. Extract features
2. Learn “visual dictionary” (Offline)
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”



BoVW Pipeline: Visual Overview

1. Extract features
2. Learn “visual dictionary” (Offline)
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”





BoVW Pipeline: Visual Overview

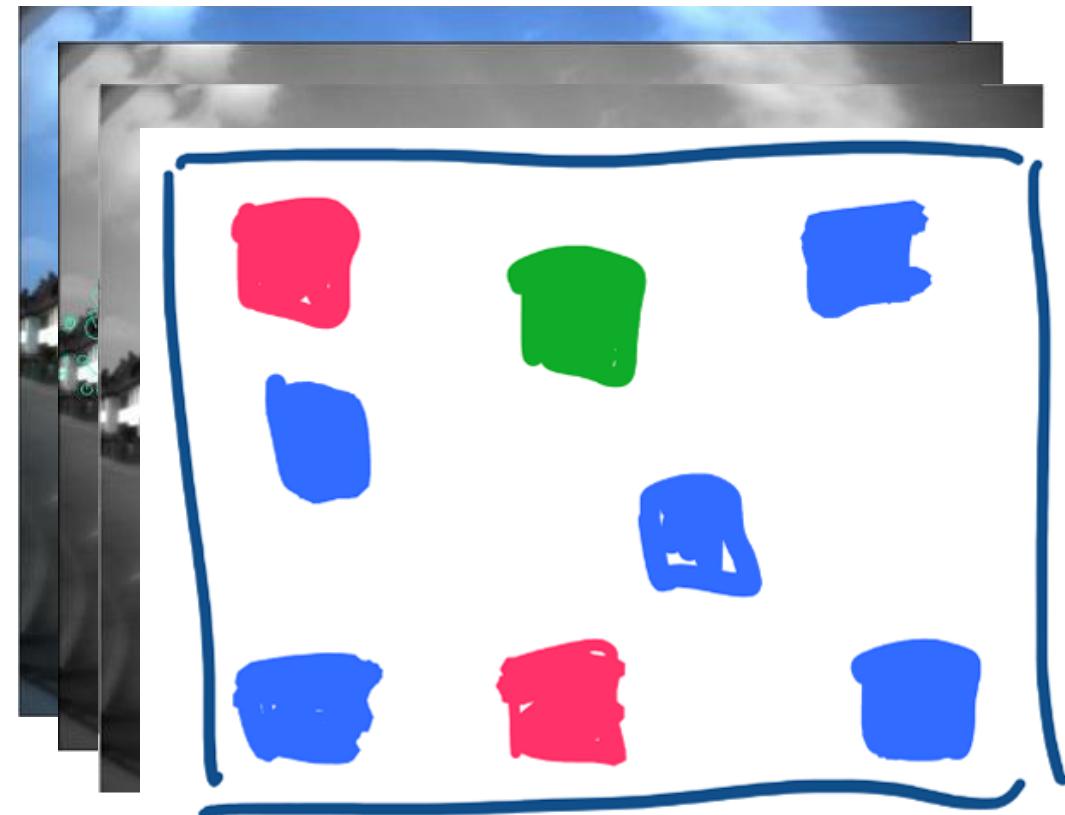
1. Extract features
2. Learn “visual dictionary”
3. **Quantize features using visual vocabulary**
4. Represent images by frequencies of “visual words”





BoVW Pipeline: Visual Overview

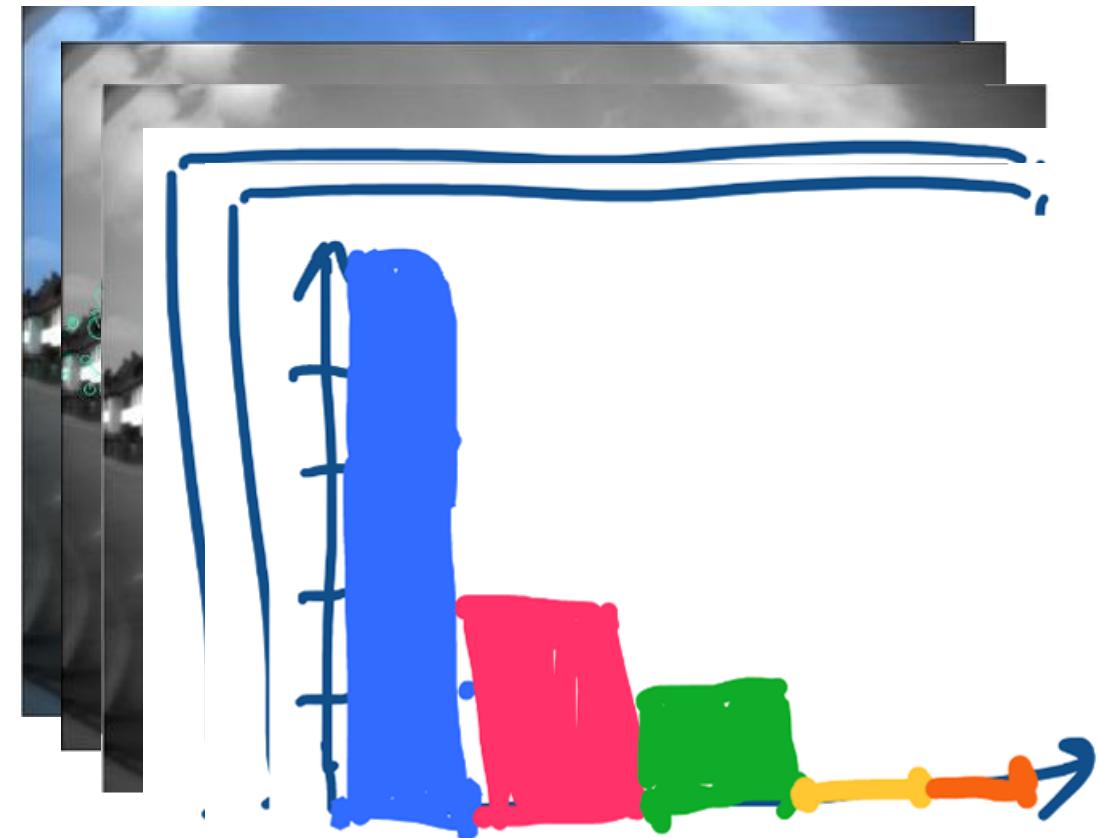
1. Extract features
2. Learn “visual dictionary”
- 3. Quantize features using visual vocabulary**
4. Represent images by frequencies of “visual words”





BoVW Pipeline: Visual Overview

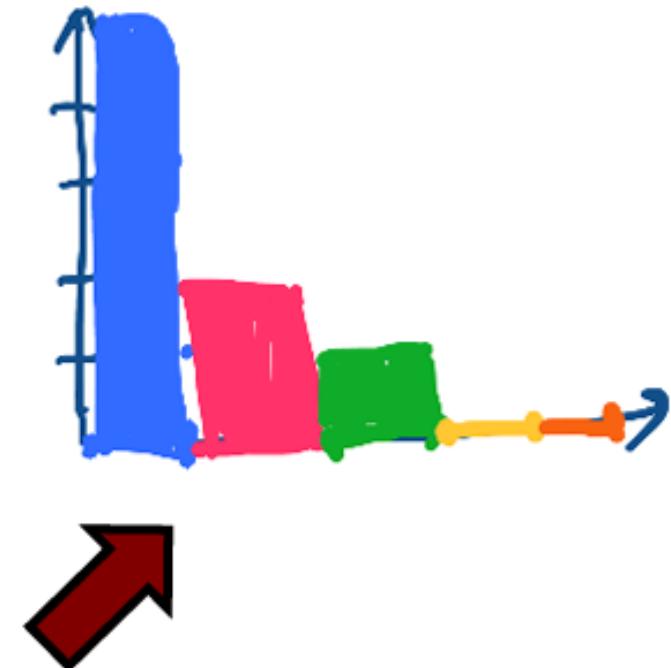
1. Extract features
2. Learn “visual dictionary”
3. Quantize features using visual vocabulary
- 4. Represent images by frequencies of “visual words”**





BoVW: Dictionary

- A dictionary defines the list of words that are considered
- The dictionary defines the x-axes of all the word occurrence histograms
- The dictionary must remain fixed

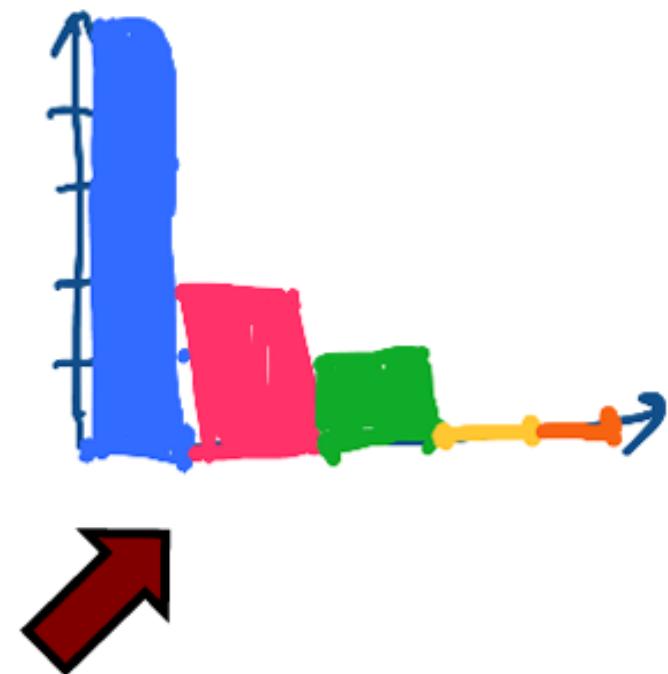




BoVW: Dictionary

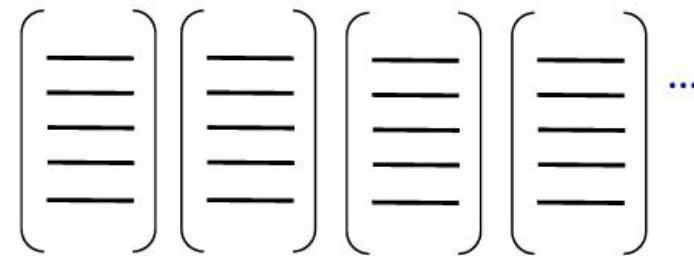
- A dictionary defines the list of words that are considered
- The dictionary defines the x-axes of all the word occurrence histograms
- The dictionary must remain fixed

The dictionary is typically learned from data. How can we do that?

