# Deep Learning in Visual Recognition

Yadong Mu
Machine Intelligence Lab
Institute of Computer Science & Technology
Peking University

Email: myd@pku.edu.cn
Website: www.muyadong.com

# Outline

- **Visual Recognition: Task Definition and Challenges**

- **Bag-of-words Models**

- **Improving BoW Models**

- **Why DL suddenly works? (AlexNet, 2012)**

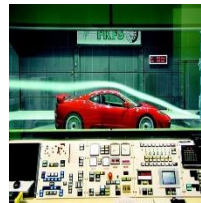- **Can it go deeper? (ResNet, 2015)**

# Visual Recognition / Image Classification

- Give a binary label to indicate whether an object is present

**Does this image contain a car?**



**cars**



**not cars**

# Image Classification

- 1<sup>st</sup> challenge: **semantic gap**

# Image Classification

- 1$^{st}$ challenge: **semantic gap**



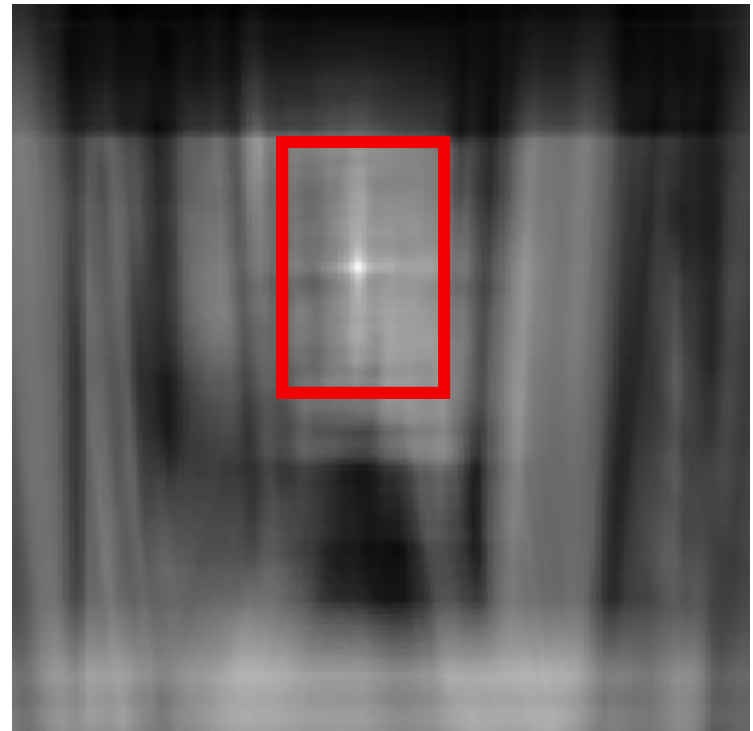**To the computer, images are essentially a 3-D array**

# Image Classification

- 1ˢᵗ challenge: **semantic gap**

This is a chair
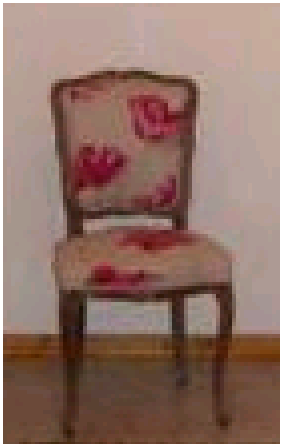
Find the chair in this image

