

Computer Vision

DATA.ML.300, 5 study credits

Esa Rahtu
Laboratory of Signal Processing, Tampere University

What is content based image retrieval

What we would like to be able to do?

- Query by example
- Given:
 - An example **query image** illustrating the users needs
 - A very large dataset of images
- Task:
 - **Rank** all images in the dataset according to relevance to the query



Difference to classification

Query: This chair



Results from dataset classified as “chair”

Classification

Difference to classification

Query: This chair



Results from dataset ranked by similarity to the query

Retrieval

Instance level retrieval

Query



Retrieved results



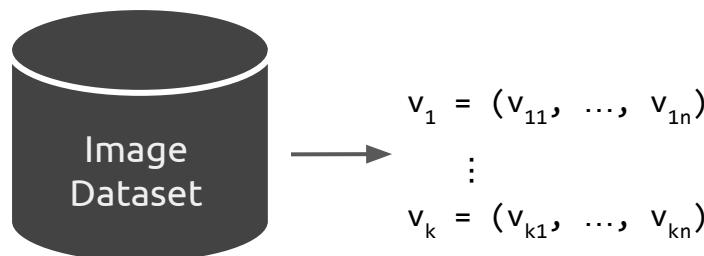
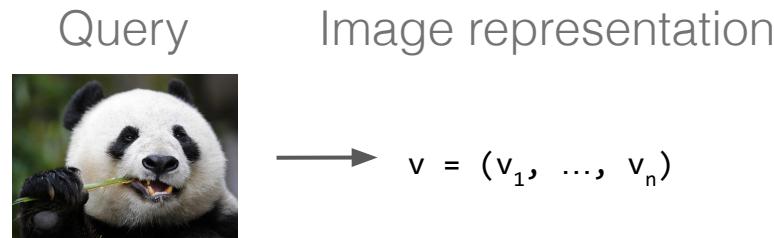
The retrieval pipeline

The retrieval pipeline

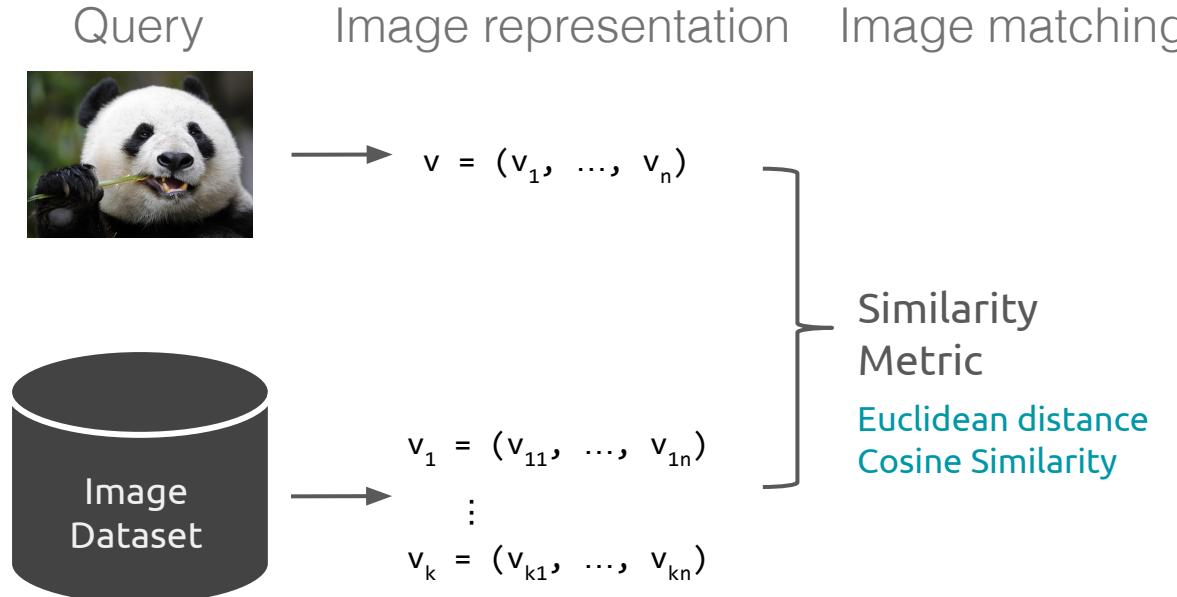
Query



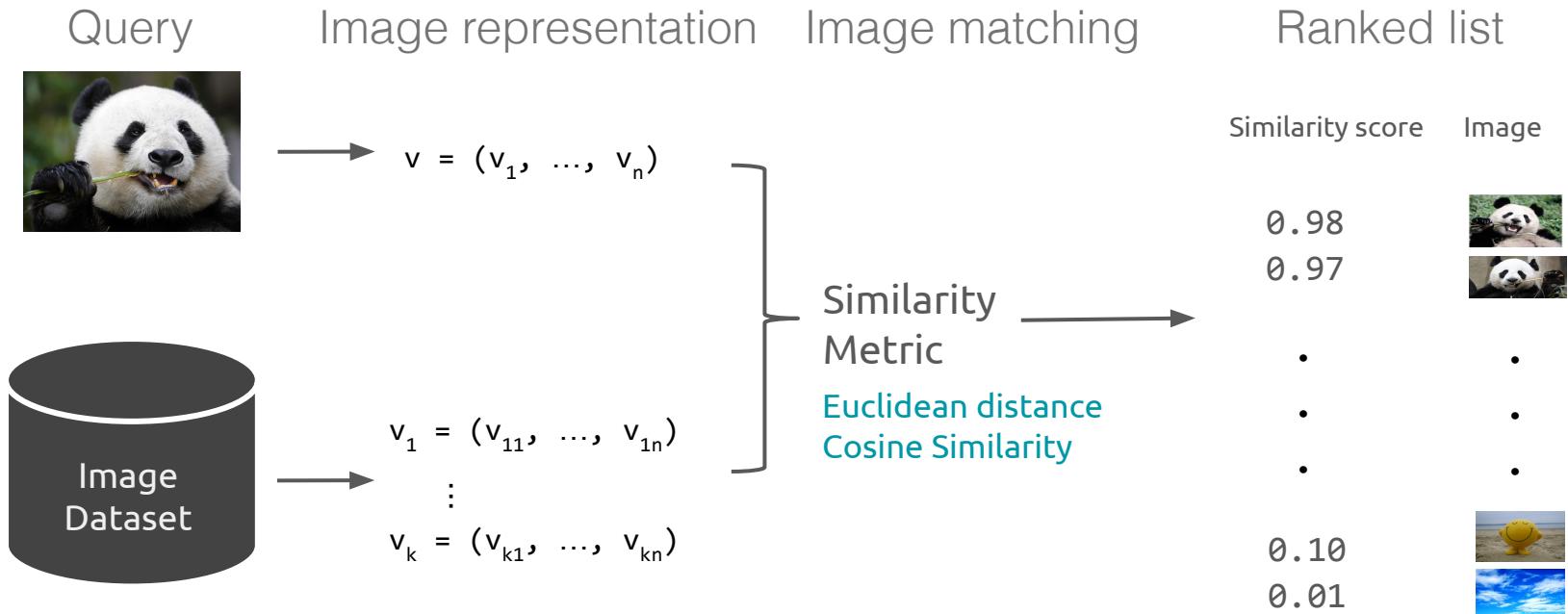
The retrieval pipeline



The retrieval pipeline

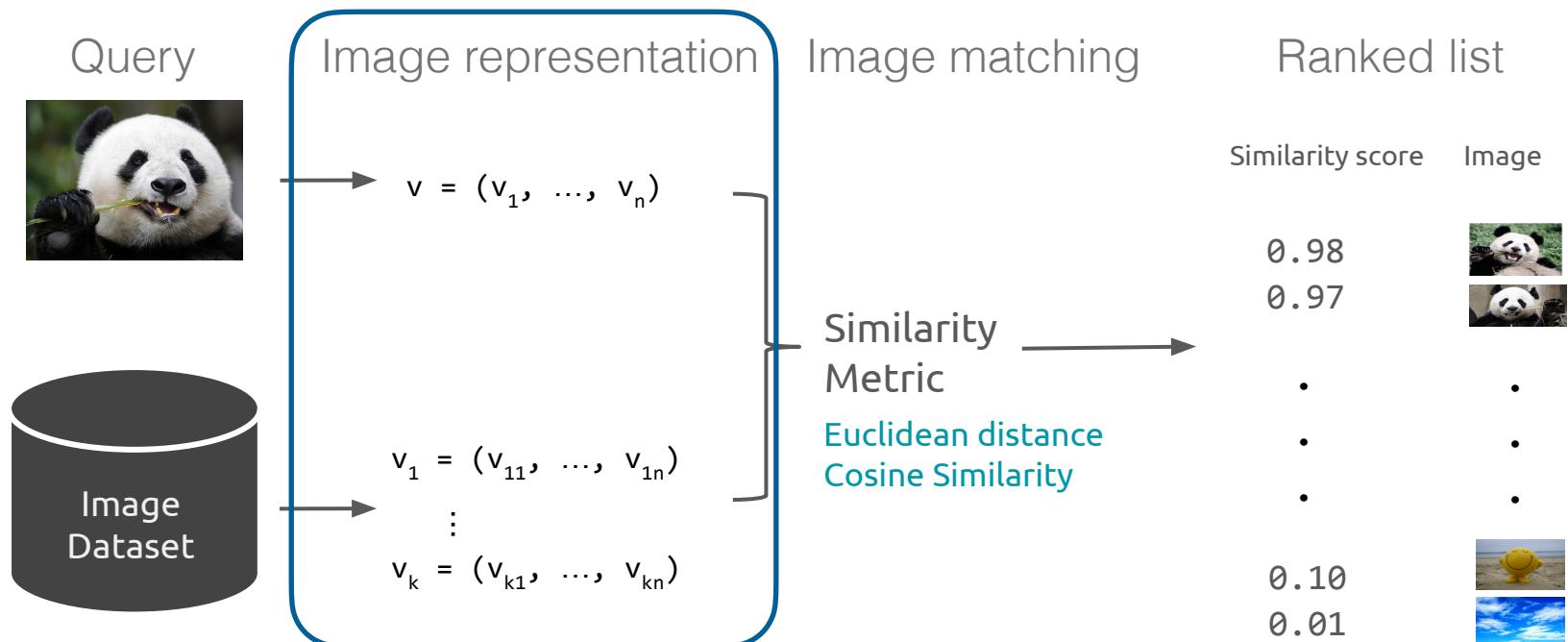


The retrieval pipeline



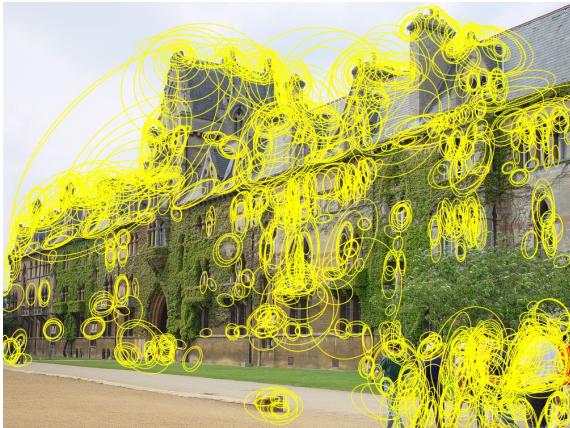
The image representation

The retrieval pipeline

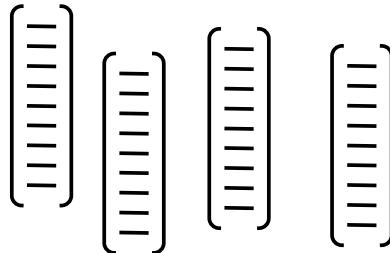


Classical approach

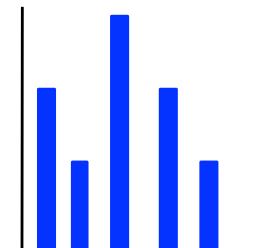
Detect variable number
of local features (e.g. SIFT)