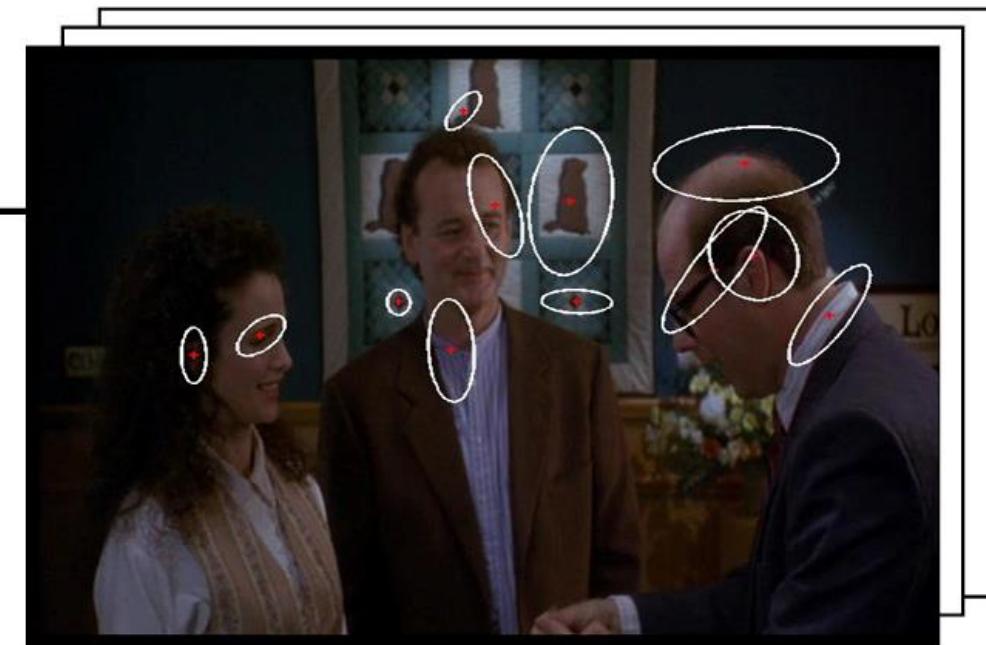




BoVW: Learn Visual Dictionary

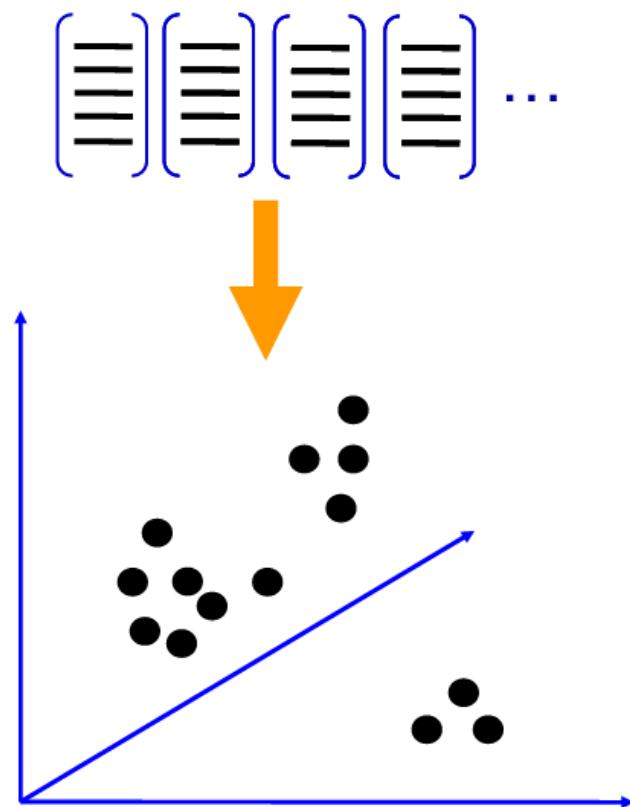


Compute SIFT Descriptor



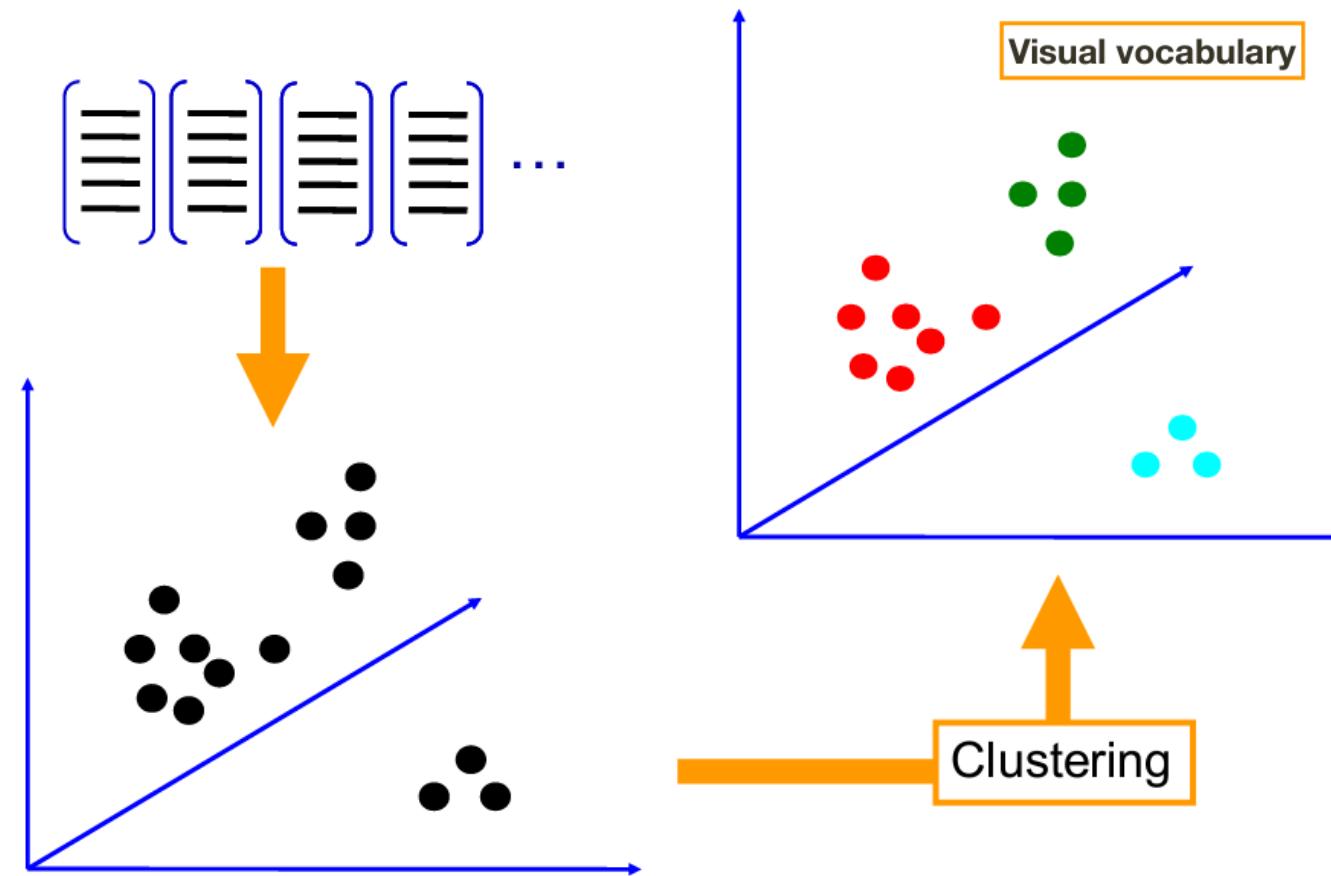


BoVW: Learn Visual Dictionary





BoVW: Learn Visual Dictionary





BoVW: K-Means Cluster Recap

- Want to minimize sum of squared Euclidean distances between points x_i and their nearest cluster centers m_k
- Algorithm
 - Randomly initialize K cluster centers
 - Iterate until convergence:
 - Assign each data point to the nearest center
 - Recompute each cluster center as the mean of all points assigned to it

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\text{point } i \text{ in cluster } k} (x_i - m_k)^2$$

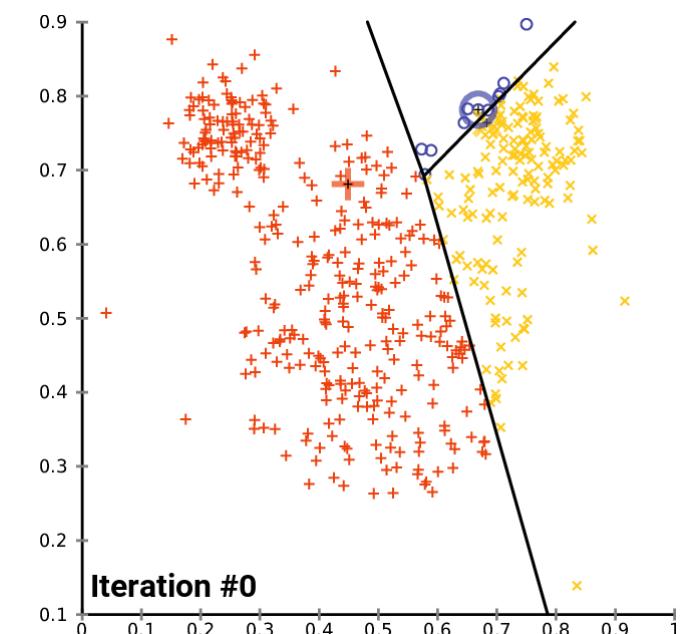


Image Credit: Wikipedia



BoVW: K-Means Cluster Recap

- Want to minimize sum of squared Euclidean distances between points x_i and their nearest cluster centers m_k
- Algorithm
 - Randomly initialize K cluster centers
 - Iterate until convergence:
 - Assign each data point to the nearest center
 - Recompute each cluster center as the mean of all points assigned to it

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\text{point } i \text{ in cluster } k} (x_i - m_k)^2$$