Output of normalized correlation

# Image Classification

- Another example: adversarial training



$$x$$

$y =$"panda"
w/ $57.7\%$
confidence

$$+ .007 \times$$

$$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

"nematode"
w/ $8.2\%$
confidence

$$=$$

$$\boldsymbol{x} + \epsilon \, \text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
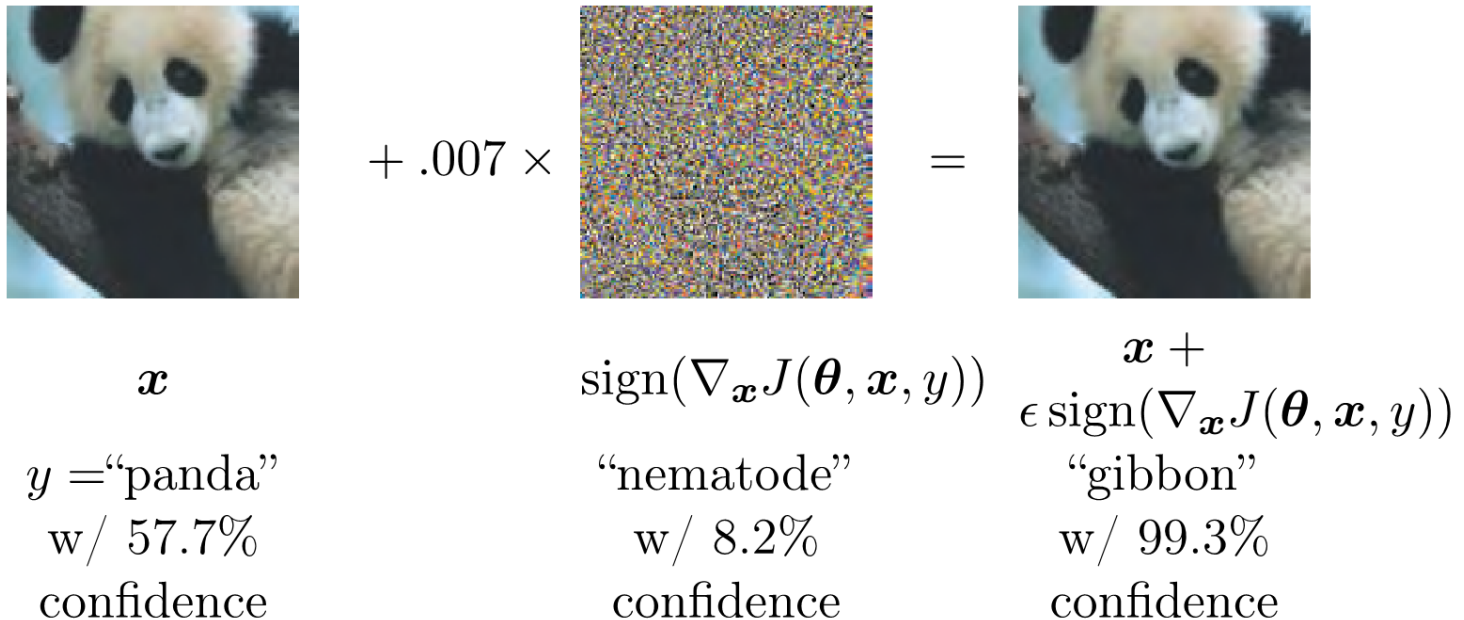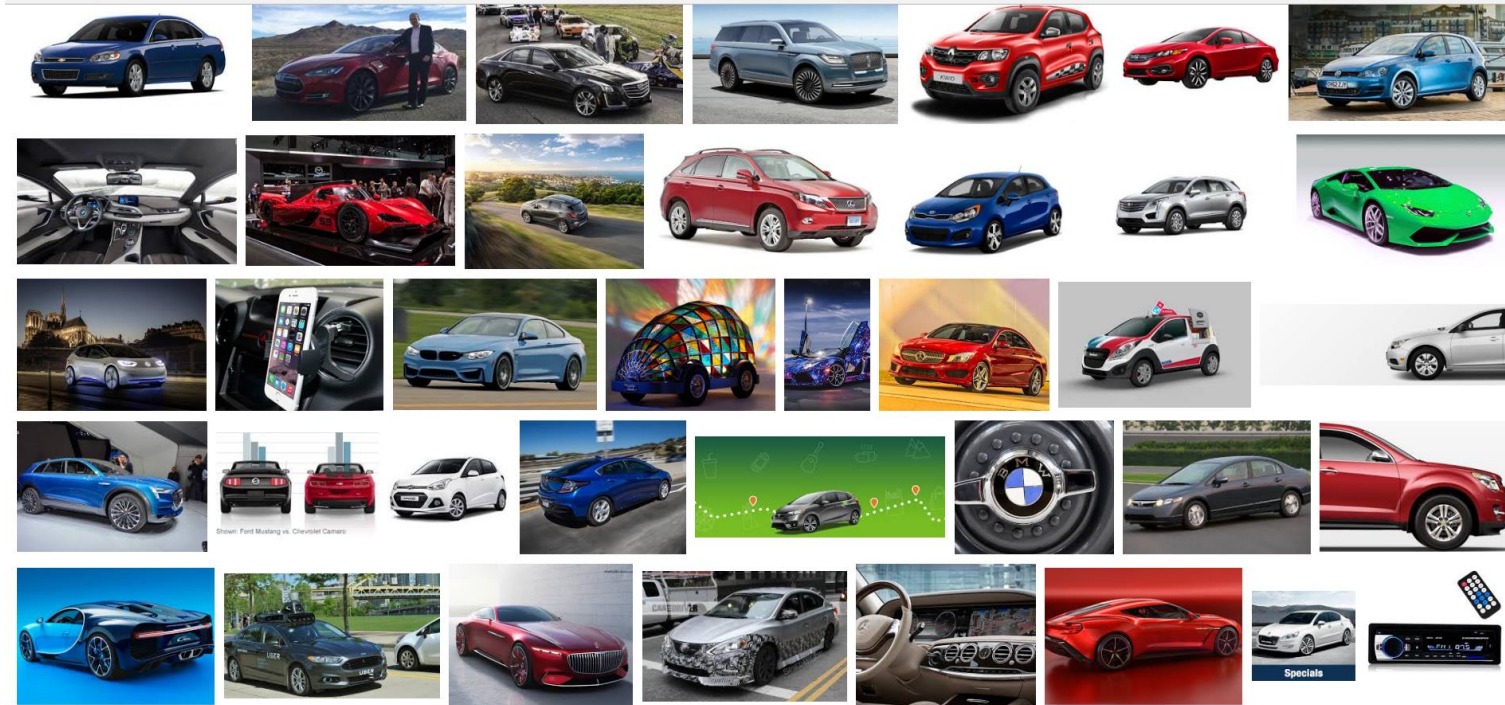
"gibbon"
w/ $99.3\%$
confidence

Figure 7.8: A demonstration of adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet's classification of the image. Reproduced with permission from Goodfellow et al. (2014b).

# Image Classification

- 2nd challenge: visual variations

# Image Classification

- 3rd challenge: multiple labels
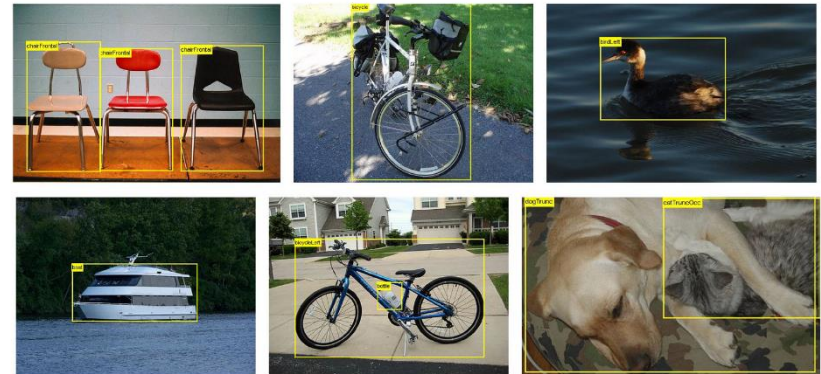
**What kind of objects do you see in this image?**

# Image Classification @ CV Community

- Pascal VOC Challenges 2005-2012

**Table 1** The VOC classes

| Vehicles | Household | Animals | Other |
|----------|-----------|---------|-------|
| Aeroplane | Bottle | Bird | Person |
| Bicycle | Chair | Cat | |
| Boat | Dining table | Cow | |
| Bus | Potted plant | Dog | |
| Car | Sofa | Horse | |
| Motorbike | TV/Monitor | Sheep | |
| Train | | | |

