

# Recognition and Learning

Dr. Sander Ali Khawaja,

Assistant Professor, Department of Telecommunication Engineering  
Faculty of Engineering and Technology, University of Sindh, Pakistan

Senior Member, IEEE – Member, ACM

<https://sander-ali.github.io>

# Computer Vision & Image Processing



# Image Features and Categorization



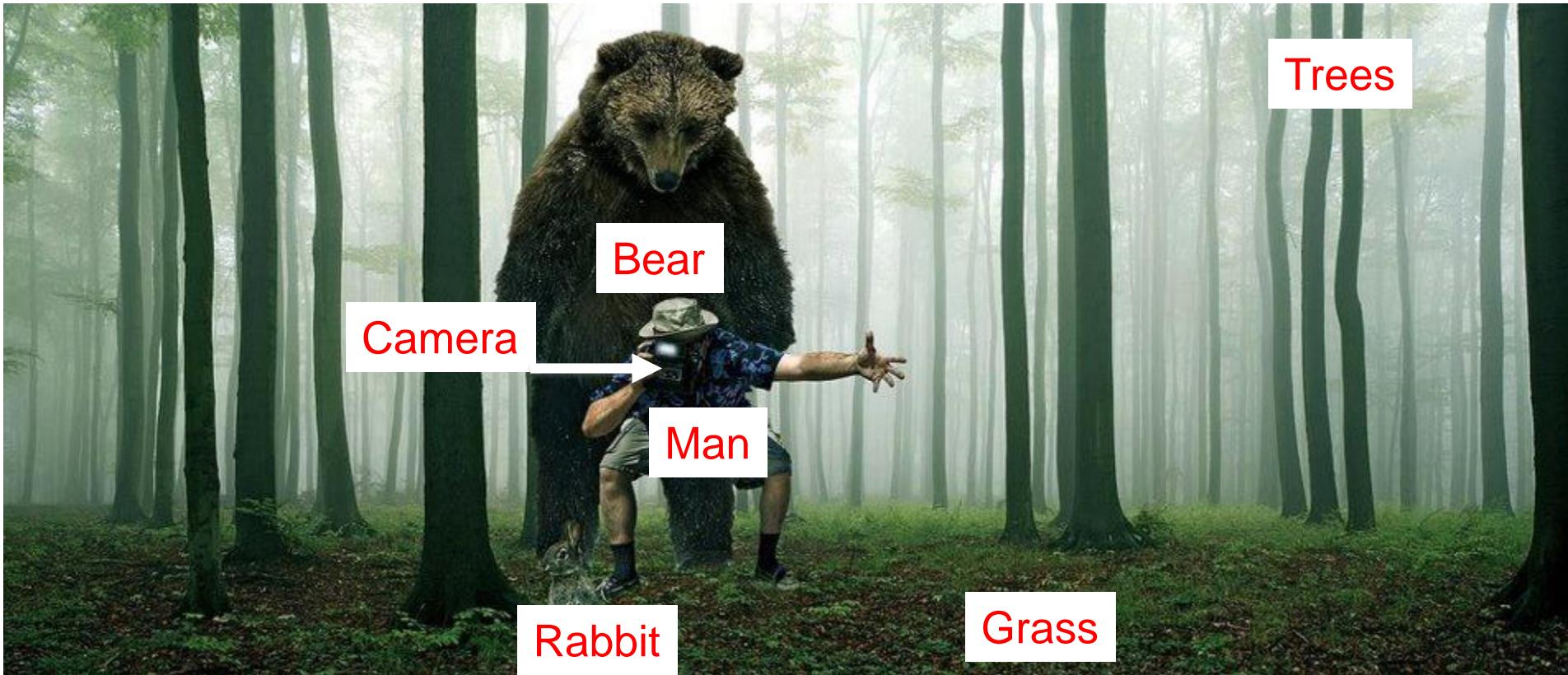
# Recognition and Learning

- Image Features and Categorization
- Convolutional Neural Networks
- Object Detection
- Part and Pixel Labeling
- Action Recognition
- Vision and Language





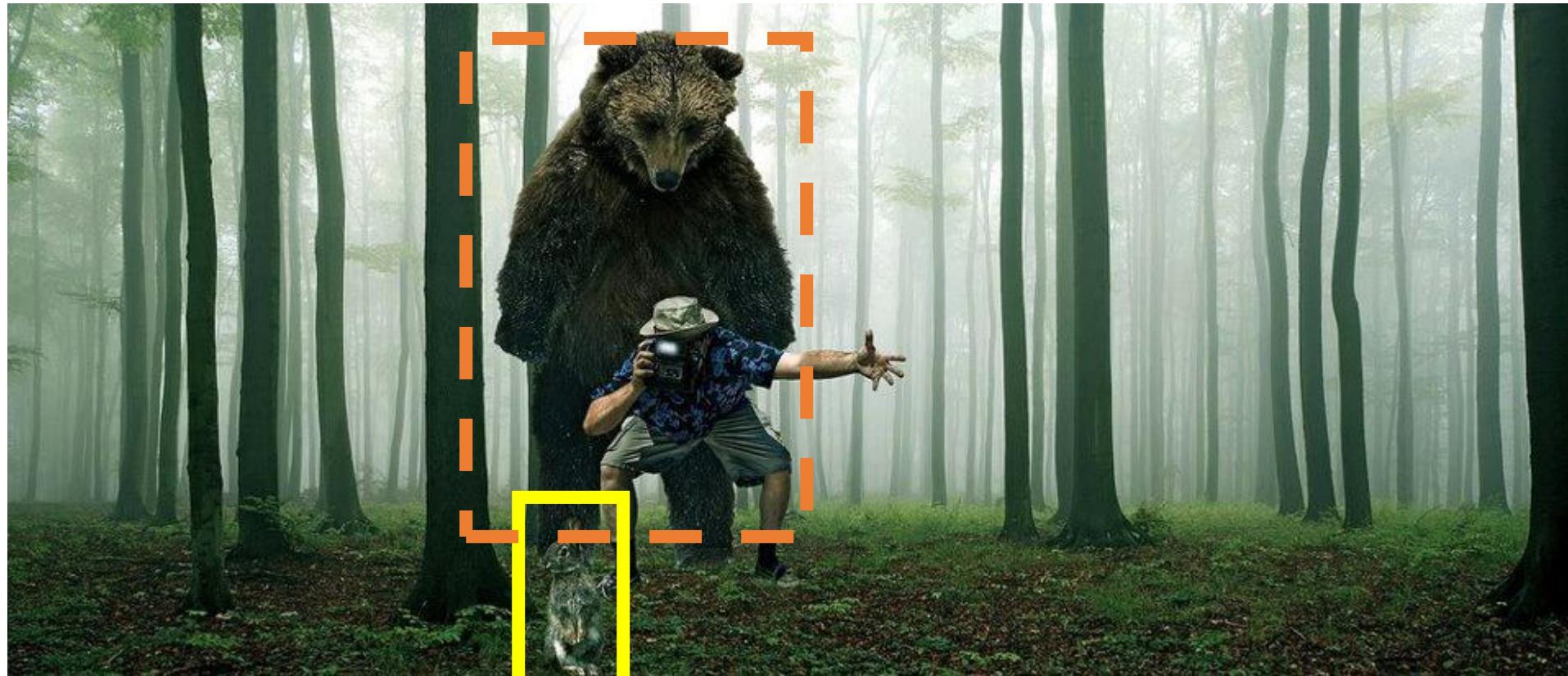
# What do you see in this image?



## Forest

Dr. Sander Ali Khowaja

# Describe, Predict, or Interact with the object based on visual cues



Is it **dangerous**?

How **fast** does it run?

Is it **alive**?

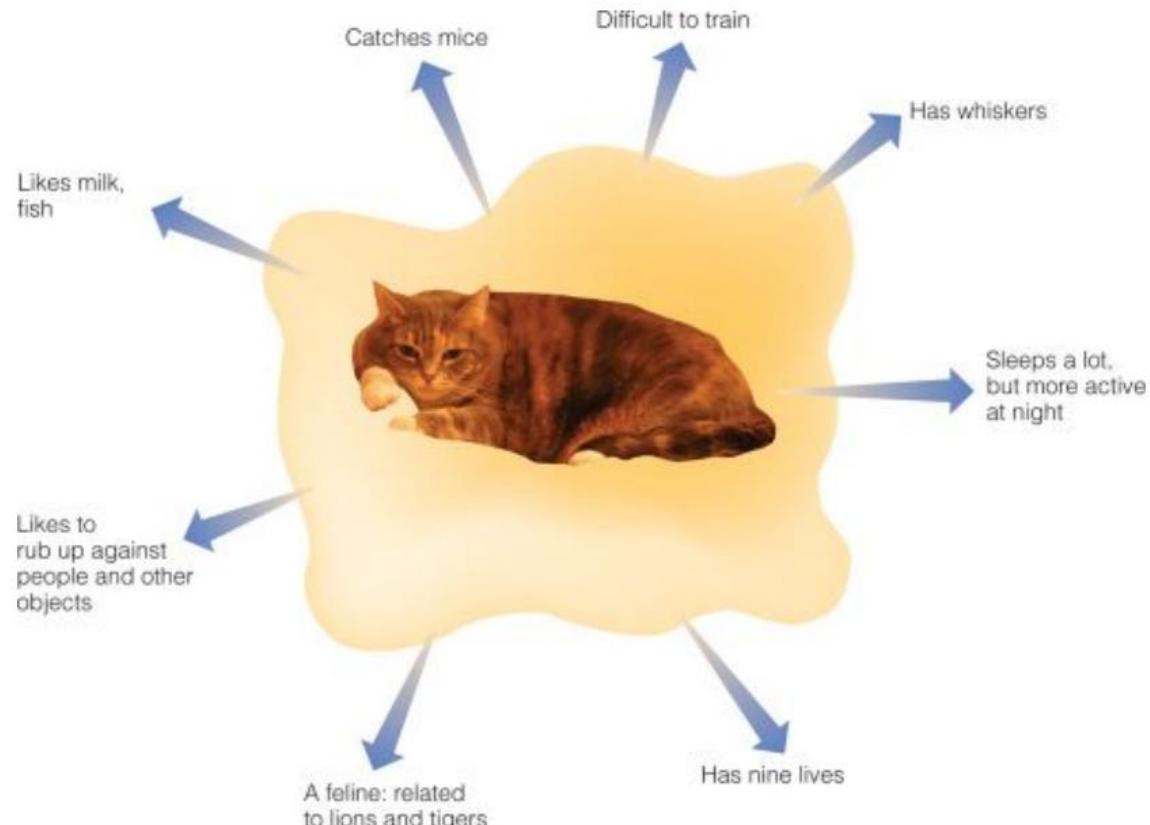
Does it have a **tail**?

Is it **soft**?

Can I **poke** with it?

# Why do we care about categories?

- From an object's category, we can make predictions about its behavior in the future, beyond of what is immediately perceived.
- Pointers to knowledge
  - Help to understand individual cases not previously encountered
- Communication



Dr. Sander Ali Khowaja

# Theory of categorization

How do we determine if something is a member of a particular category?

- Definitional approach
- Prototype approach
- Exemplar approach



# Definitional approach: classical view of categories

- Plato & Aristotle

- Categories are defined by a list of properties shared by all elements in a category
- Category membership is binary
- Every member in the category is equal

## The Categories (Aristotle)



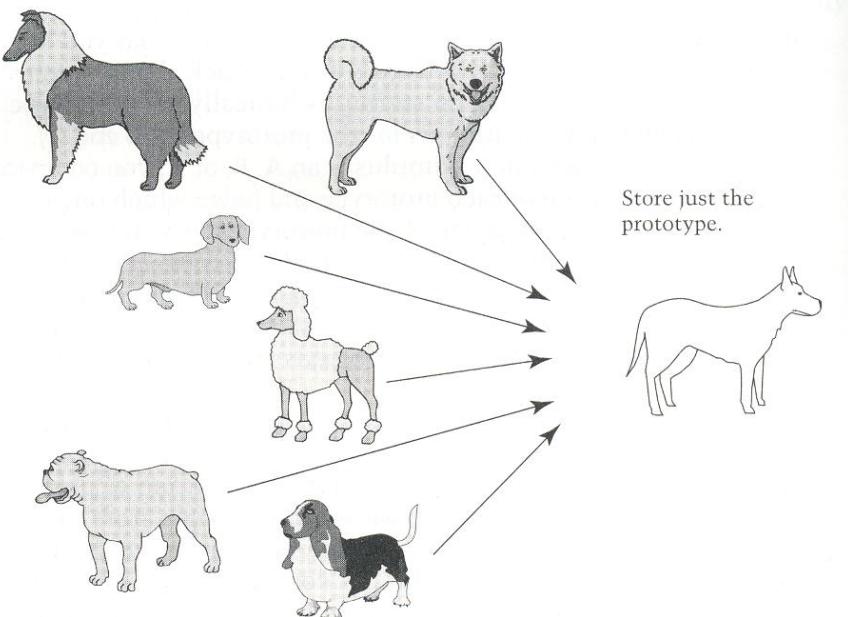
Aristotle by Francesco Hayez

Slide Credit: A. A. Efros

Dr. Sander Ali Khowaja

# Prototype or sum of exemplars?

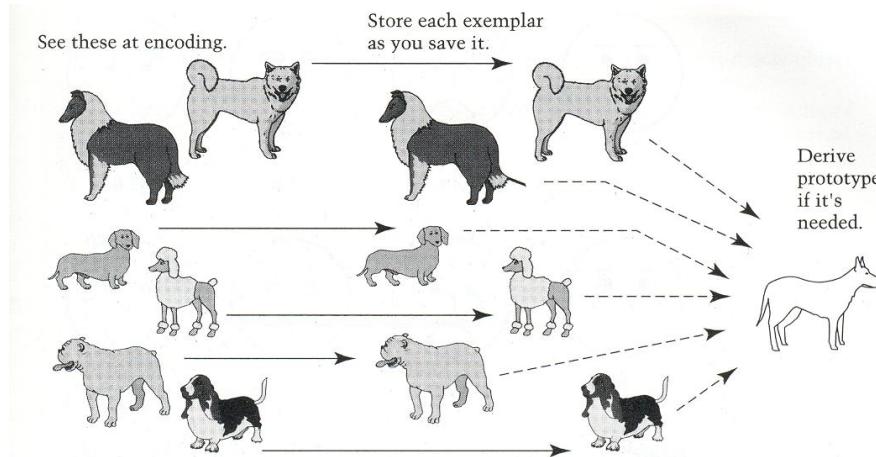
## Prototype Model



**Figure 7.3.** Schematic of the prototype model. Although many exemplars are seen, only the prototype is stored. The prototype is updated continually to incorporate more experience with new exemplars.

Category judgments are made by comparing a new exemplar to the prototype.

## Exemplars Model



**Figure 7.4.** Schematic of the exemplar model. As each exemplar is seen, it is encoded into memory. A prototype is abstracted only when it is needed, for example, when a new exemplar must be categorized.

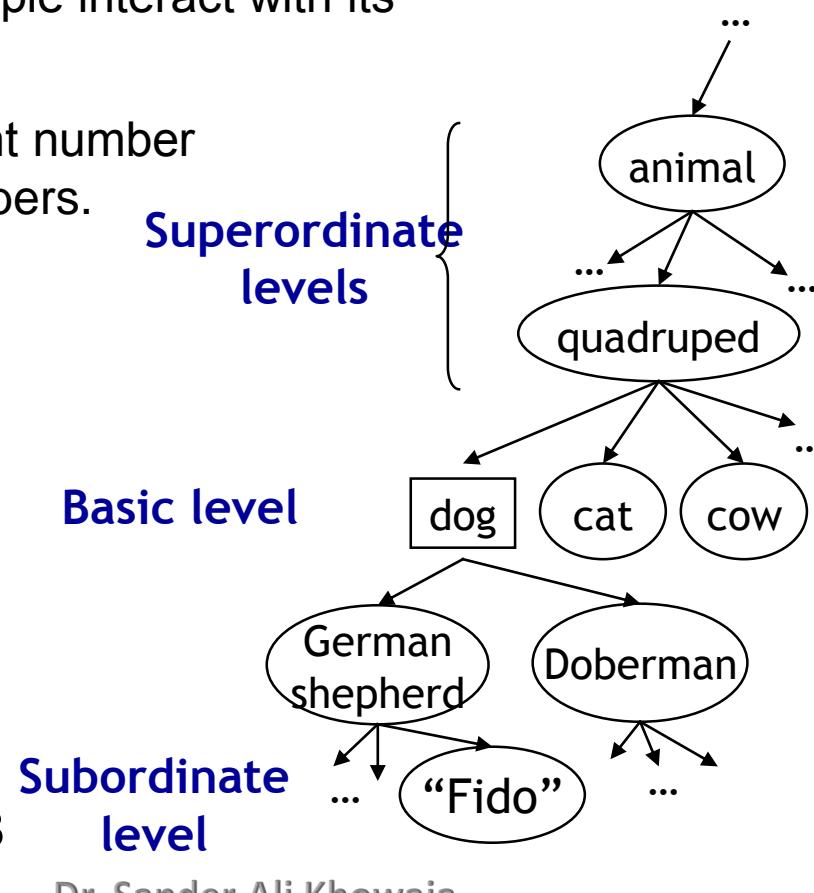
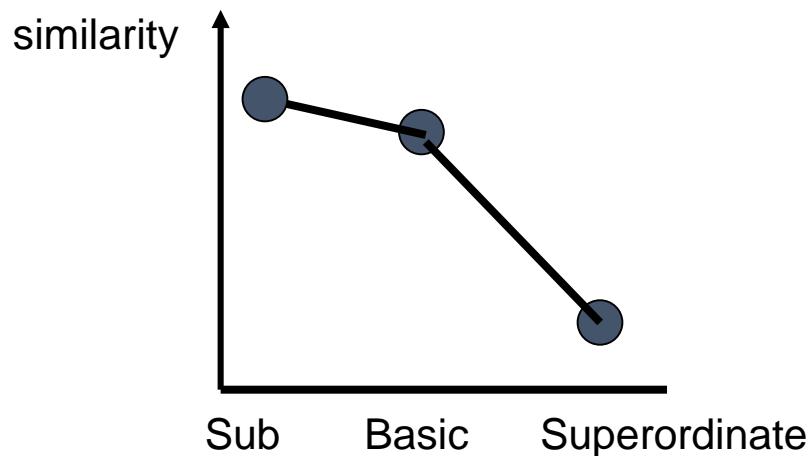
Category judgments are made by comparing a new exemplar to all the old exemplars of a category or to the exemplar that is the most appropriate

Slide Credit: Torralba

# Levels of categorization [Rosch 70s]

Definition of Basic Level:

- **Similar shape:** Basic level categories are the highest-level category for which their members have similar shapes.
- **Similar motor interactions:** ... for which people interact with its members using similar motor sequences.
- **Common attributes:** ... there are a significant number of attributes in common between pairs of members.



[Rosch et al. Principle of categorization, 1978](#)

11

Dr. Sander Ali Khowaja



# Image Categorization

- Cat vs Dog



# Image Categorization

- Object recognition



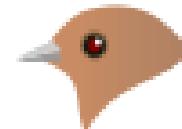
Caltech 101 Average Object Images

Dr. Sander Ali Khowaja

# Image Categorization



Generalist



Insect catching



Grain eating



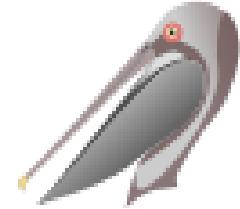
Coniferous-seed eating



Nectar feeding



Chiseling



Dip netting



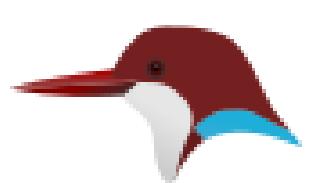
Surface skimming



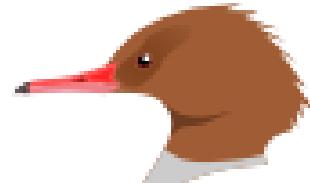
Scything



Probing



Aerial fishing



Pursuit fishing



Scavenging



Raptorial



Filter feeding

Visipedia Project

# Image Categorization

- Place recognition



spare bedroom

teenage bedroom

romantic bedroom



darkest forest path

wintering forest path

greener forest path



wooded kitchen

messy kitchen

stylish kitchen



rocky coast

misty coast

sunny coast

Places Database [[Zhou et al. NIPS 2014](#)]

# Image Categorization

- Visual font recognition



[Chen et al. CVPR 2014]

Dr. Sander Ali Khowaja

# Image Categorization

- Dating historical photos



1940



1953



1966



1977

[Palermo et al. ECCV 2012]

# Image Categorization

- Image style recognition



HDR



Macro



Baroque



Roccoco



Vintage



Noir



Northern Renaissance



Cubism



Minimal



Hazy



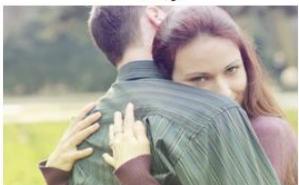
Impressionism



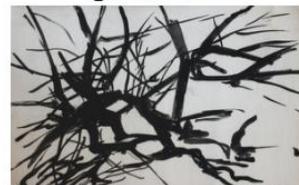
Post-Impressionism



Long Exposure



Romantic



Abs. Expressionism



Color Field Painting

Flickr Style: 80K images covering 20 styles.

Wikipaintings: 85K images for 25 art genres.