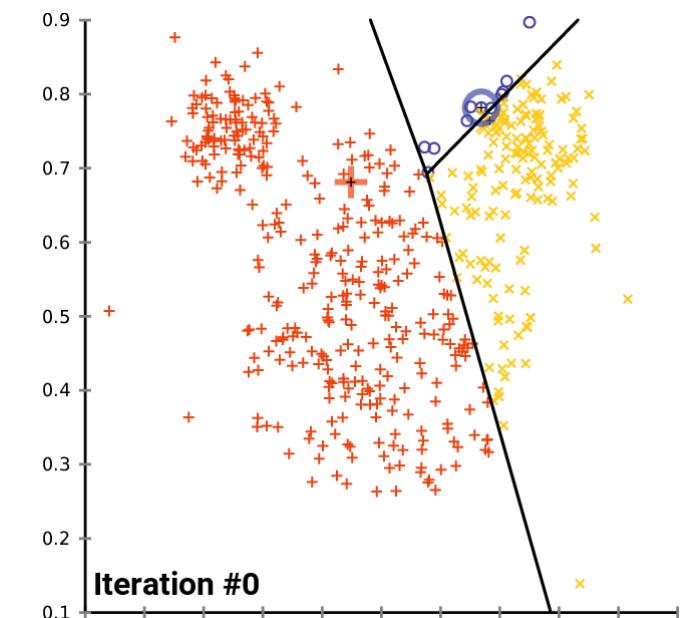
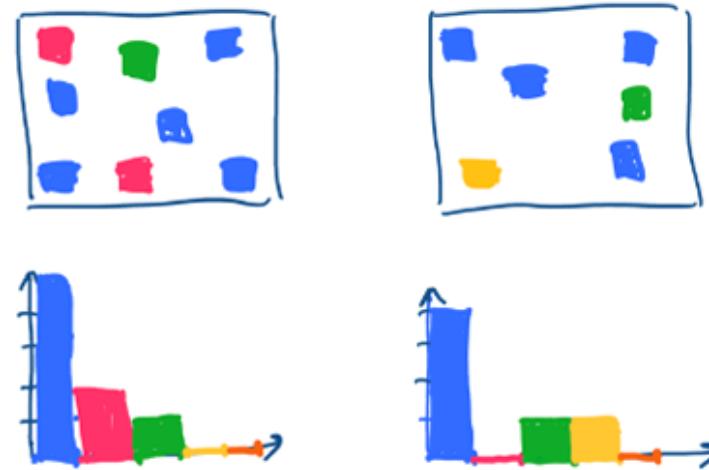
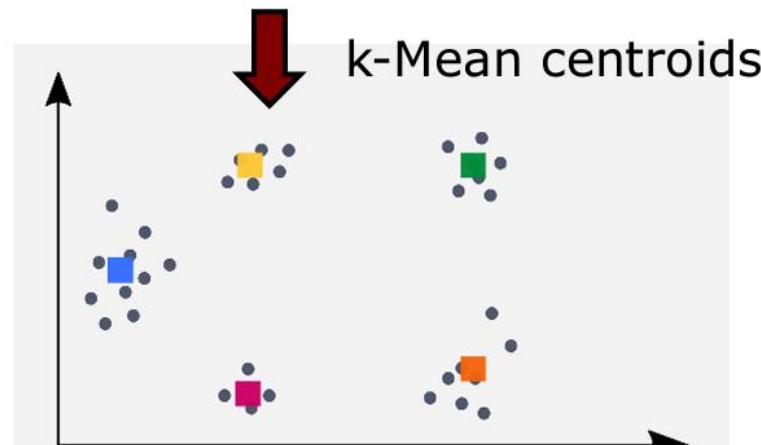


Image Credit: Wikipedia

We use k-means to compute the dictionary of visual words.



BoVW: Represent Images as Visual Words



Every image turns into a histogram!



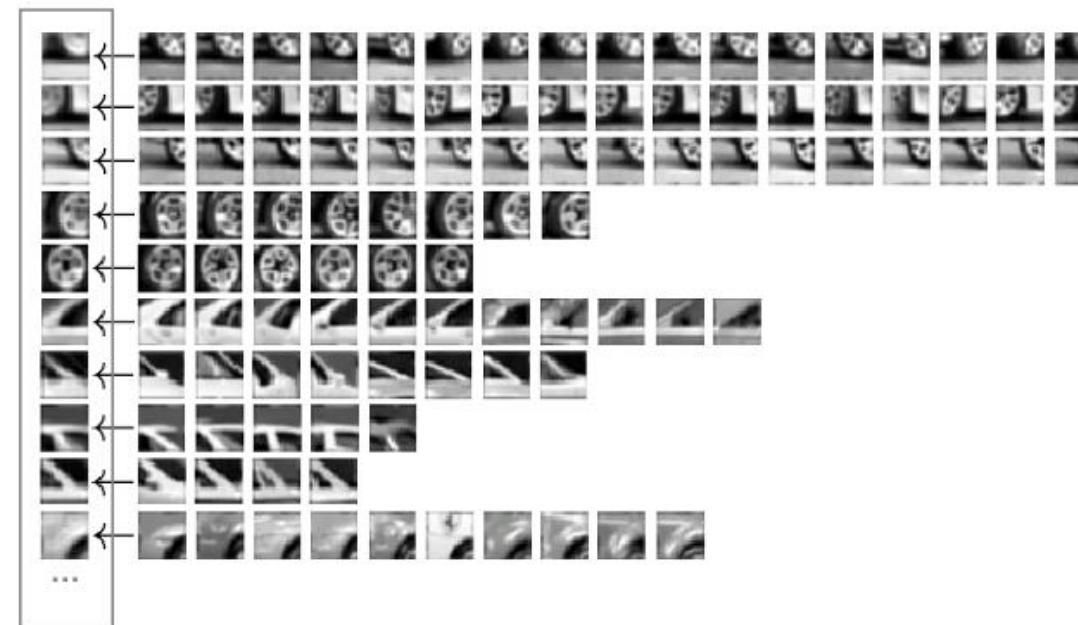
NYU

TANDON SCHOOL
OF ENGINEERING

Visual Place Recognition (cfeng@nyu.edu)

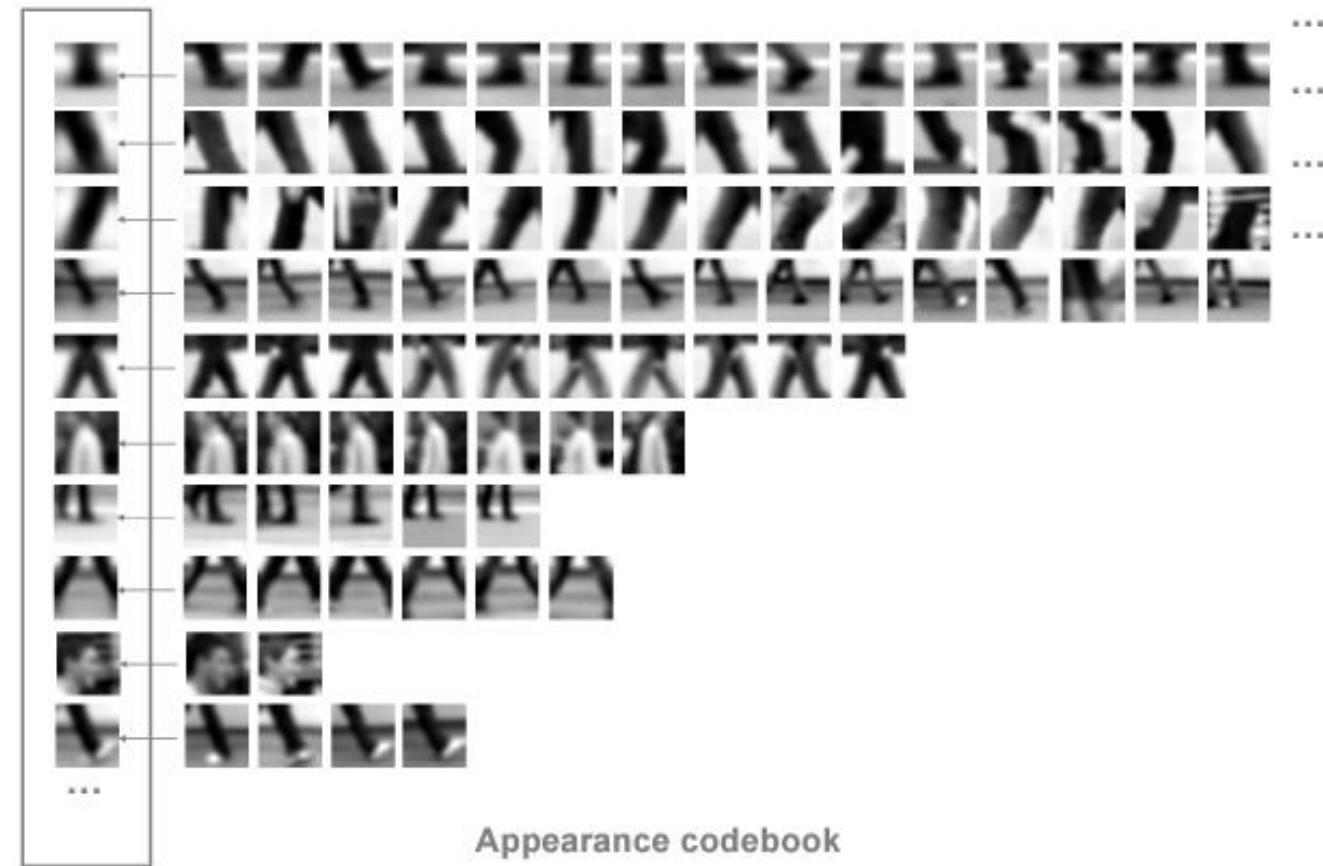


BoVW: Visual Dictionary Example





BoVW: Visual Dictionary Example



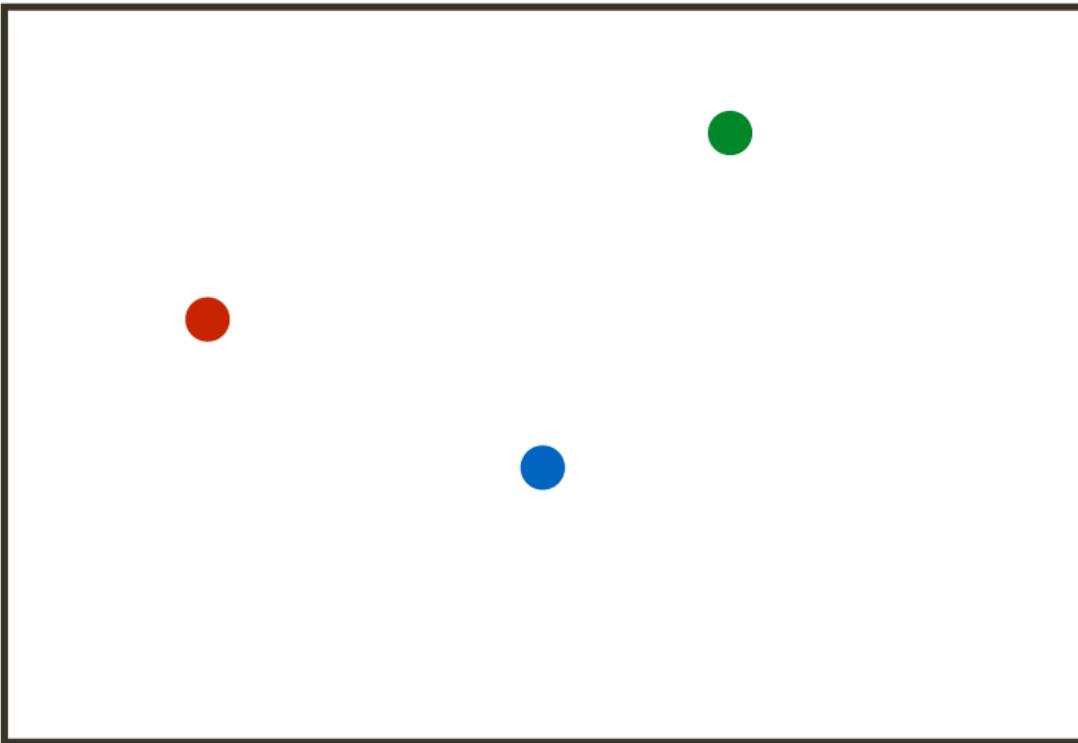
BoVW: Summary

- Compact summary of the image content
- Largely invariant to viewpoint changes and deformations
- Ignores the spatial arrangement
- Unclear how to choose optimal size of the dictionary
 - Too small: Words not representative of all image regions
 - Too large: Over-fitting