- Large Scale Visual Recognition Challenge (ILSVRC) 2010-2016

### 1000 synsets for Object classification/localization

kit fox, Vulpes macrotis

English setter

Australian terrier

grey whale, gray whale, devilfish, Eschrichtius gibbosus, Eschrichtius robustus

lesser panda, red panda, panda, bear cat, cat bear, Ailurus fulgens

Egyptian cat

ibex, Capra ibex

Persian cat

cougar, puma, catamount, mountain lion, painter, panther, Felis concolor

gazelle

porcupine, hedgehog

sea lion

**IMAGENET**

14,197,122 images, 21841 synsets indexed

Explore  Download  Challenges  Publications  CoolStuff  About

Not logged in. Login | Signup

ImageNet is an image database organized according to the WordNet hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. Currently we have an average of over five hundred images per node. We hope ImageNet will become a useful resource for researchers, educators, students and all of you who share our passion for pictures.
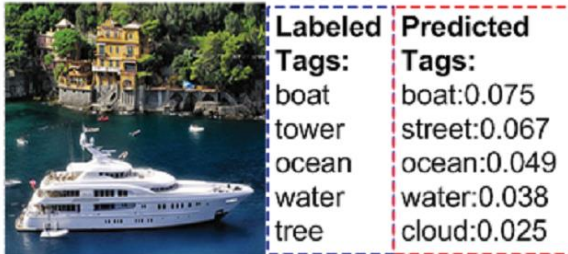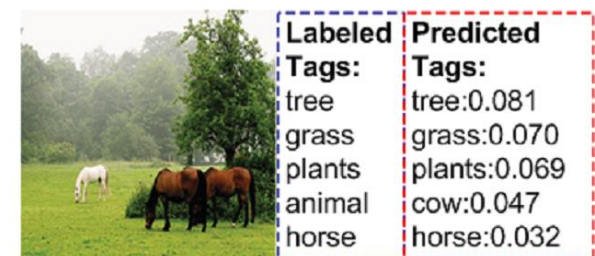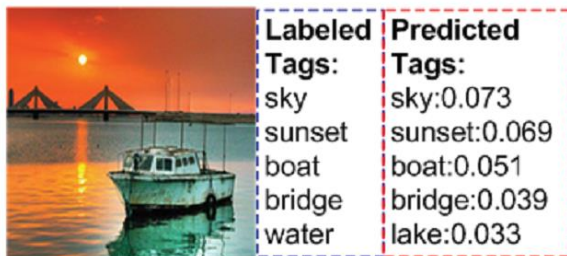Click here to learn more about ImageNet, Click here to join the ImageNet mailing list.

http://www.image-net.org/

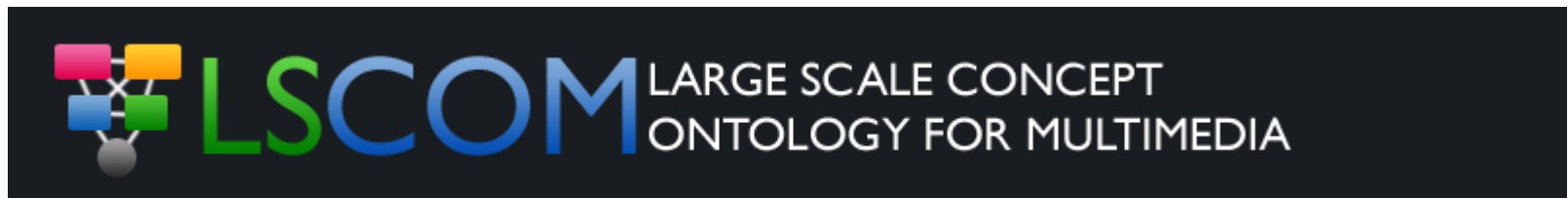http://www.image-net.org/challenges/LSVRC/2016/index

# Image Classification @ Multimedia Community

- Image tagging: YFCC100M



- Concept detection: TRECVID SIN task

**Ontology design -> Image collecting -> Image annotation -> Concept detection**
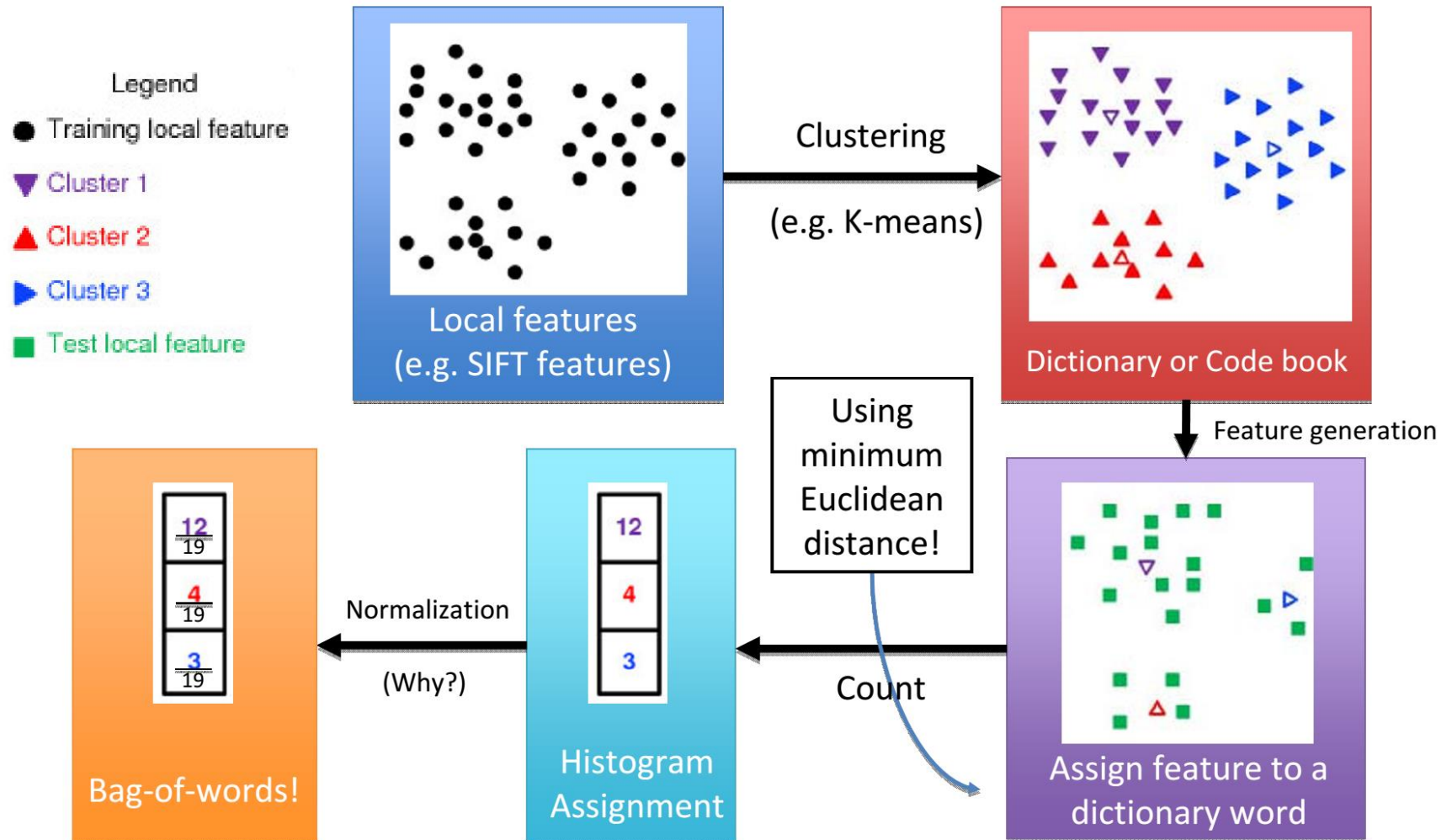
**Object** → **Bag of 'words'**