[[Karayev et al. BMVC 2014](#)]

Dr. Sander Ali Khowaja

# Region Categorization

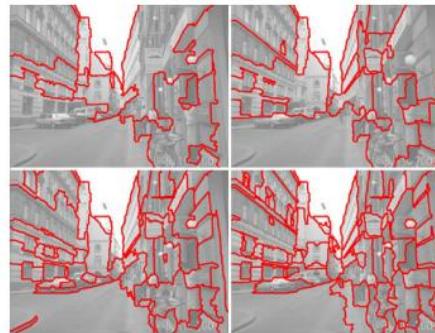
- Layout prediction



Input



Superpixels



Multiple Segmentations



Surface Layout



a

Assign regions to orientation  
Geometric context [[Hoiem et al. IJCV 2007](#)]

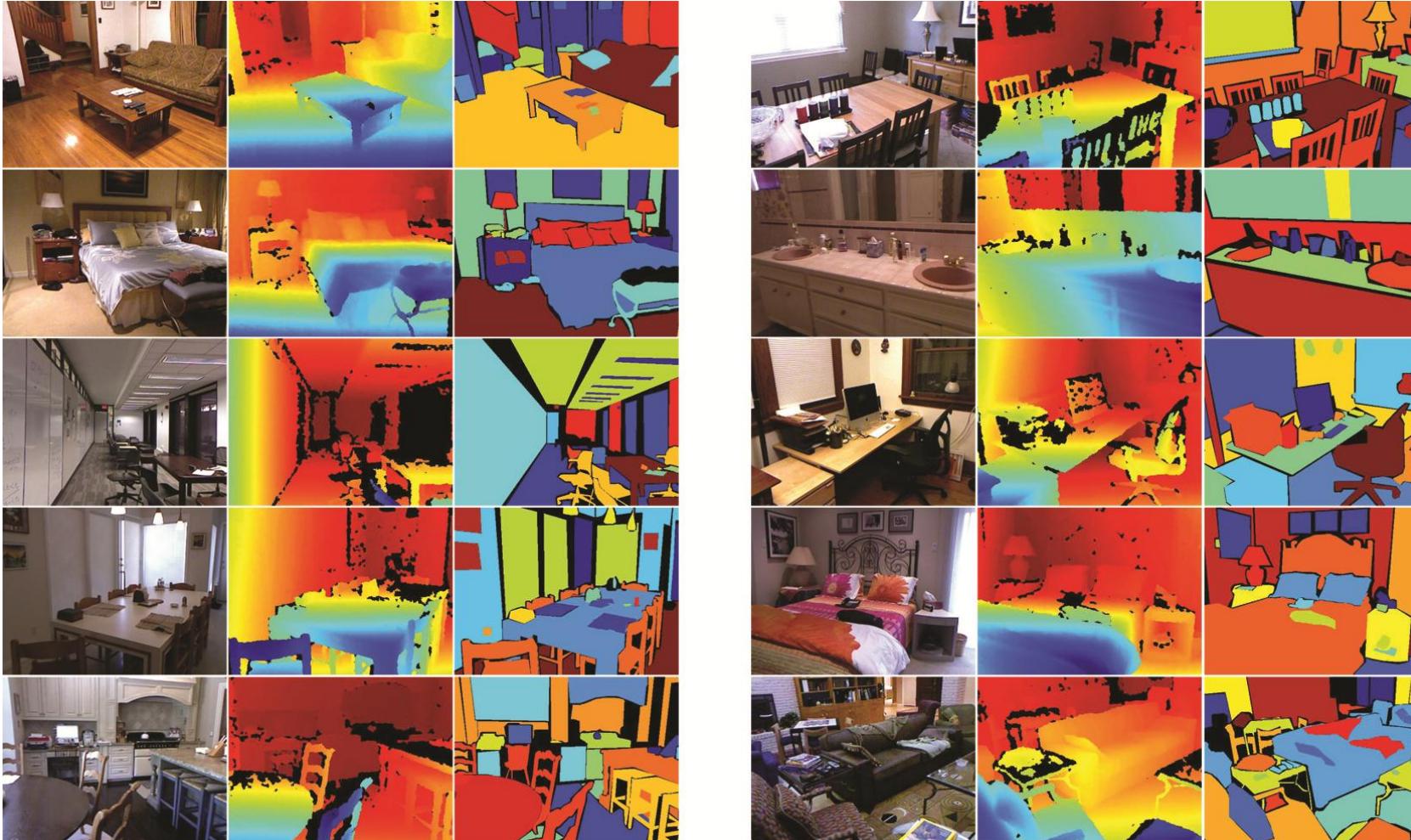


c d

Assign regions to depth  
Make3D [[Saxena et al. PAMI 2008](#)]

# Region Categorization

- Semantic segmentation from RGBD images

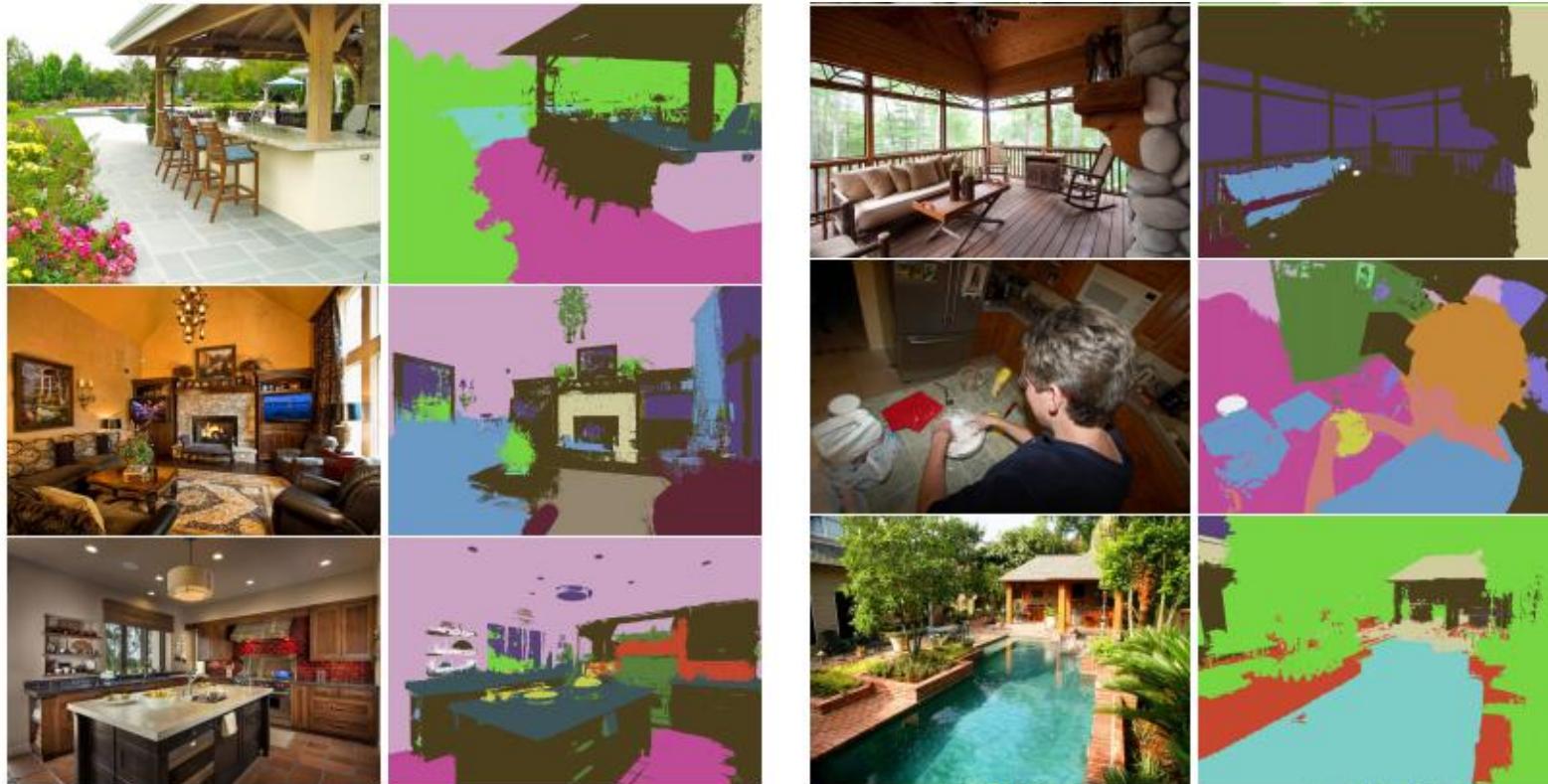


[Silberman et al. ECCV 2012]

Dr. Sander Ali Khowaja

# Region Categorization

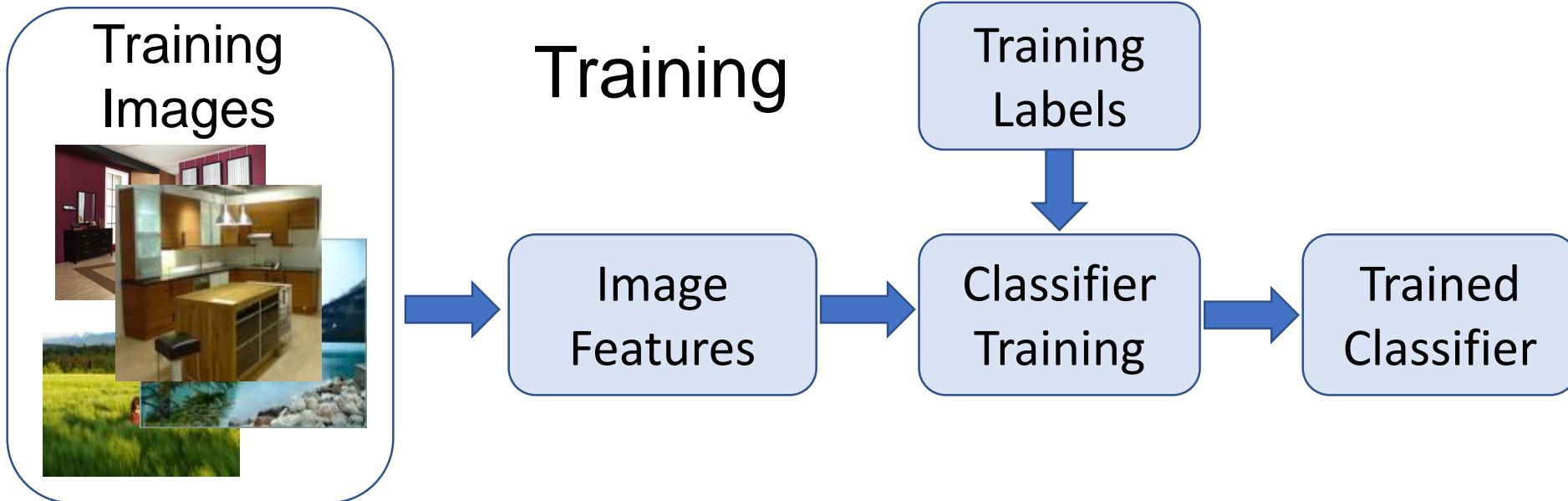
- Material recognition



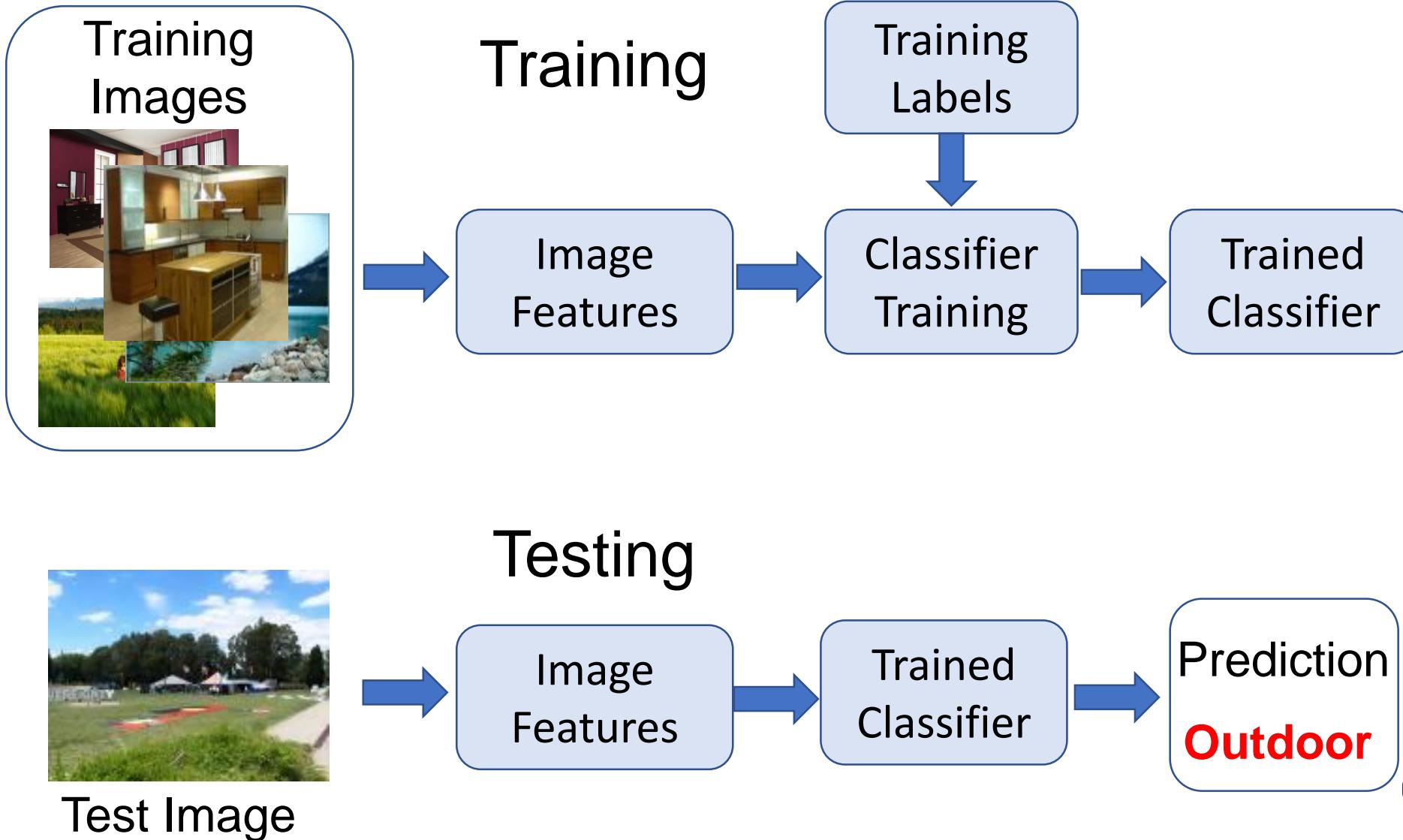
[Bell et al. CVPR 2015]

Dr. Sander Ali Khowaja

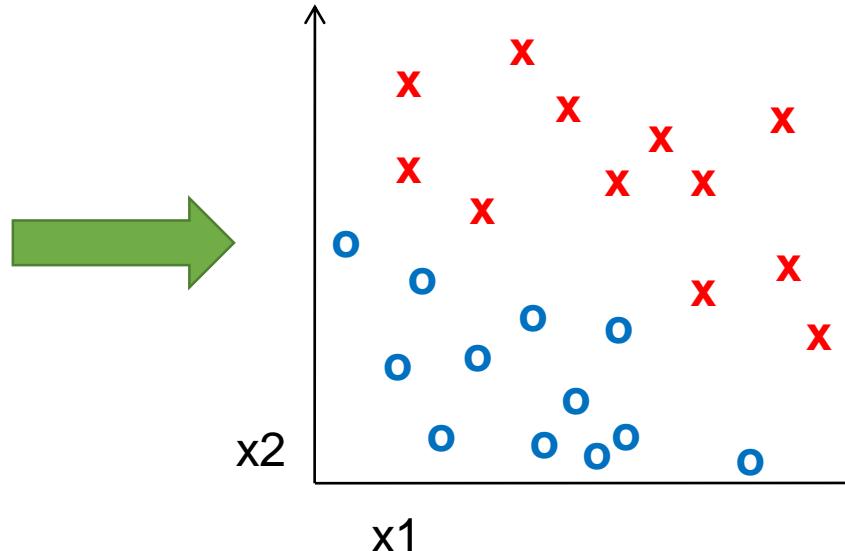
# Training Phase



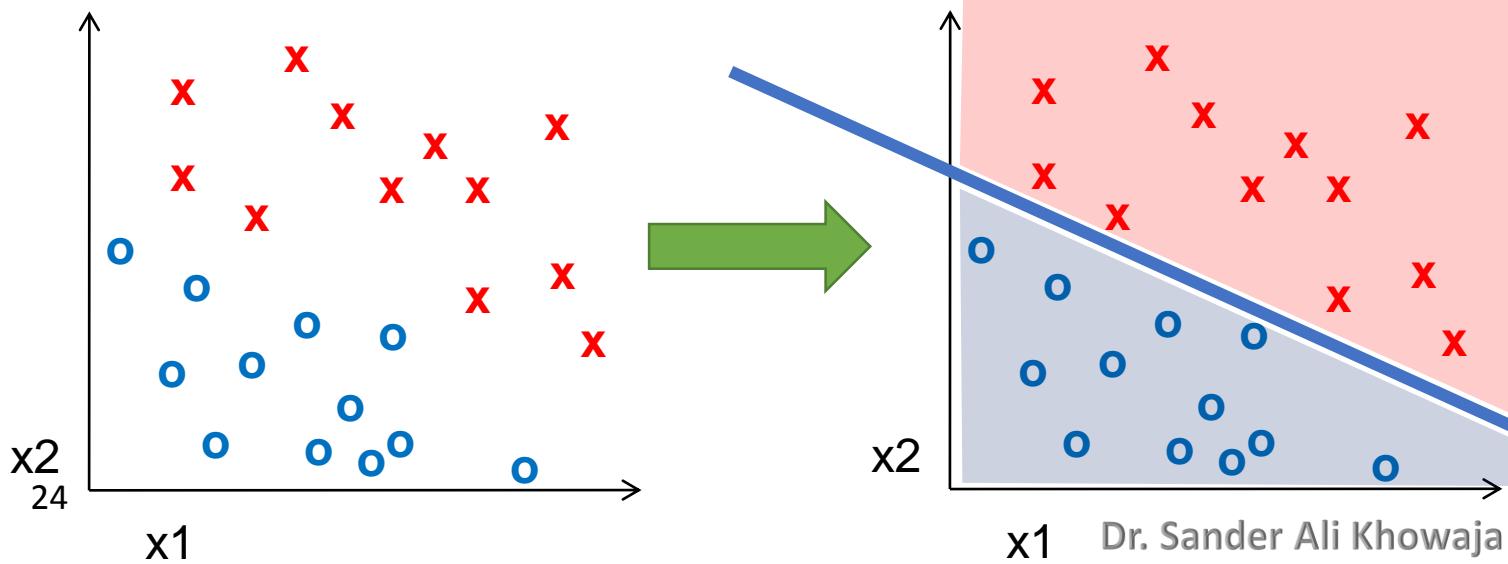
# Testing Phase



- **Image features:** map images to feature space



- **Classifiers:** map feature space to label space



Dr. Sander Ali Khawaja



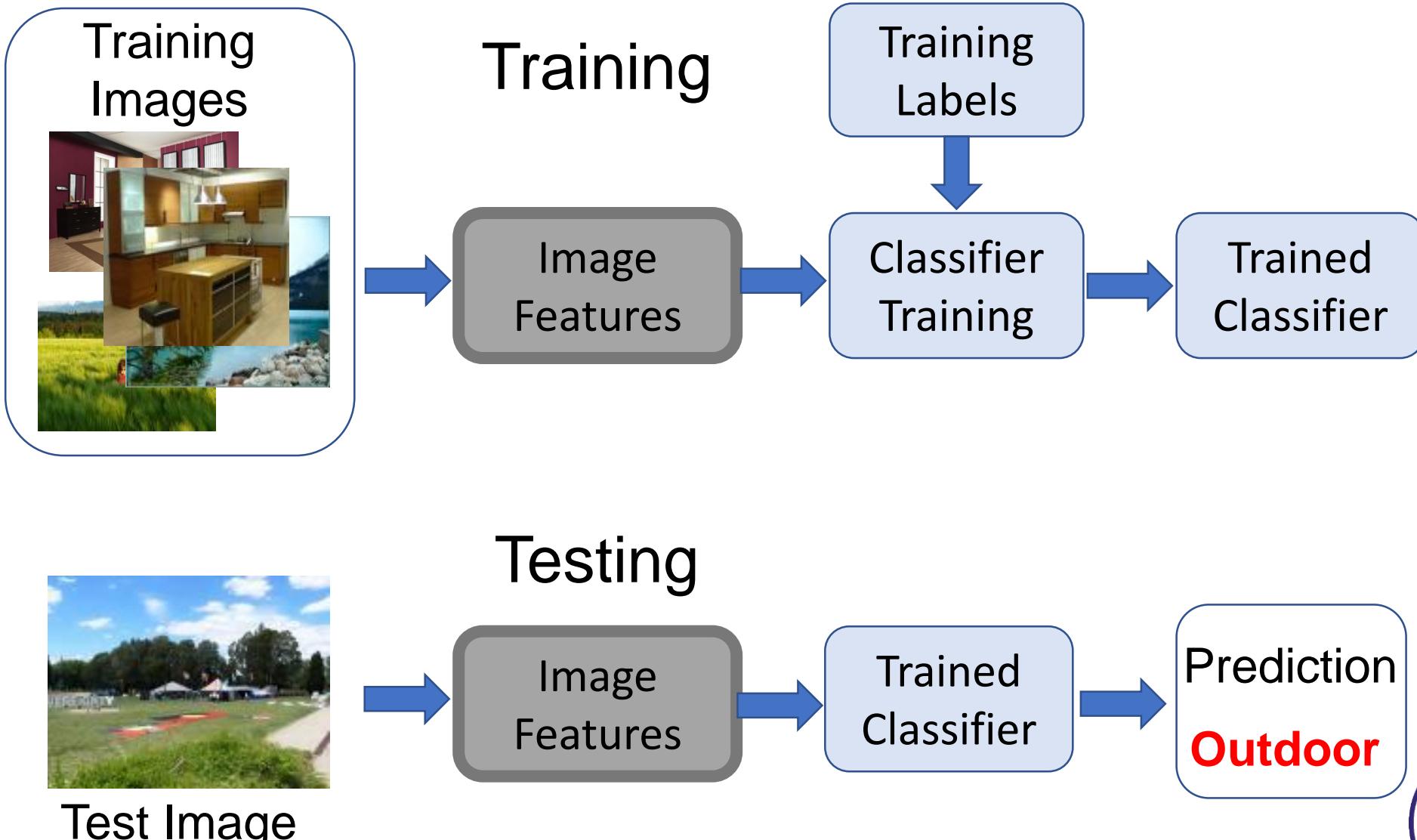
# Different types of classification

- **Exemplar-based:** transfer category labels from examples with most similar features
  - What similarity function? What parameters?
- **Linear classifier:** confidence in positive label is a weighted sum of features
  - What are the weights?
- **Non-linear classifier:** predictions based on more complex function of features
  - What form does the classifier take? Parameters?
- **Generative classifier:** assign to the label that best explains the features (makes features most likely)
  - What is the probability function and its parameters?

Note: You can always fully design the classifier by hand, but usually this is too difficult. Typical solution: learn from training examples.



# Testing Phase



# What are good features for

- recognizing a beach?



# What are good features for...

- recognizing cloth fabric?



# What are good features for...

- recognizing a mug?



# What are the right features?

Depend on what you want to know!

- Object: shape
  - Local shape info, shading, shadows, texture
- Scene : geometric layout
  - linear perspective, gradients, line segments
- Material properties: albedo, feel, hardness
  - Color, texture
- Action: motion
  - Optical flow, tracked points



# General principles of representation

- **Coverage**
  - Ensure that all relevant info is captured
- **Concision**
  - Minimize number of features without sacrificing coverage
- **Directness**
  - Ideal features are independently useful for prediction



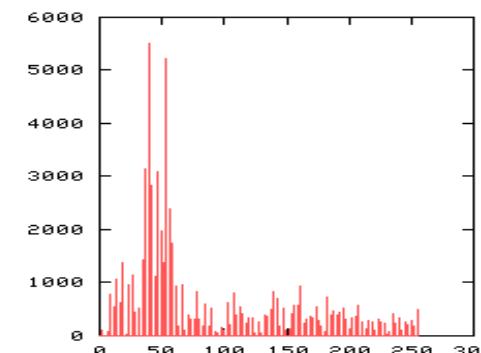
# Image Representations

- Templates
  - Intensity, gradients, etc.
- Histograms
  - Color, texture, SIFT descriptors, etc.
- Average of features

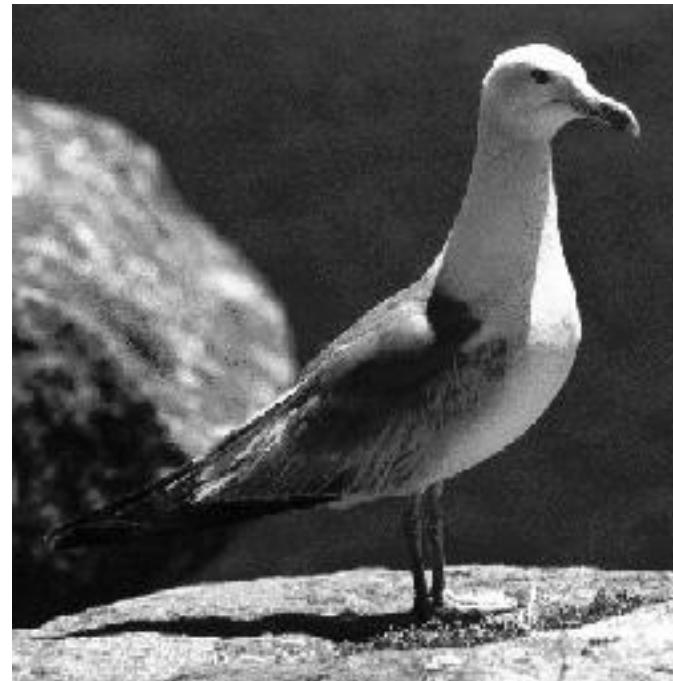
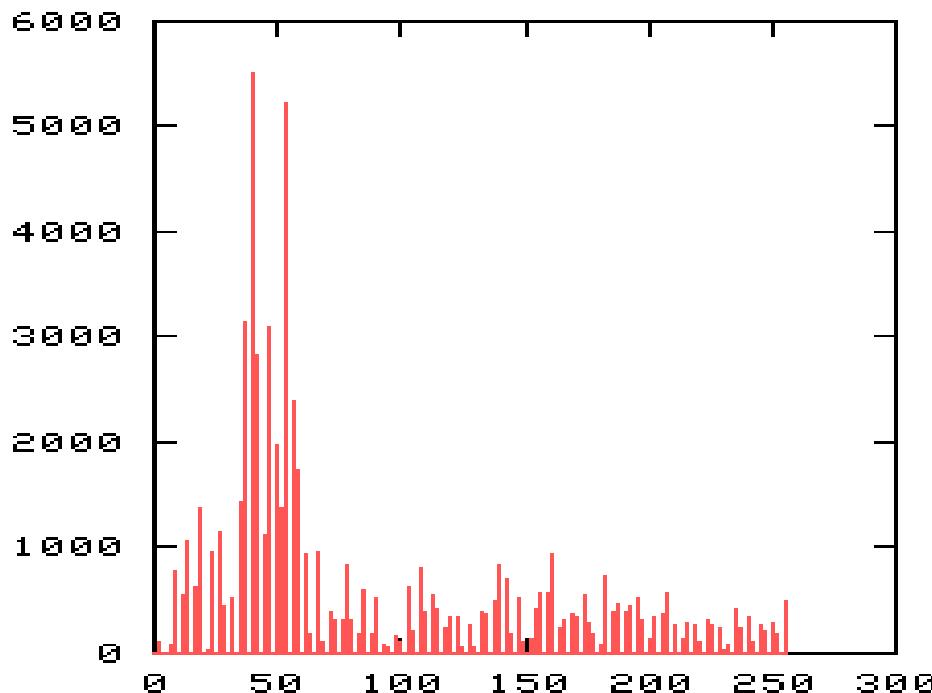