Compute descriptor
(e.g. SIFT)

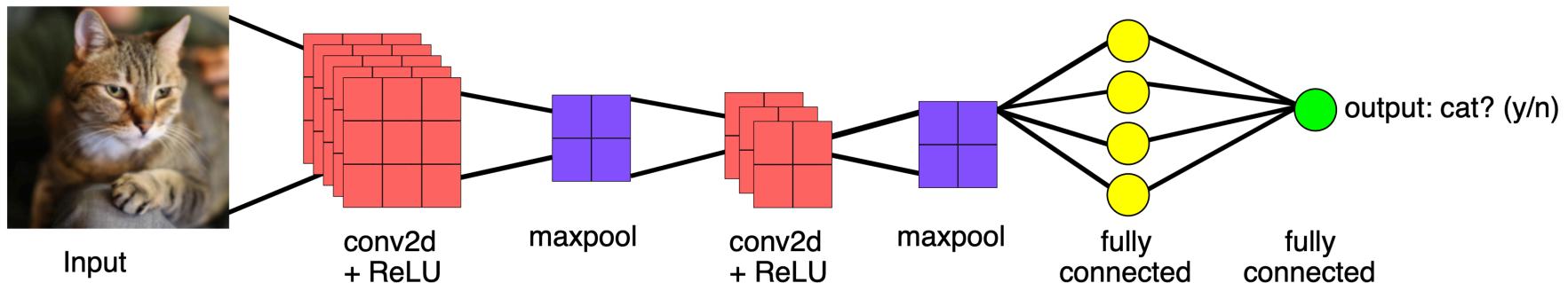


Quantise to form an
image level descriptor



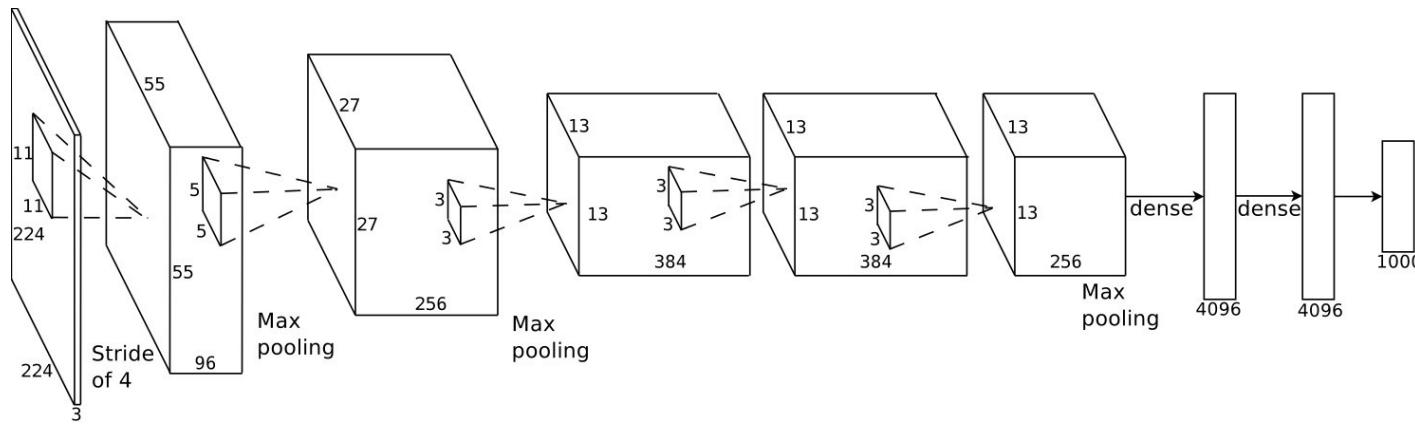
CNNs and retrieval

- CNNs are state-of-the-art in classification
- Could they be used for retrieval?



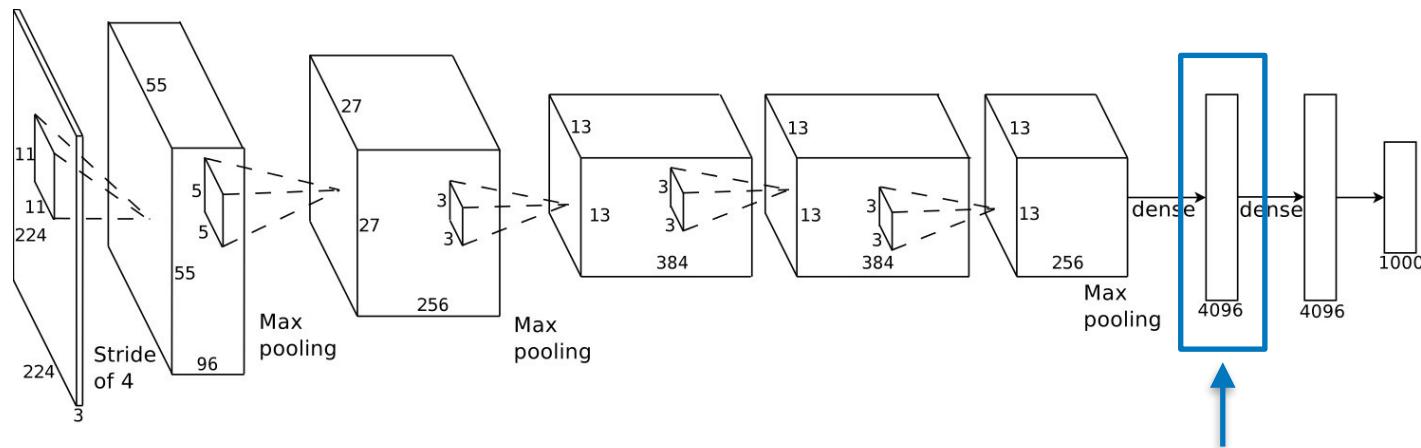
Using off-the-self CNN representations

image



Using off-the-self CNN representations

image



FC layers as global
feature representation

Using off-the-self CNN representations

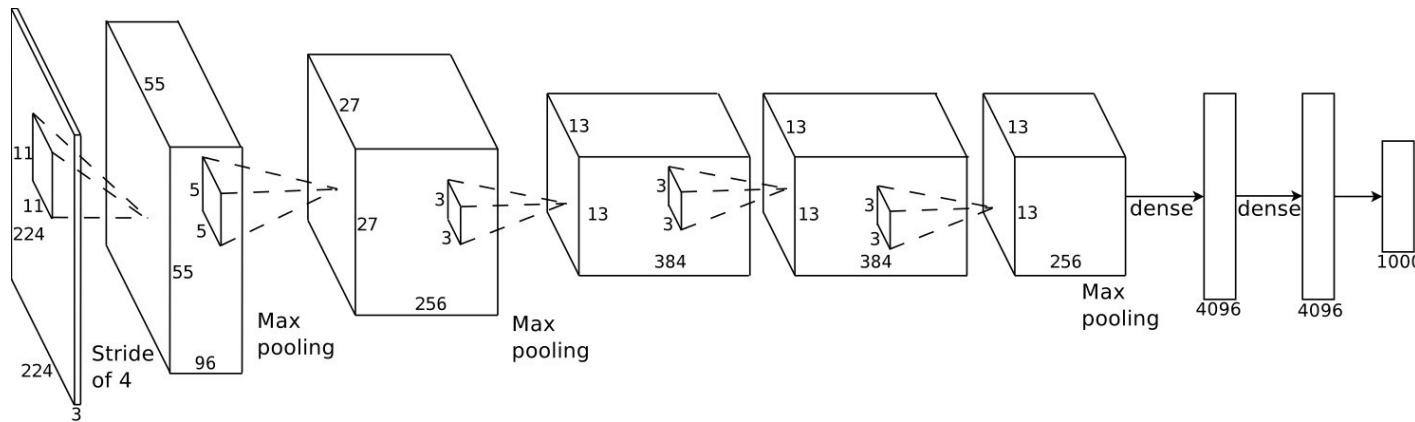
- Babenko et al. [1]
 - FC7 layer features (4096D)
 - Euclidean distance
 - Slightly outperform traditional SIFT baseline after fine-tuning
- Razavian et al. [2]
 - Extracts features from several sub windows (sliding window)
 - Excellent results, but computationally impractically heavy

[1] Babenko et al, Neural codes for image retrieval, 2014, CVPR

[2] Razavian et al., CNN features off-the-shelf: an astounding baseline for recognition, 2014, CVPR

Using off-the-self CNN representations

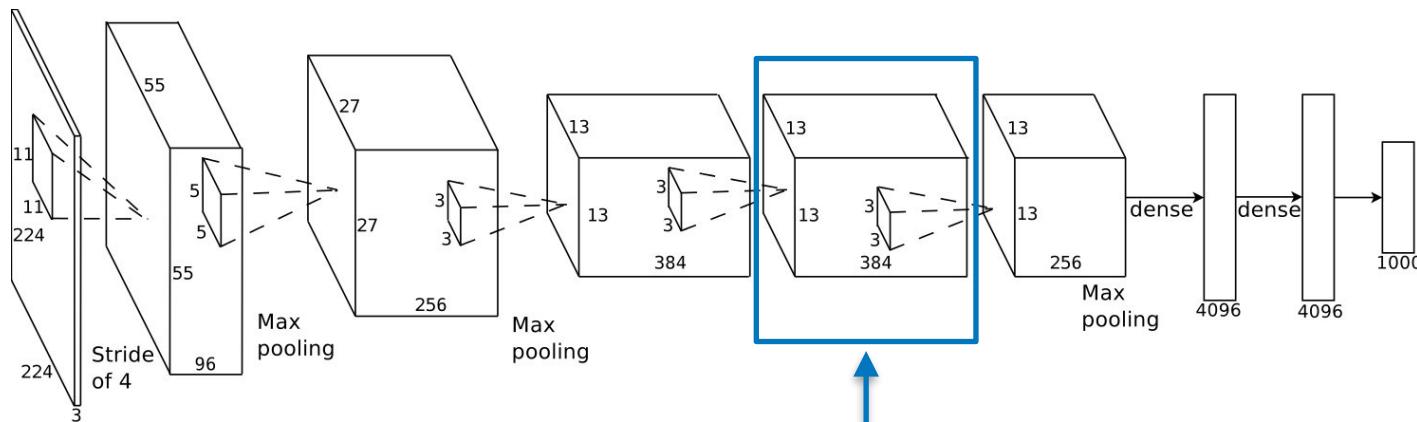
image



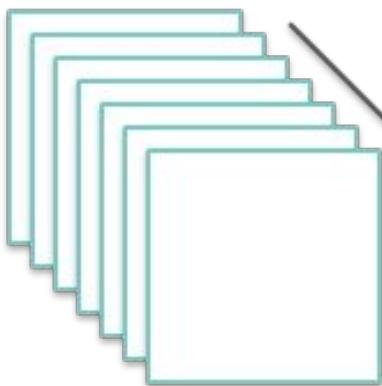
Could we obtain spatial information without explicit sliding window?