# Pipeline of Visual Bag-of-Words Model



Legend
- Training local feature
- Cluster 1
- Cluster 2
- Cluster 3
- Test local feature

Local features
(e.g. SIFT features)

Clustering

(e.g. K-means)

Dictionary or Code book

Feature generation

Using minimum Euclidean distance!

Assign feature to a dictionary word

Count

Histogram Assignment

$\frac{12}{19}$

$\frac{4}{19}$

$\frac{3}{19}$

Bag-of-words!

Normalization

(Why?)

# Local Features – the "Words"

- SIFT (Scale-invariant feature transform)



**http://www.vlfeat.org/overview/sift.html**

- And its variants!

- SIFT (sift)
- HueSIFT (huesift)
- HSV-SIFT (hsvsift)
- OpponentSIFT (opponentsift)

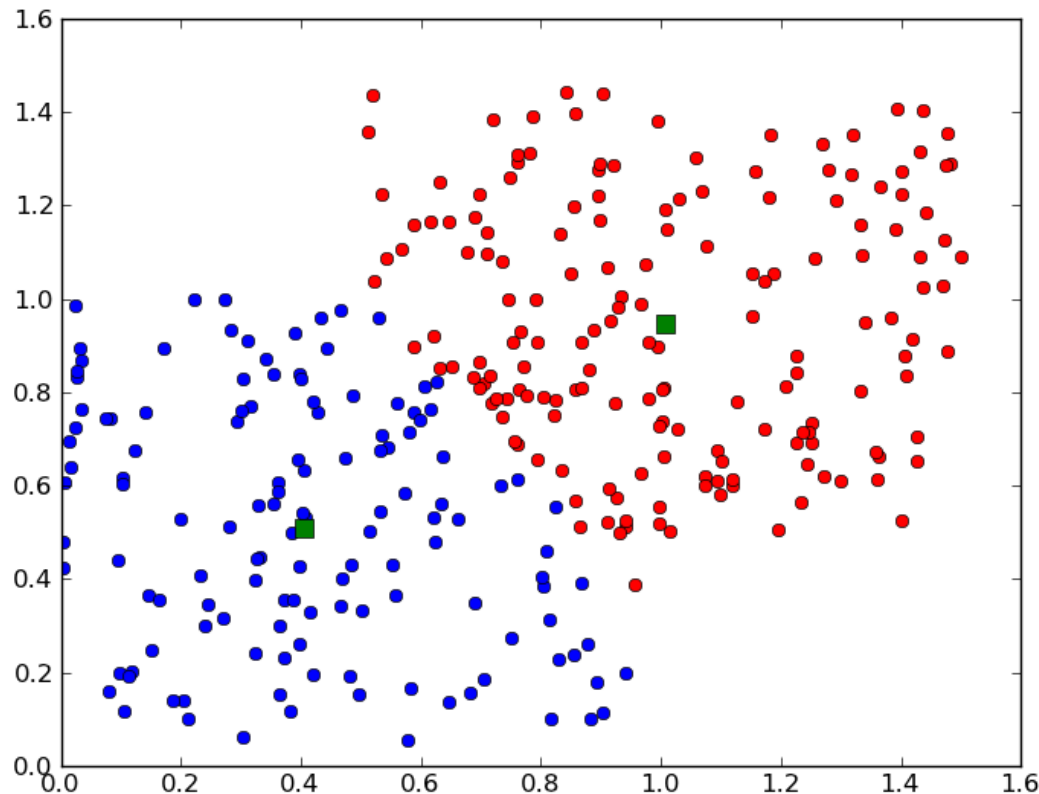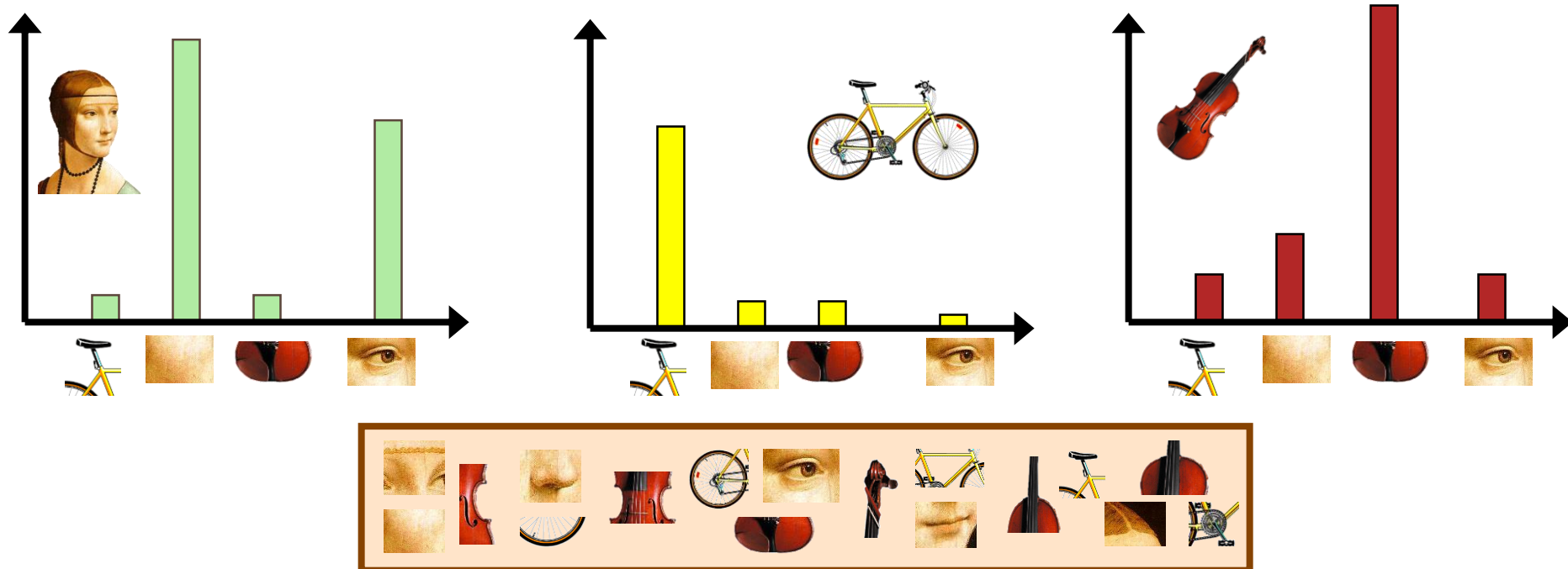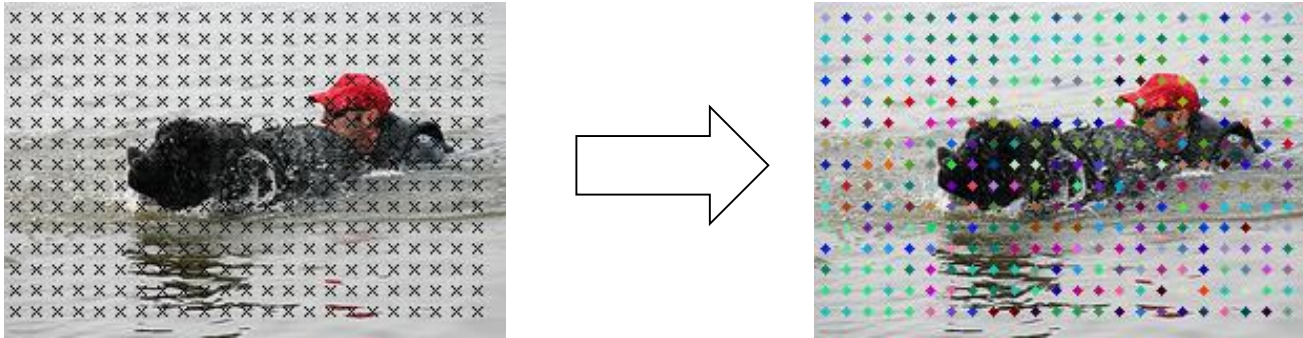- *rg*SIFT (rgsift)
- C-SIFT (csift)
- RGB-SIFT(rgbsift)

