## Slide 1

*Vision par Ordinateur :*
*Indexation et recherche d'images*

CNN-based descriptors

---

## Slide 2

*Milestones of instance retrieval*



Video Google
*Sivic and Zisserman*

Hierarchical K-Means
*Stewénius and Nistér*

Approximate K-Means
*Philbin et al.*

Improved FV
*Perronnin et al.*

VLAD
*Jégou et al.*

CNN off-the-shelf
*Razavian et al.*

R-MAC
*Tolias et al.*

"The end of the early years"
*Smeulders et al.*

Hamming Embedding
*Jégou et al.*

CNN for ImageNet
*Krizhevsky et al.*

Neural codes
*Babenko et al.*

VLAD-CNN
*Ng et al.*

SIFT-based

CNN-based

« SIFT Meets CNN: A Decade Survey of Instance Retrieval »  Liang Zheng, Yi Yang, and Qi Tian, 2015

2

---

## Slide 3

*General pipeline of SIFT- and CNN-based retrieval models*



3

---

## Slide 4

*AlexNet*



- 3-channel RGB input, 224 x 224
- 8 layers : 5 conv + 3 fc
- ReLU follows each convolutional and fully connected layer
- Data augmentation, dropout
- Stochastic gradient descent with momentum

4

---

## Slide 5

*Pre-trained CNNs*

- AlexNet pre-trained on ImageNet for classification
- last fully connected layer (fc6): global descriptor of dimension
  $$k = 4096$$
- nearest neighbors in ImageNet according to Euclidean distance



Krizhevsky, Sutskever, Hinton. NIPS 2012. Imagenet Classification with Deep Convolutional Neural Networks.

5

---

## Slide 6

Neural codes for image retrieval

- fine-tuning by softmax on 672 classes of 200k landmark photos
- outperforms VLAD and Fisher vectors on standard retrieval benchmarks, but still inferior to SIFT local descriptors



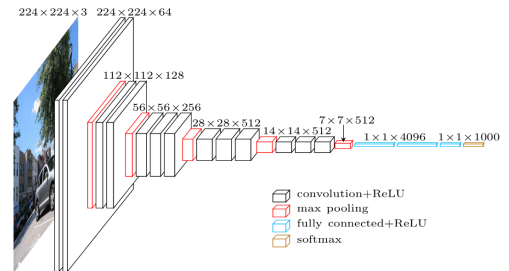Babenko, Slesarev, Chigorin, Lempitsky. ECCV 2014. Neural codes for Image Retrieval.

6

## CNN features Off-the-shelf

- For each image, extract multiple sub-patches of different sizes at different locations

- For each extracted sub-patch, its CNN representation is the L2 normalized output of the first fully connected layer (dim=4096)

- PCA dimensionality reduction → whitening →L2 renormalization (500-D)

A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "Cnn features off-the-shelf: an astounding baseline for recognition," in CVPR Workshops, 2014

## VGG-16



$224 \times 224 \times 3$  $224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$  $14 \times 14 \times 512$  $7 \times 7 \times 512$

$1 \times 1 \times 4096$  $1 \times 1 \times 1000$

- convolution+ReLU
- max pooling
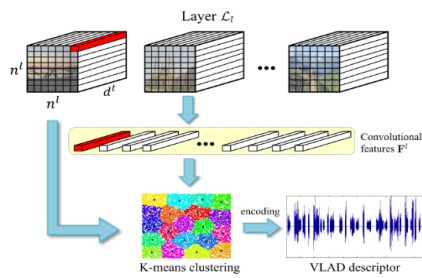- fully connected+ReLU
- softmax

- Depth increased up to 19 layers,
- Kernel sizes reduced to 3, strides to 1

Simonyan and Zisserman 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition.