Using off-the-self CNN representations

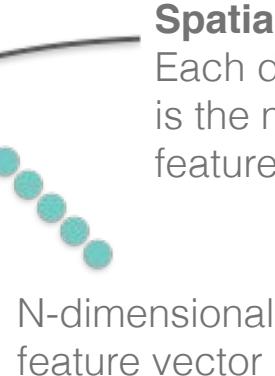
image



Obtaining descriptors from conv layers



N feature maps
dimensions H, W



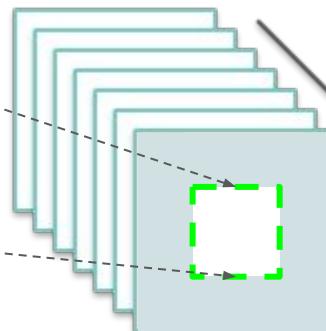
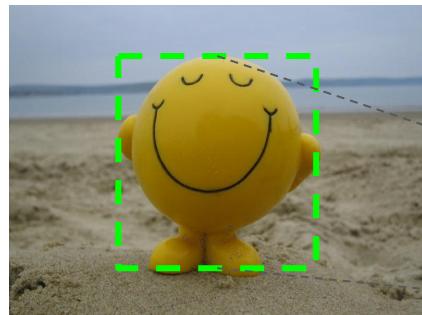
Spatial max-pooling:

Each dimension of the final vector is the max value of the corresponding feature map

N-dimensional
feature vector

Obtaining descriptors from conv layers

- Pooling features allow to describe specific parts of an image



N feature maps
Dimensions H, W

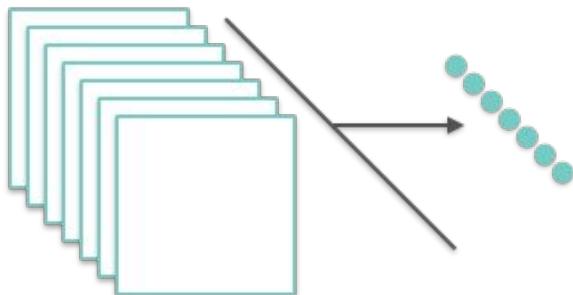
Spatial max-pooling:

Each dimension of the final vector is the max value of the corresponding feature map

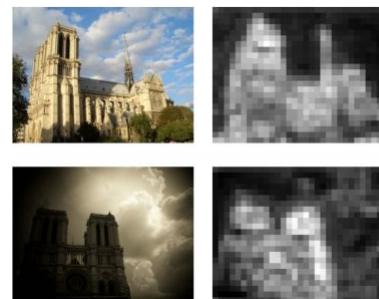
N -dimensional
feature vector

Obtaining descriptors from conv layers

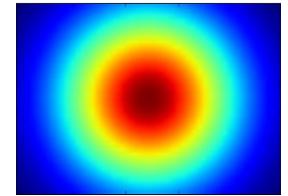
Sum/max pooling of a conv layer



Apply spatial weighting on the features before pooling them



Weighting based on
feature ‘strength’ [2]



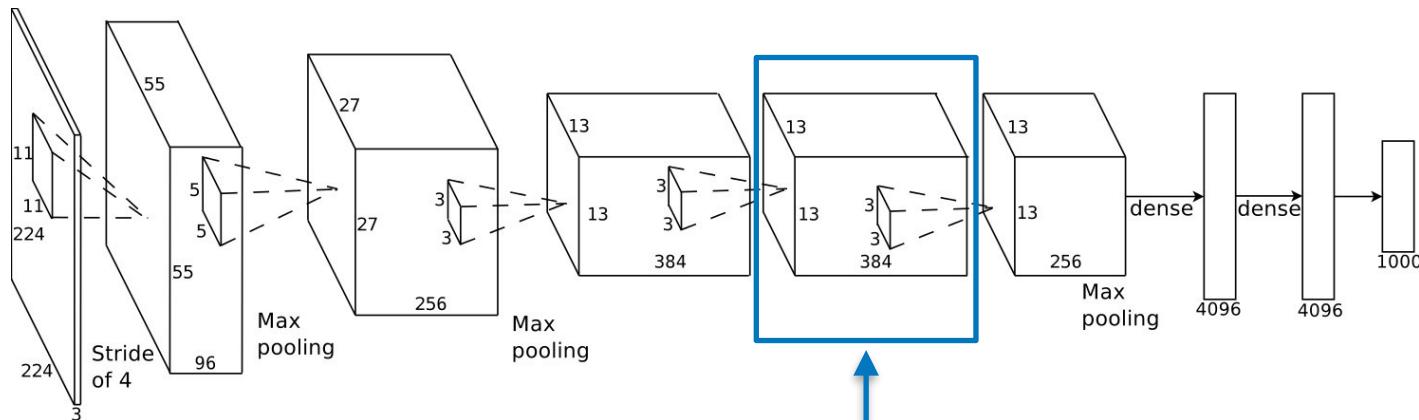
Weighting based on
distance to center [1]

[1] Babenko et al, Aggregating local deep features for image retrieval, 2015, ICCV

[2] Kalantidis et al., Cross-dimensional weighting for aggregated deep convolutional features, 2016, ECCV

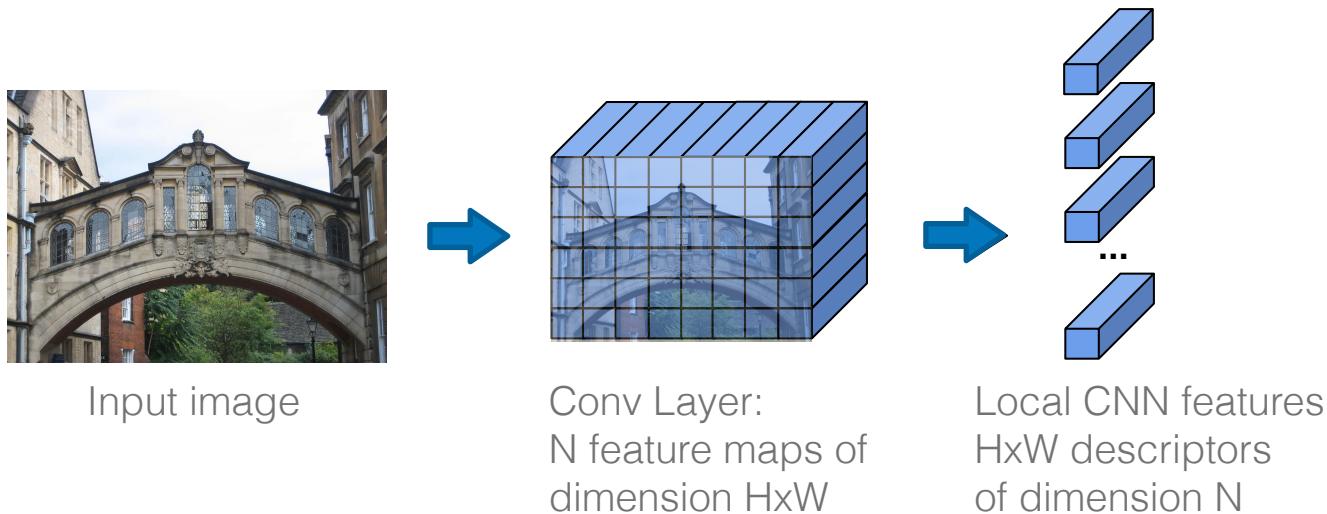
Obtaining descriptors from conv layers

image



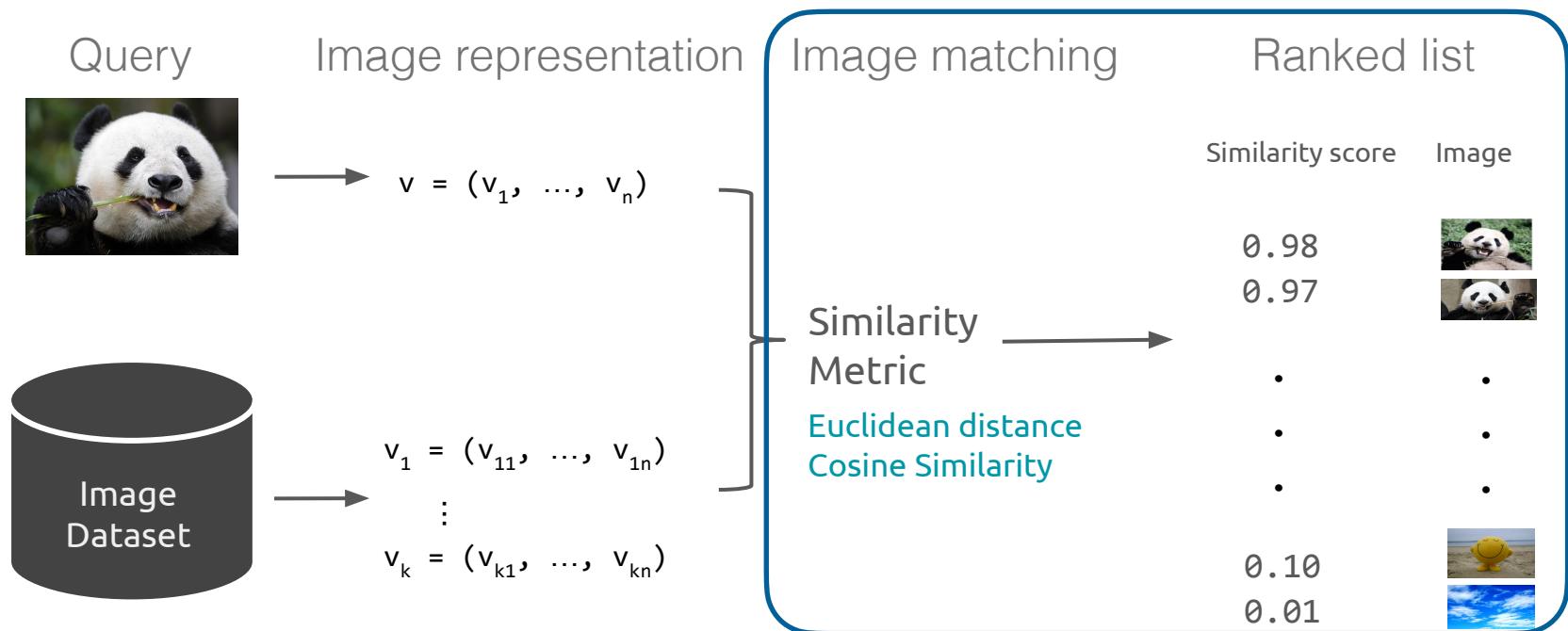
Treating as local
feature representations