Image  
Intensity

Gradient  
template



# Image Representations: histograms

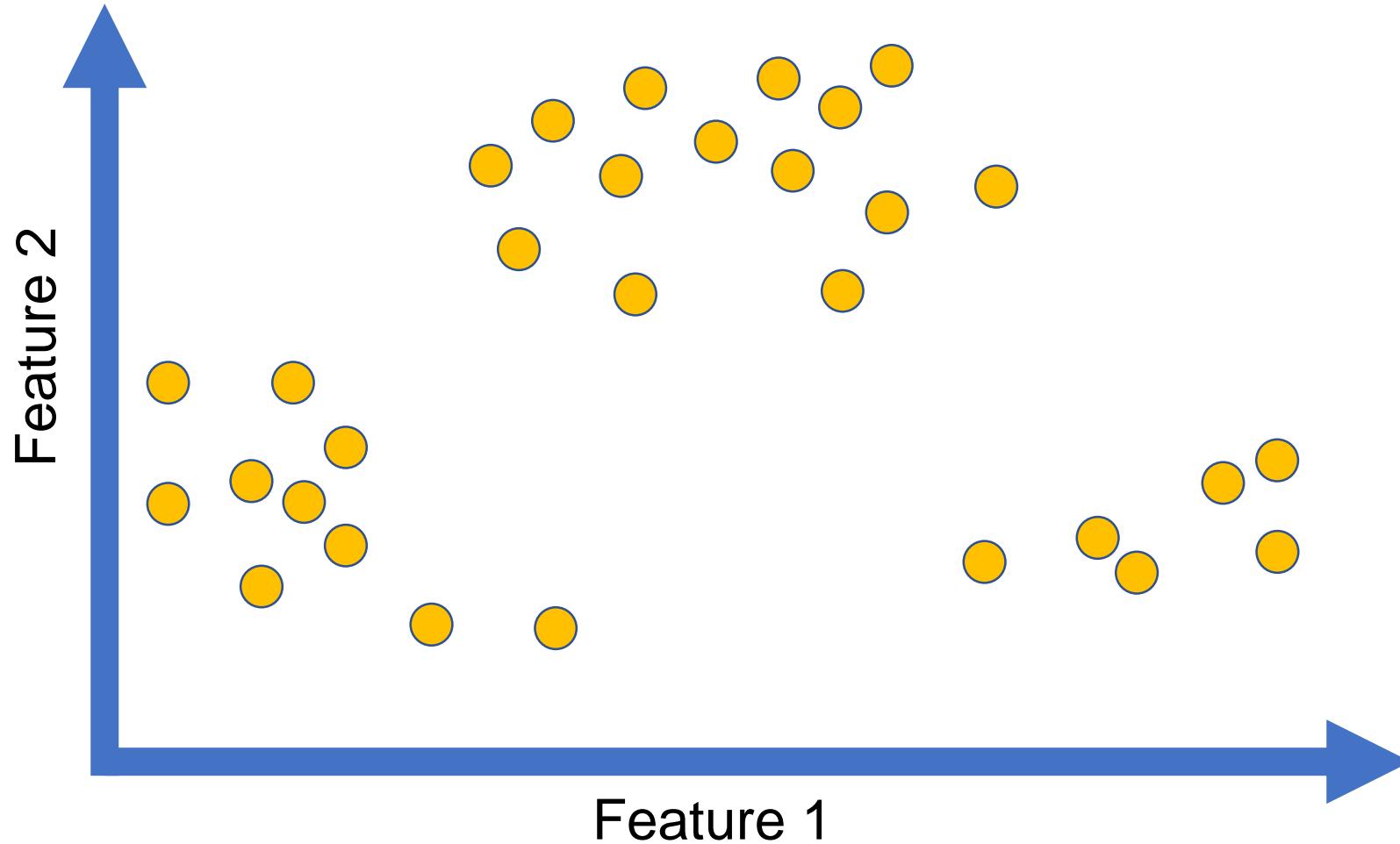


## Global histogram

- Represent distribution of features
  - Color, texture, depth, ...

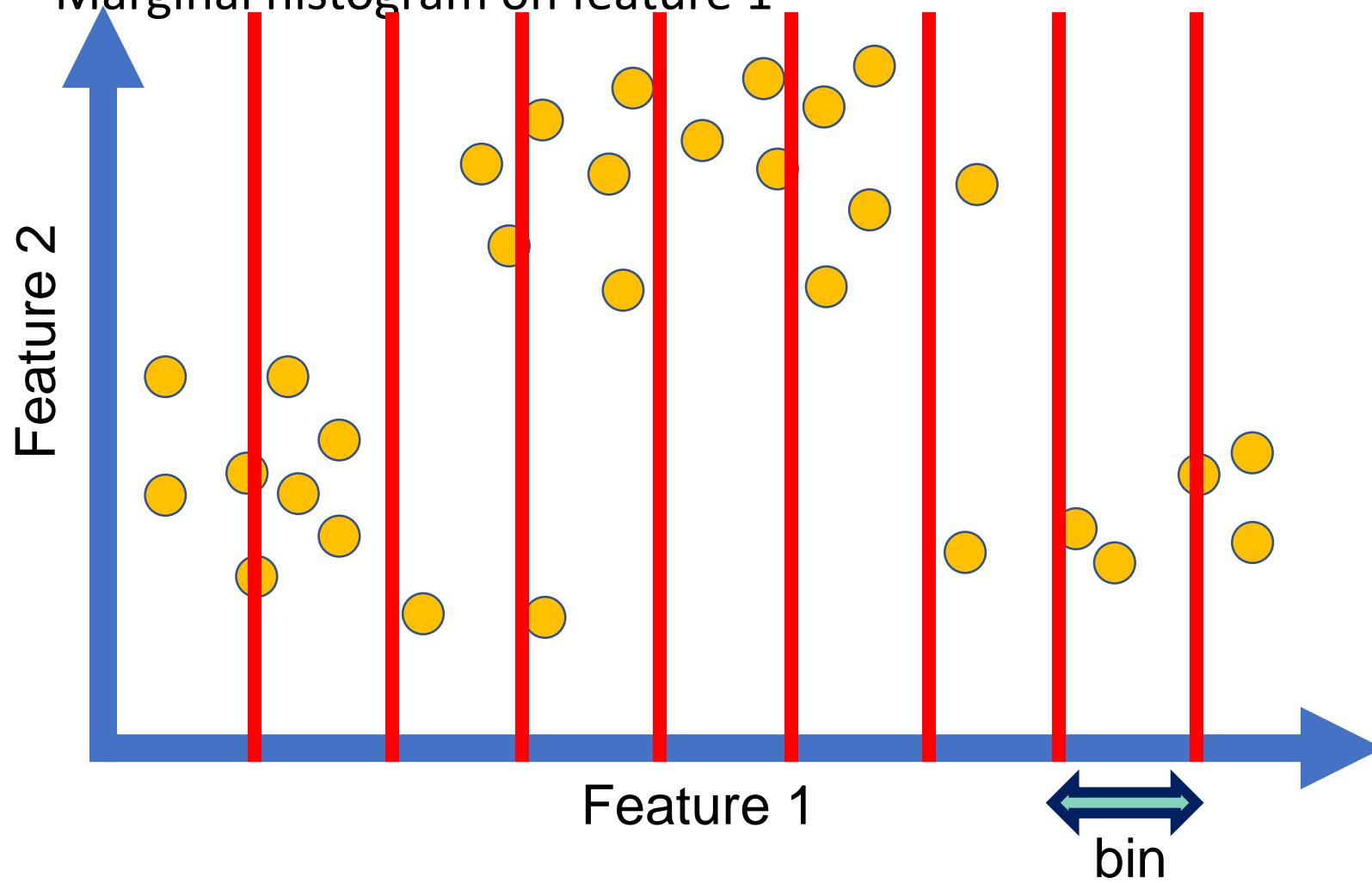
# Image Representations: Histograms

- Data samples in 2D



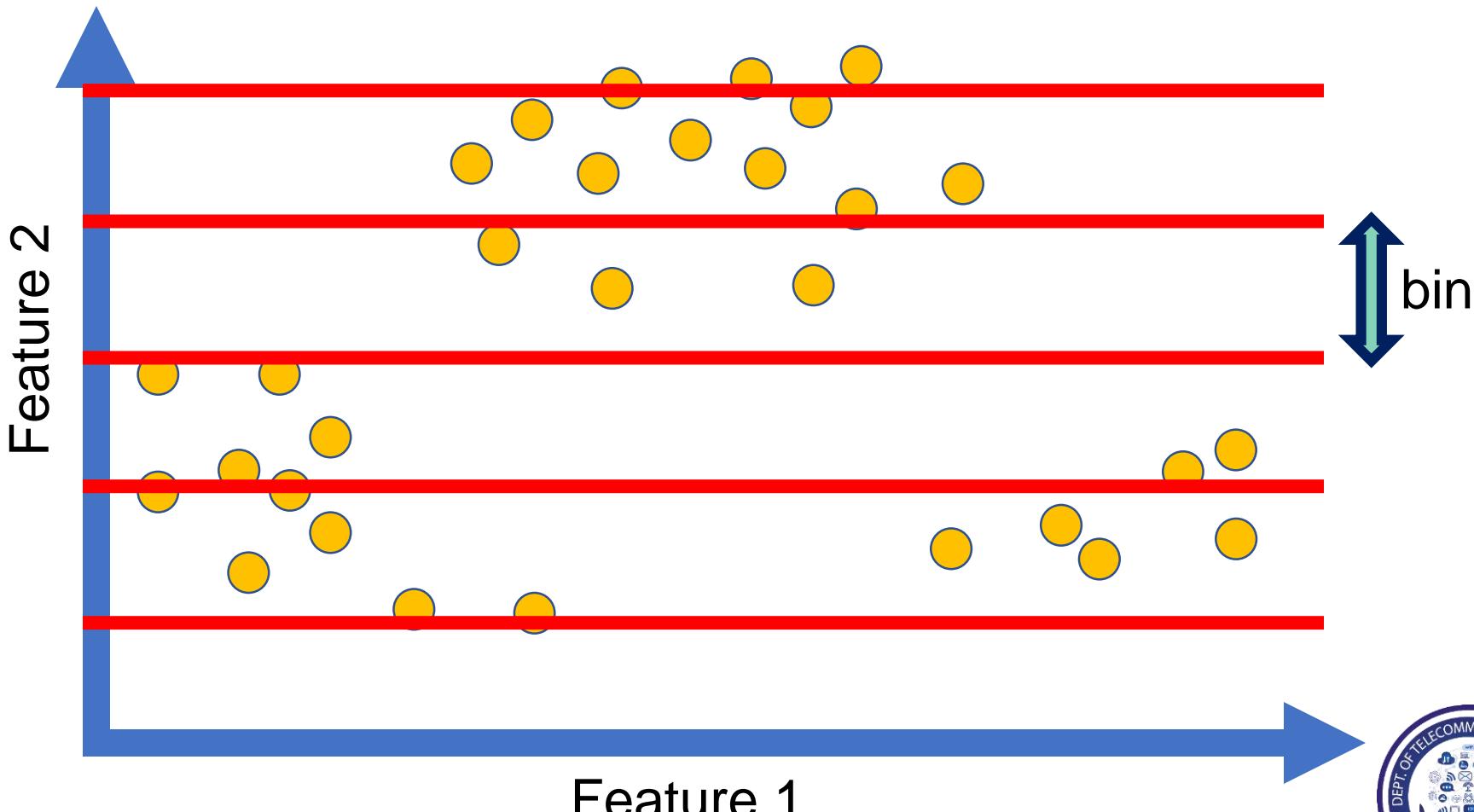
# Image Representations: Histograms

- Probability or count of data in each bin
- Marginal histogram on feature 1



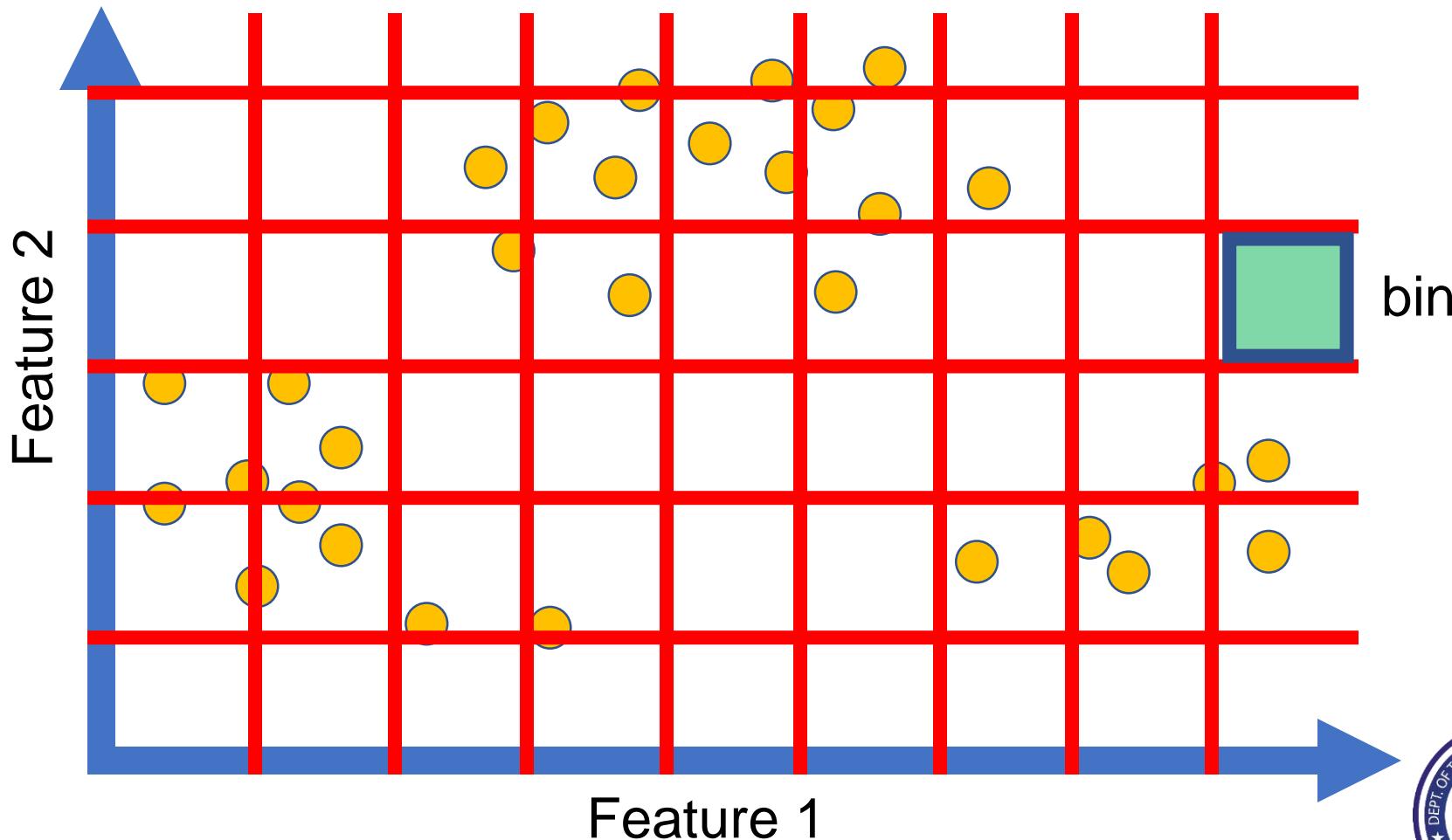
# Image representations: histograms

- Marginal histogram on feature 2

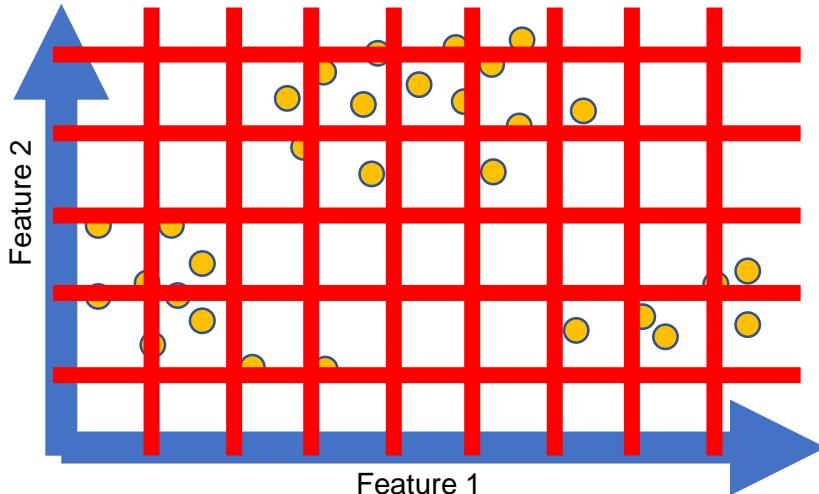


# Image representations: histograms

- Joint histogram

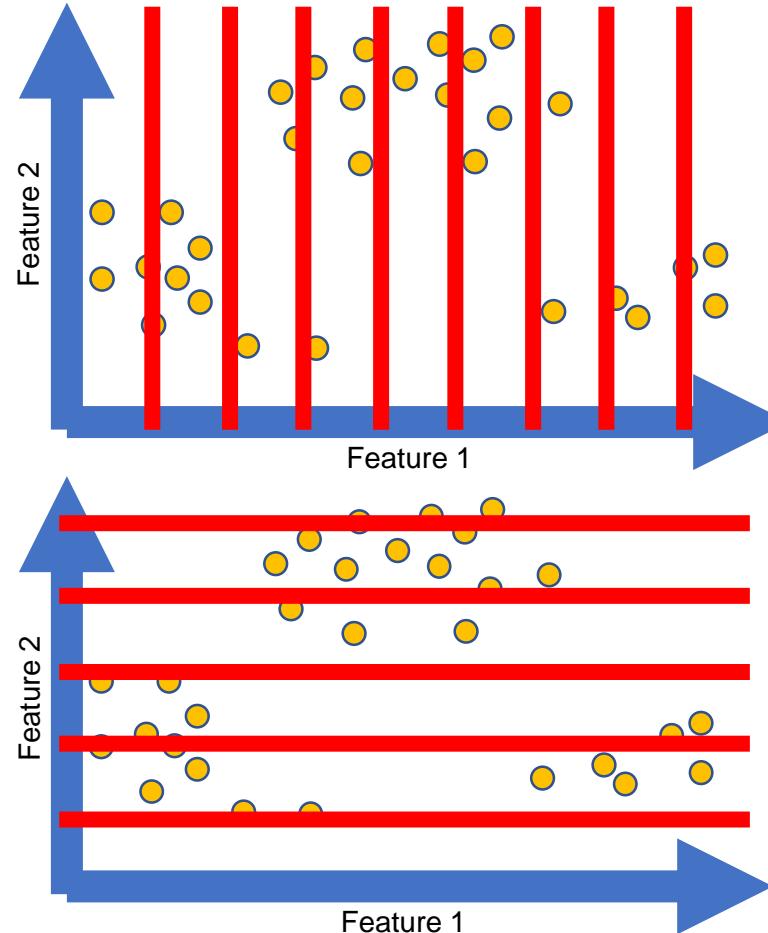


# Modeling multi-dimensional data



## Joint histogram

- Requires lots of data
- Loss of resolution to avoid empty bins

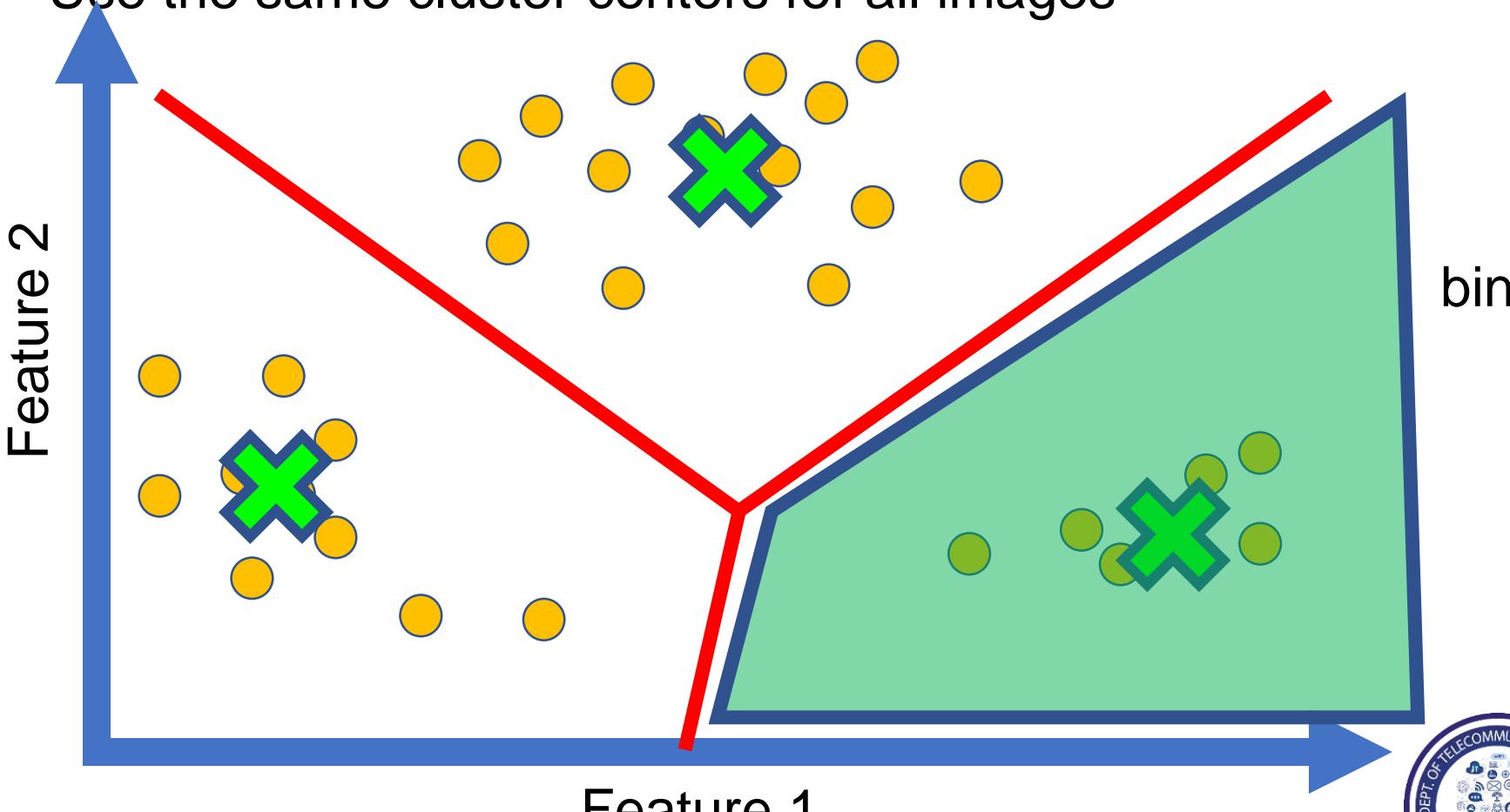


## Marginal histogram

- Requires independent features
- More data/bin than joint histogram

# Modeling multi-dimensional data

- Clustering
- Use the same cluster centers for all images



# Computing histogram distance

- Histogram intersection

$$\text{histint}(h_i, h_j) = 1 - \sum_{m=1}^K \min(h_i(m), h_j(m))$$

- Chi-squared Histogram matching distance

$$\chi^2(h_i, h_j) = \frac{1}{2} \sum_{m=1}^K \frac{[h_i(m) - h_j(m)]^2}{h_i(m) + h_j(m)}$$