## VLAD-CNN



Layer $\mathcal{L}_t$

Convolutional features $\mathbf{F}^l$

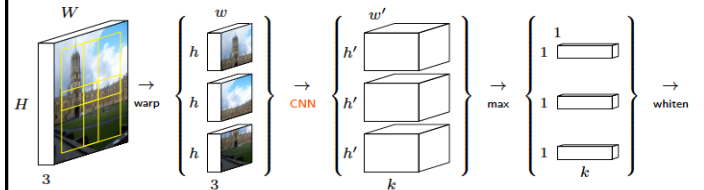K-means clustering          encoding          VLAD descriptor

- Consider different layers of VGG-16
- For each layer, VLAD encoding (k=100)
- L2-normalization, PCS-whitening (128-D)

J. Ng, F. Yang, and L. Davis, "Exploiting local features from deep networks for image retrieval," CVPR Workshops, 2015.
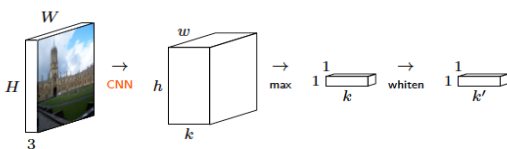
## Regional CNN features



$\rightarrow$ warp   $\rightarrow$ CNN   $\rightarrow$ max   $\rightarrow$ whiten

- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, warped into w x h = 227 x 227
- each region yields a w' x h' x k = 36 x 36 x 256 dimensional feature at the last convolutional layer of AlexNet
- global spatial max-pooling
- L2-normalization, PCA-whitening of each descriptor

Razavian, Sullivan, Maki and Carlsson 2015. Visual Instance Retrieval with Deep Convolutional Networks.

## Global max-pooling (MAC)



$\rightarrow$ CNN   $\rightarrow$ max   $\rightarrow$ whiten

- VGG-16 last convolutional layer, k = 512
- global spatial max-pooling
- L2-normalization, PCA-whitening, L2-normalization
- MAC: maximum activation of convolutions

Tolias, Sicre and Jegou. ICLR 2016. Particular Object Retrieval with Integral Max-Pooling of CNN Activations.
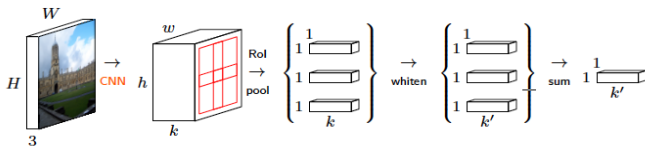
## Global max-pooling (MAC)



- receptive fields of 5 components of MAC vectors that contribute most to image similarity
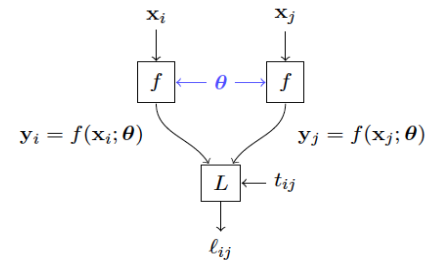
## Regional max-pooling (R-MAC)



- VGG-16 last convolutional layer, k = 512
- fixed multiscale overlapping regions, spatial max-pooling
- L2-normalization, PCA-whitening, L2-normalization
- sum-pooling over all descriptors, L2-normalization

Tolias, Sicre and Jegou. ICLR 2016. Particular Object Retrieval with Integral Max-Pooling of CNN Activations.
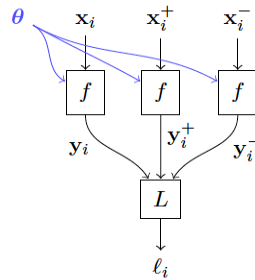
## Siamese architecture



$$\mathbf{y}_i = f(\mathbf{x}_i; \boldsymbol{\theta}) \qquad \mathbf{y}_j = f(\mathbf{x}_j; \boldsymbol{\theta})$$

- an input sample is a pair $(x_i, x_j)$
- both $x_i$, $x_j$ go through the same function f with shared parameters θ
- Contrastive loss $l_{ij}$ is measured on output pair $(y_i, y_j)$ and target $t_{ij}$

Chopra, Hadsell, Lecun, CVPR 2005. Learning a Similarity Metric Discriminatively, with Application to Face Verification.

## Triplet architecture

- an input sample is a triplet
  $(x_i, x^+_i, x^-_i)$

- $x_i$, $x^+_i$, $x^-_i$ go through the same function f with shared parameters

- loss $l_i$ measured on output triplet $(y_i, y^+_i, y^-_i)$



Wang, Song, Leung, Rosenberg, Wang, Philbin, Chen, Wu. CVPR 2014. Learning Fine-Grained Image Similarity with Deep Ranking.