VLAD (Vector of Locally Aggregated Descriptors)

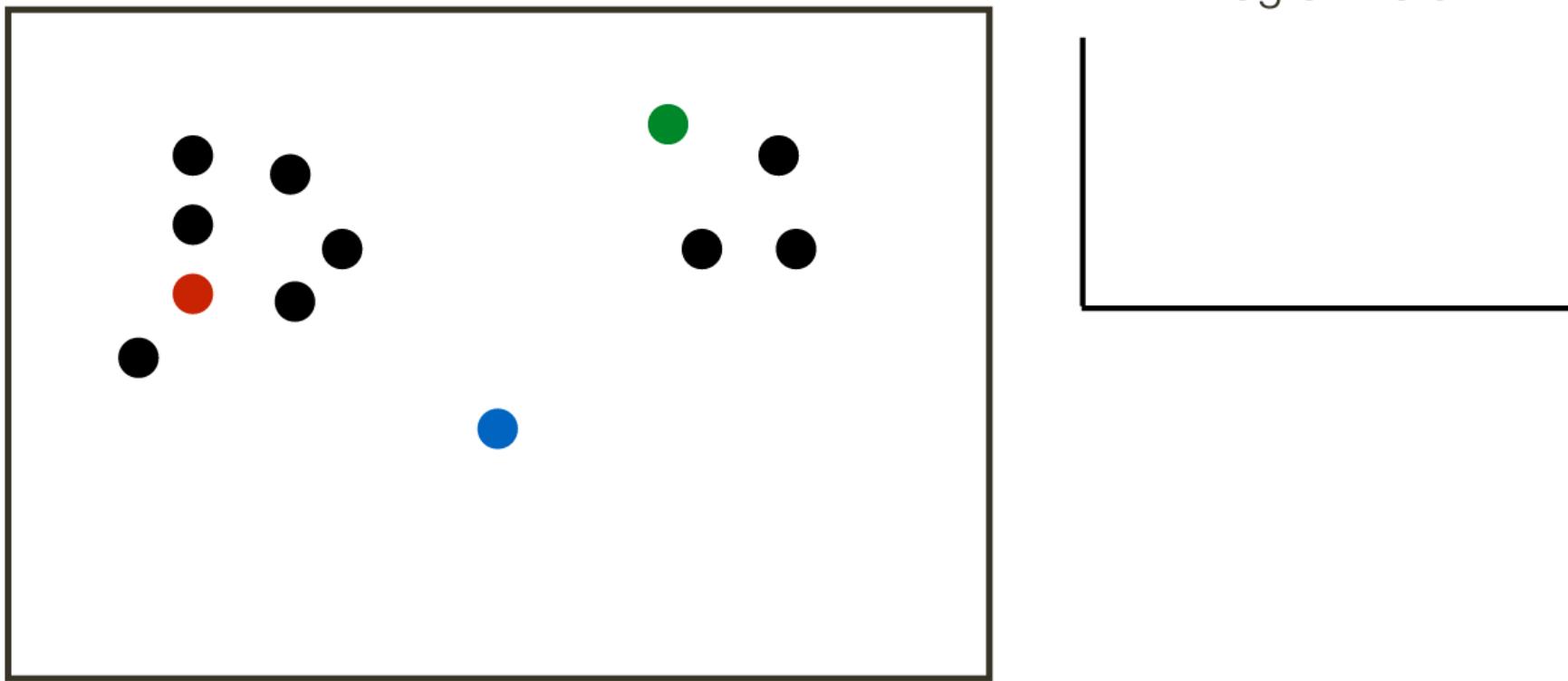
- There are more advanced ways to ‘count’ visual words than incrementing its histogram bin
- For example, it might be useful to describe how local descriptors are quantized to their visual words
- In the VLAD representation, instead of incrementing the histogram bin by one, we increment it by the residual vector $x - c(x)$