- Earth mover's distance  
(Cross-bin similarity measure)
  - minimal cost paid to transform one distribution into the other

[Rubner et al. [The Earth Mover's Distance as a Metric for Image Retrieval](#), IJCV 2000]

Dr. Sander Ali Khowaja



# Histograms: implementation issues

- Quantization

- Grids: fast but applicable only with few dimensions
- Clustering: slower but can quantize data in higher dimensions



Few Bins

Need less data

Coarser representation

Many Bins

Need more data

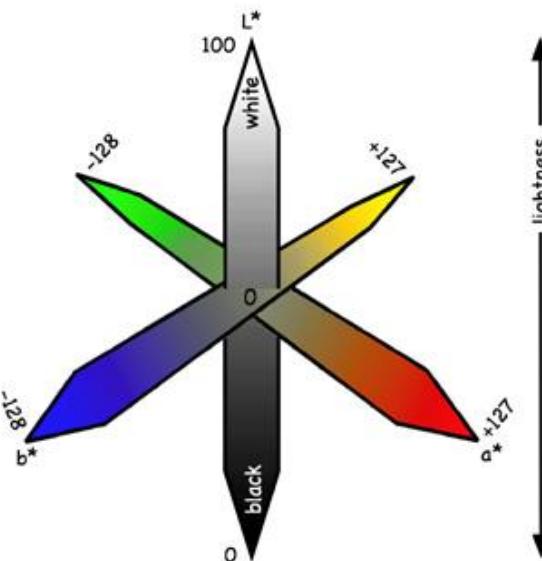
Finer representation

- Matching

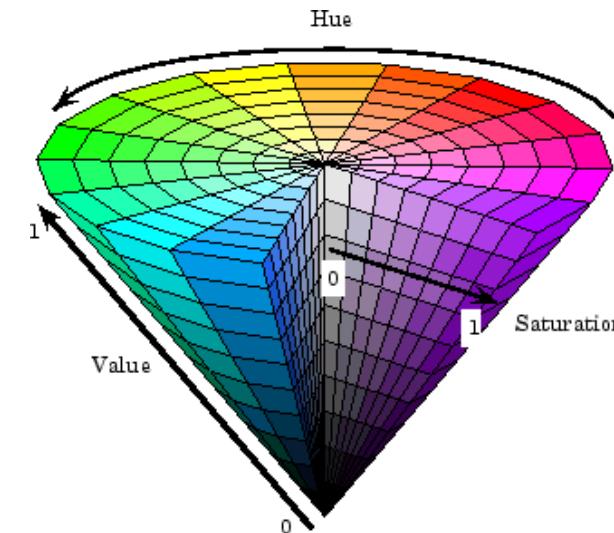
- Histogram intersection or Euclidean may be faster
- Chi-squared often works better
- Earth mover's distance is good for when nearby bins represent similar values

# What kind of things do we compute histograms of?

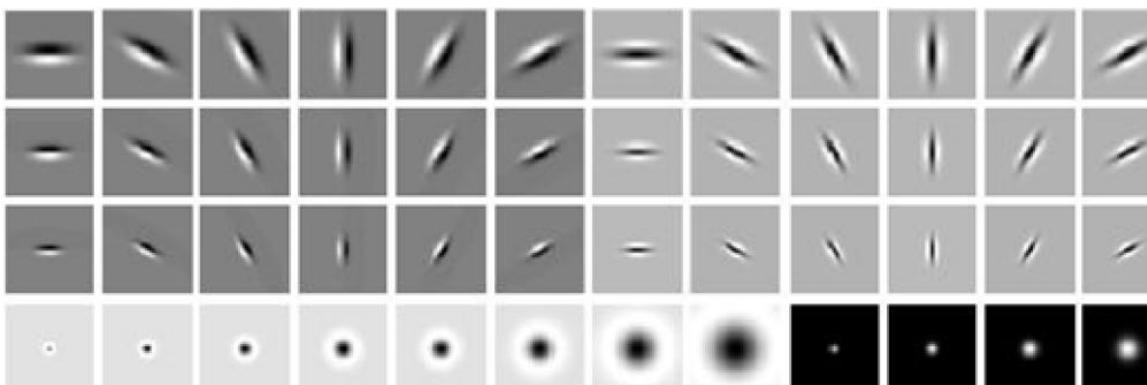
- Color



L\*a\*b\* color space



HSV color space

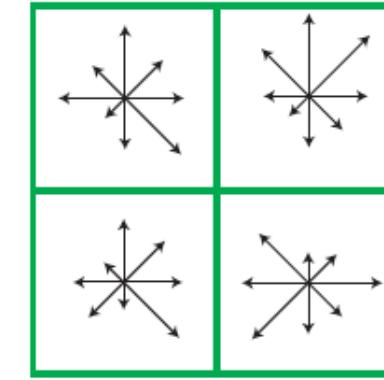
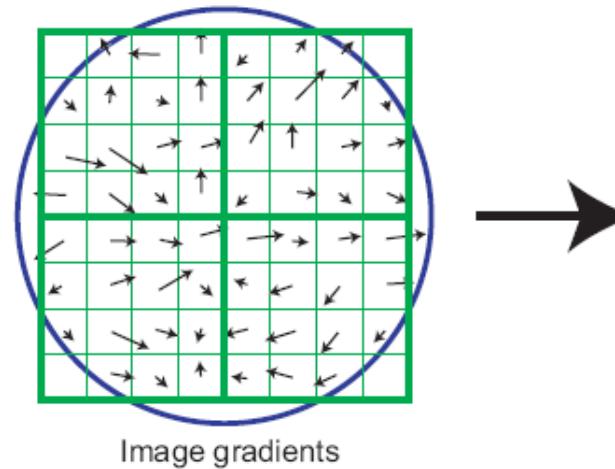


Dr. Sander Ali Knowaja



# What kind of things do we compute histograms of?

- Histograms of descriptors



SIFT – [Lowe IJCV 2004]

- “Bag of 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 our brain via our eyes. For a long time it was believed that the retinal image was processed directly in the visual centers in the brain. In 1960, however, a movie shot by David Hubel and Torsten Wiesel showed that the visual system is much more complex than previously thought. Following the path of the optic nerve from the eye to the various centers of the cerebral cortex, Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a top-down 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.

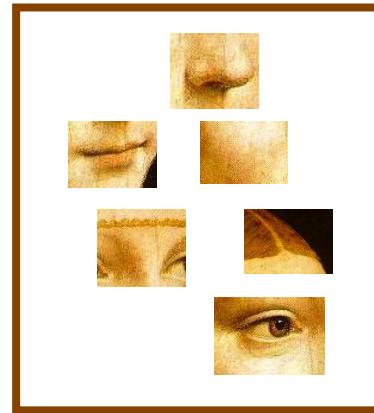
**sensory, brain,  
visual, perception,  
retinal, cerebral cortex,  
eye, cell, optical  
nerve, image  
Hubel, Wiesel**

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 \$600bn last year. The surplus is forecast to rise to \$660bn. That is likely to annoy the US, which is pressuring China's leaders to let the yuan appreciate. The US says the yuan is undervalued and that the Chinese government needs to allow the market to determine its value. It also needs to increase its demand so that it can buy more from the country. China has been buying up large amounts of yuan against the dollar and other currencies and permitted it to trade within a narrow band. But the US wants the yuan to be allowed to trade 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.

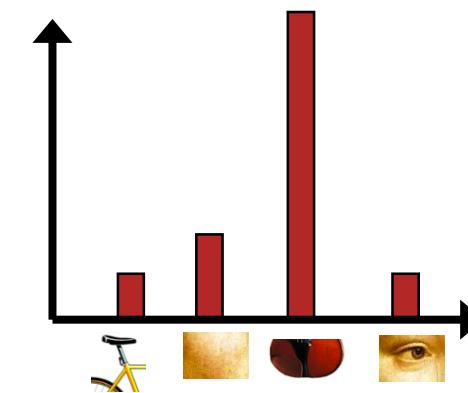
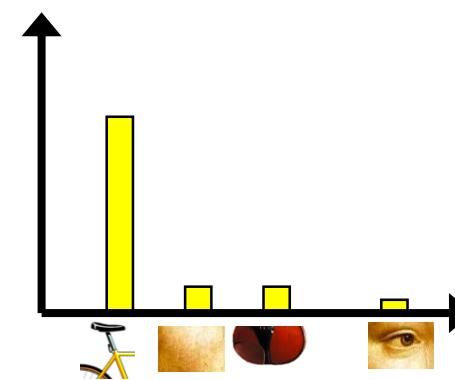
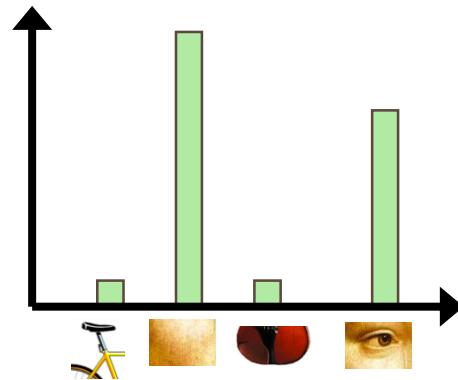
**China, trade,  
surplus, commerce,  
exports, imports, US,  
yuan, bank, domestic,  
foreign, increase,  
trade, value**

# Bag of visual words

- Image patches



- BoW histogram



- Codewords

# Image categorization with bag of words

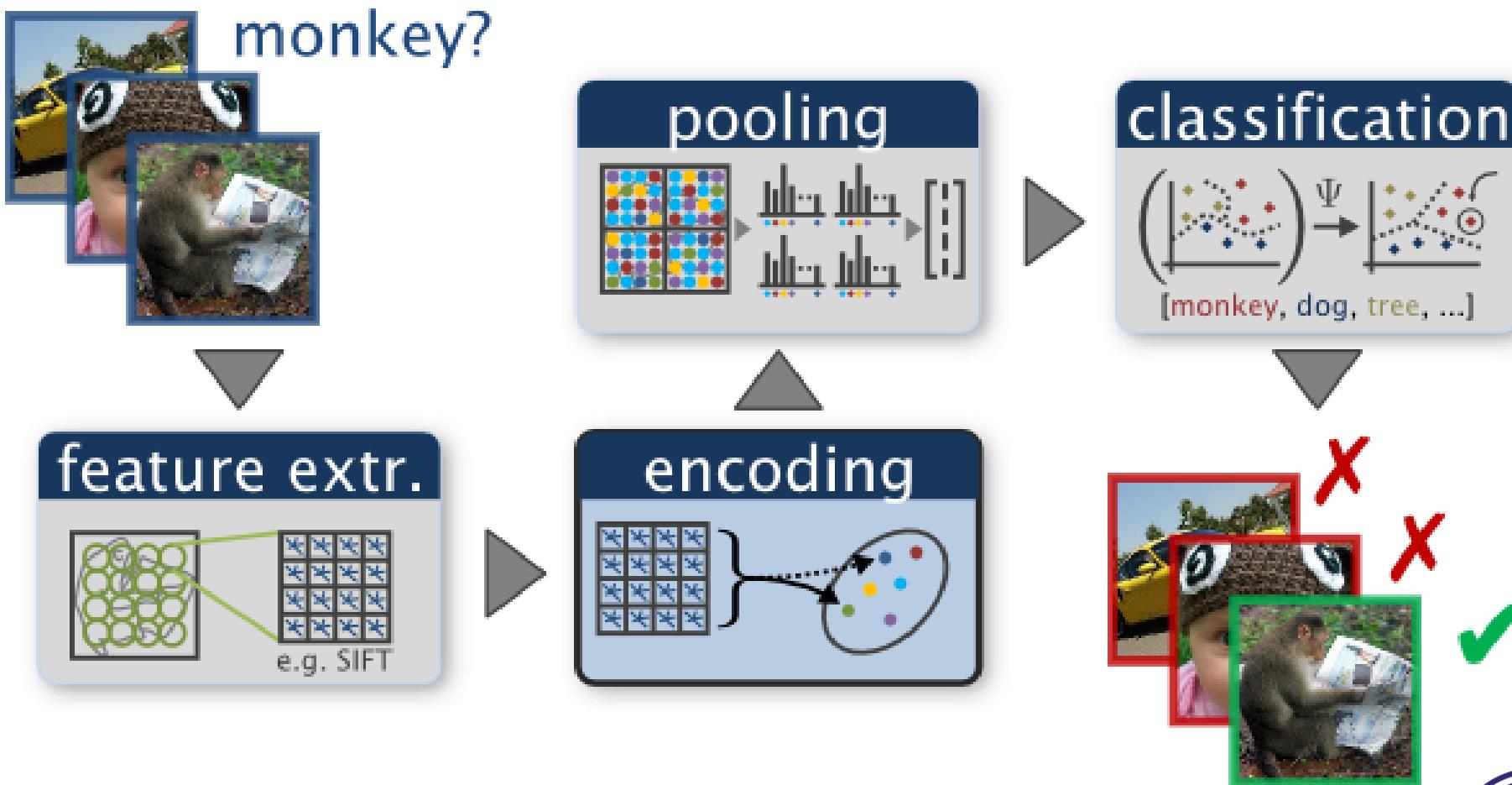
## Training

1. Extract keypoints and descriptors for all training images
2. Cluster descriptors
3. Quantize descriptors using cluster centers to get “visual words”
4. Represent each image by normalized counts of “visual words”
5. Train classifier on labeled examples using histogram values as features

## Testing

1. Extract keypoints/descriptors and quantize into visual words
2. Compute visual word histogram
3. Compute label or confidence using classifier

# Bag of visual words image classification

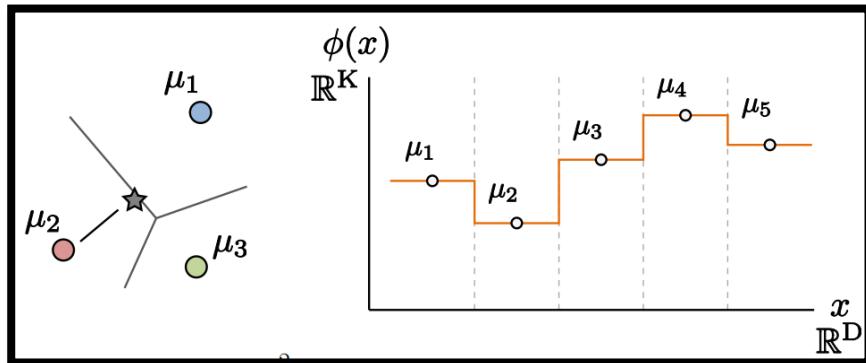


[Chatfield et al. BMVC 2011]  
Dr. Sander Ali Khawaja

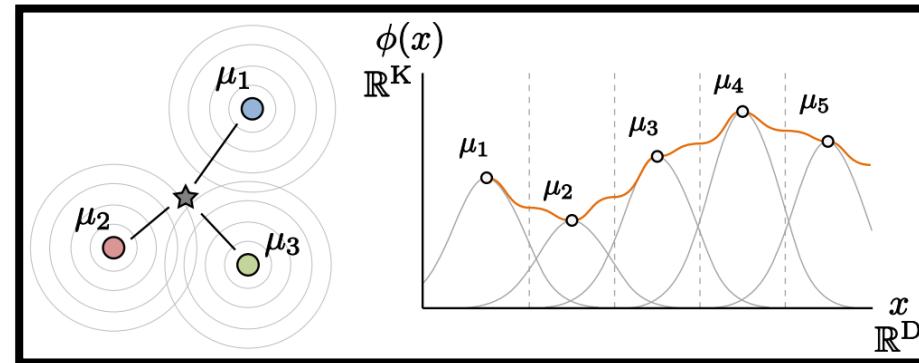


# Feature encoding

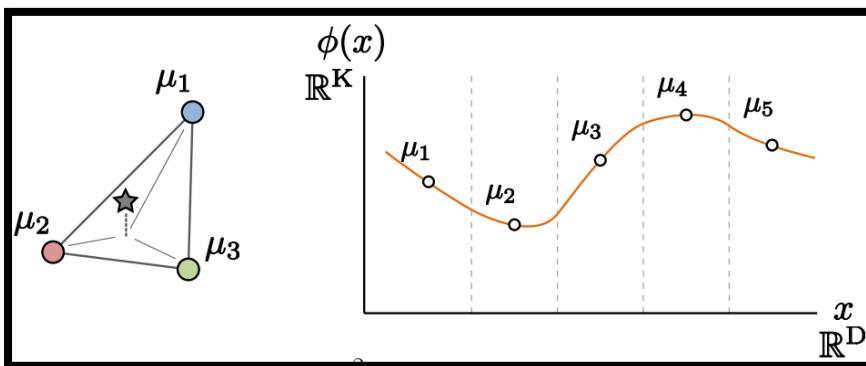
- Hard/soft assignment to clusters



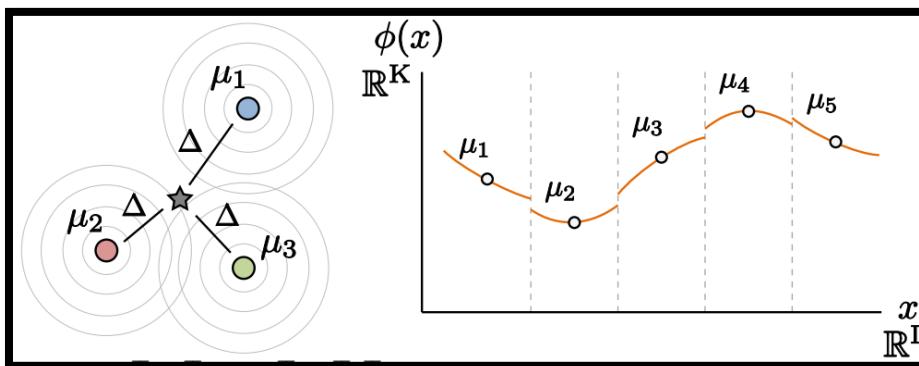
Histogram encoding



Kernel codebook encoding



Locality constrained encoding



Fisher encoding

# Fisher vector encoding

- Fit Gaussian Mixture Models

$$\Theta = (\mu_k, \Sigma_k, \pi_k : k = 1, \dots, K)$$

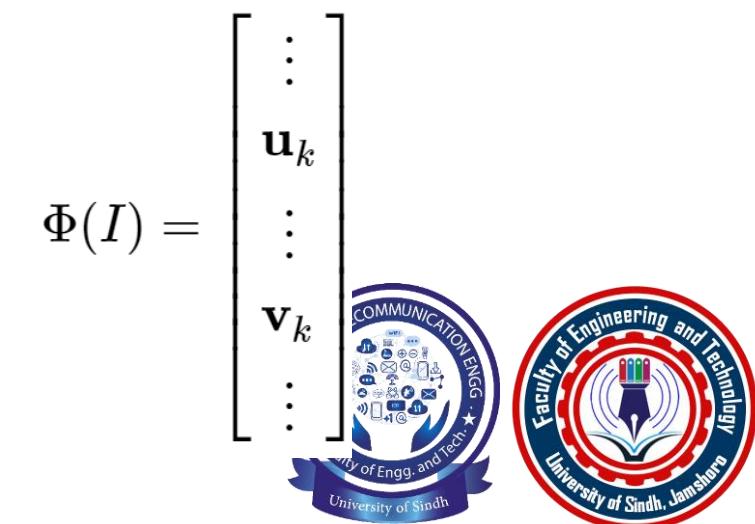
- Posterior probability

$$q_{ik} = \frac{\exp\left[-\frac{1}{2}(\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k)\right]}{\sum_{t=1}^K \exp\left[-\frac{1}{2}(\mathbf{x}_i - \mu_t)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_t)\right]}$$

- First and second order differences to cluster k

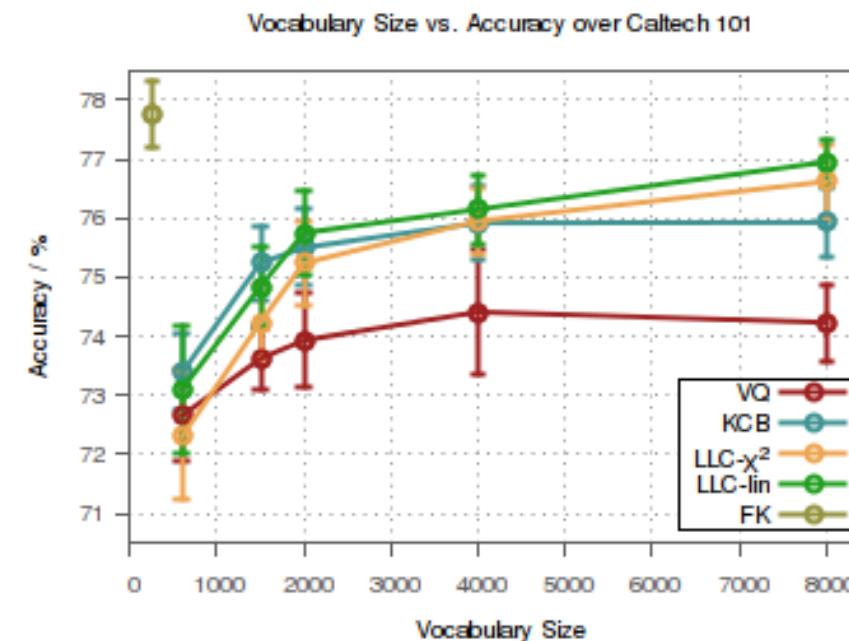
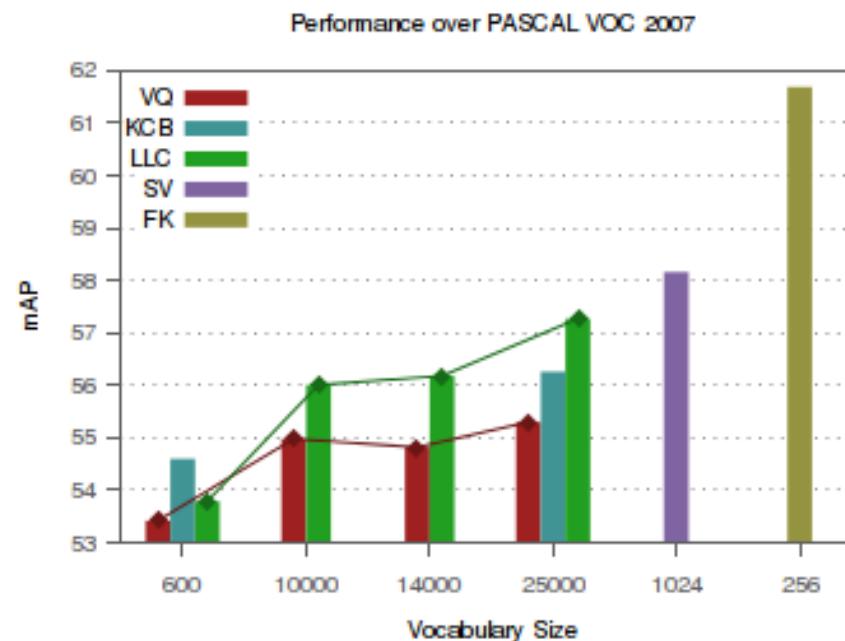
$$u_{jk} = \frac{1}{N\sqrt{\pi_k}} \sum_{i=1}^N q_{ik} \frac{x_{ji} - \mu_{jk}}{\sigma_{jk}},$$

$$v_{jk} = \frac{1}{N\sqrt{2\pi_k}} \sum_{i=1}^N q_{ik} \left[ \left( \frac{x_{ji} - \mu_{jk}}{\sigma_{jk}} \right)^2 - 1 \right]$$