Obtaining descriptors from conv layers



Obtaining similarity and ranking

The retrieval pipeline



Calculating similarities

- Euclidean distance or cosine similarity between feature vectors

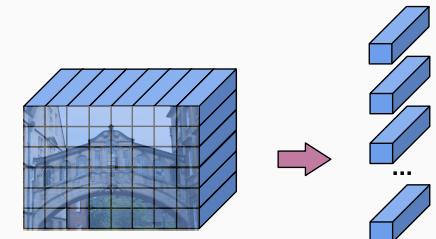
$$d_{L2}(x, a) = \sqrt{\sum_{i=1}^n (x_i - a_i)^2} \quad \text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

Calculating similarities

- Euclidean distance or cosine similarity between feature vectors

$$d_{L2}(x, a) = \sqrt{\sum_{i=1}^n (x_i - a_i)^2} \quad \text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

- Exhaustive evaluation unfeasible for large datasets
- What if we have multiple descriptors per image?



Idea: Inverted file index

- For text documents, an efficient way to find all **pages** that contain a specific **word** is to use an index
- We want to find all **images** in which a **feature** occurs
- To apply the idea in practice, we need to map the features to “**visual words**”

Index

A

Abel, Micah 18
Adkins, Melinda 40, 55
Adnrick, Kayla 20
Akers, William 8, 51
Algebra, Advanced 34
Allen, Shyla 8
Allshouse, Damon 32
Anderson, Cassie 38, 59
Andrick, Kayla 20
Asterino, Thomas 66
Audia, Sidney 16

B

Baker, Catherine 36
Baker, Savannah 44, 69
Bartimus, Jerad 40
Basketball, Junior Varsity Boys' 50, 51
Basketball, Junior Varsity Girls' 50
Basketball, Varsity Boys' 48, 49, 52, 53
Basketball, Varsity Girls' 52
Baylor, Nicholas 14
Bell, Charles 22, 23
Benedum, Anastasia 16, 56
Bennett, Brittany 34, 50, 57
Bennett, Coach Greg 50
Bennett, Dylan 6

Brooks, Orry 34, 51
Brooks, Susan 46
Brown, Benjamin 14
Brown, Delante 6
Brown, Nikki 42
Brumage, Brett 30
Brumage, Brittany 28
Brummage, Lindsey 18
Burks, Dylan 10
Burns, Tamika 10
Burton, Cherise 40, 52, 58, 67
Burton, Coach Ed 52
Burton, Julia 57
Burton, Linda 60
Burton, Lucille 18, 46
Butcher, Jack 34, 35
Butler, Alexis 14
Byrne, Devin 24

C

Canfield, Marisa 32
Carnes, Sara 14, 55
Carpenter, Phoenix 8
Centeno, Ashley 26
Cheerleaders, Junior Varsity 54, 55
Cheerleaders, Varsity 55
Clay, Jamie 6
Clevenger, Clayton 30
Club, Girls' 56, 57
Clutter, Jacob 20, 21
Cogar, Tyler 38
Cole, Anthony 10
Cole, Marcus 12

Denoon, Christian 30
Dent, Ashton 44
Dent, Bailee 12
des Mond, Adrian 68
Devalt, Jack 36, 37, 51
di Rosa, Rene 68
Dick, Jacob 8
Dick, Victoria 16
Dixon, Bethany 26
Donini, Jasmine 28, 29
Du Toit, Kim 68
Duckworth, Caleb 6
Duckworth, Mariah 10
Duckworth, Melysa 36
Duckworth, Trenton 6
Dukich, Mikayla 32, 44
Duskey, David 26

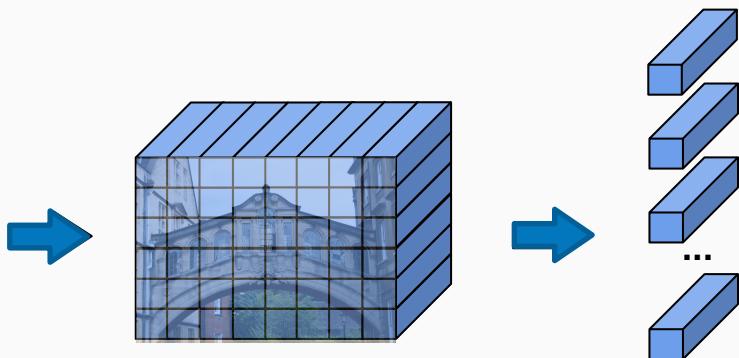
E

Echols, LC 34, 51
Eckles, Carson 30
Eddy, Brady 14
Efaw, Adrianna 26
Efaw, Breanna 36, 50, 54
Efaw, Coach Susie 54
Efaw, Danielle 14
Efaw, Hannah 42, 55, 57, 69
Elliott, Thomas 12
Eubank, Jacob 14, 15
Evan, Izaiha 28, 48, 49
Evans, Mason 32
Evanson, Jamie 42, 69
Evanson, Jenna 3, 5, 55, 69
Evanson, Jennifer 41