VLAD: Example



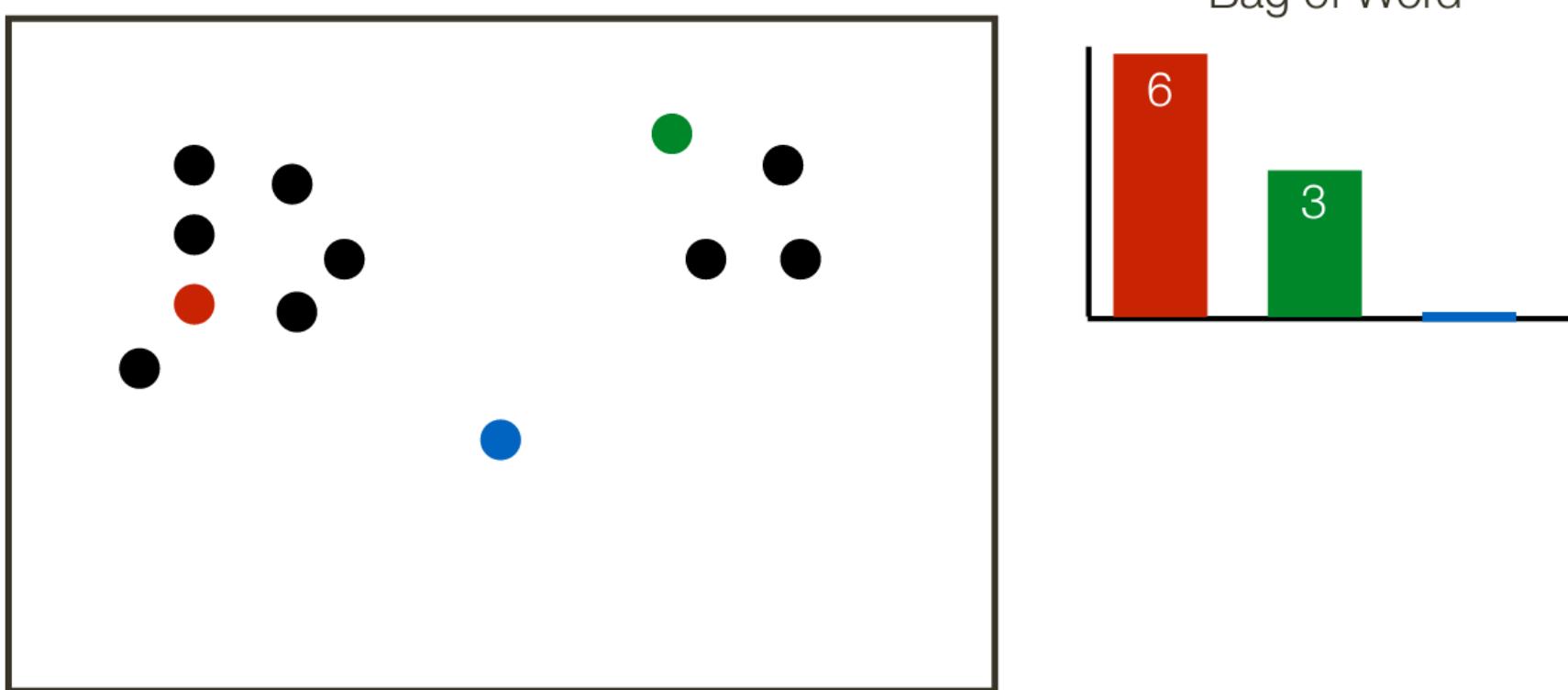


VLAD: Example



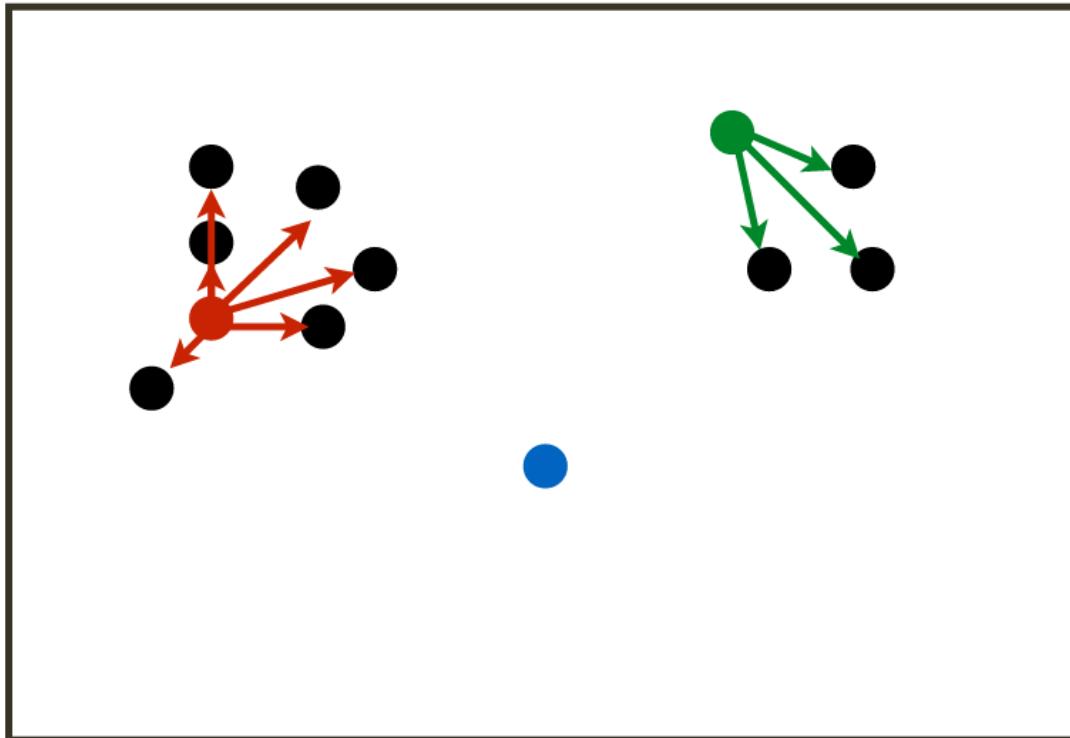


VLAD: Example





VLAD: Example



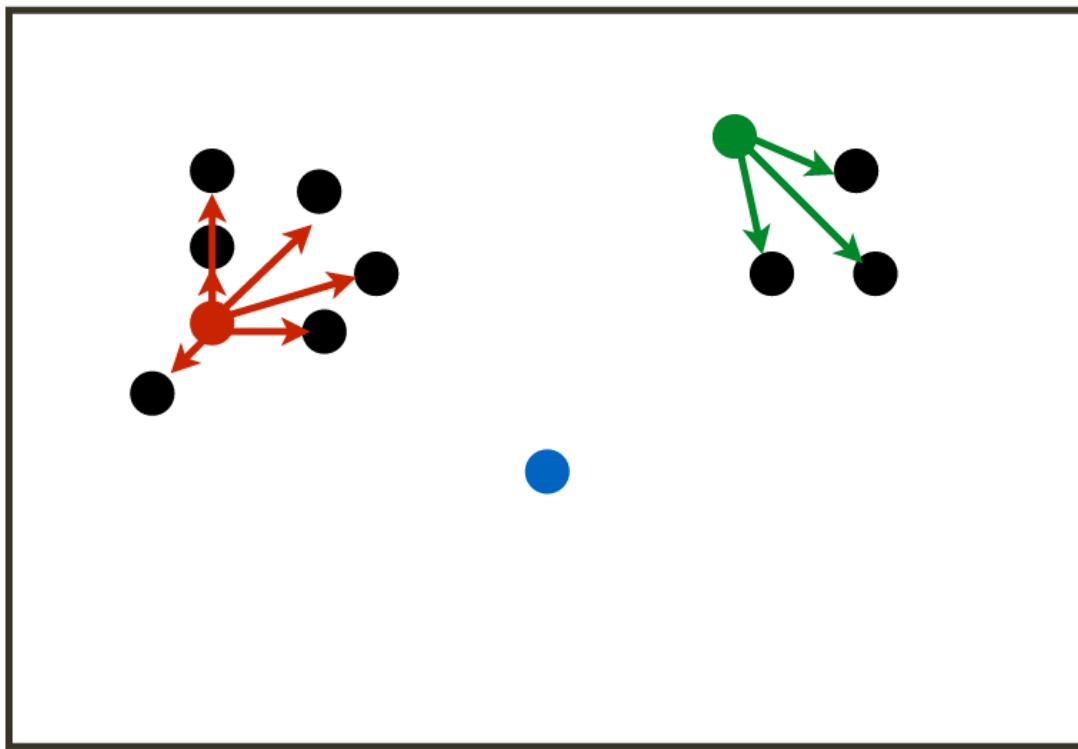
Bag of Word



VLAD



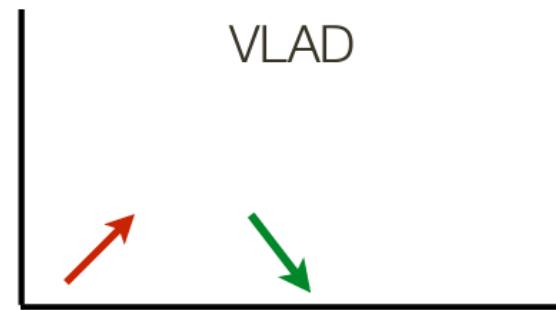
VLAD: Example



Bag of Word



VLAD





BoVW

Sparse (with a large vocab.)

- ▶ Inverted file indexing
- ▶ Size per image ~=
#features x 8 byte
(4 bytes for the index and 4 bytes
for the weight)
E.g. $4000 \times 8 = 32K$ byte
- + Can provide matches

VLAD (FV)

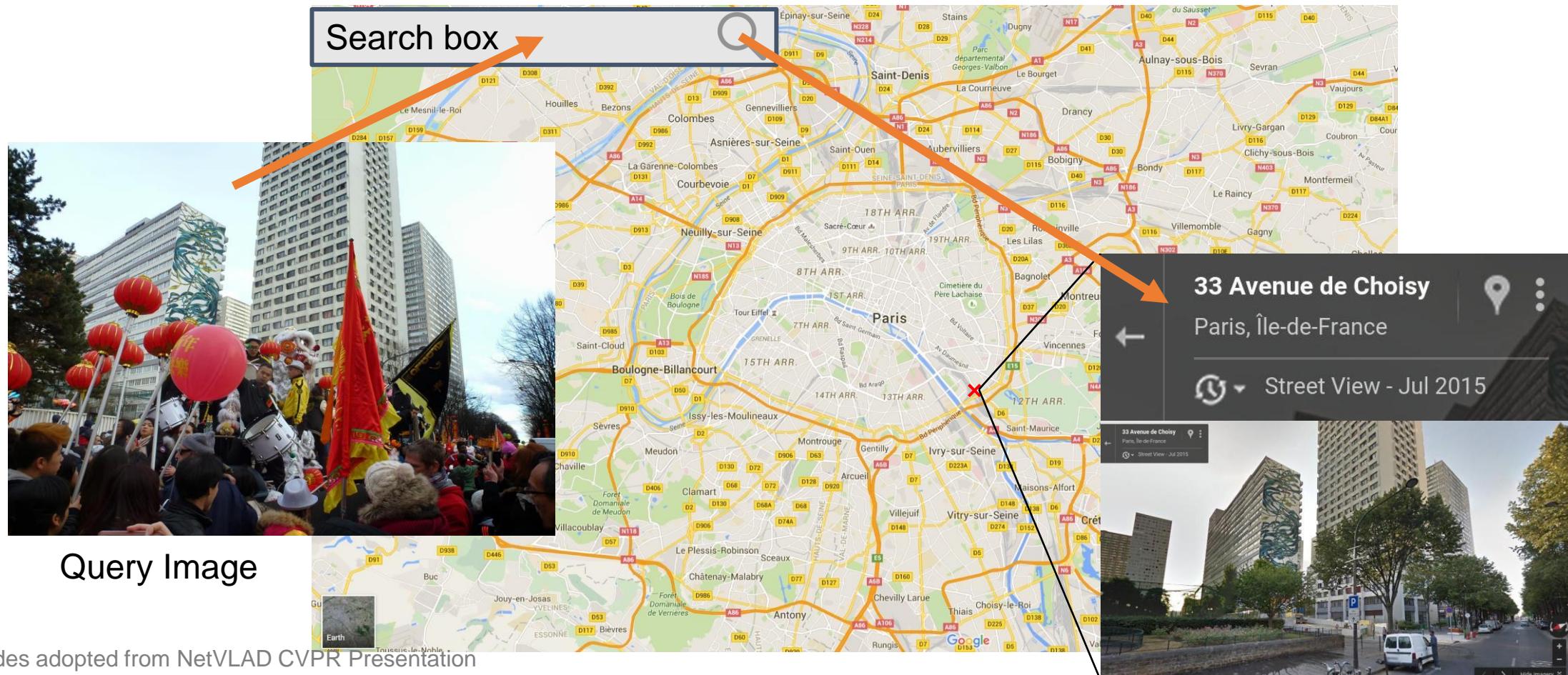
Dense

- ▶ ANN, product quantization
- ▶ Size per image ~=
#centroids x desc dim. x 4 byte
(4 bytes (single precision) per
dimension}
E.g. $256 \times 128 \times 4 = 65K$ byte
- + No extra memory requirement
to use densely detected features



NetVLAD: CNN architecture for weakly supervised place recognition

- Goal: Visual Place Recognition



NetVLAD: Difficulty in VPR

- **Lighting changes: Different time of day / year**
- Changes in camera viewpoint
- Occluders and ambiguous objects: People, cars, trees, pavement...
- Big data: World-scale localization



NetVLAD: Difficulty in VPR

- Lighting changes: Different time of day / year
- **Changes in camera viewpoint**
- Occluders and ambiguous objects: People, cars, trees, pavement...
- Big data: World-scale localization



NetVLAD: Difficulty in VPR

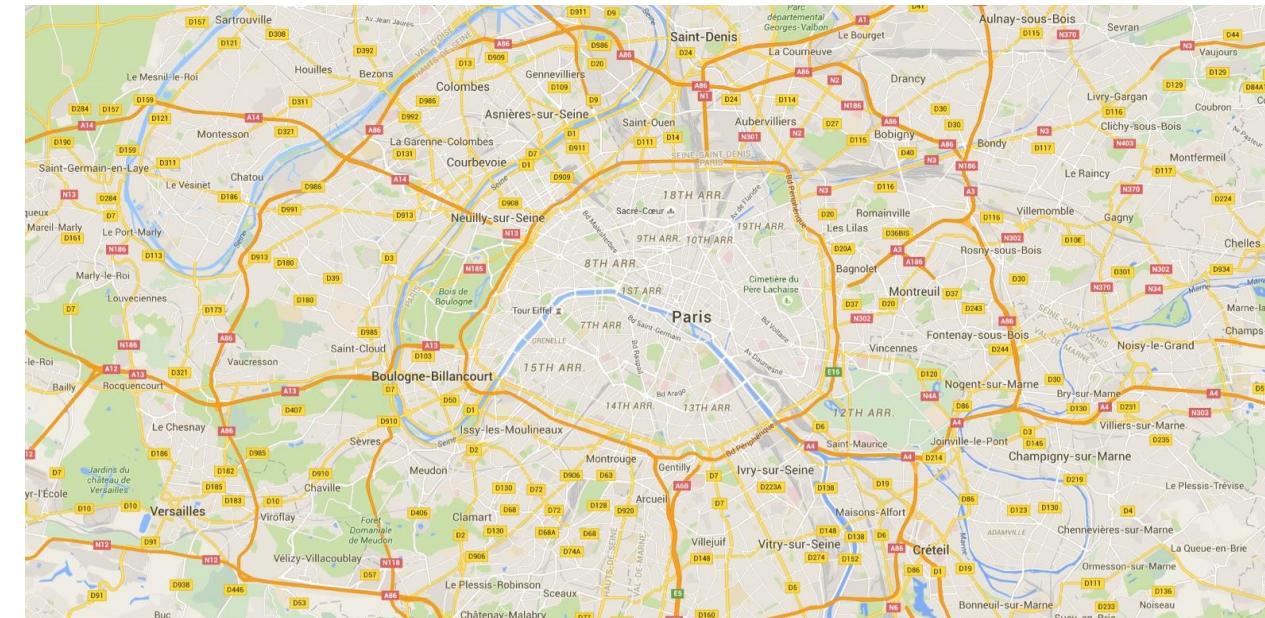
- Lighting changes: Different time of day / year
- Changes in camera viewpoint
- **Occluders and ambiguous objects: People, cars, trees, pavement...**
- Big data: World-scale localization



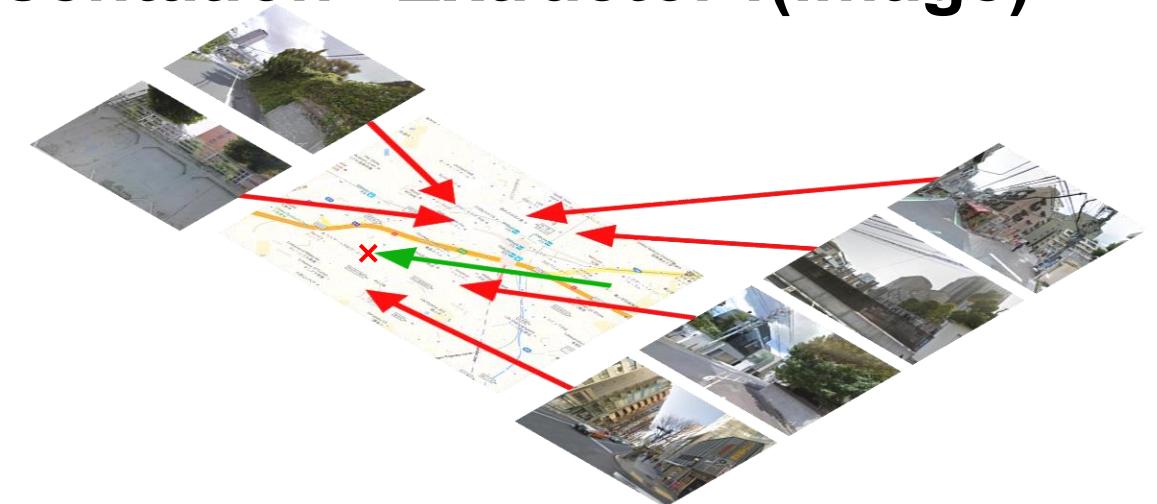


NetVLAD: Difficulty in VPR

- Lighting changes: Different time of day / year
- Changes in camera viewpoint
- Occluders and ambiguous objects: People, cars, trees, pavement...
- **Big data: World-scale localization**