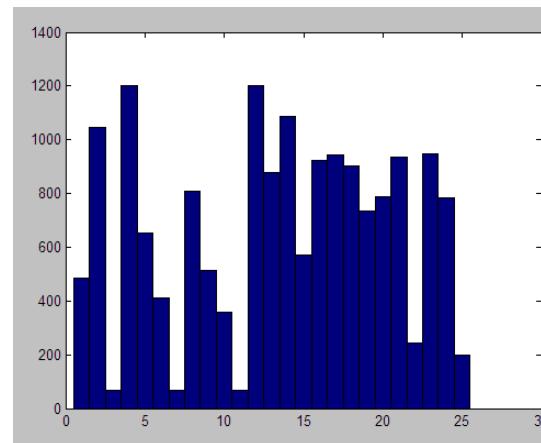


# Performance comparisons

- Fisher vector encoding outperforms others
- Higher-order statistics helps

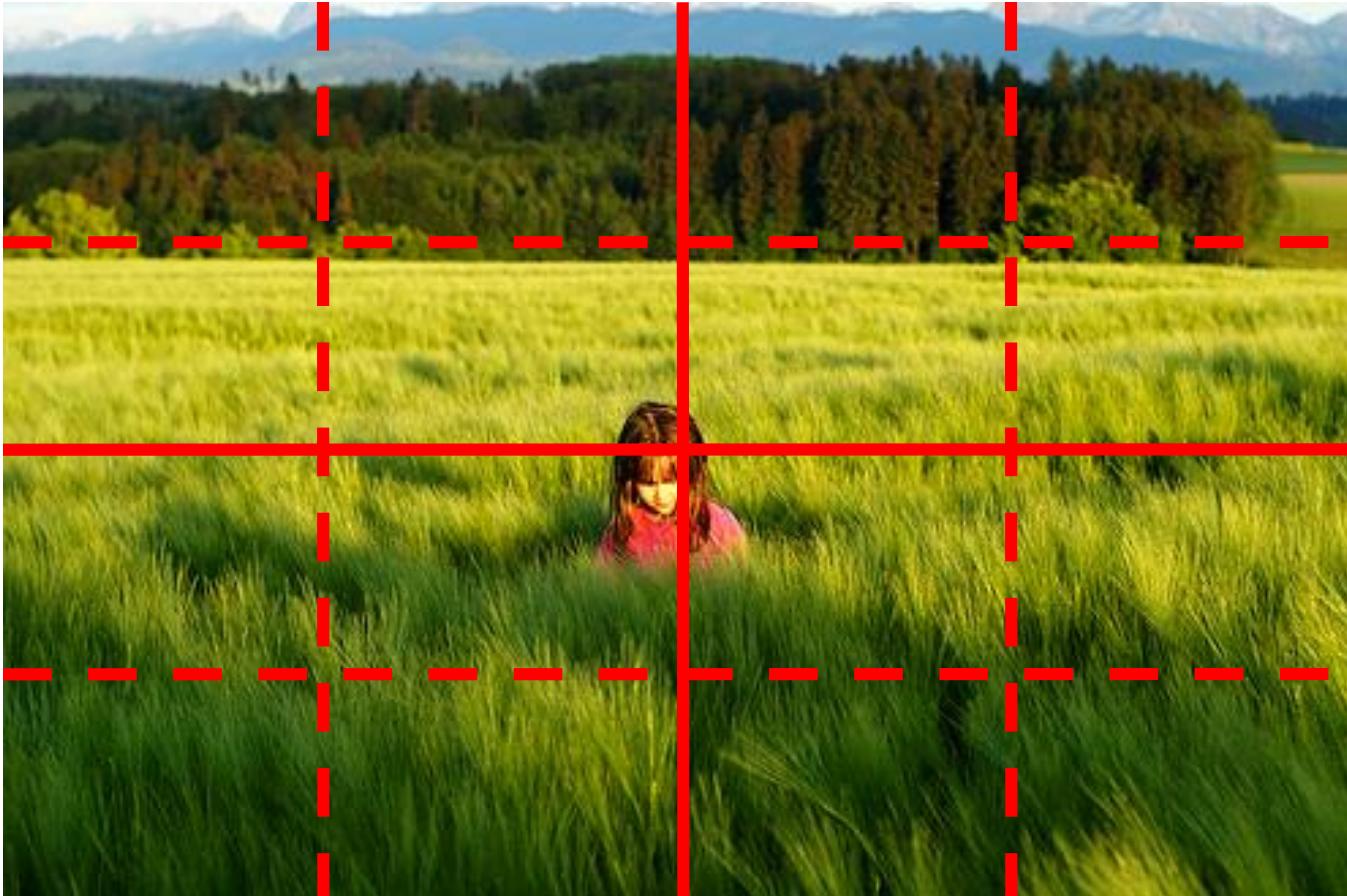


# But what about spatial layout?



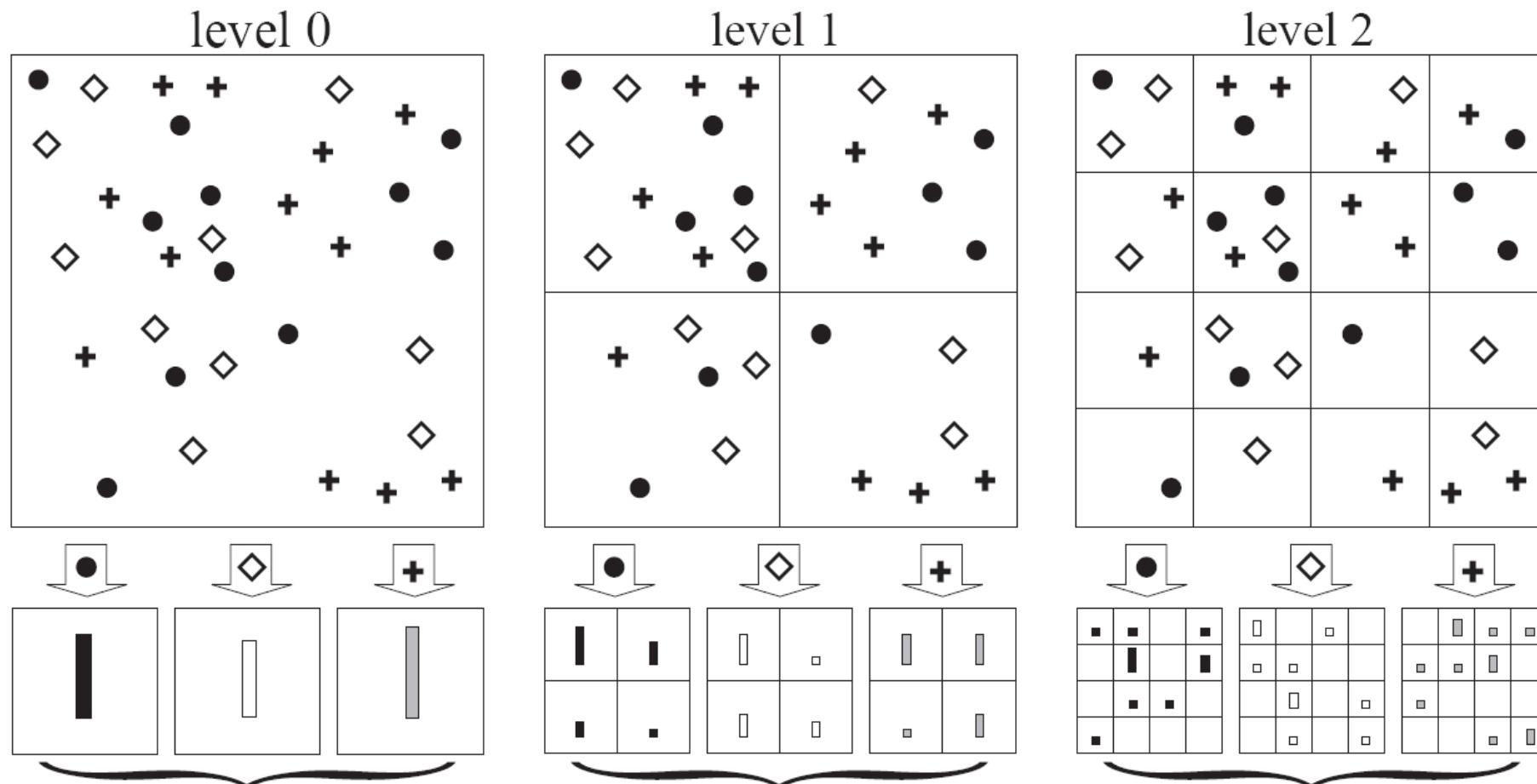
All of these images have the same color histogram

# Spatial pyramid



Compute histogram in each spatial bin

# Spatial pyramid



High number of features – PCA to reduce dimensionality

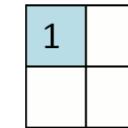
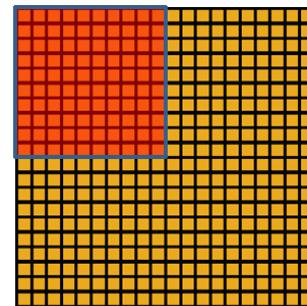
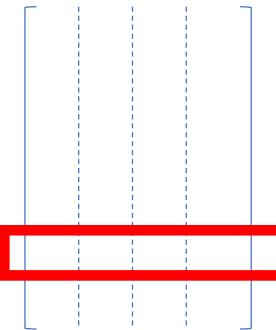
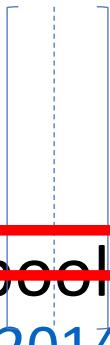
[Lazebnik et al. CVPR 2006]

Dr. Sander Ali Khowaja



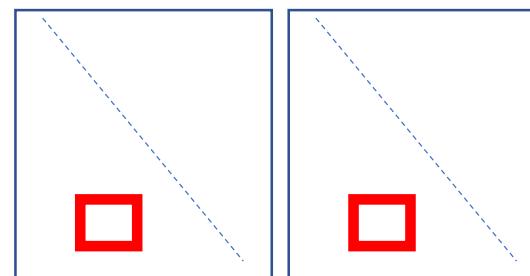
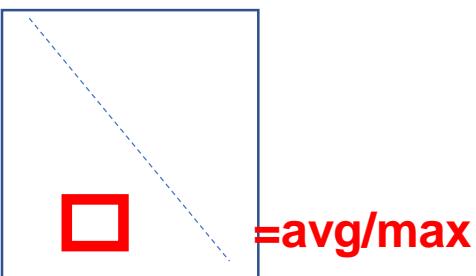
# Pooling

- Average/max pooling



Convolved      Pooled  
feature      feature

Source: Unsupervised Feature  
Learning and Deep Learning

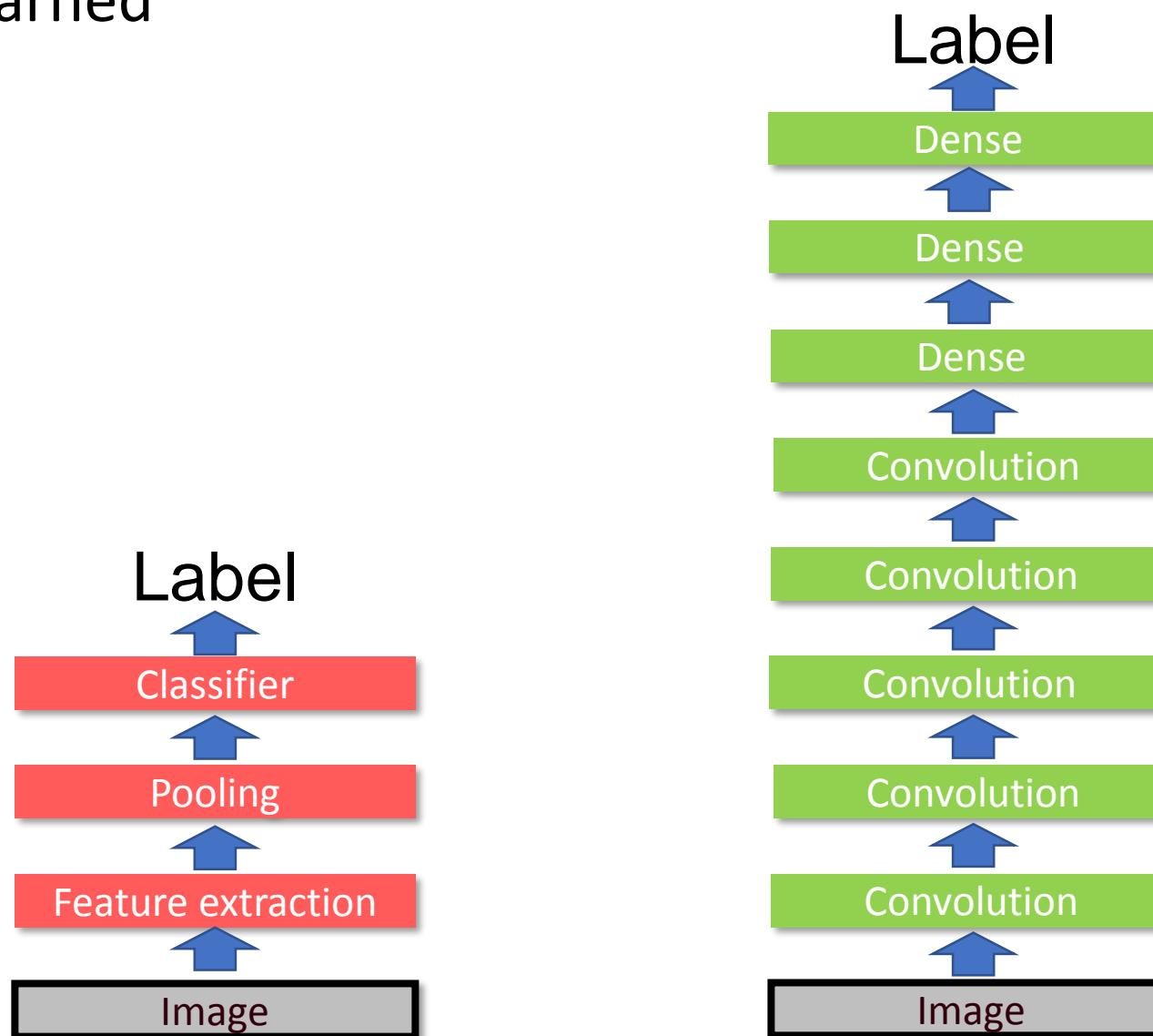


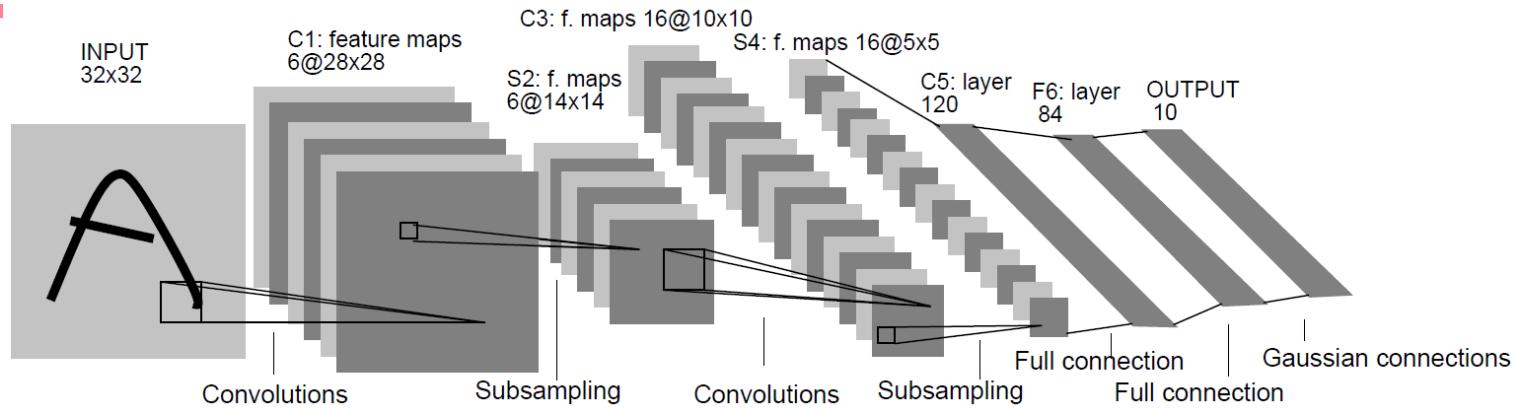
$$\mathbf{G}_{avg}(R_j) = \frac{1}{|F_{R_j}|} \sum_{i:(\mathbf{f}_i \in R_j)} \mathbf{x}_i \cdot \mathbf{x}_i^\top$$

$$\mathbf{G}_{max}(R_j) = \max_{i:(\mathbf{f}_i \in R_j)} \mathbf{x}_i \cdot \mathbf{x}_i^\top$$

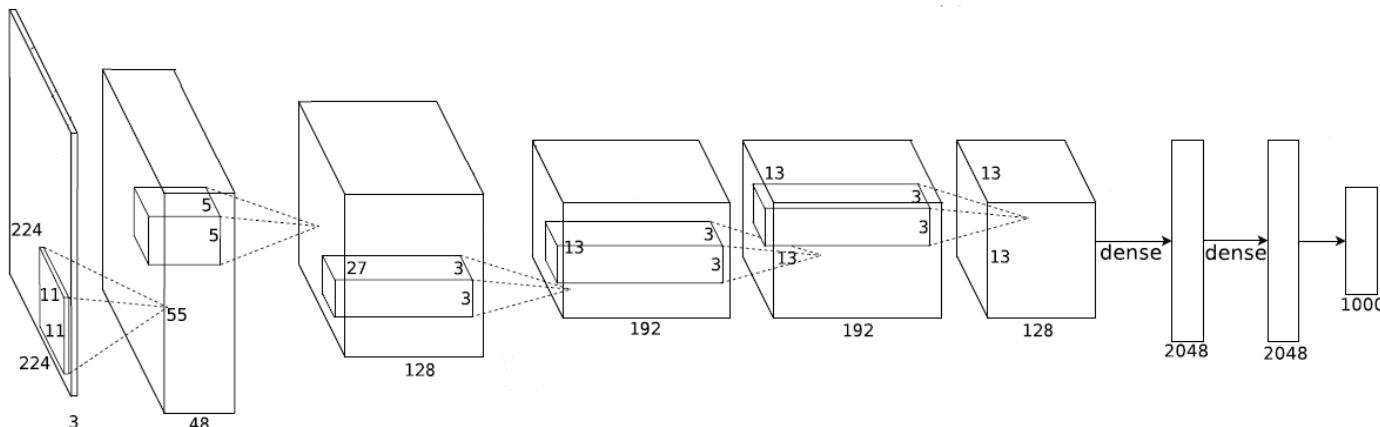
# Shallow vs. deep learning

- Engineered vs. learned features

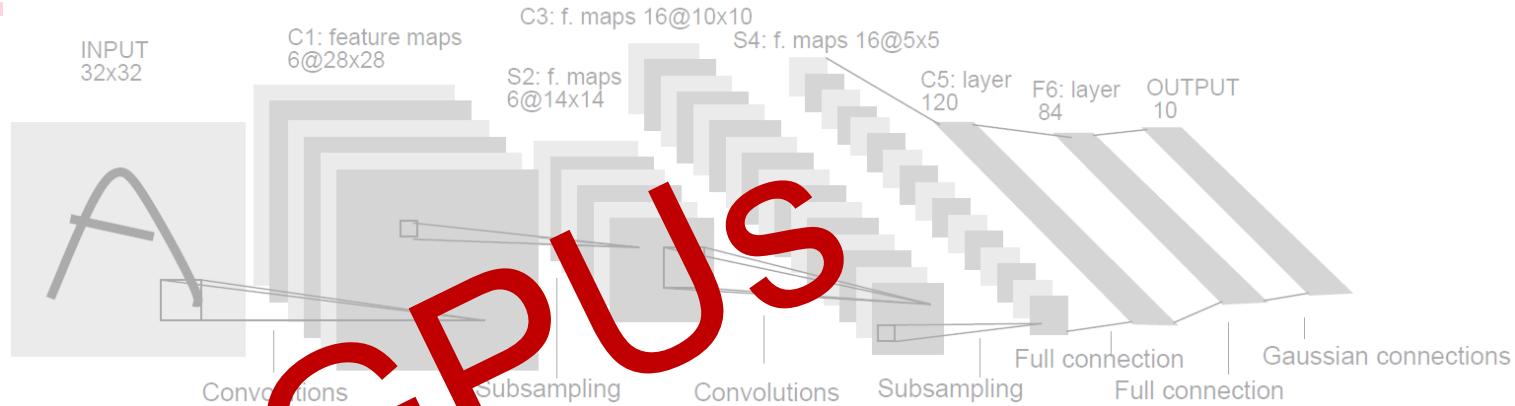




**Gradient-Based Learning Applied to Document Recognition**, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, **1998**

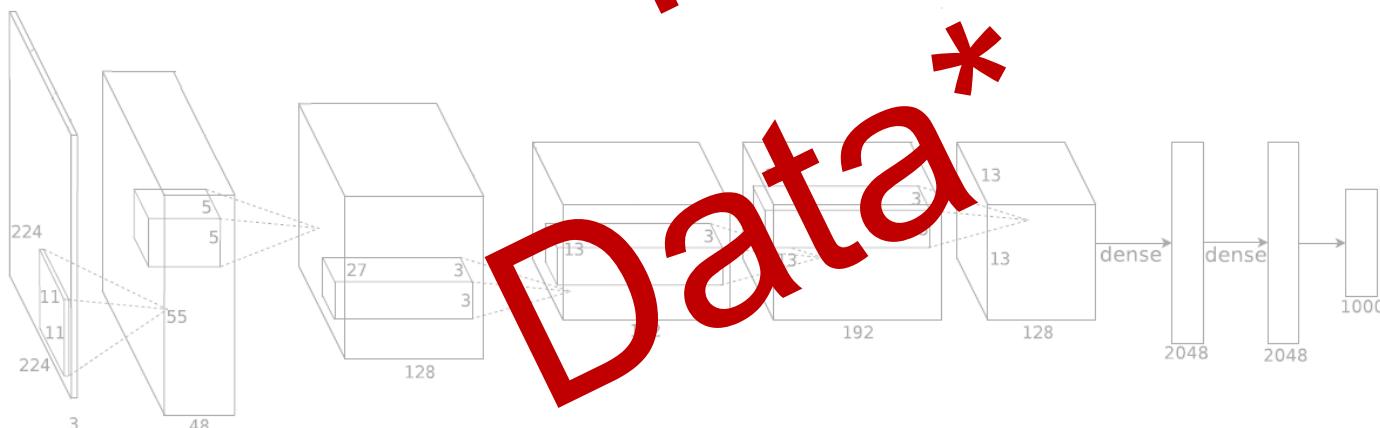


**Imagenet Classification with Deep Convolutional Neural Networks**, Krizhevsky, Sutskever, and Hinton, NIPS **2012**



Gradient-Based Learning Applied to Document  
Recognition, LeCun, Bottou, Bengio and Haffner, Proc. of  
the IEEE, 1998