Obtaining visual words



Input image



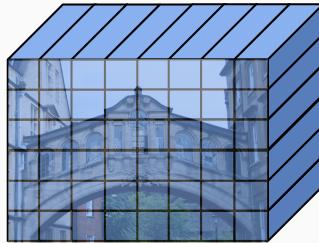
Conv Layer

Local CNN
descriptors of
dimension N

Obtaining visual words



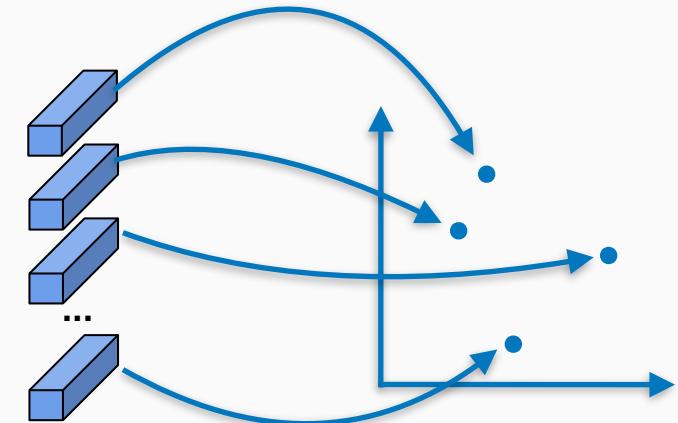
Input image



Conv Layer

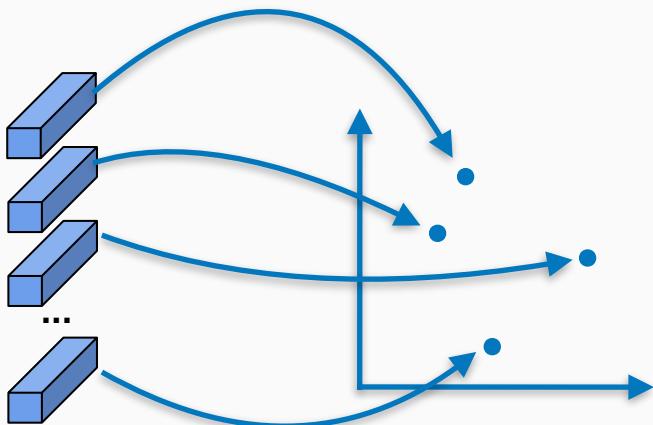


Local CNN
descriptors of
dimension N



Represent as points
in N dimensional
feature space

Obtaining visual words



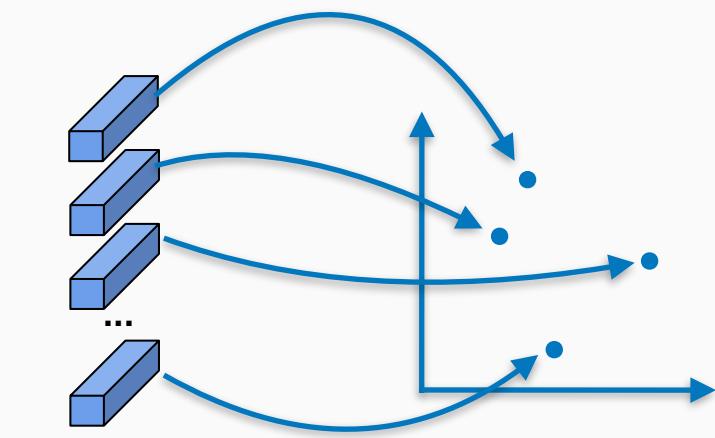
Local CNN
descriptors of
dimension N

Represent as points
in N dimensional
feature space



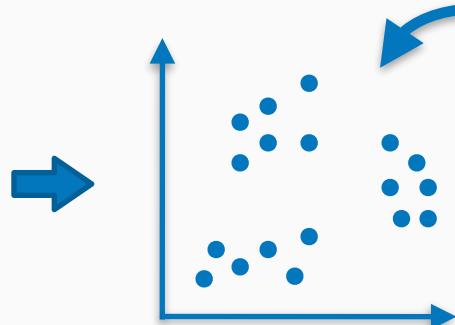
Training dataset

Obtaining visual words



Local CNN
descriptors of
dimension N

Represent as points
in N dimensional
feature space

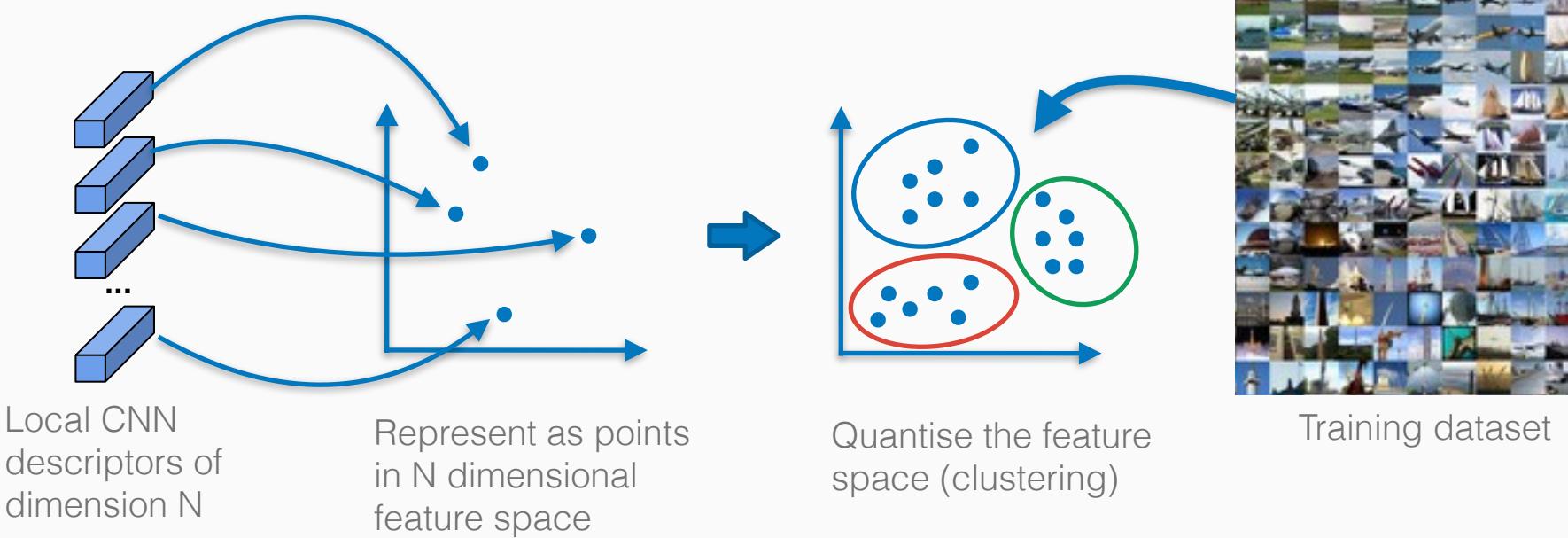


Collect descriptor from
the entire training set

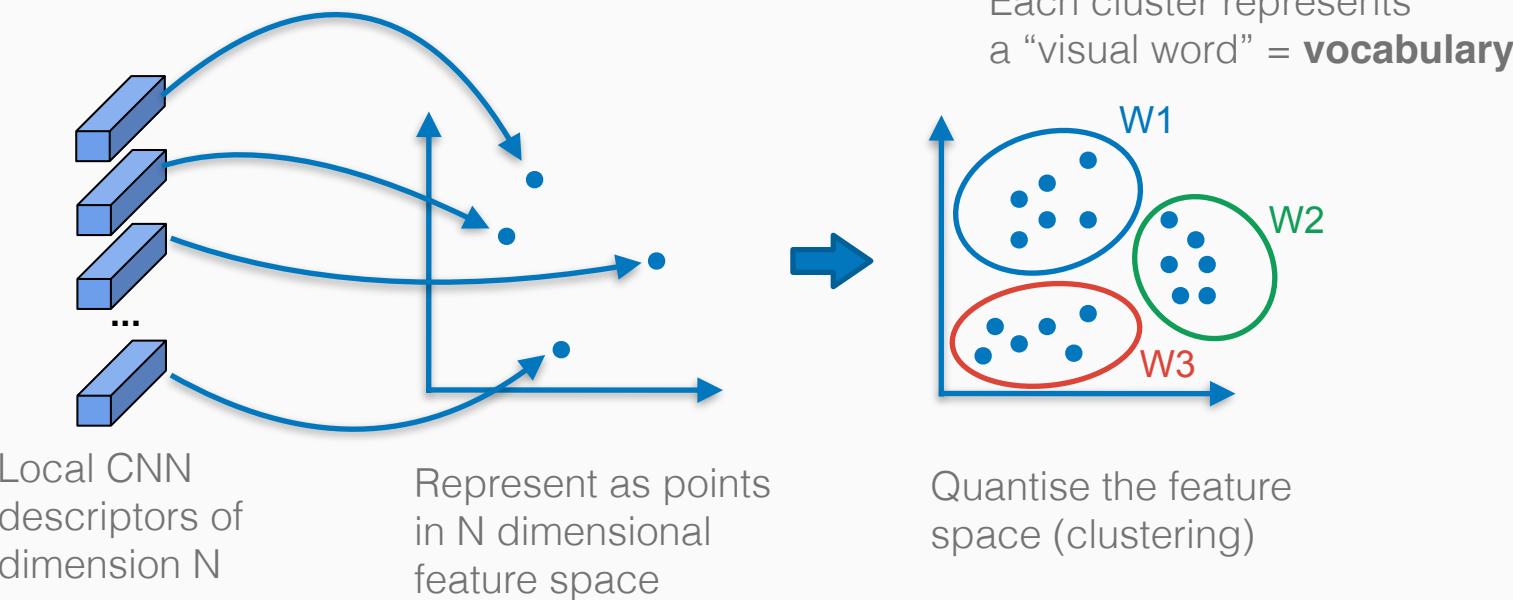


Training dataset

Obtaining visual words



Obtaining visual words



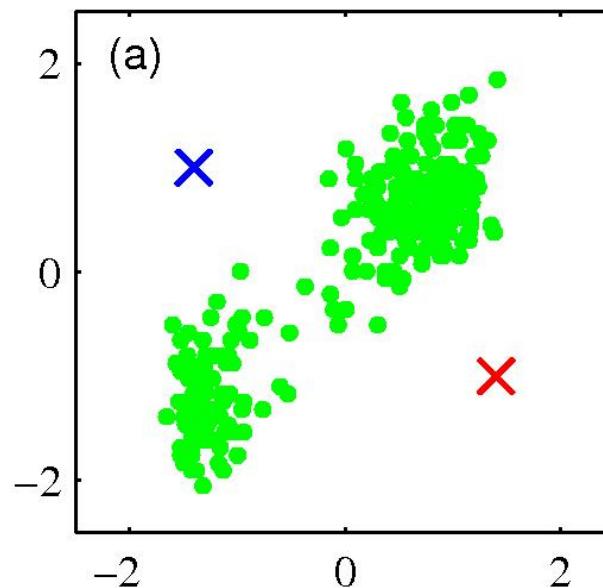
K-mean clustering

- Idea is to minimise sum of Euclidean distances between points x_i and their nearest cluster centres m_k

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\substack{\text{point } i \text{ in} \\ \text{cluster } k}} (\mathbf{x}_i - \mathbf{m}_k)^2$$

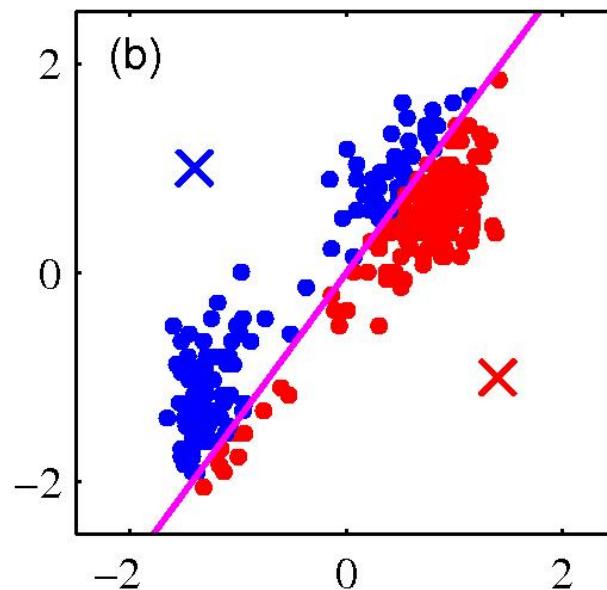
- Algorithm:
 1. Randomly initialise K cluster centres
 2. Assign all points to nearest cluster center
 3. Recompute cluster centres as the mean of all points assigned to it
 4. Return to 2 until convergence

K-means clustering: example



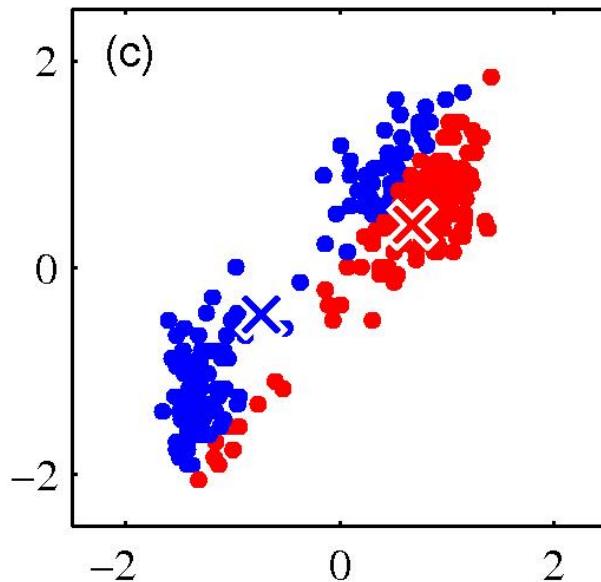
Initial state
2-dimensional space
 $K = 2$

K-means clustering: example



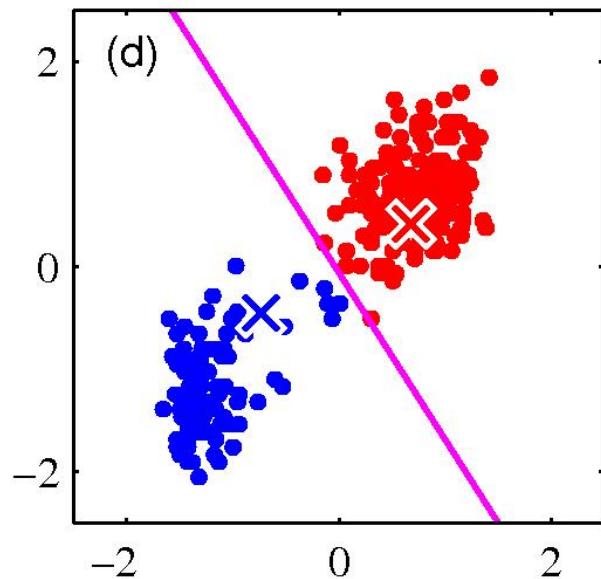
Result after step 3.
(assign data points to
nearest cluster centre)