**http://koen.me/research/colordescriptors**

**Dense SIFT**

# Visual Feature Clustering – the "Dictionary"

- Iterate between "assignment" and "center-updating" steps
  - Each sample is assigned to the most similar center
  - Update centers by averaging all assigned samples

# Local Feature Quantization

- Each SIFT descriptor is quantized into a visual word using the nearest cluster center.
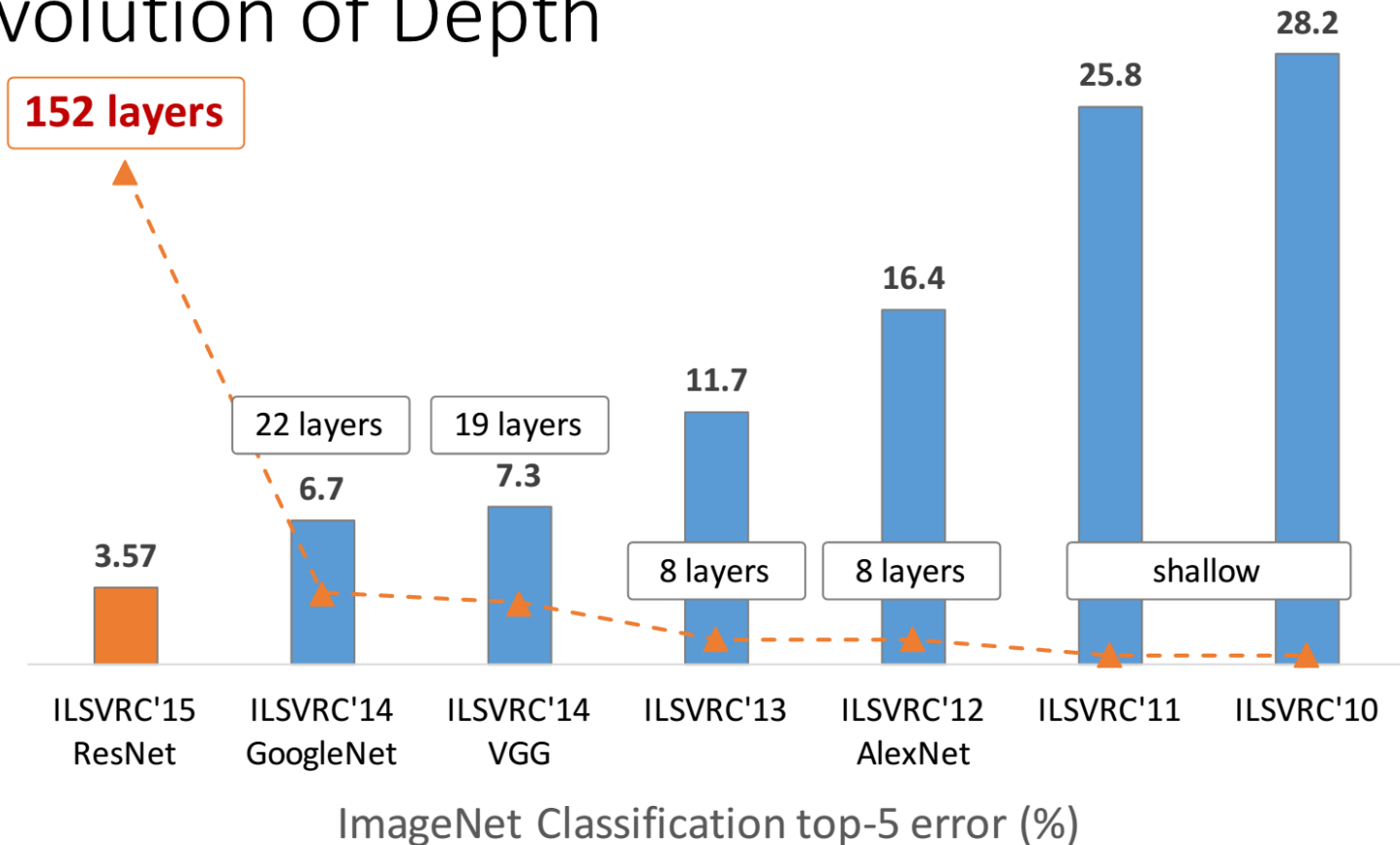
# How can we improve the BoW model