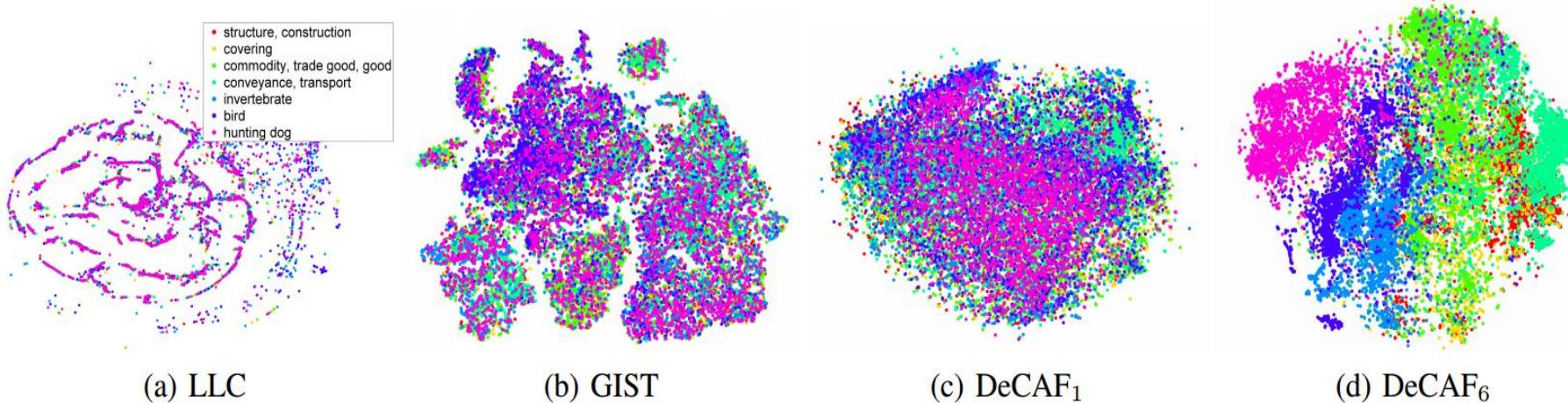
+



Imagenet Clas:  
Networks, Kriz

\* Rectified activations and dropout

# Convolutional activation features



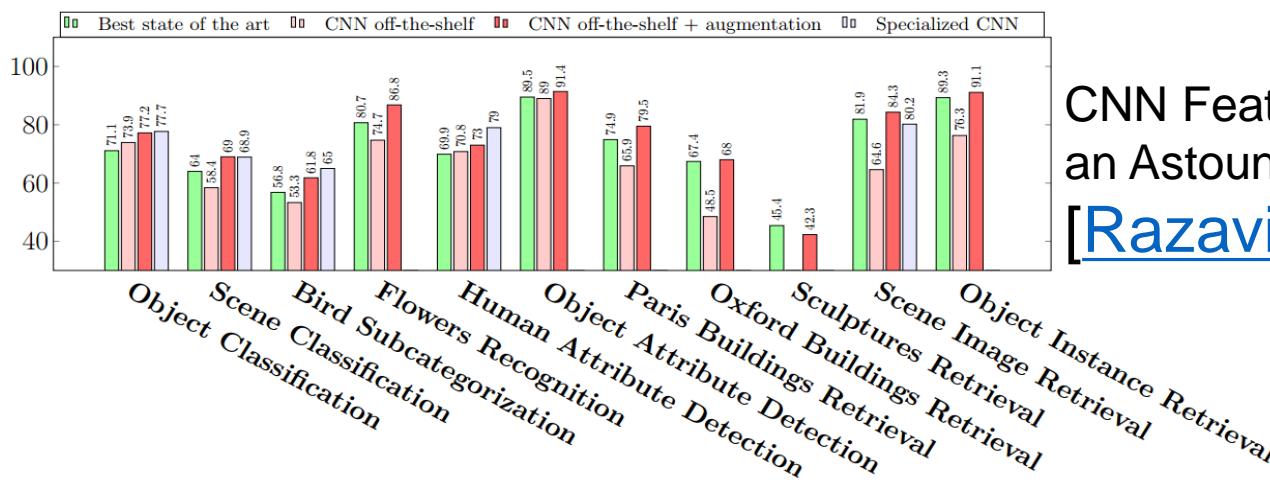
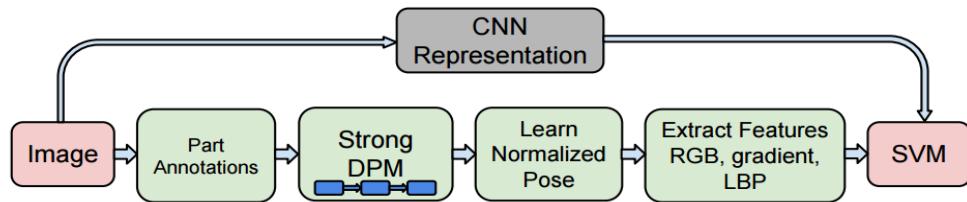
(a) LLC

(b) GIST

(c) DeCAF<sub>1</sub>

(d) DeCAF<sub>6</sub>

[Donahue et al. ICML 2013]

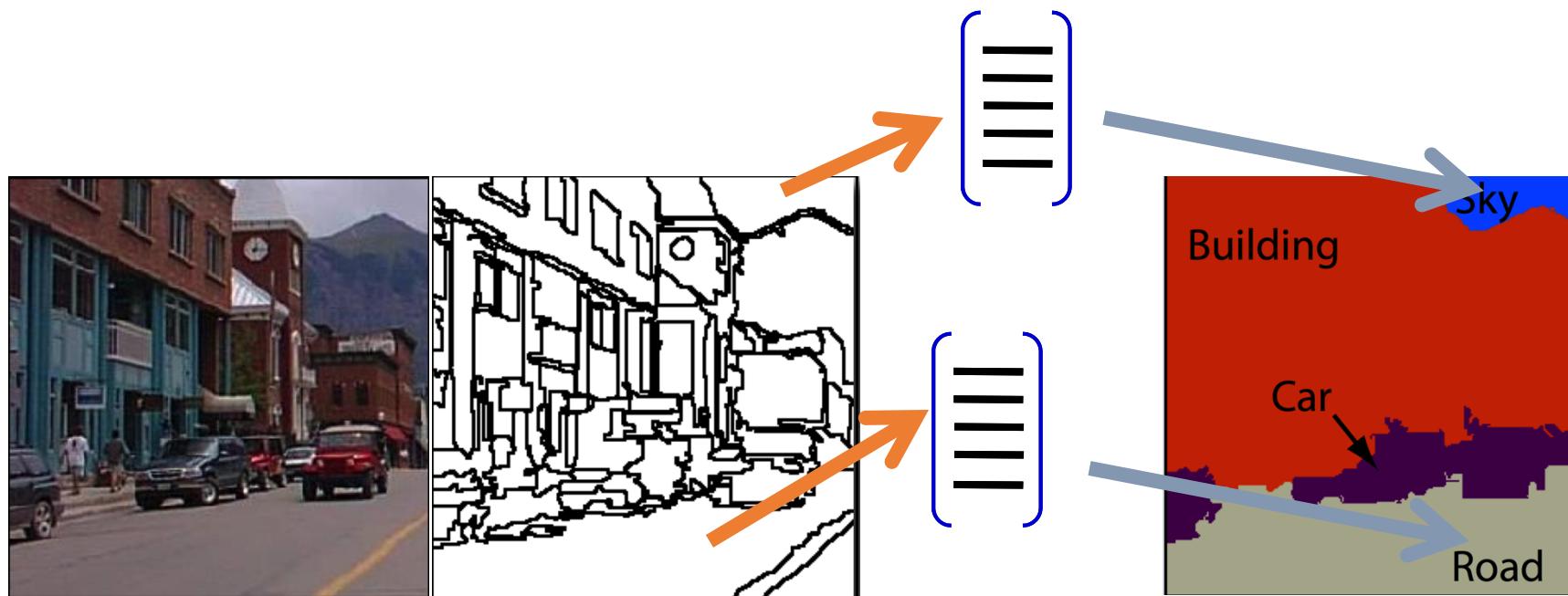


CNN Features off-the-shelf:  
an Astounding Baseline for Recognition  
[Razavian et al. 2014]



# Region representation

- Segment the image into superpixels
- Use features to represent each image segment



# Region representation

- Color, texture, BoW
  - Only computed within the local region
- Shape of regions
- Position in the image



# Working with regions

- Spatial support is important – multiple segmentation



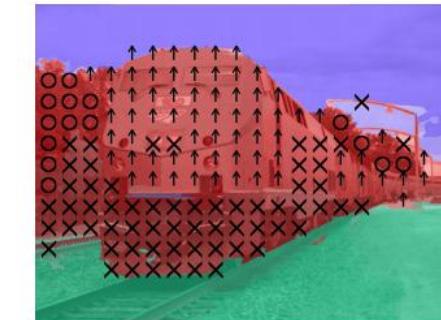
(a) Input



(b) Superpixels



(c) Multiple Hypotheses



(d) Geometric Labels

Geometric context [[Hoiem et al. ICCV 2005](#)]

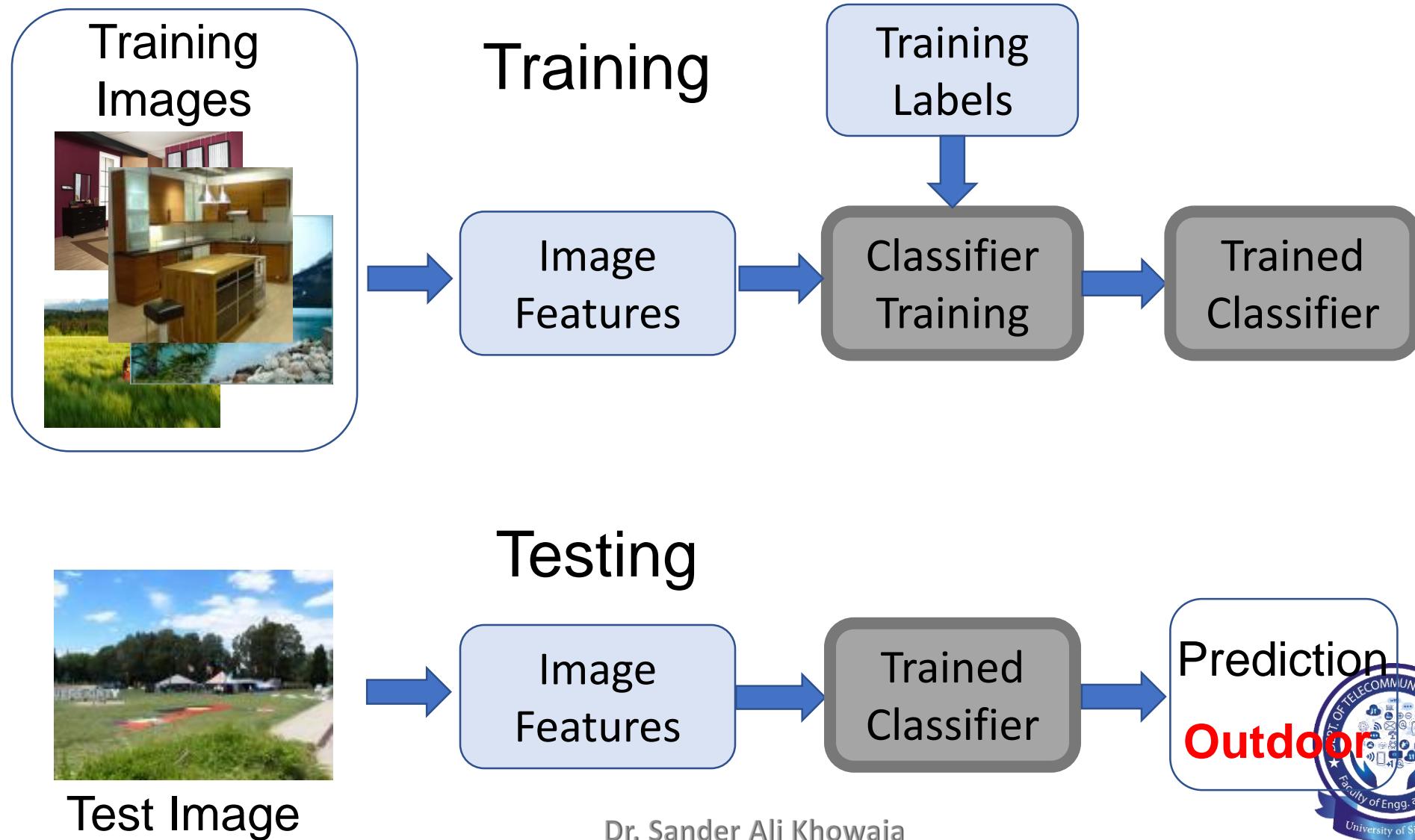
- Spatial consistency – MRF smoothing

# Things to remember

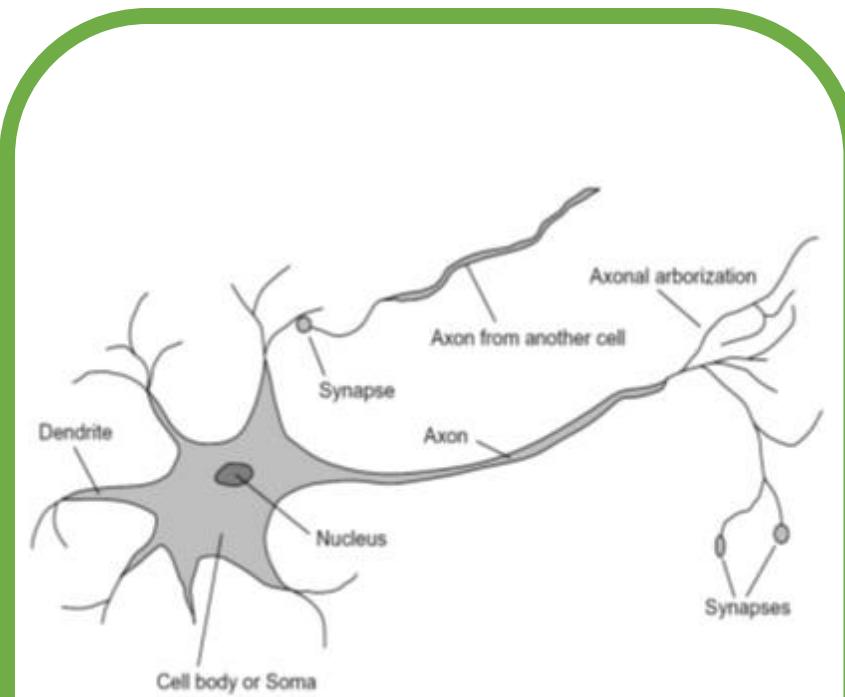
- Visual categorization help transfer knowledge
- Image features
  - Coverage, concision, directness
  - Color, gradients, textures, motion, descriptors
  - Histogram, feature encoding, and pooling
  - CNN as features
- Image/region categorization



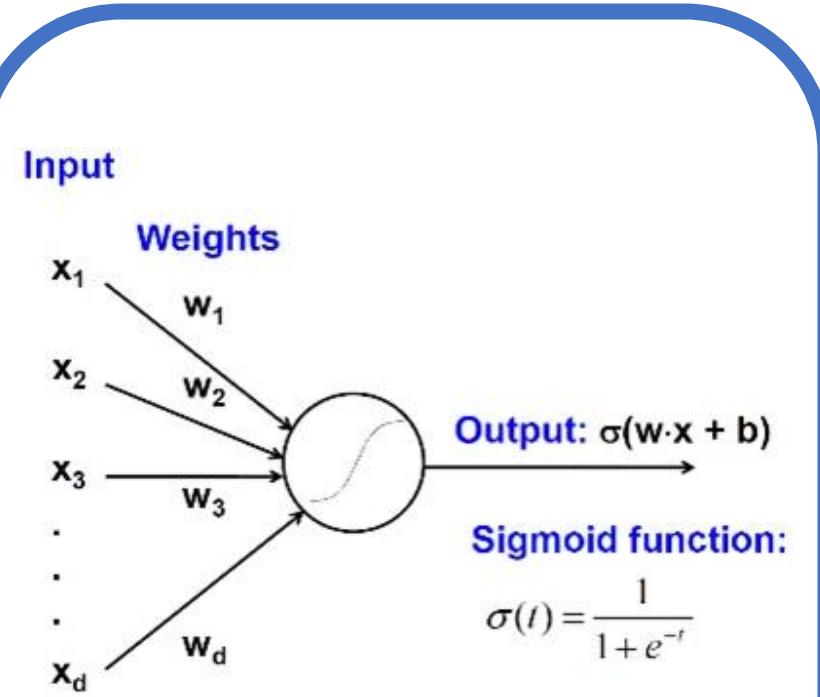
# Next lecture - Classifiers



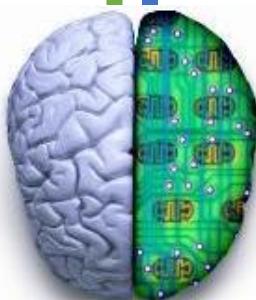
# Biological neuron and Perceptrons



# A biological neuron



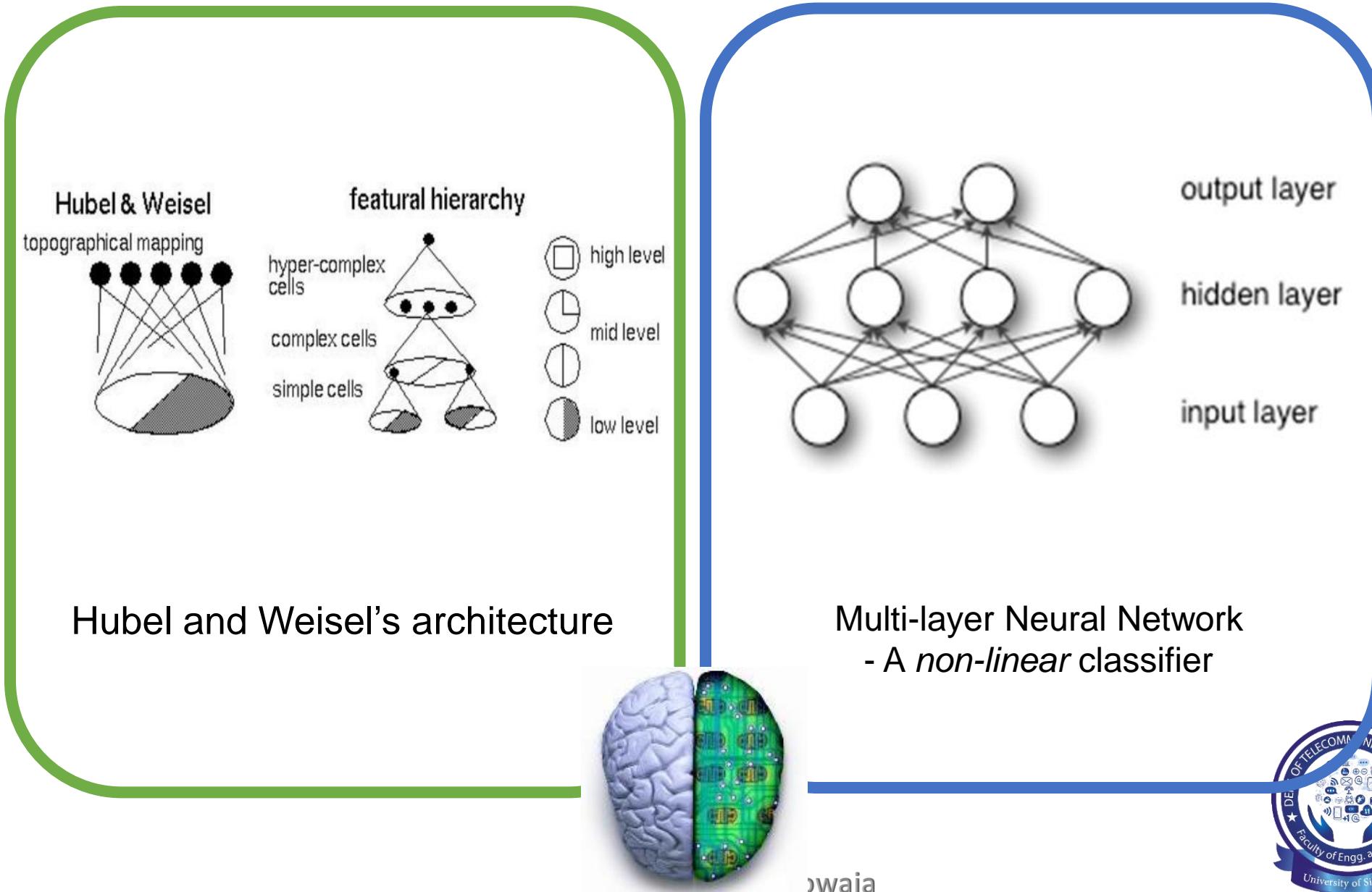
## An artificial neuron (Perceptron) - a linear classifier



)waja



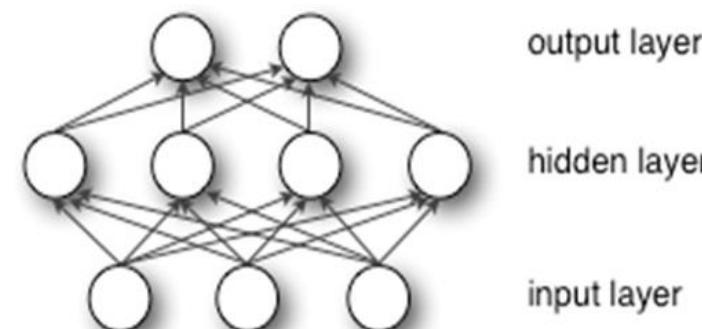
# Hubel/Wiesel Architecture and Multi-layer Neural Network



# Multi-layer Neural Network

- A non-linear classifier
- **Training:** find network weights  $\mathbf{w}$  to minimize the error between true training labels  $y_i$  and estimated labels  $f_{\mathbf{w}}(\mathbf{x}_i)$

- Minimization can be done if  $f$  is differentiable
- This training method is called **back-propagation**

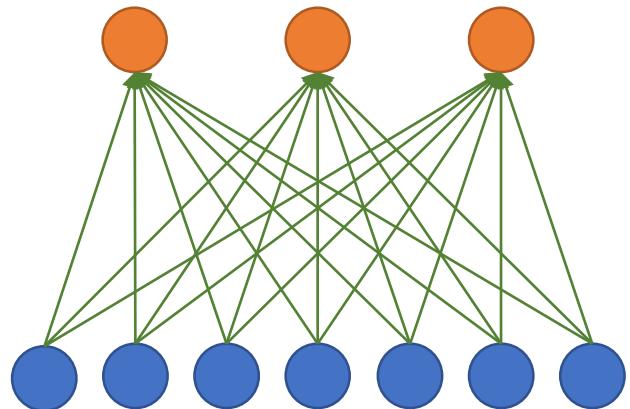


# Convolutional Neural Networks

- Also known as CNN, ConvNet, DCN
- CNN = a multi-layer neural network with
  1. Local connectivity
  2. Weight sharing



# CNN: Local Connectivity

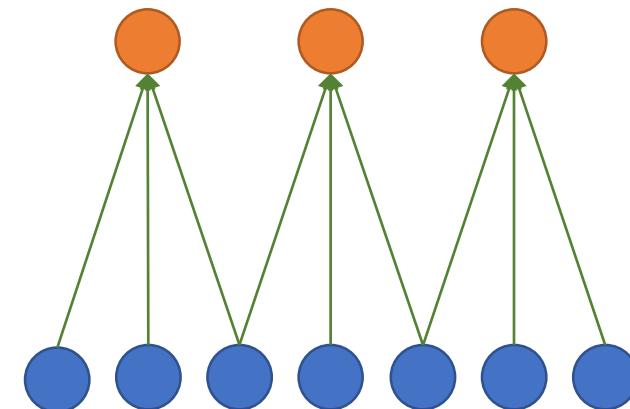


Hidden layer

Input layer

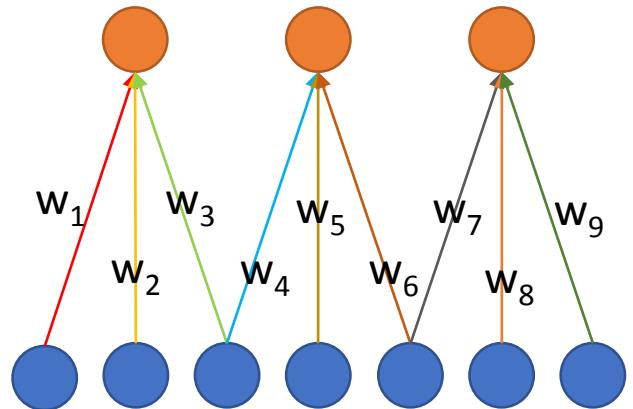
## Global connectivity

- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
  - Global connectivity:  $3 \times 7 = 21$
  - Local connectivity:  $3 \times 3 = 9$



## Local connectivity

# CNN: Weight Sharing

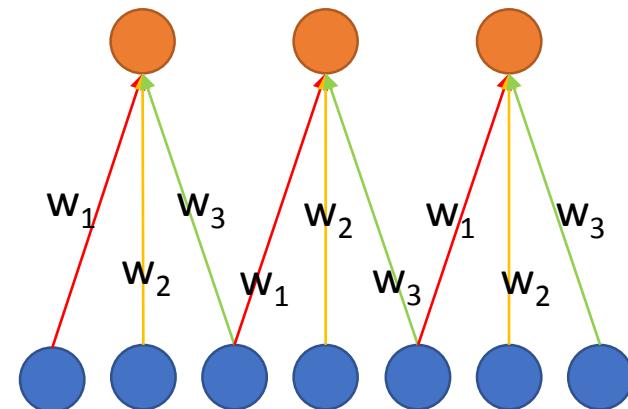


**Without weight sharing**

- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
  - Without weight sharing:  $3 \times 3 = 9$
  - With weight sharing :  $3 \times 1 = 3$