K-means clustering: example



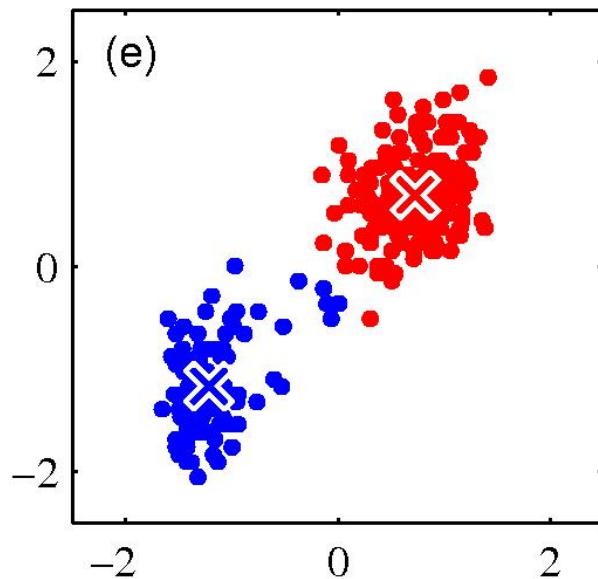
Result after step 4.
(move cluster centres
to the means of data points)

K-means clustering: example



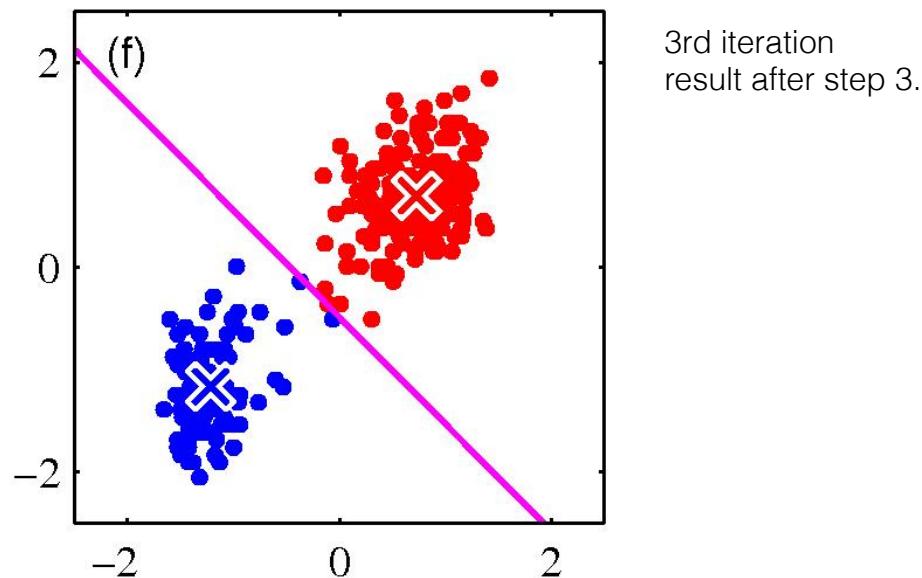
2nd iteration
result after step 3.

K-means clustering: example

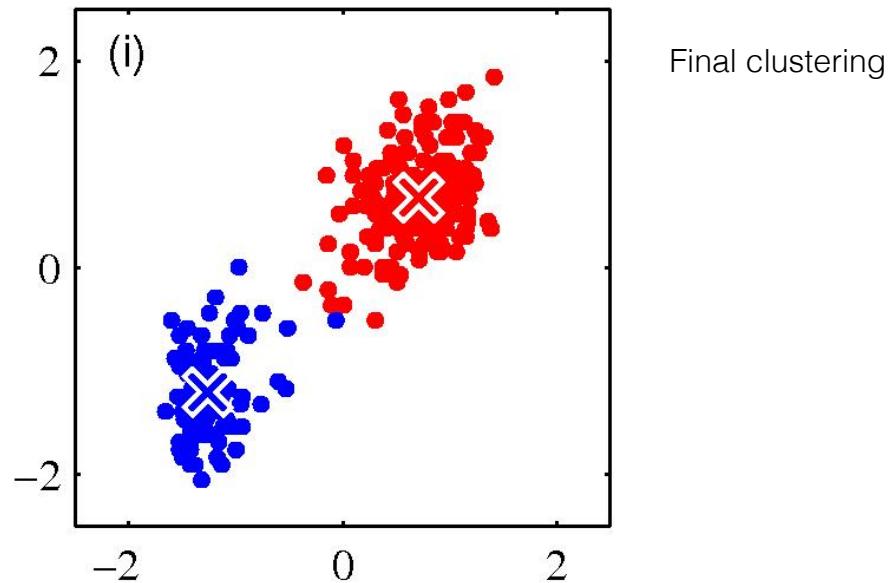


2nd iteration
result after step 4.

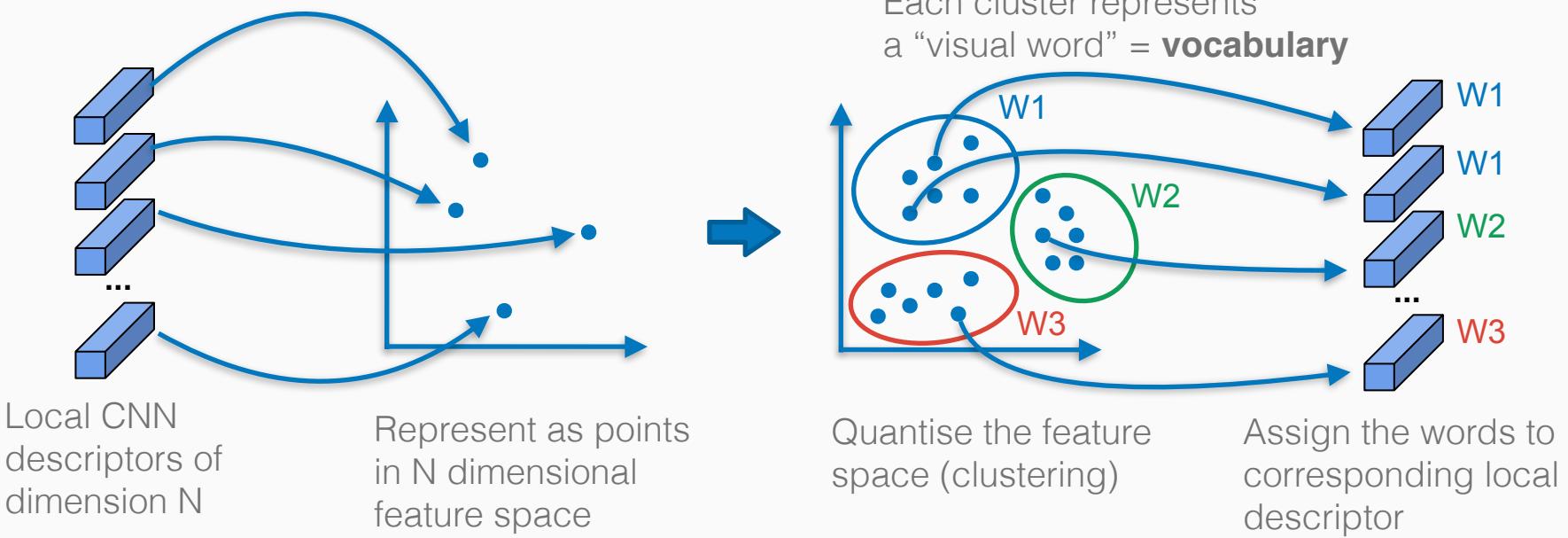
K-means clustering: example



K-means clustering: example

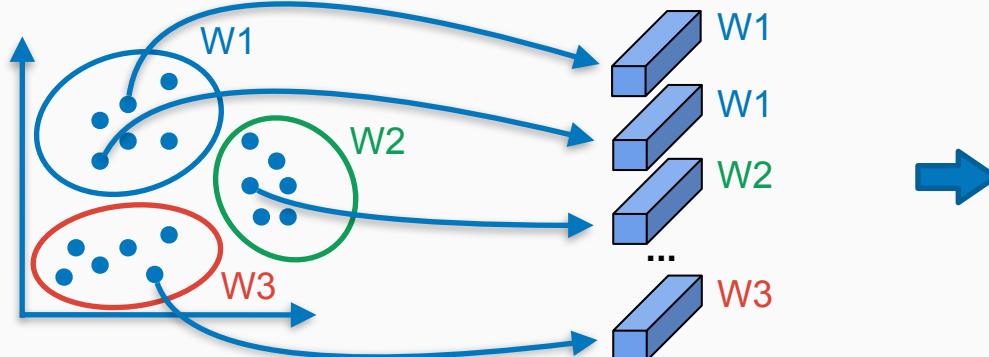


Obtaining visual words



Forming inverted index

Each cluster represents
a “visual word” = **vocabulary**



Quantise the feature
space (clustering)

Assign the words to
corresponding local
descriptor

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	
5	2,3
6	
7	3
8	
...	...

Visual words in retrieval (pipeline)

- Training phase
 - Compute local (CNN) descriptors from the **training** images
 - Form a **vocabulary** by clustering the local descriptors from the training set
 - Assign visual words to each **database image** using vocabulary and form an index
- Query phase
 - At query time, obtain the local descriptors and visual words for the **query image**
 - **Search** matching images from the index

Example

Database images



Example

Database images



→ W1 W1 W3 ...



→ W2 W5 W1 ...



→ W7 W5 W5 ...

Example

Database images



W1 W1 W3 ...



W2 W5 W1 ...



W7 W5 W5 ...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	
5	2,3
6	
7	3
8	
...	...

Example

New query image



→ W5 ...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	
5	2,3
6	
7	3
8	
...	...

Example

New query image



W5 ...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	
5	2,3
6	
7	3
8	
...	...

Example

New query image



W5

...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	
5	2,3
6	
7	3
8	
...	...



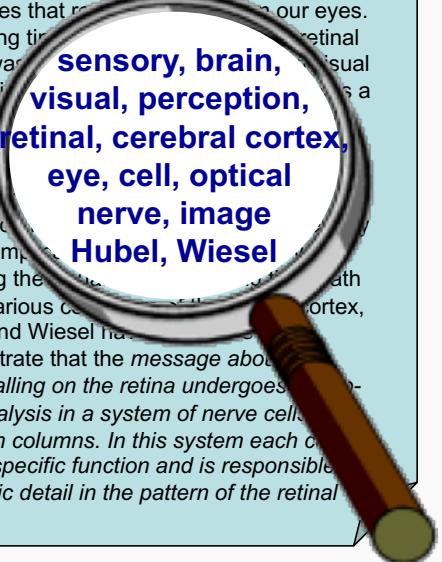
Inverted file index

- Inverted file index is efficient only if it is **sparse**
- If most pages/images contain most words, then indexing gives no advantage over exhaustive comparison to each dataset image
- Problem solved?
- But how to summarise and compare the content of an entire image?