- **(Locality): Spatial Pyramid Matching & Pyramid Match Kernel**

- **(Aggregation): VLAD, Fisher Vector, Soft Assignment**

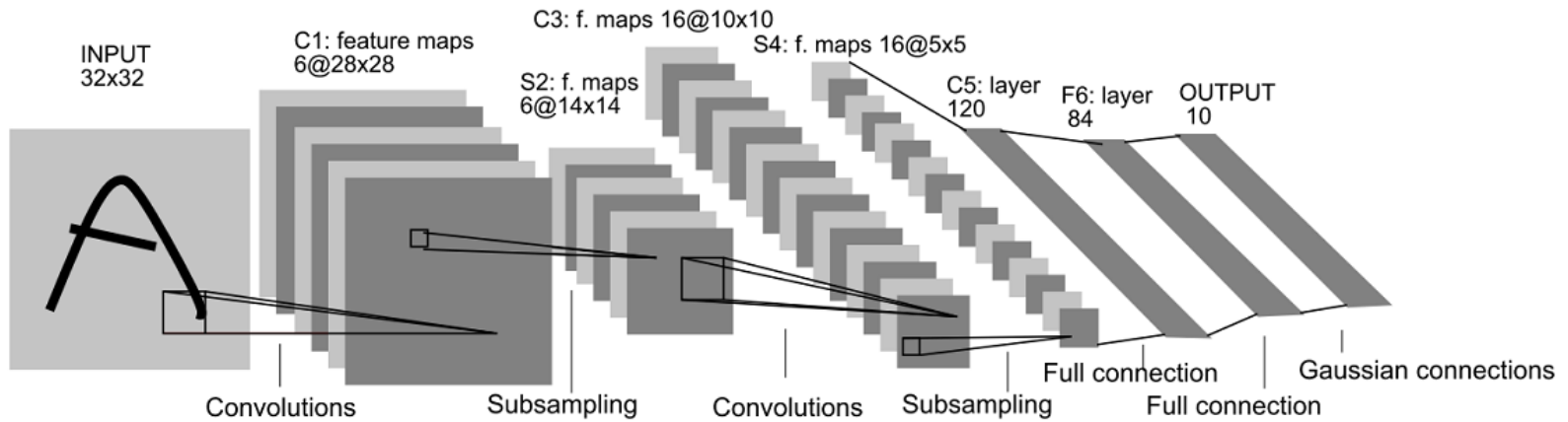- **(Dictionary): Vocabulary Tree, Sparse Coding**

# AlexNet

- Named after Alex Krizhevsky, proposed in 2012



Revolution of Depth

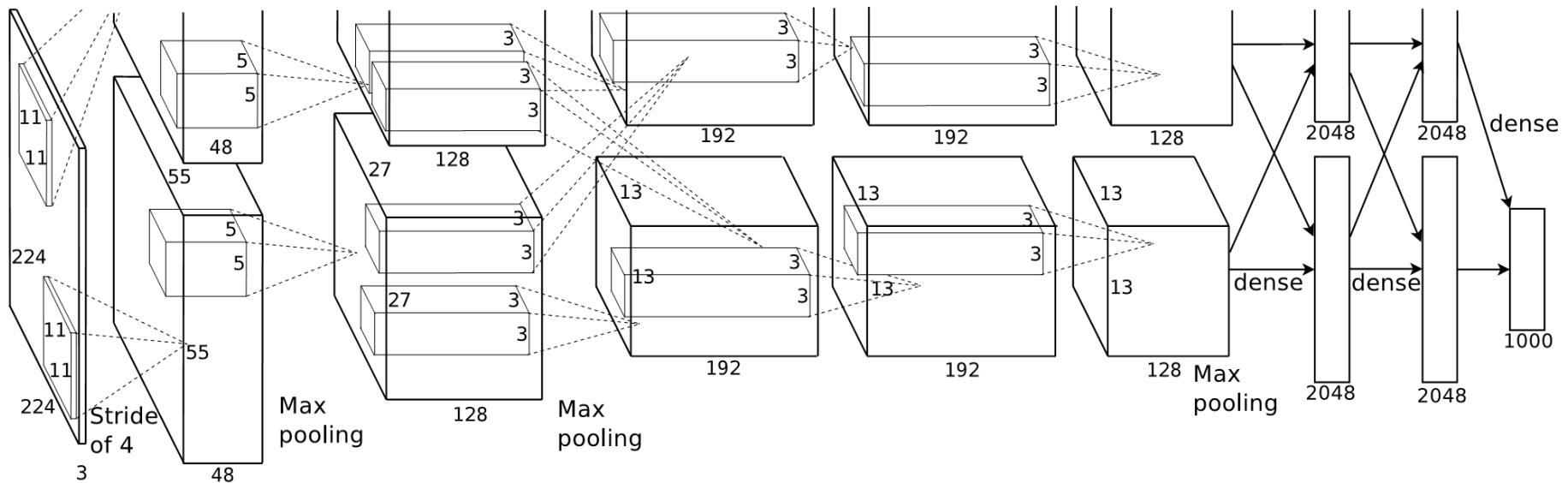ImageNet Classification top-5 error (%)

# LeNet-5



- Input: 32x32 pixel image. Largest character is 20x20
  (All important info should be in the center of the receptive field of the highest level feature detectors)

- Cx: Convolutional layer

- Sx: Subsample layer

- Fx: Fully connected layer

- Black and White pixel values are normalized:
    E.g. White = -0.1, Black =1.175 (Mean of pixels = 0, Std of pixels =1)

14

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, november 1998.