NetVLAD: Design an “Image Representation” Extractor $f(\text{image})$

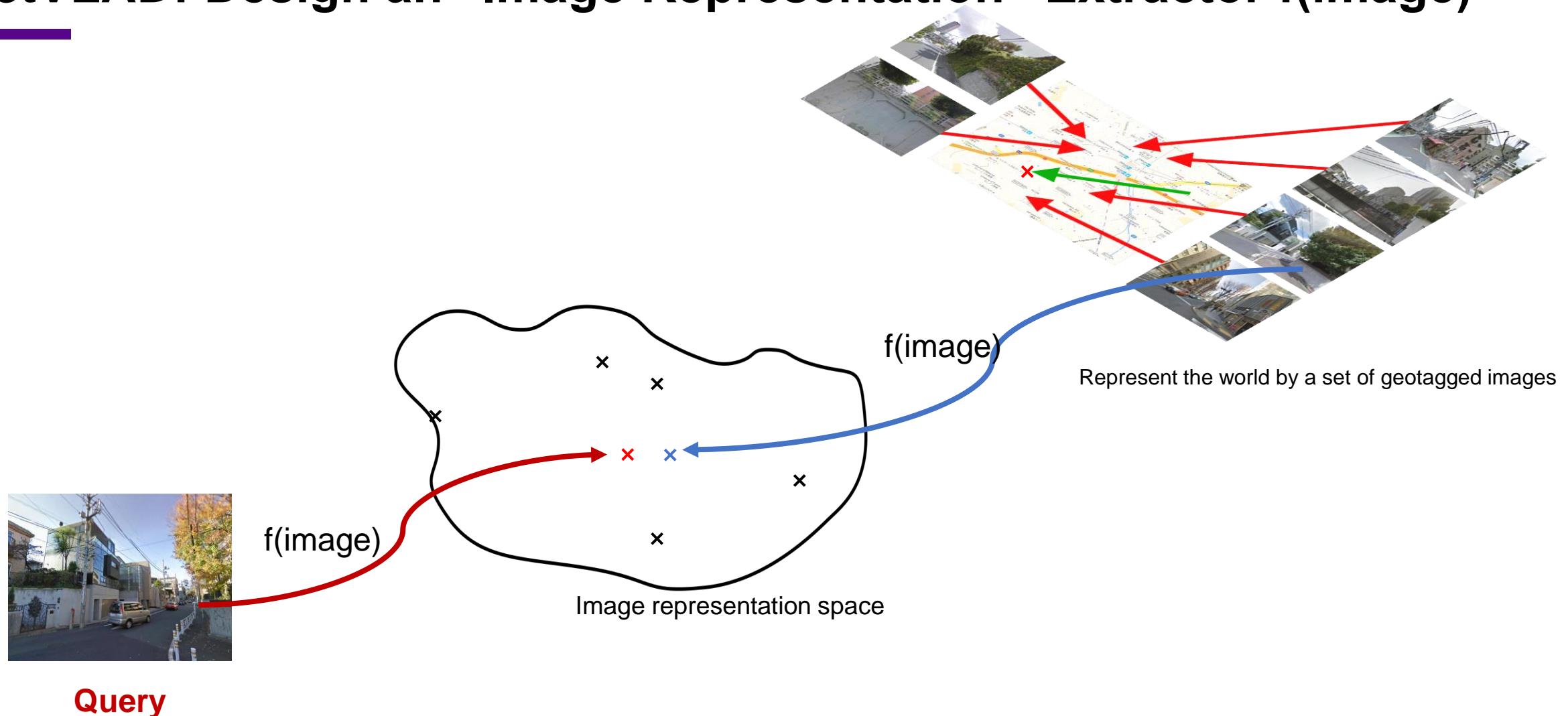


Represent the world by a set of geotagged images

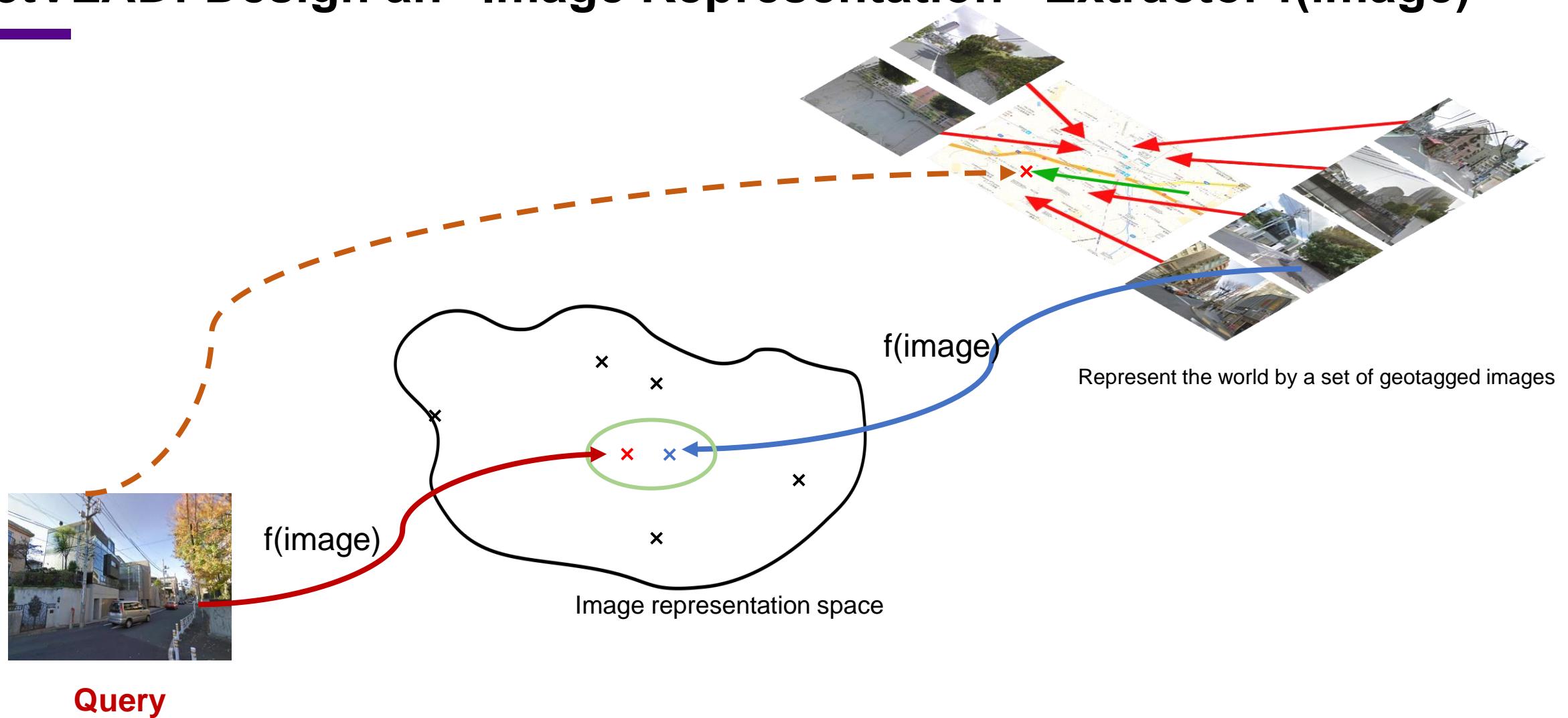


Query

NetVLAD: Design an “Image Representation” Extractor $f(\text{image})$



NetVLAD: Design an “Image Representation” Extractor $f(\text{image})$



NetVLAD: Can we apply CNNs to place recognition?

Questions:

1. (Model) What is a good CNN architecture?
2. (Data) How to get the lots of annotated training data?
3. (Loss) What is the appropriate loss for end-to-end training?

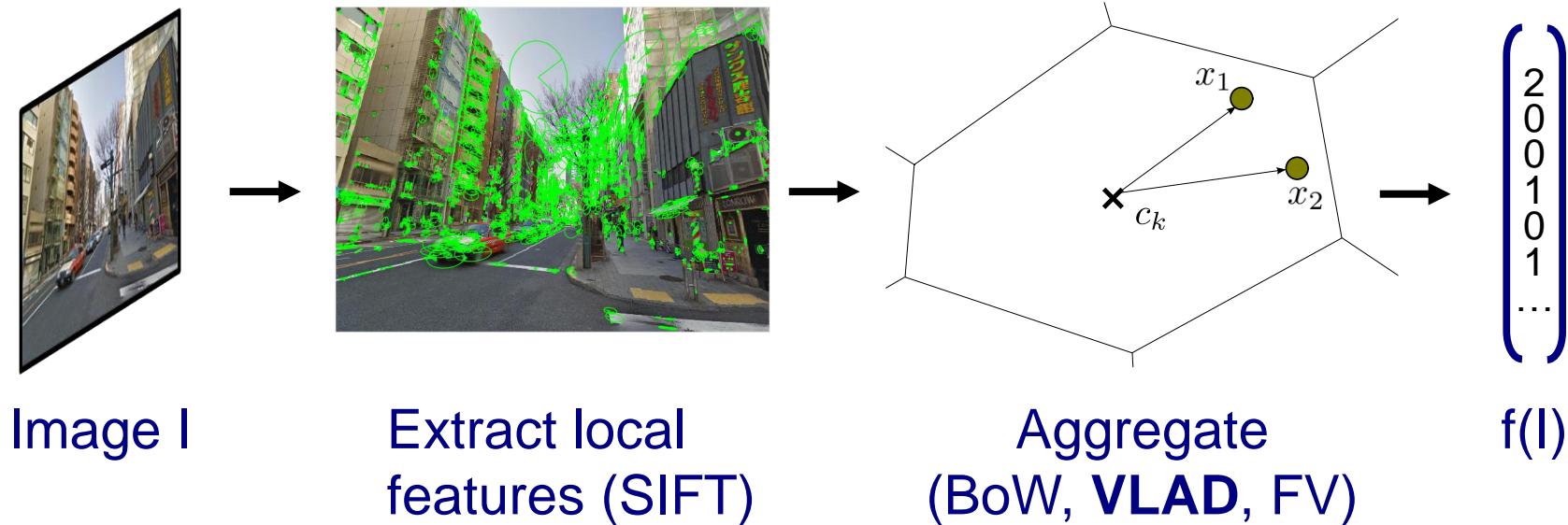
NetVLAD: Can we apply CNNs to place recognition?