Comparing entire images with visual words

Analogy to documents

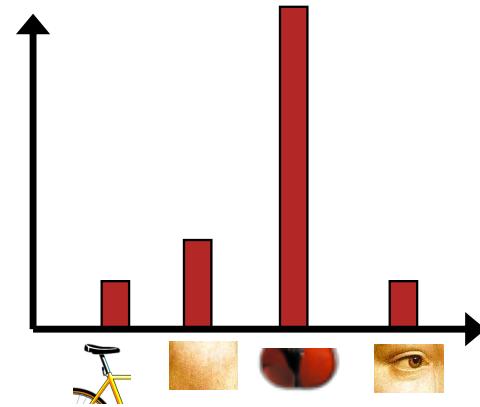
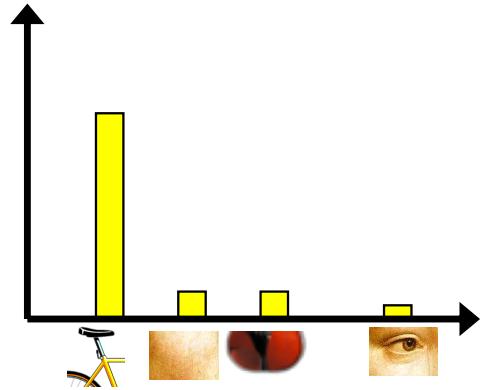
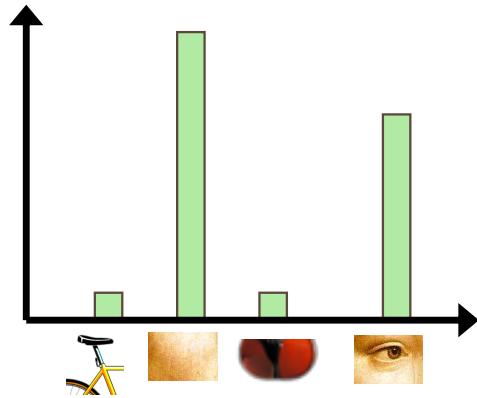
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach us through our eyes. For a long time it was believed that the retinal image was processed by the visual centers in the cerebral cortex, just as a movie screen receives an image. In 1960, Hubel and Wiesel discovered that the visual perception is more complex than that. Following the path to the various centers in the cerebral cortex, Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a two-stage analysis in a system of nerve cells stored in columns. In this system each column has its specific function and is responsible for a specific detail in the pattern of the retinal image.



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 predicted 30% increase in exports to \$750bn, compared with \$660bn. That would annoy the US, which is already annoyed by China's trade policies. The US has deliberately agreed to let the Chinese government manage the yuan is. The Chinese government also needs to encourage foreign investment so that the country can develop. China has been trying to control the yuan against the dollar since 1994 and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to float freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

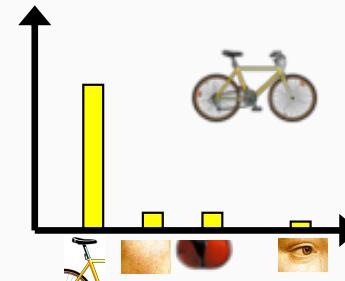
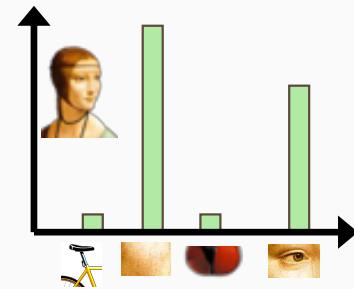


Bags of visual words



Bags of visual words

- Summarise the image using the histogram of word occurrences
- Analogous to bag of words representation used for documents

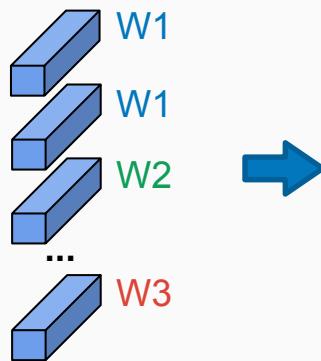


Bags of visual words

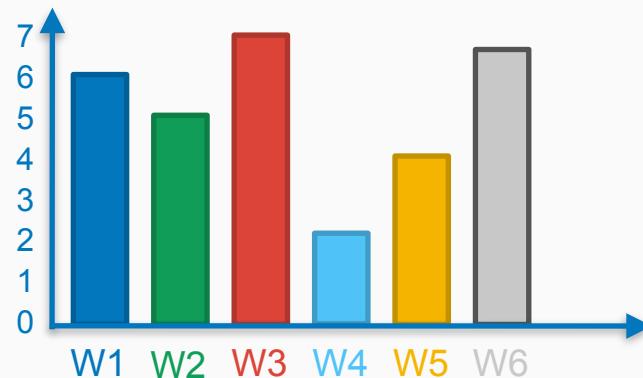


Quantised local
(CNN) descriptors

Bags of visual words

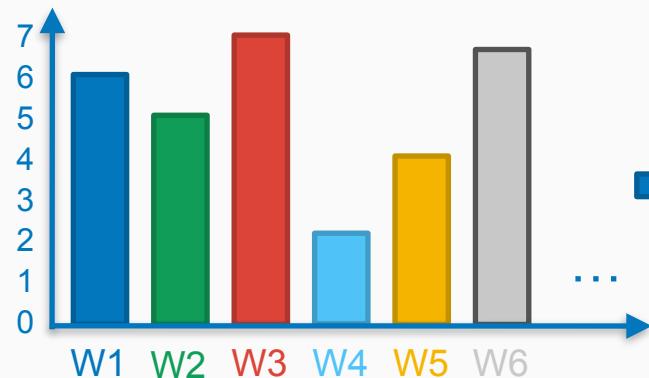
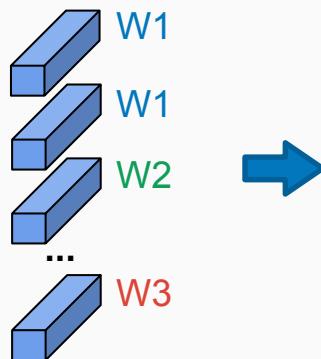


Quantised local
(CNN) descriptors



Count the occurrences of
each word in the image

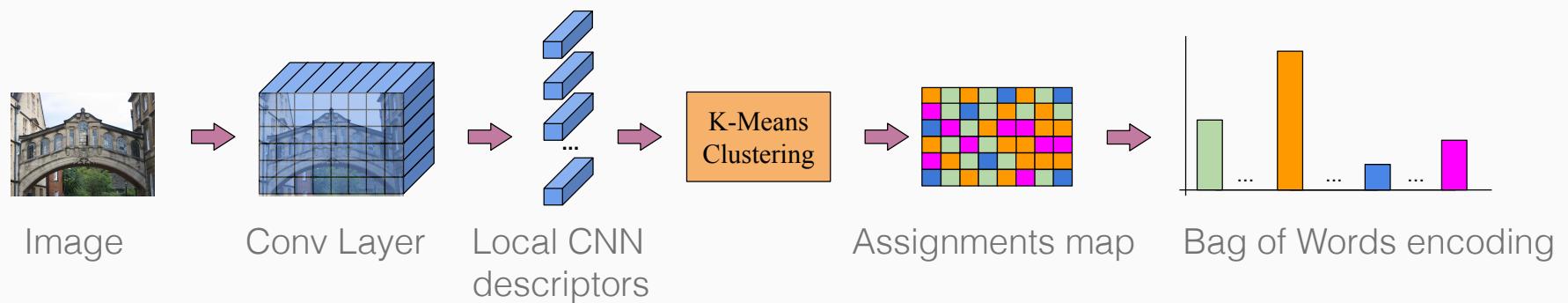
Bags of visual words



[6, 5, 7, 2, 4, 6, ...]

Quantised descriptor for the entire image content

Bags of local convolutional features pipeline



Comparing bags of words

- Measure similarity using normalised scalar product between the bags of words descriptor vectors
- Rank dataset images based on similarity to the query image

$$sim(d_j, q) = \frac{\langle d_j, q \rangle}{\|d_j\| \|q\|}$$

Inverted file index and bags of words similarity

1. Extract words from the query image



Query image

W5,W9,...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	8
5	2,3
6	
7	3
8	
9	8,3
...	...

Inverted file index and bags of words similarity

1. Extract words from the query image
2. Find relevant database images from inverted file index



Query image

W5,W9,...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	8
5	2,3
6	
7	3
8	
9	8,3
...	...



Inverted file index and bags of words similarity

1. Extract words from the query image
2. Find relevant database images from inverted file index
3. Compare word counts



Query image

W5,W9,...

Inverted index file

Word #	Image #
1	1,2
2	2
3	1
4	8
5	2,3
6	
7	3
8	
9	8,3
...	...

Three database images are shown with their corresponding word count vectors:

- The first image has vector [1, 2, 0, 0, 9, 0, ...].
- The second image has vector [0, 0, 0, 0, 1, 0, ...].
- The third image has vector [0, 0, 0, 7, 0, 0, ...].

How to evaluate retrieval results?

How to evaluate retrieval quality?



Query

Dataset size: 10 images
Relevant (total): 5 images

How to evaluate retrieval quality?



Query

Dataset size: 10 images
Relevant (total): 5 images



Results (ordered)

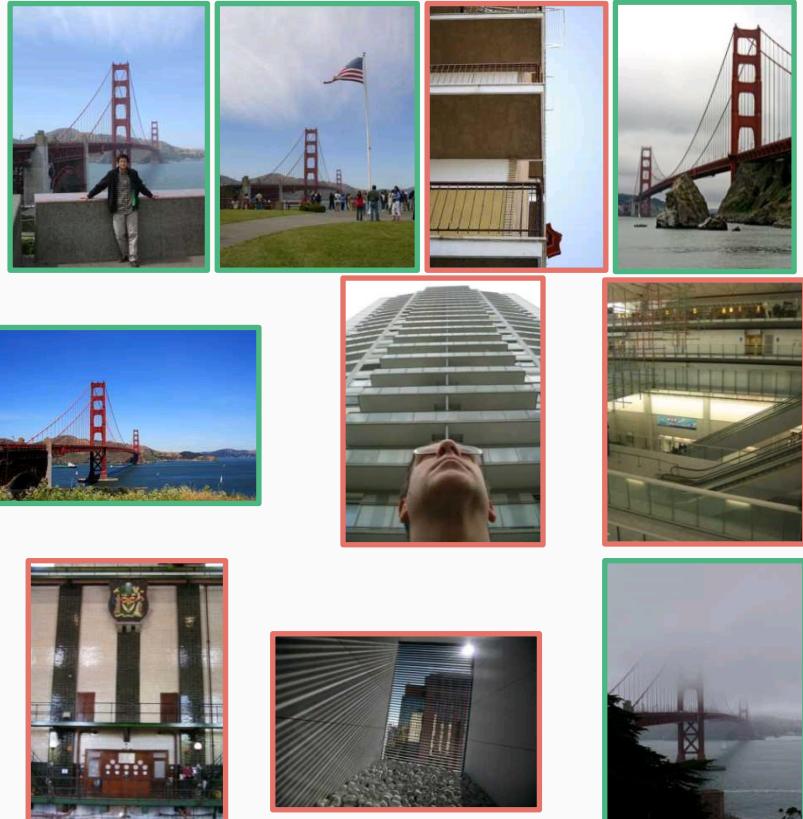
Credit: O. Chum

How to evaluate retrieval quality?