# AlexNet

- Much larger than LeNet-5
- Trained on two GTX 580 GPU
- Largest networks at its time
- Utilize multiple engineering tricks (dropout, ReLU)



Alex Krizhevsky et al., ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012

# Why DL Suddenly works?

*...It may be that the primary barriers to the success of neural networks were psychological (practitioners did not expect neural networks to work, so they did not make a serious effort to use neural networks)...*
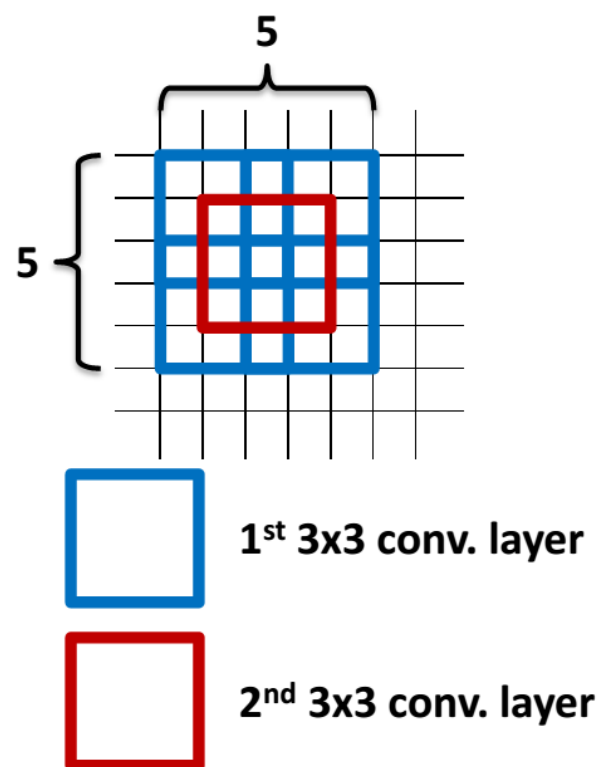
-- Goodfellow et al. "deep Learning"

# Why DL Suddenly works? – My Two Cents

- Emerging of big visual data

- GPU -> large network

- New engineering tricks (dropout, ReLU etc.)

# VGG Net

Why 3x3 layers?

- Stacked conv. layers have a large receptive field
    - two 3x3 layers – 5x5 receptive field
    - three 3x3 layers – 7x7 receptive field
- More non-linearity
- Less parameters to learn
    - ~140M per net



1st 3x3 conv. layer

2nd 3x3 conv. layer

# Network Design

**Key design choices:**

- 3x3 conv. kernels – very small
- conv. stride 1 – no loss of information

Other details:

- Rectification (ReLU) non-linearity
- 5 max-pool layers (x2 reduction)
- no normalisation
- 3 fully-connected (FC) layers