Questions:

1. **(Model) What is a good CNN architecture?**
2. **(Data) How to get the lots of annotated training data?**
3. **(Loss) What is the appropriate loss for end-to-end training?**



NetVLAD: Model



NetVLAD: Model

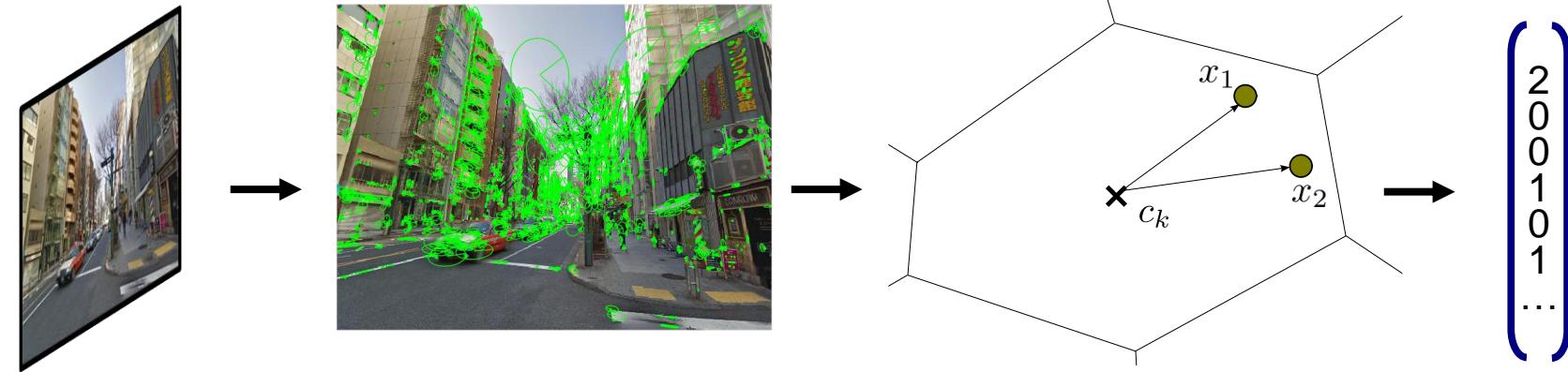
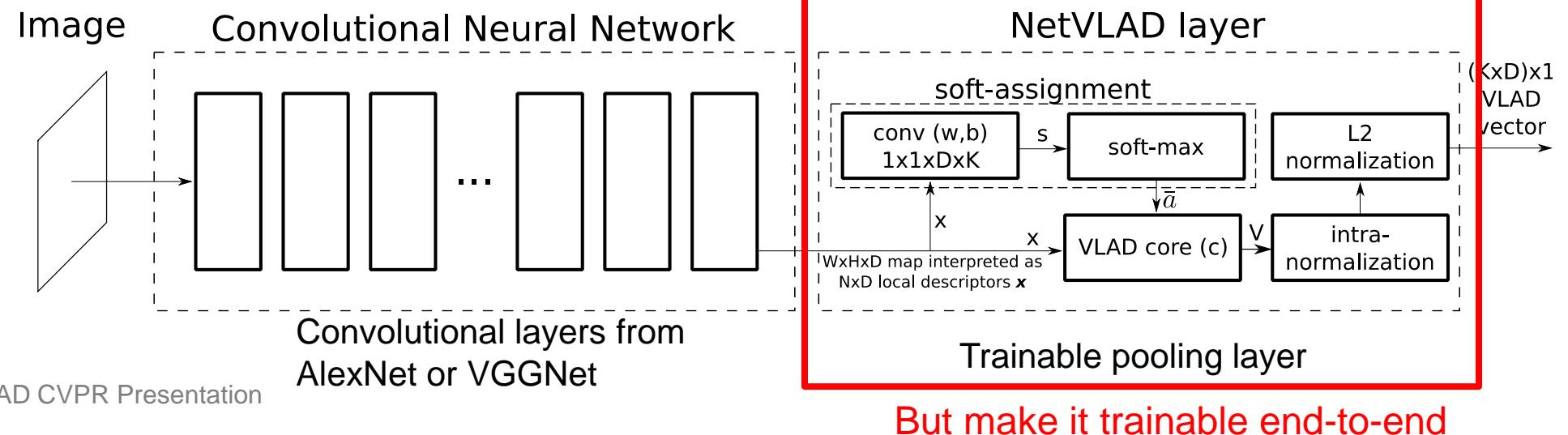


Image I

Extract local
features (SIFT)

Aggregate
(BoW, **VLAD**, FV)

$f(I)$





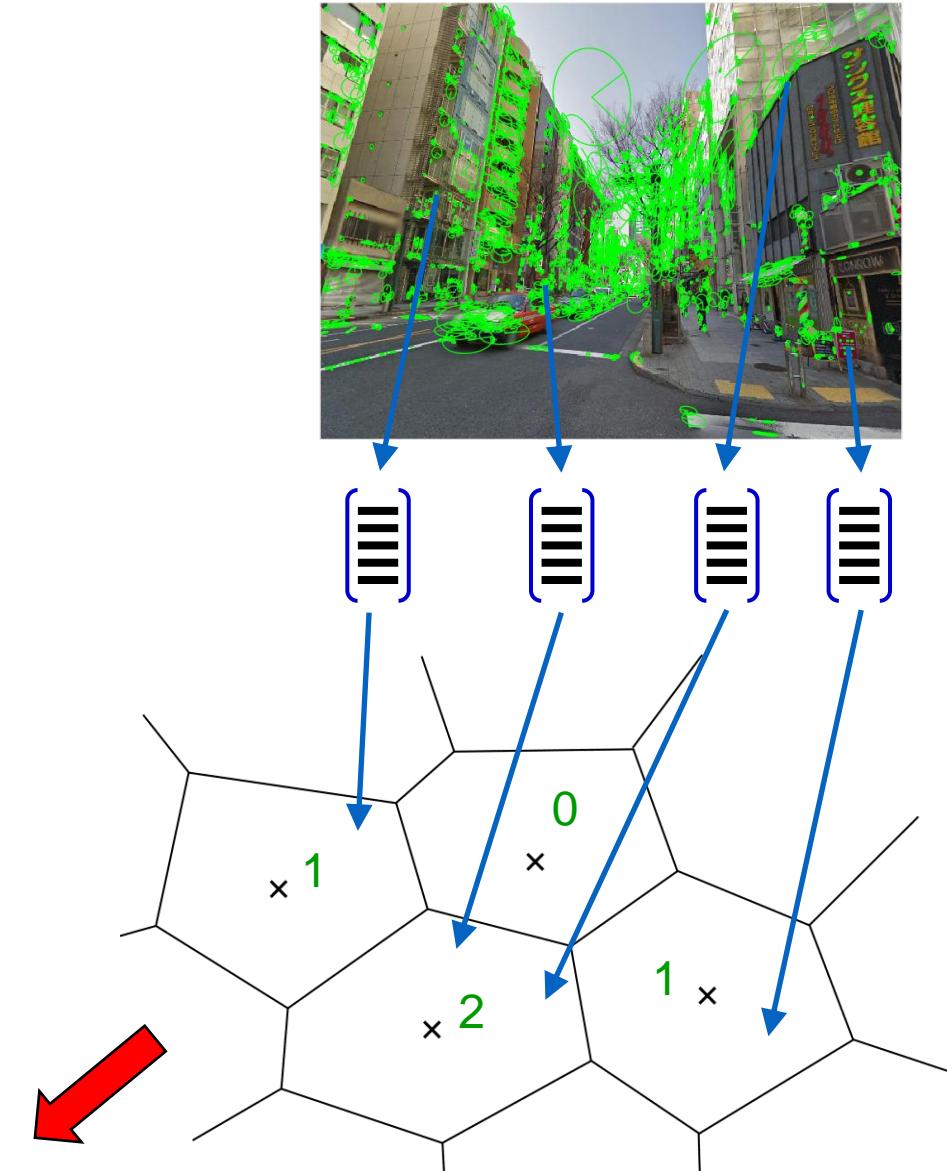
NetVLAD: BoVW Recap

0/1 assignment of descriptor i to cluster k

$$B(k) = \sum_{i=1}^N a_k(x_i)$$

Sum over all N descriptors in the image

$$B = [1, 0, 2, 1, \dots]$$





NetVLAD: VLAD Recap

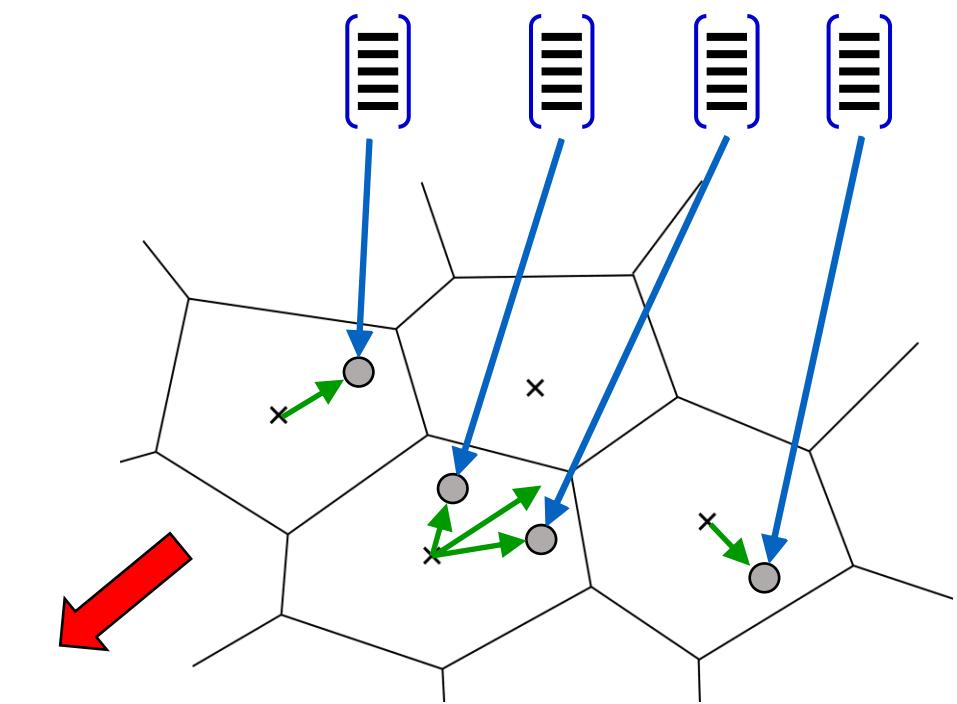
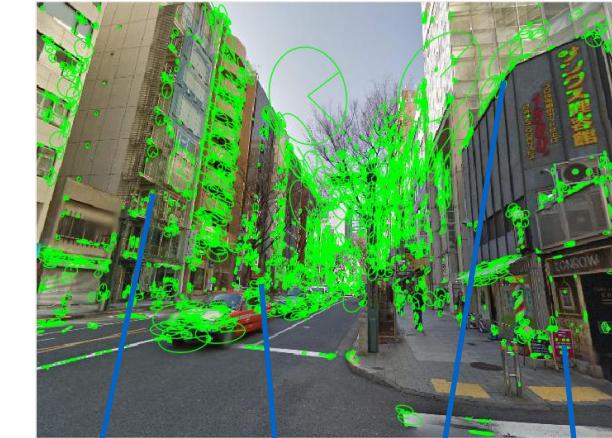
0/1 assignment of descriptor i to cluster k