Query

Dataset size: 10 images
Relevant (total): 5 images



Results (ordered)

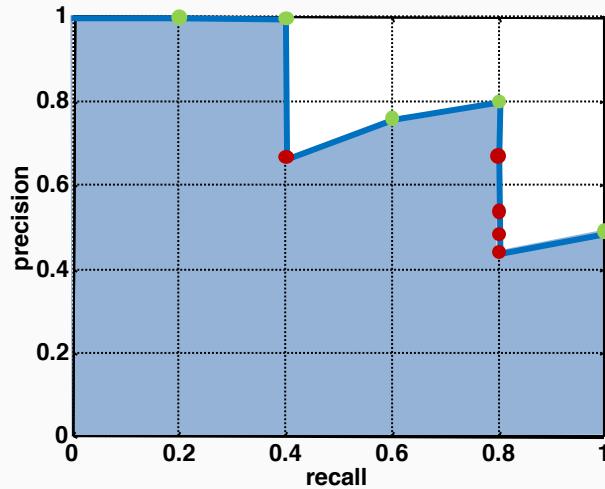
Credit: O. Chum

How to evaluate retrieval quality?

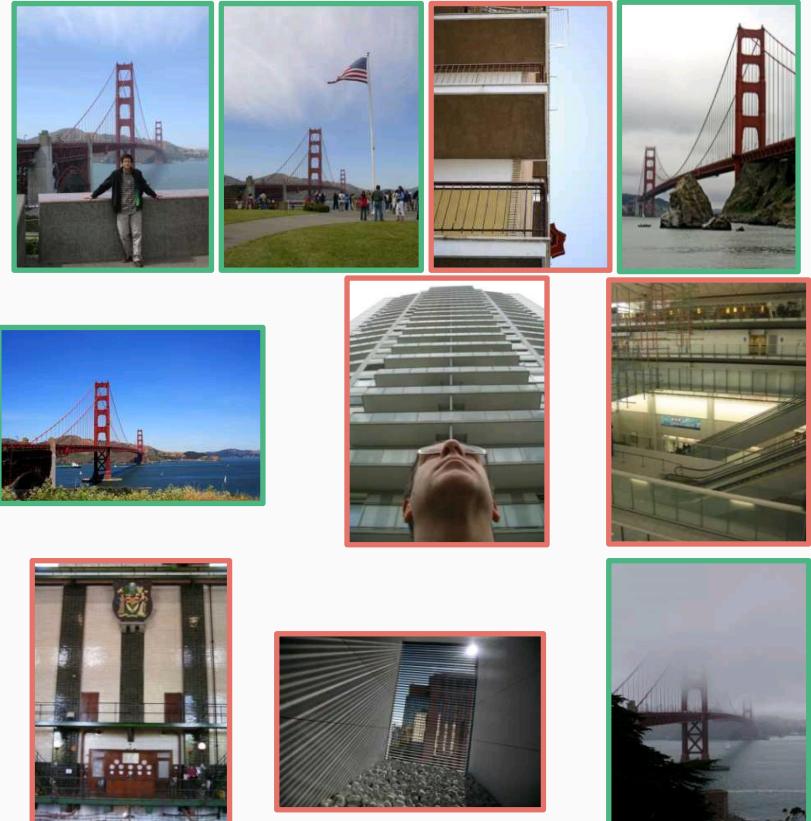


Dataset size: 10 images
Relevant (total): 5 images

Query



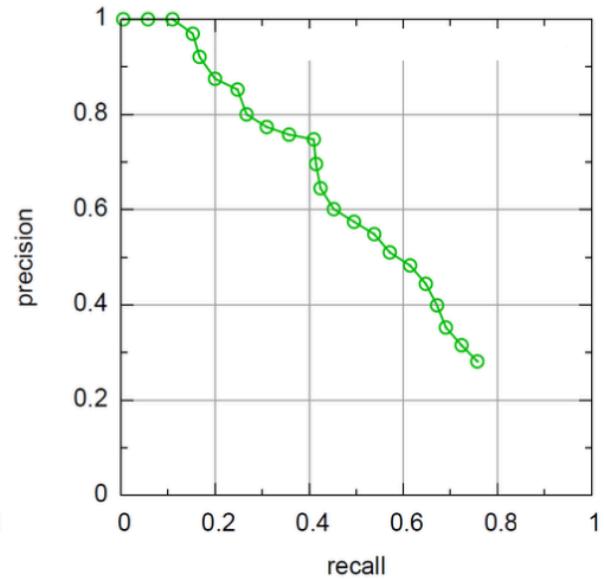
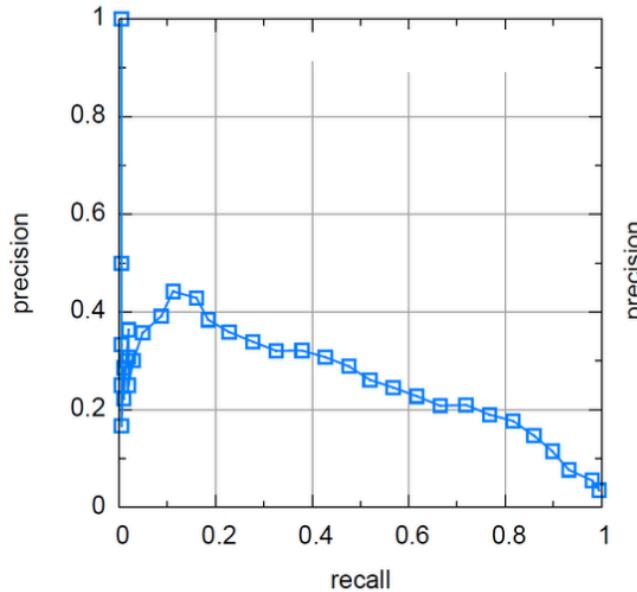
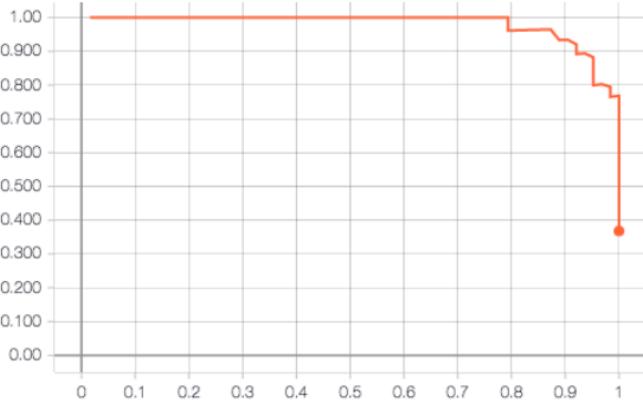
Precision = #relevant / #returned
Recall = #relevant / #total relevant



Results (ordered)

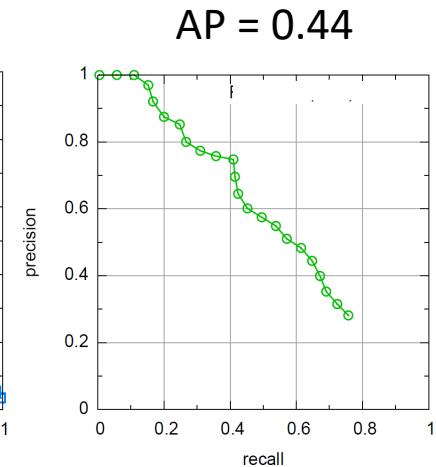
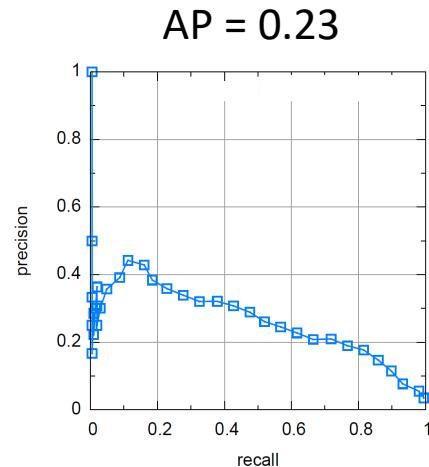
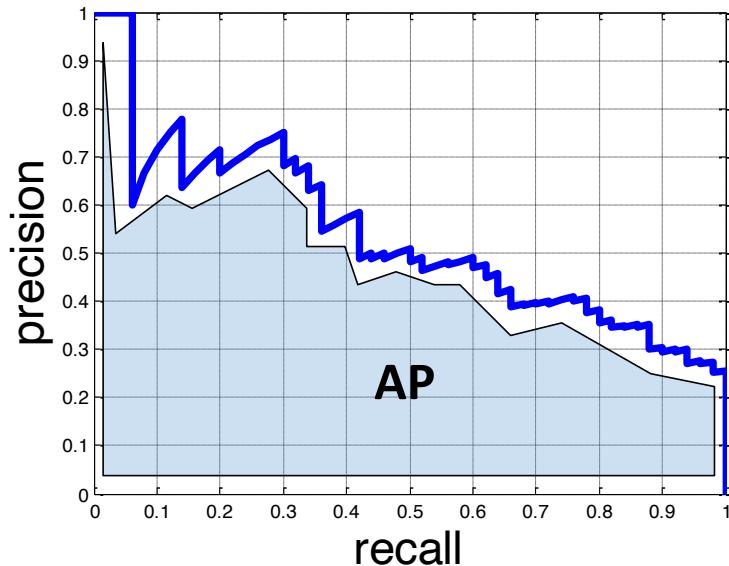
Credit: O. Chum

A pot-pourri of PR curves



Average precision

- Area under Precision-Recall curve
- Single score to assess performance



A good score requires both high recall and high precision

Credit: A. Zisserman

Things to remember

- Two phases of image retrieval:
 - Obtain image descriptors for database and query image
 - Calculate similarity metrics and rank the database images according to similarity with query
- (Pre-trained) CNNs can form powerful image descriptors for retrieval
- Descriptors can be quantised to “visual words” for efficient search
 - Learn a vocabulary from a training set by clustering
 - Summarise image by distribution of words (bag of words)
- Precompute index to enable fast search at query time (inverted index)
- Retrieval quality can be evaluated using precision vs recall plot

That's it folks!