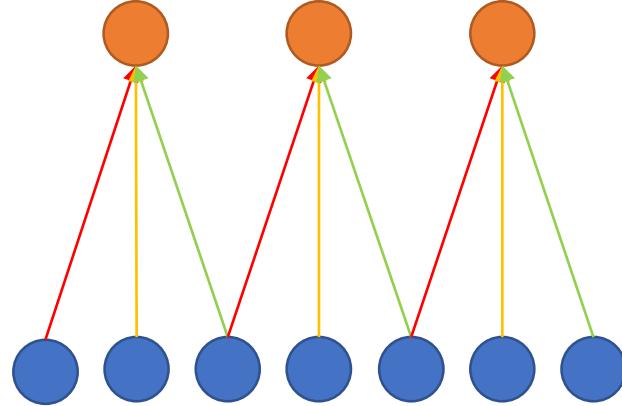
Hidden layer

Input layer

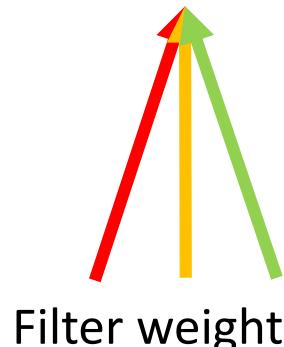


**With weight sharing**

# CNN with multiple input channels



**Single input channel**

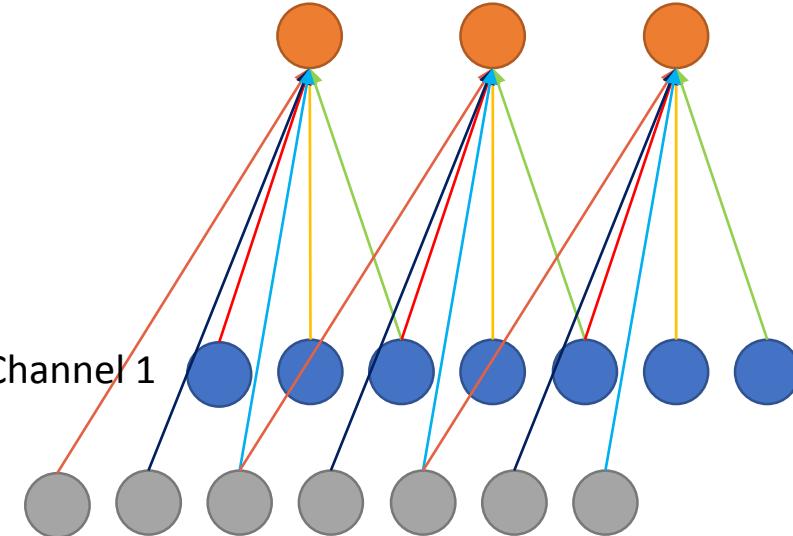


**Filter weights**

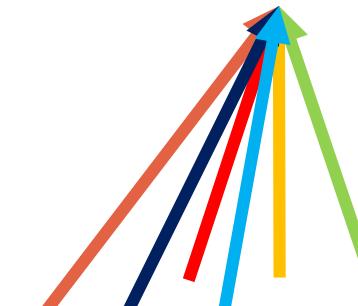
Hidden layer

Input layer

Channel 2

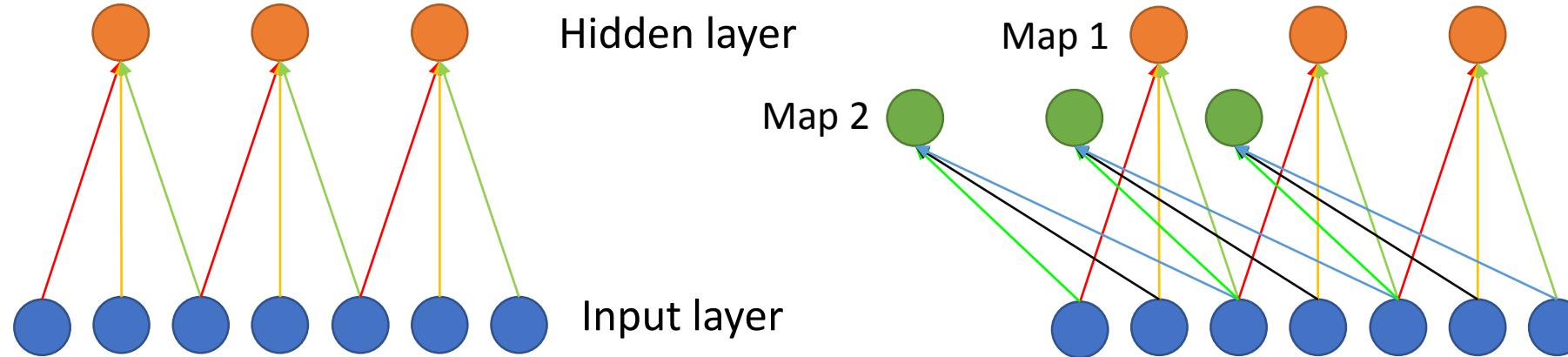


**Multiple input channels**

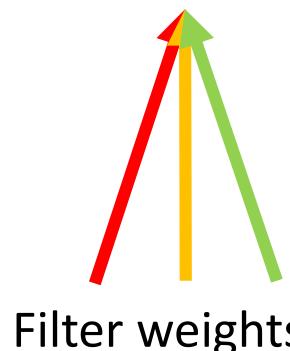


**Filter weights**

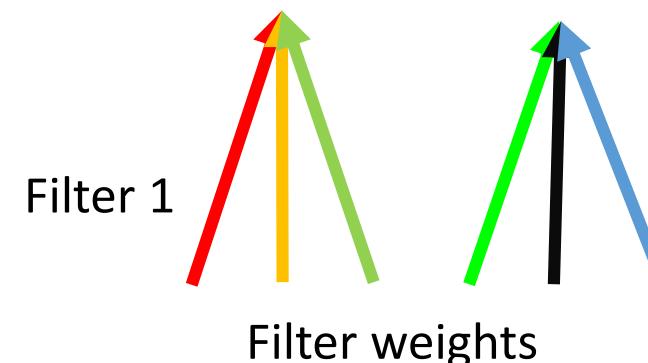
# CNN with multiple output maps



**Single output map**

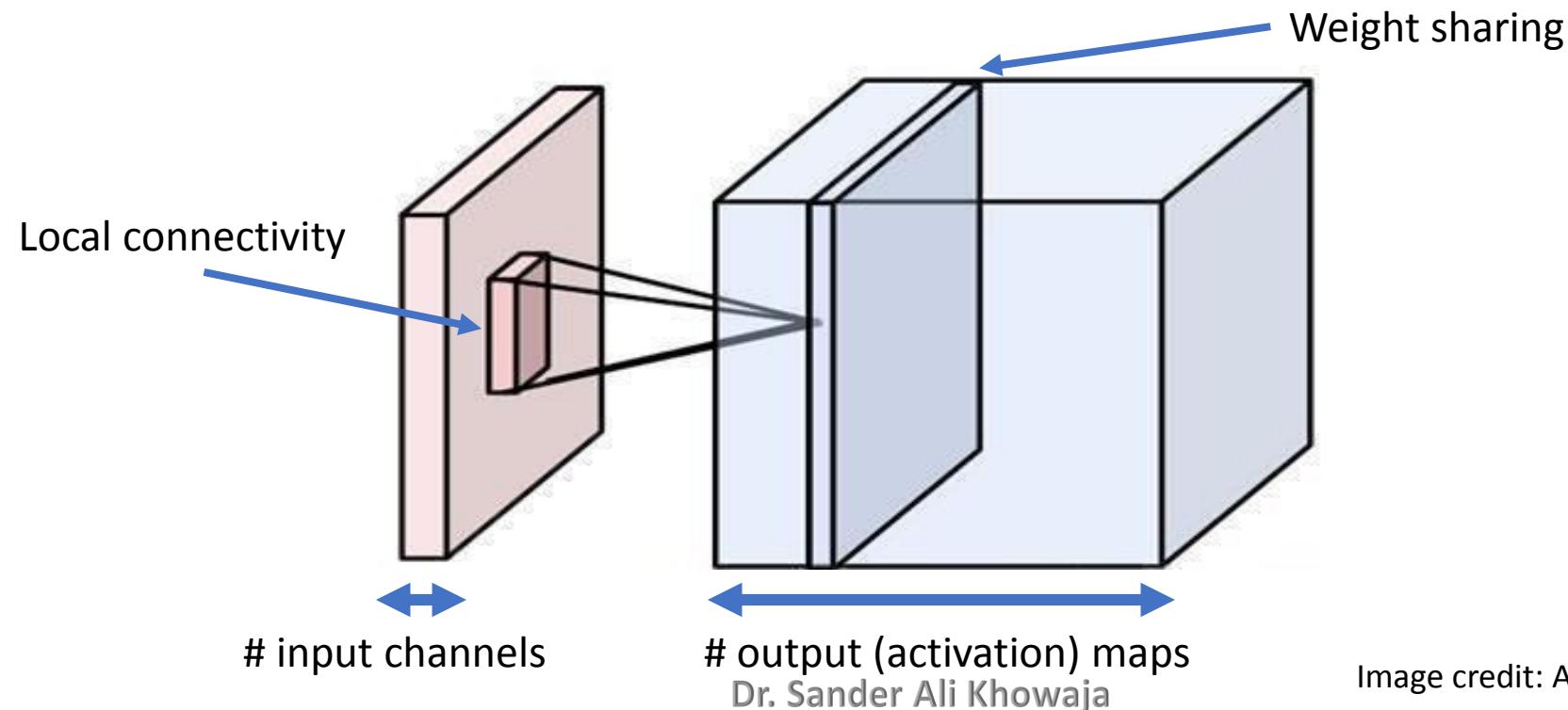


**Multiple output maps**



# Putting them together

- Local connectivity
- Weight sharing
- Handling multiple input channels
- Handling multiple output maps

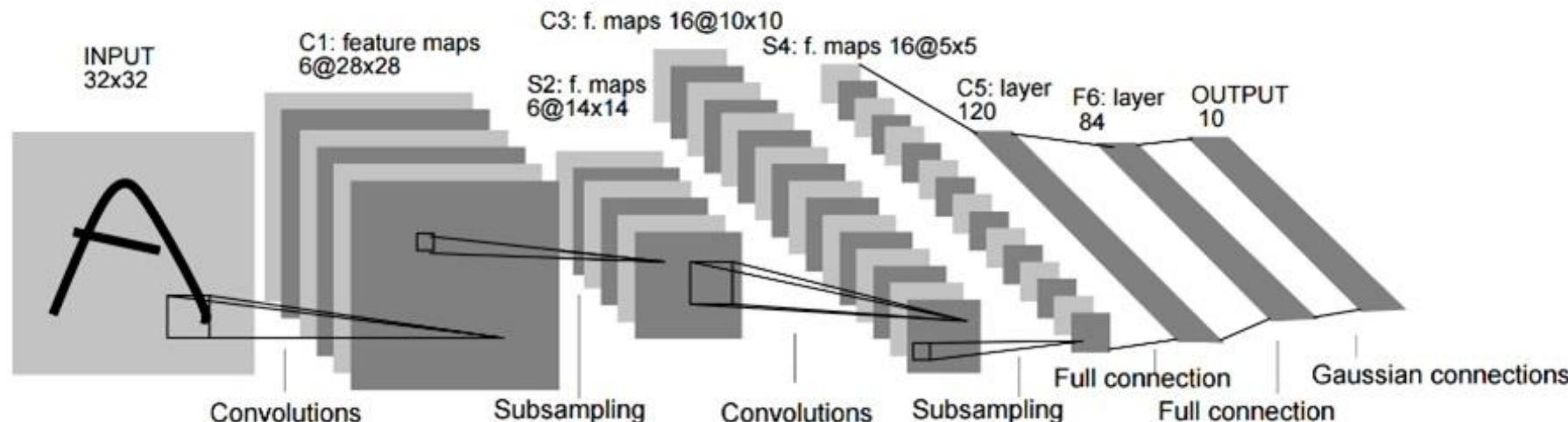


Dr. Sander Ali Khowaja

Image credit: A. Karpathy



# LeNet [LeCun et al. 1998]



LeNet-1 from 1993

Gradient-based learning applied to document recognition [[LeCun, Bottou, Bengio, Haffner 1998](#)]

Dr. Sander Ali Khawaja

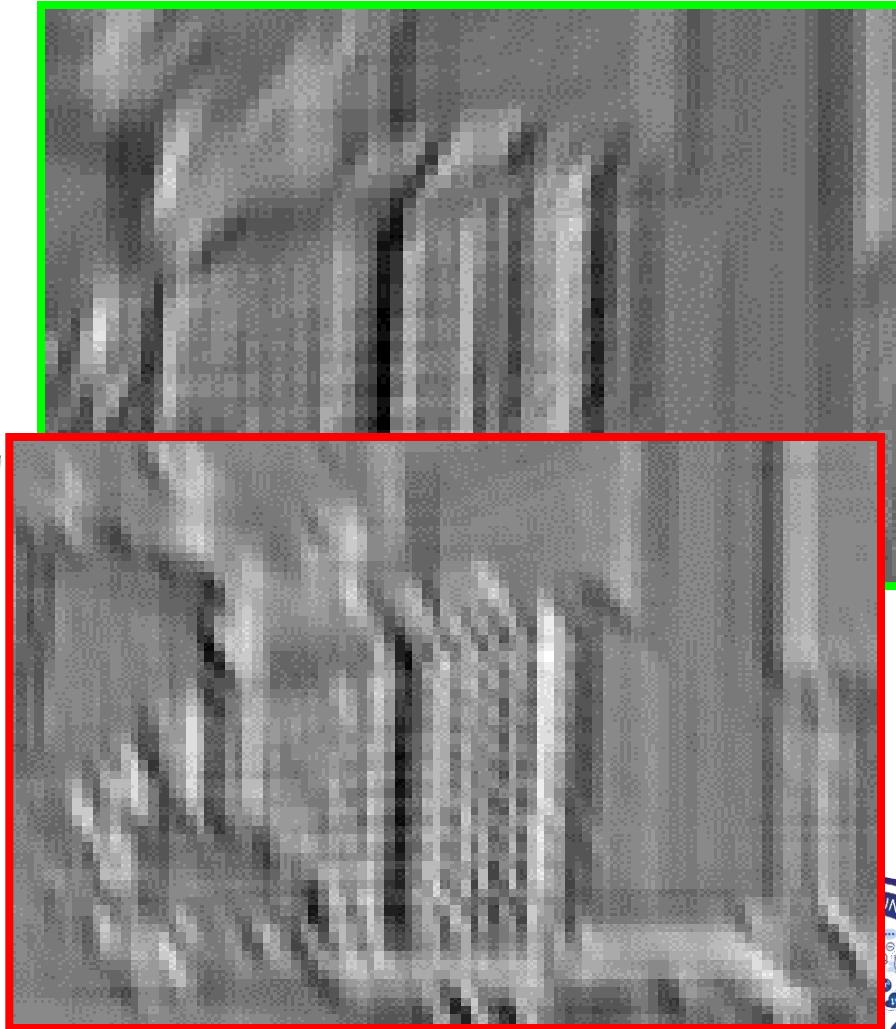


# What is a Convolution?

- Weighted moving sum



Input



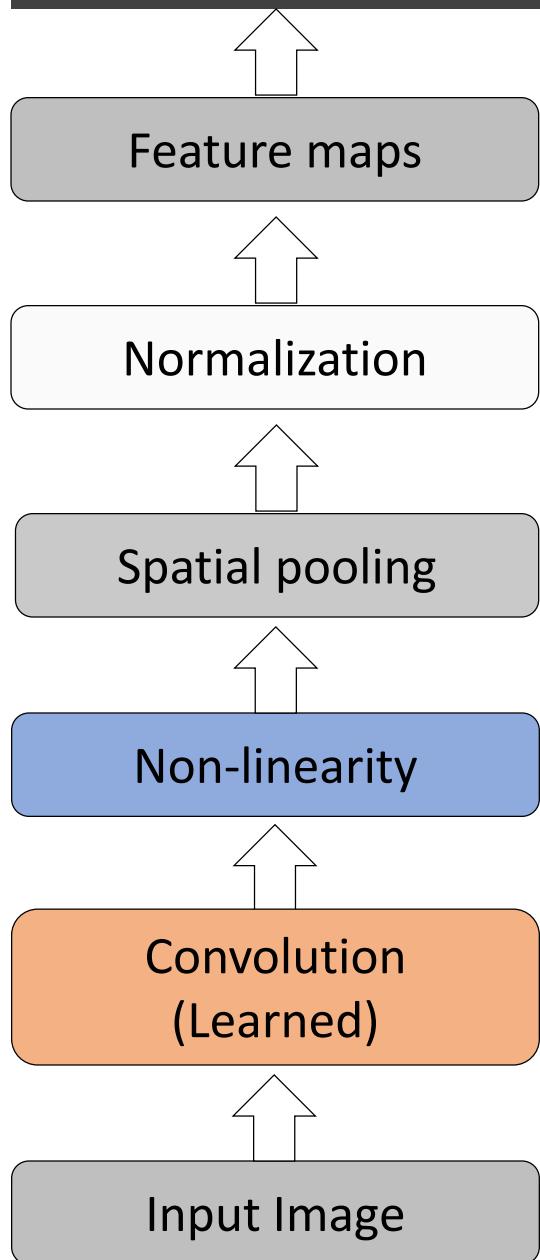
Feature Activation Map

Dr. Sander Ali Khowaja

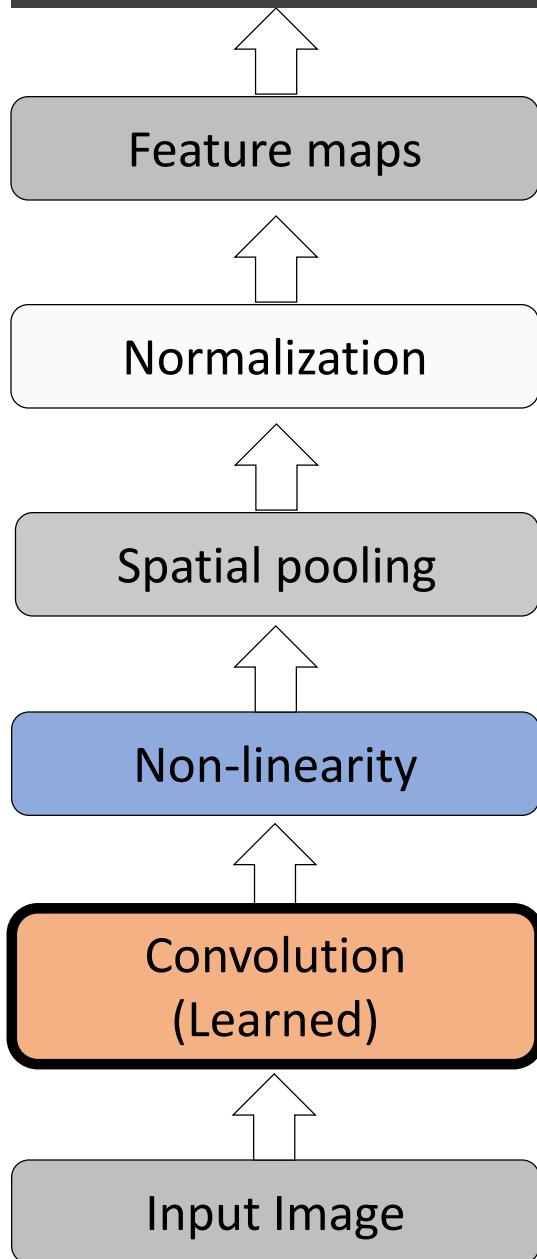
slide credit: S. Lazebnik



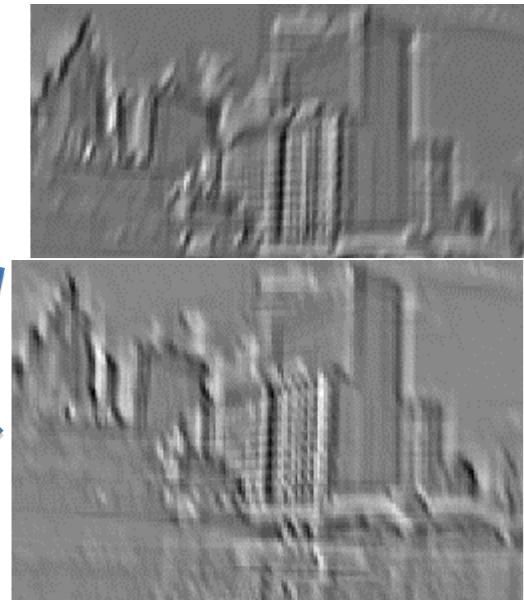
# Convolutional Neural Networks



# Convolutional Neural Networks



Dr. Sander Ali Khowaja

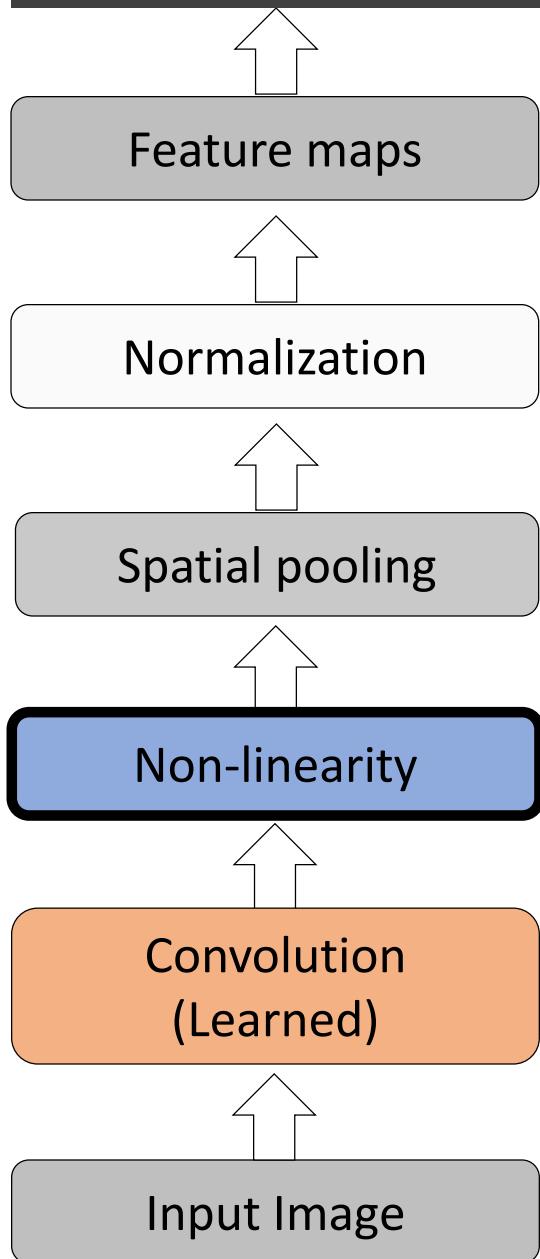


Feature Map

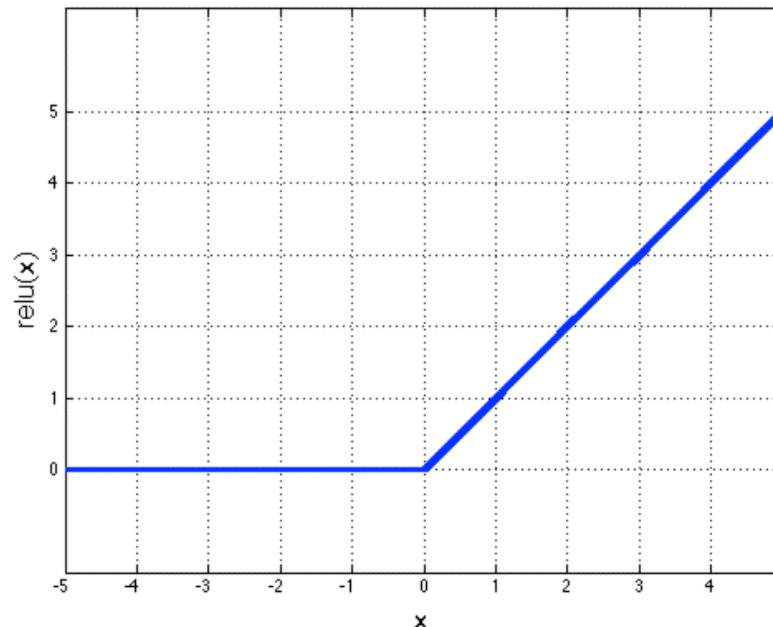
slide credit: S. Lazebnik



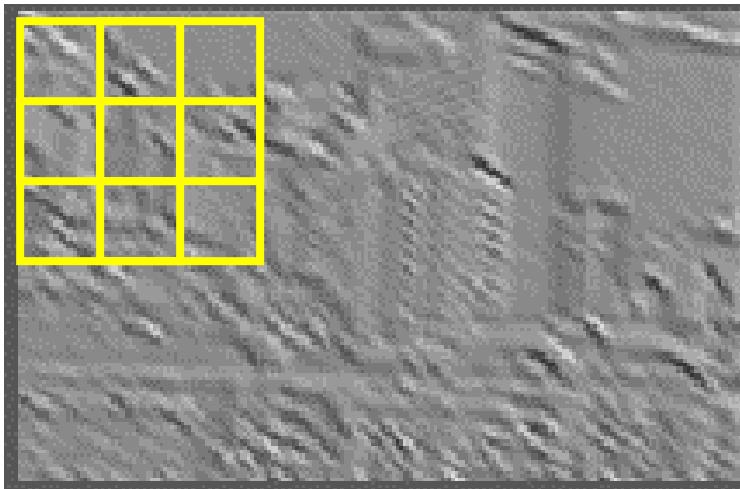
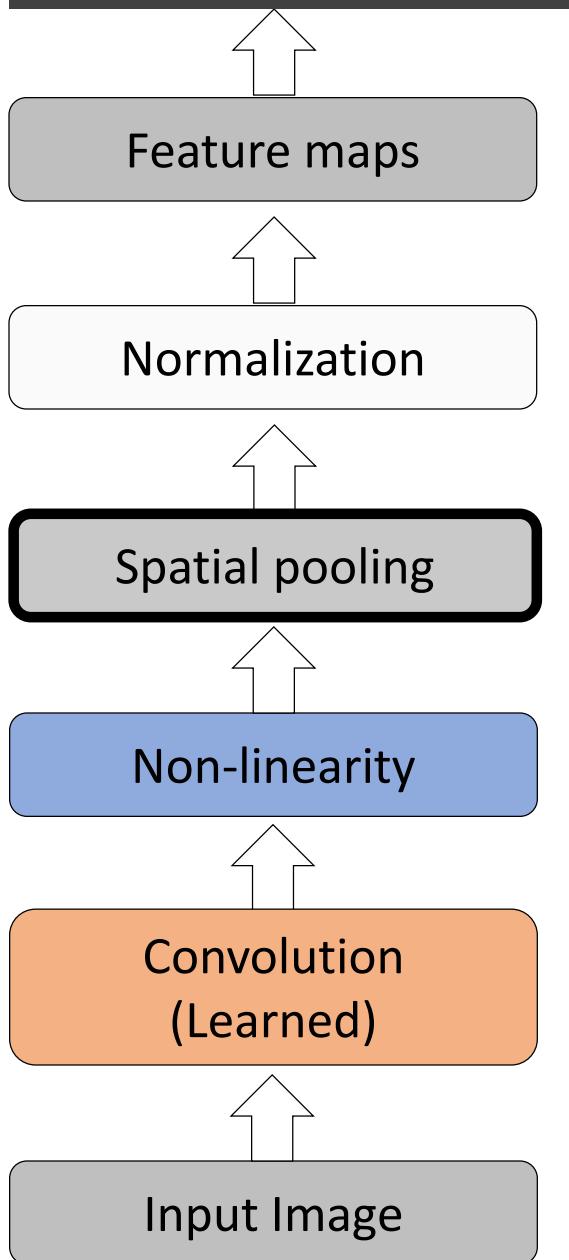
# Convolutional Neural Networks