$$V(:, k) = \sum_{i=1}^N a_k(x_i)(x_i - c_k)$$

Residual vector

Sum over all N descriptors in the image

$$V = [\nearrow, \dots, \nearrow, \searrow, \dots]$$



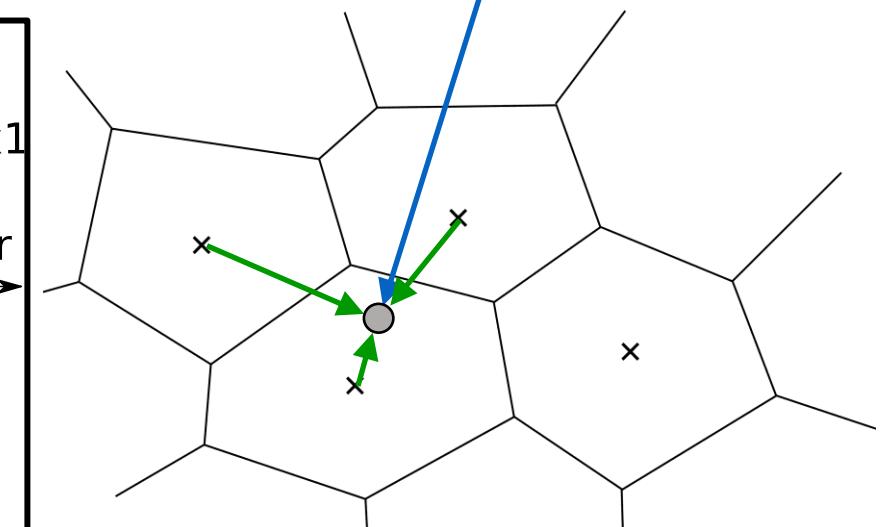
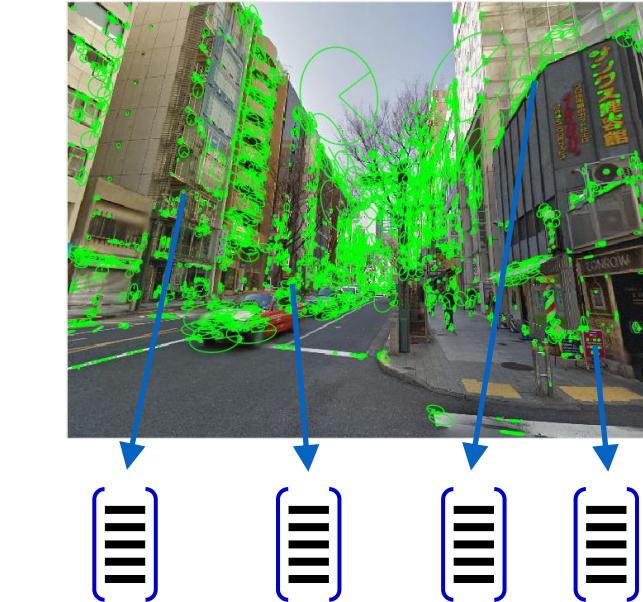
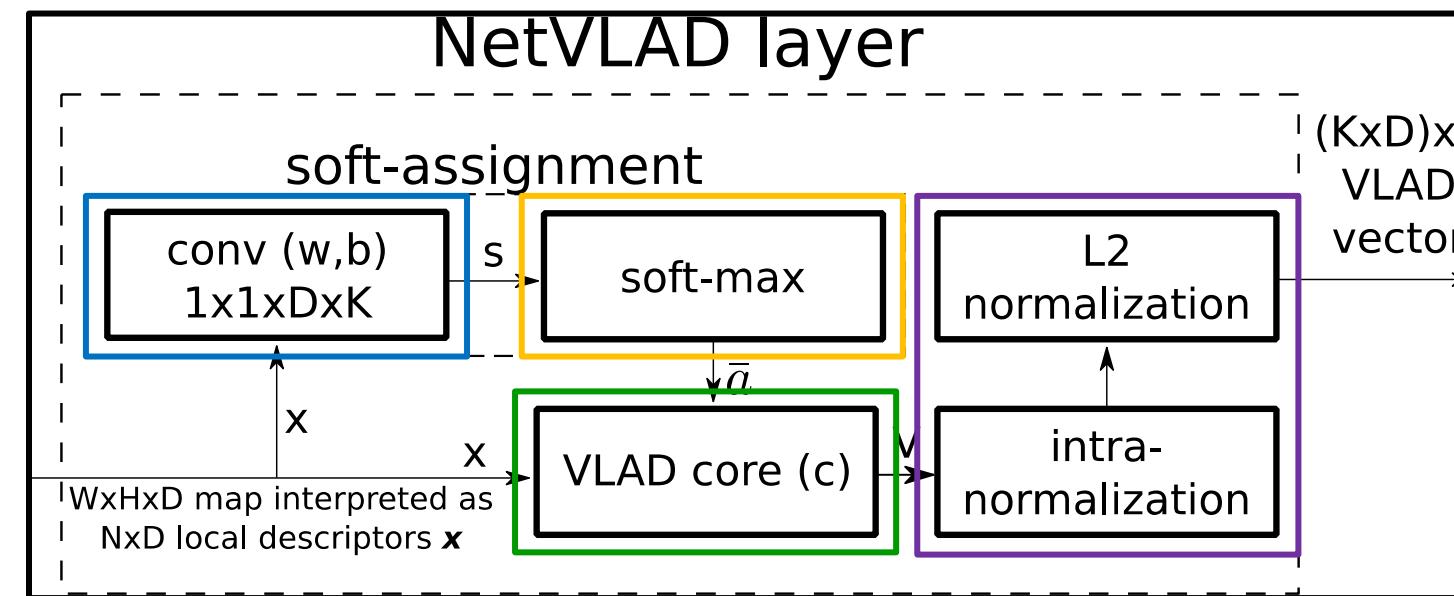


NetVLAD: NetVLAD Modification

$$V(:, k) = \sum_{i=1}^N \frac{e^{w_k^T x_i + b_k}}{\sum' e^{w_{k'}^T x_i + b_{k'}}} (x_i - c_k)$$

Soft-max → \bar{a}

Residual vector

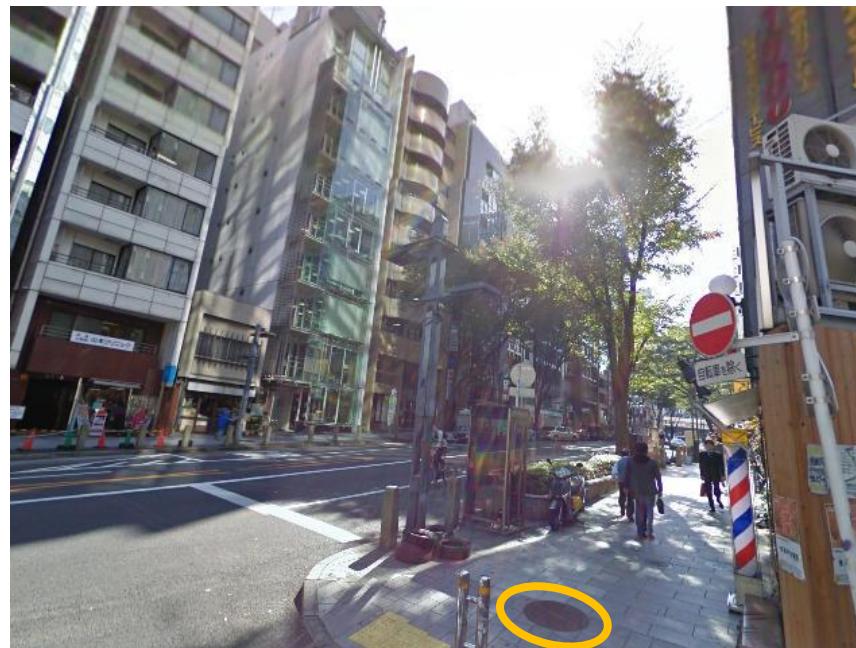


NetVLAD: Can we apply CNNs to place recognition?

Questions:

1. (Model) What is a good CNN architecture?
2. (Data) How to get the lots of annotated training data?
3. (Loss) What is the appropriate loss for end-to-end training?

NetVLAD: Google Street View Time Machine



same locations at different times and seasons

NetVLAD: Google Street View Time Machine

Tokyo 24/7 [Torii et al. 15]

Database: 76k images from Street View

Queries: 315 images from mobile phone cameras



NetVLAD: Can we apply CNNs to place recognition?

Questions:

1. (Model) What is a good CNN architecture?
2. (Data) How to get the lots of annotated training data?
3. **(Loss) What is the appropriate loss for end-to-end training?**



NetVLAD: Weakly Supervised Ranking Loss

$$L_\theta = \sum_j l \left(\min_i d_\theta^2(q, p_i^q) + m - d_\theta^2(q, n_j^q) \right)$$

hinge loss

margin

Sum over negatives

Distance to the best potential positive

Distance to the negative

NetVLAD: Results

Query



Top result

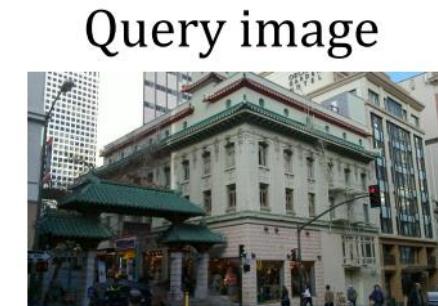


Green: Correct
Red: Incorrect

Code and trained networks available online:

<http://www.di.ens.fr/willow/research/netvlad/> (the link is in the paper)

Summary: Visual Place Recognition Pipeline



Database images



Offline

1. Feature detection & description
2. Training visual vocabulary
3. Image description

4. Feature detection & description
5. Image description
6. Initial ranking
7. Re-ranking with geometric verification

Posenet: Convolutional networks for real-time 6-DOF camera relocalization

Convolutional networks
for real-time 6-DOF
camera relocalization

Alex Kendall, Matthew Grimes, Roberto Cipolla

References

- Szeliski 2022
 - Section 6.3
- Forsyth & Ponce 2011
 - Section 15.1, Chapter 17
- <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>
- Ren, Shaoqing, Kaiming He, Ross Girshick, and Jian Sun. "Faster r-cnn: Towards real-time object detection with region proposal networks." arXiv preprint arXiv:1506.01497 (2015).
- Redmon, Joseph, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 779-788. 2016.