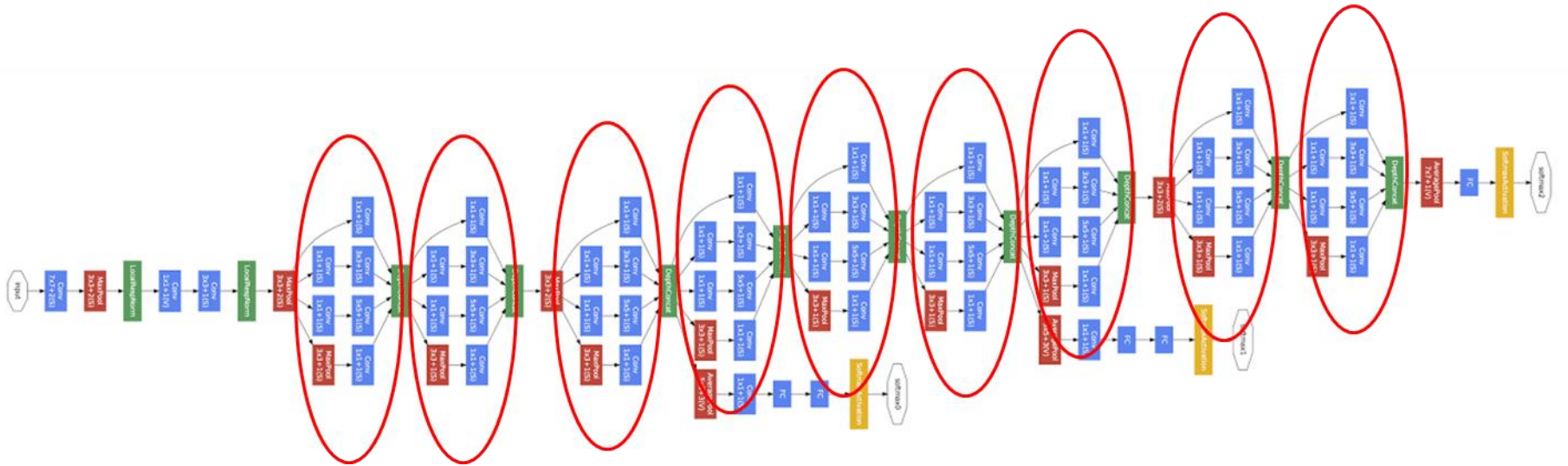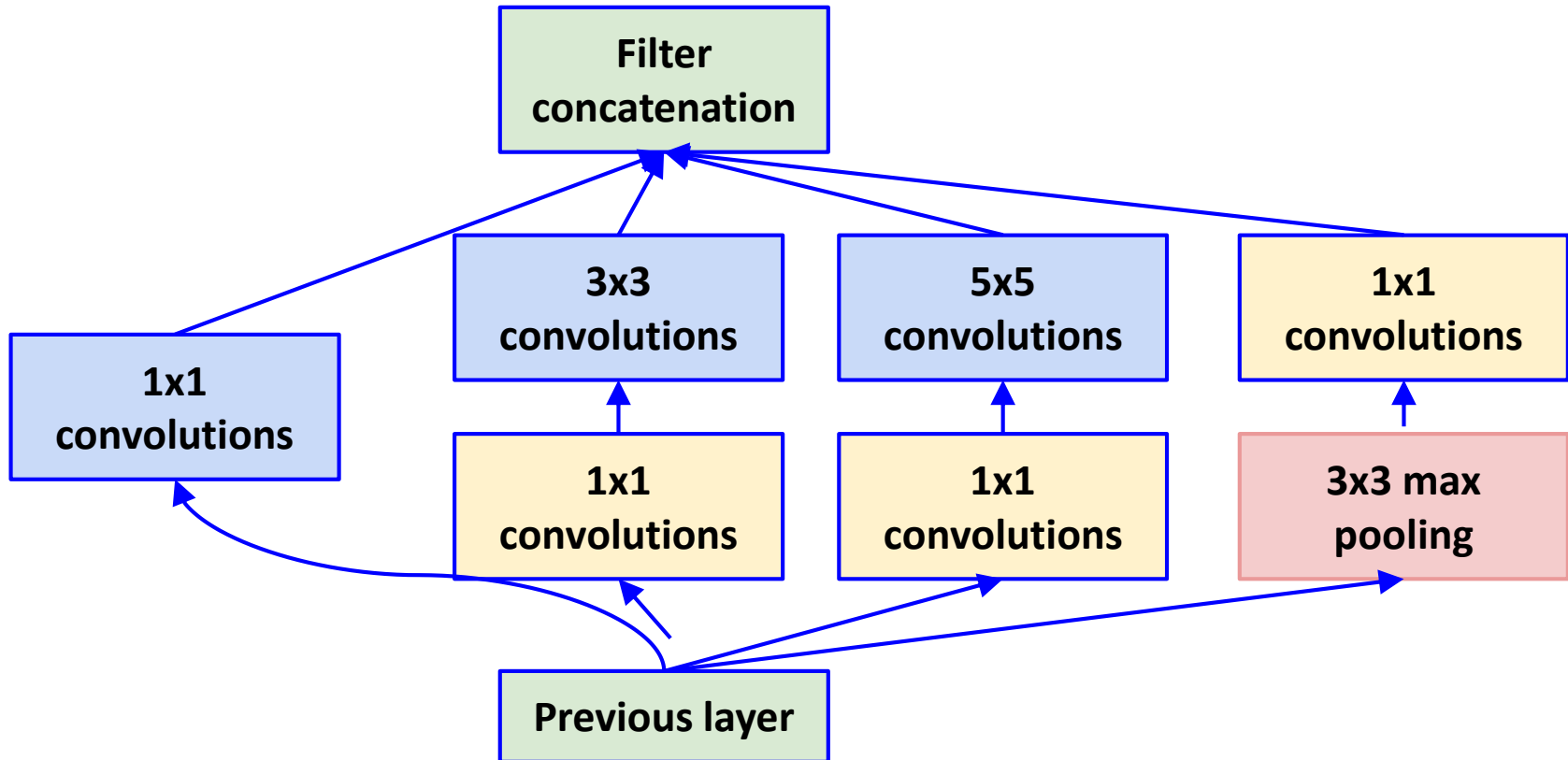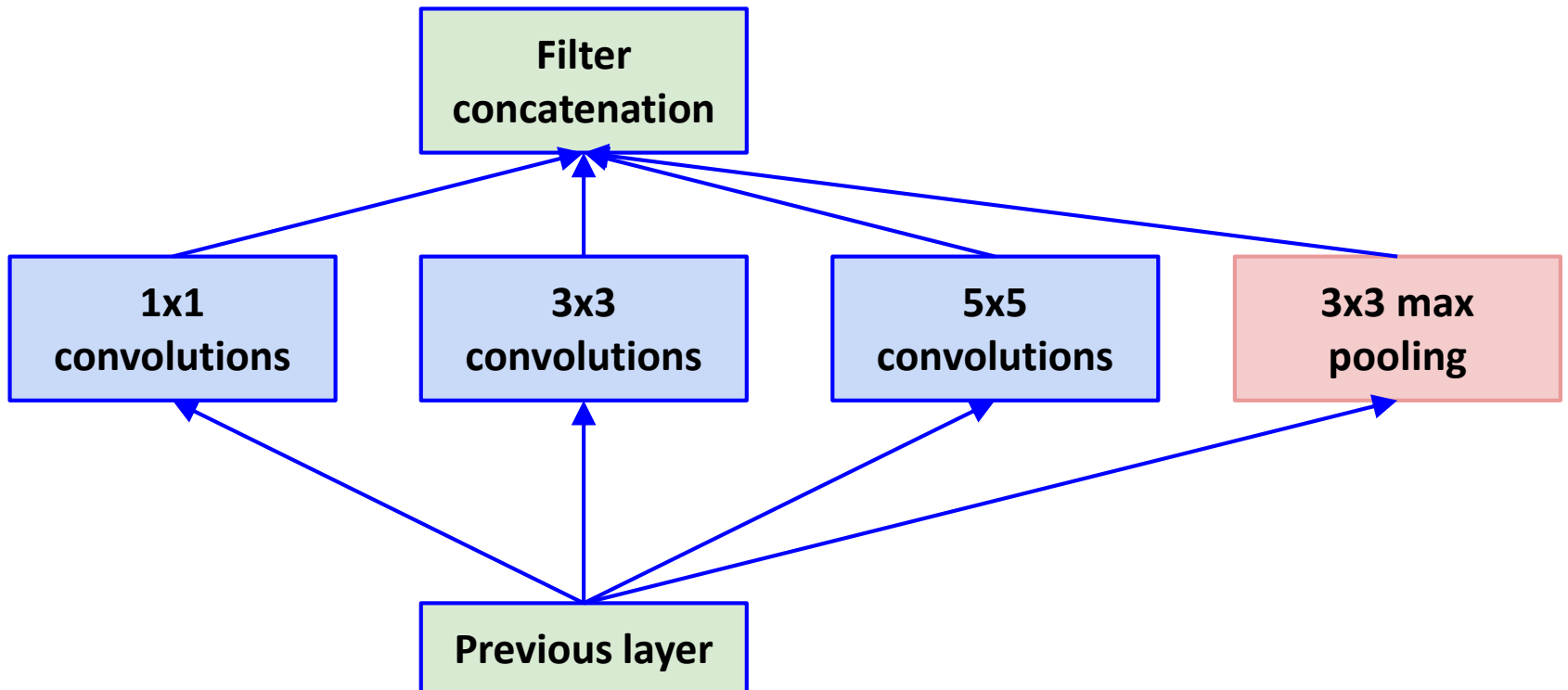| image |
| :---: |
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| FC-4096 |
| FC-4096 |
| FC-1000 |
| softmax |

# GoogLeNet

**Inception**

**Network in a network in a network…**

Convolution
Pooling
Softmax
Other

# Inception module

# Naive idea (does not work!)

# ResNet

- See He Kaiming's ICML tutorial