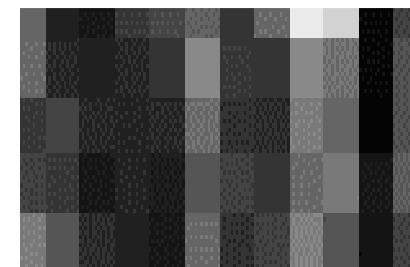
Rectified Linear Unit (ReLU)



# Convolutional Neural Networks



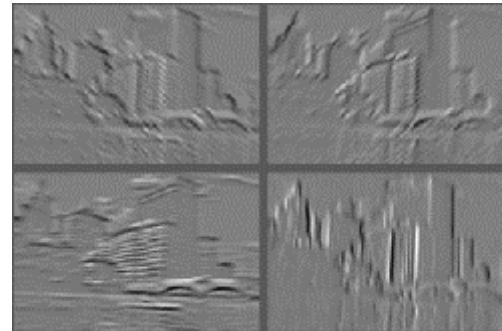
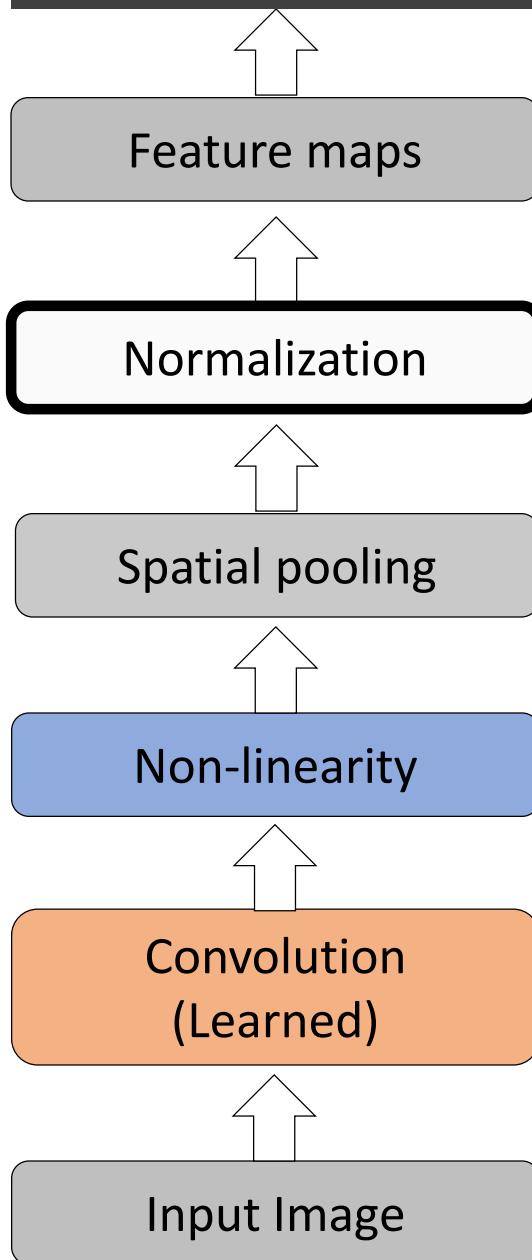
Max pooling



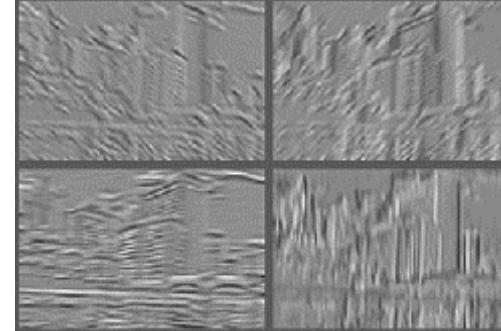
Max-pooling: a non-linear down-sampling

Provide *translation invariance*

# Convolutional Neural Networks

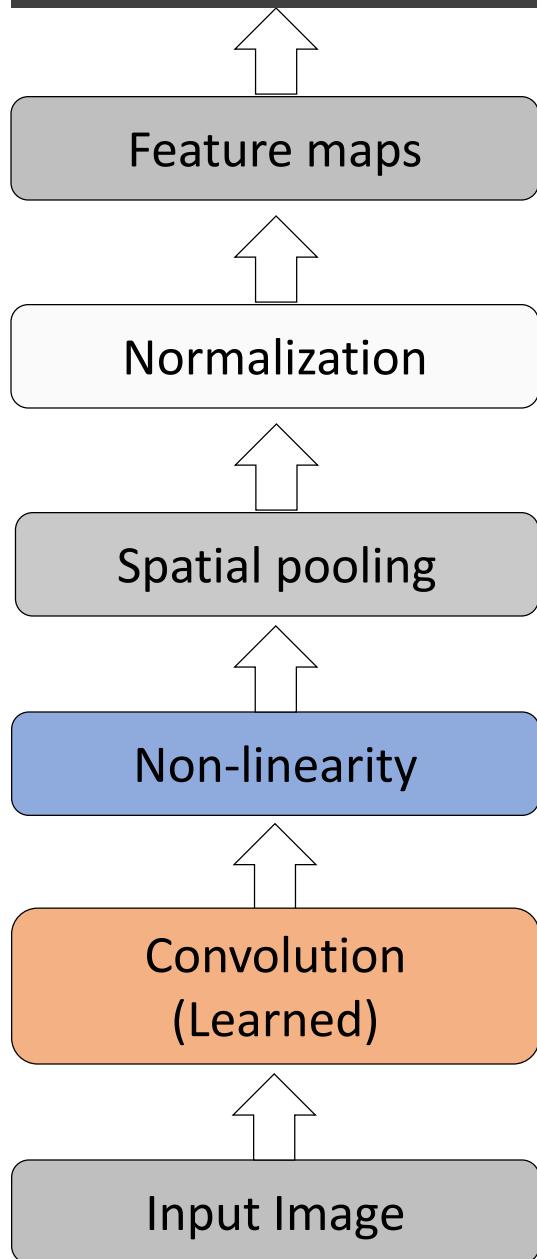


Feature Maps



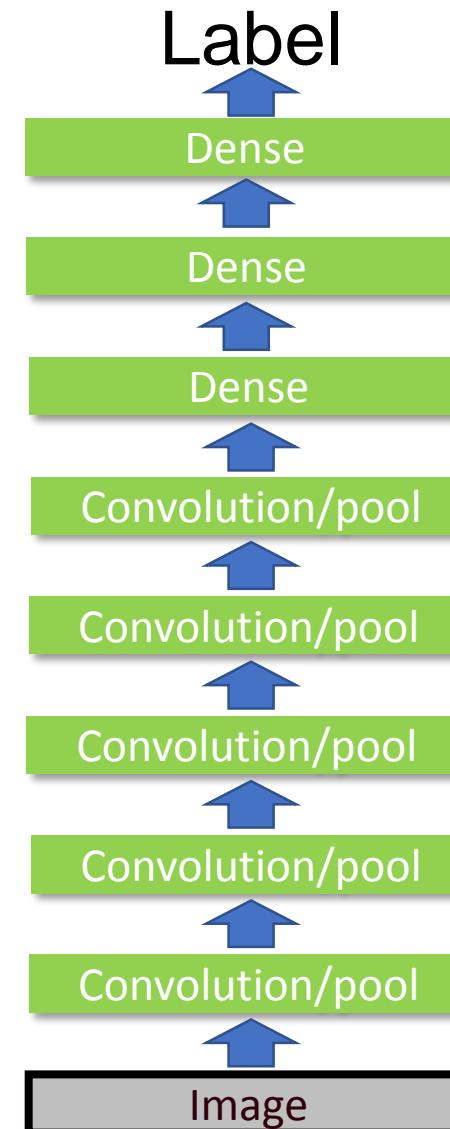
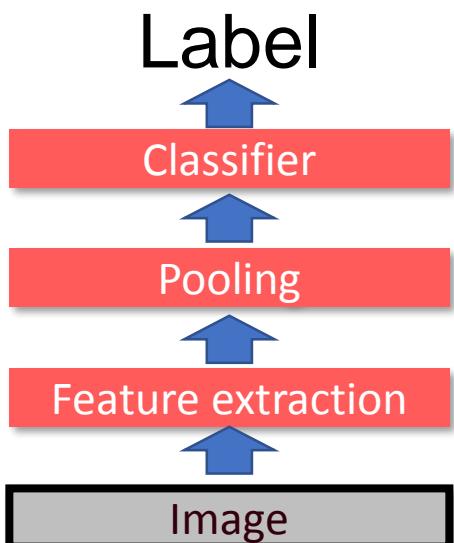
Feature Maps  
After Contrast  
Normalization

# Convolutional Neural Networks

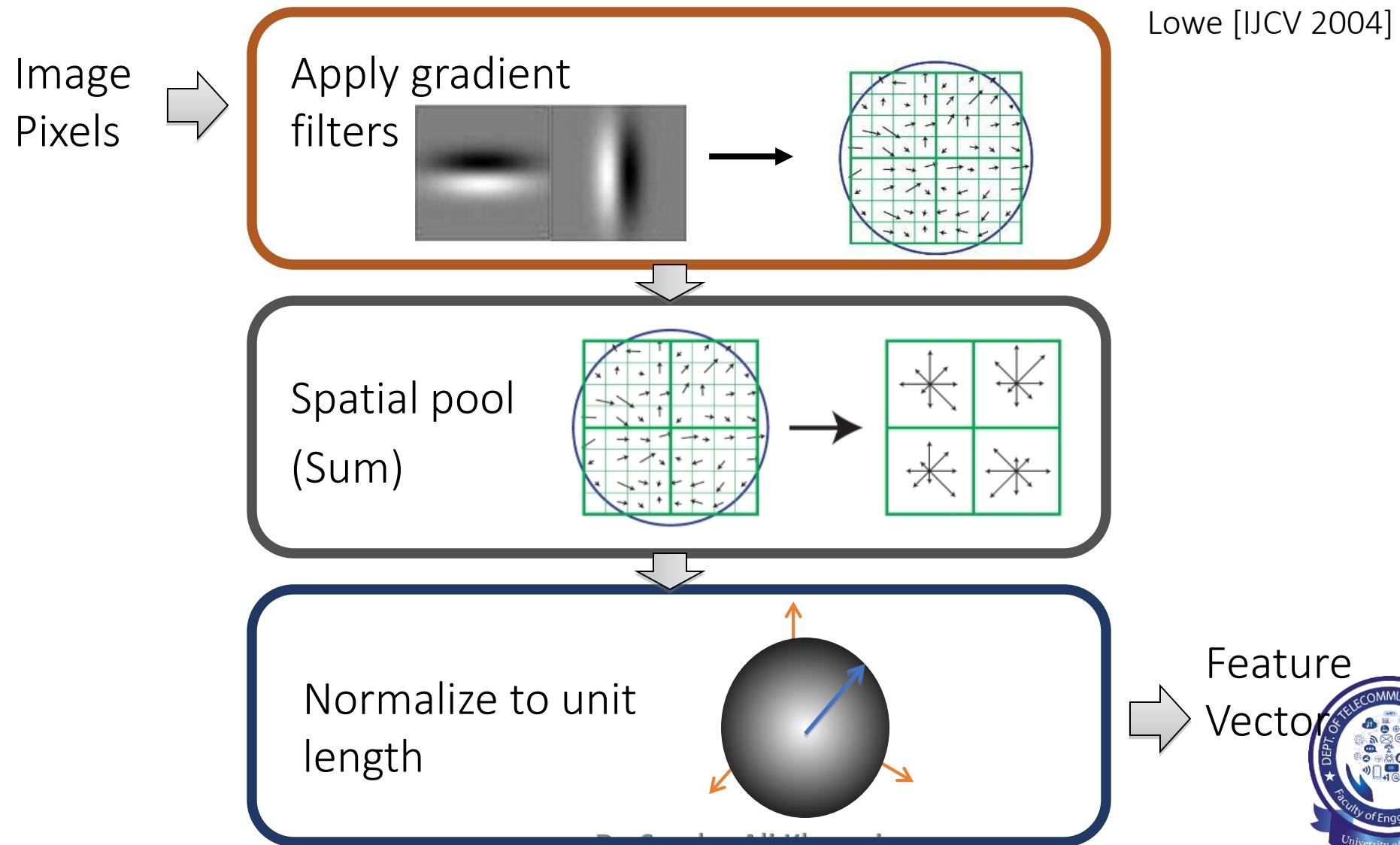


# Engineered vs. learned features

Convolutional filters are trained in a supervised manner by back-propagating classification error



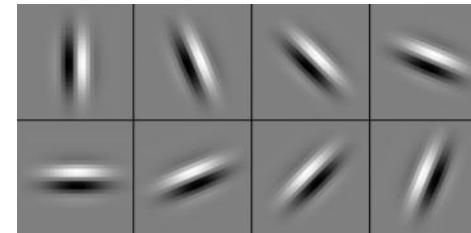
# SIFT Descriptor



# SIFT Descriptor

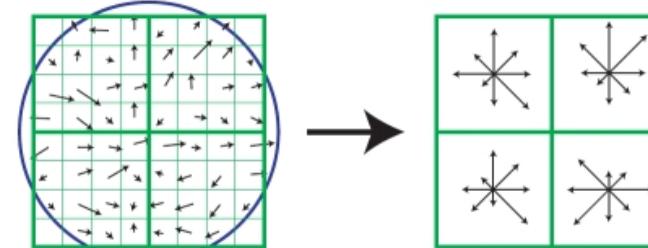
Image  
Pixels

Apply  
oriented filters

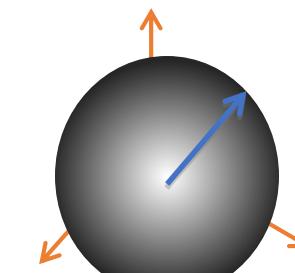


Lowe [IJCV 2004]

Spatial pool  
(Sum)



Normalize to unit  
length



Feature  
Vector

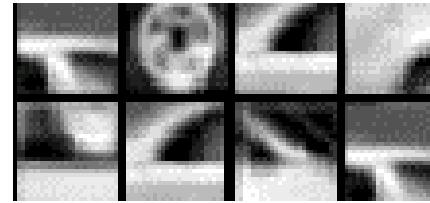
slide credit: R. Fergus



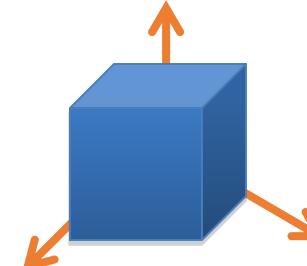
# Spatial Pyramid Matching

SIFT  
Features

Filter with  
Visual Words



Max



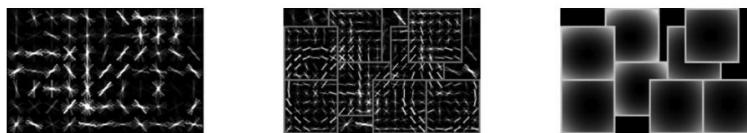
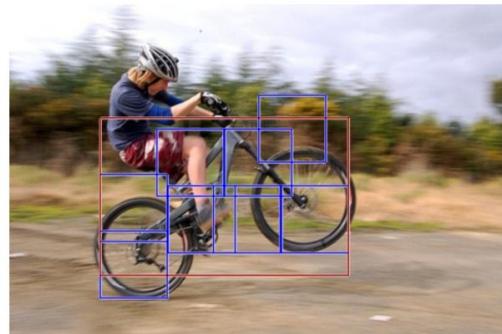
Multi-scale  
spatial pool  
(Sum)



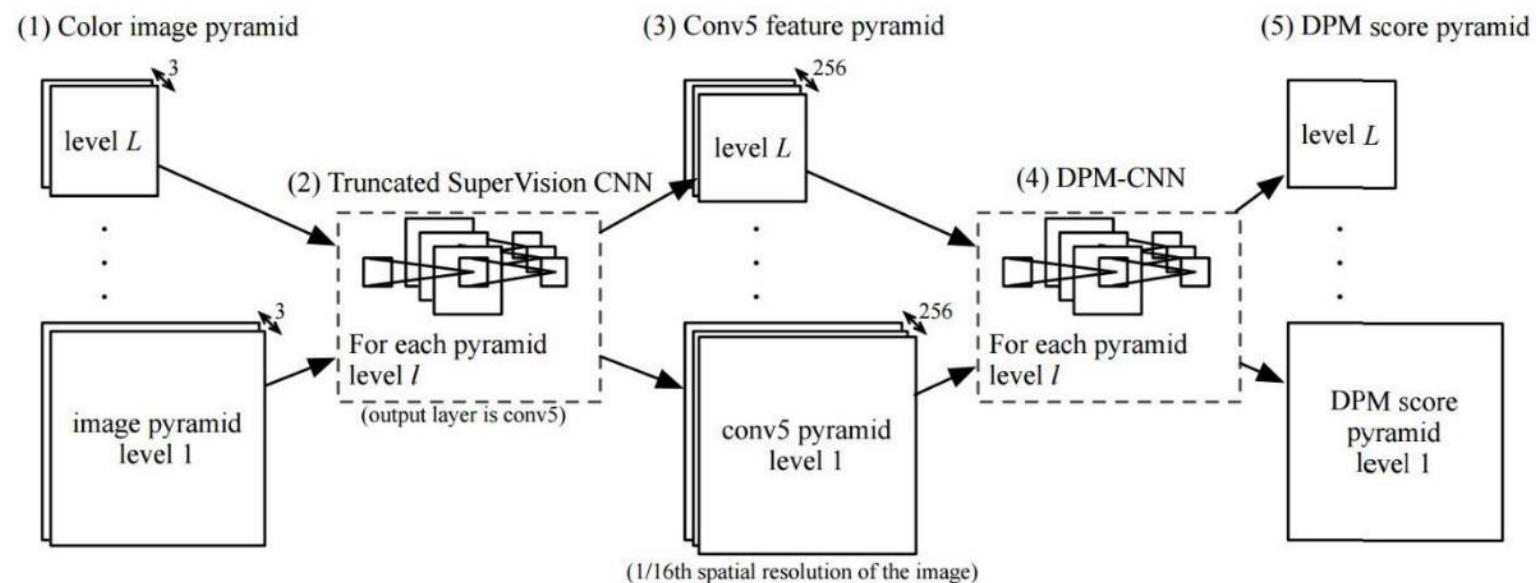
Lazebnik,  
Schmid,  
Ponce  
[CVPR 2006]



# Deformable Part Model



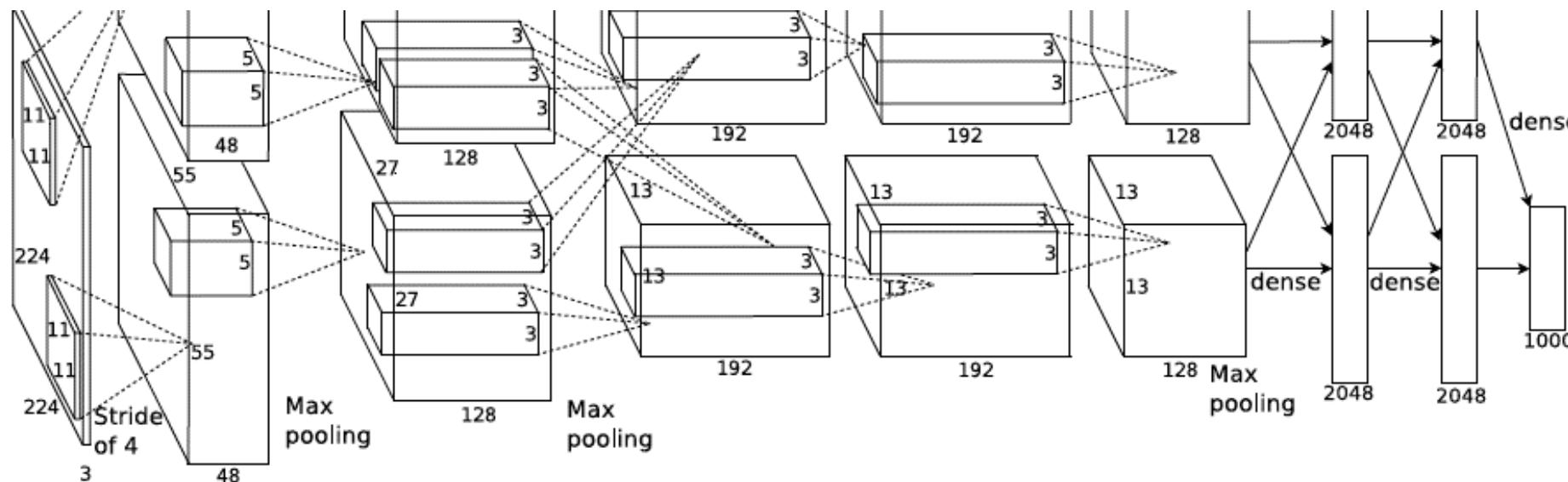
**DeepPyramid DPM**



Deformable Part Models are Convolutional Neural Networks [Girshick et al. CVPR 15]  
Dr. Sander Ali Khawaja



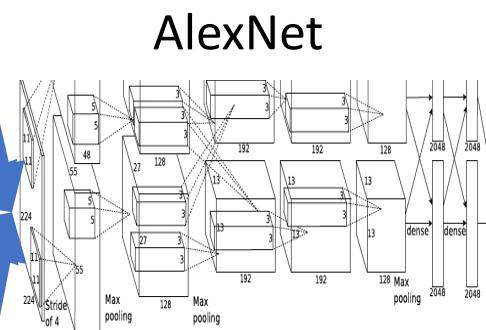
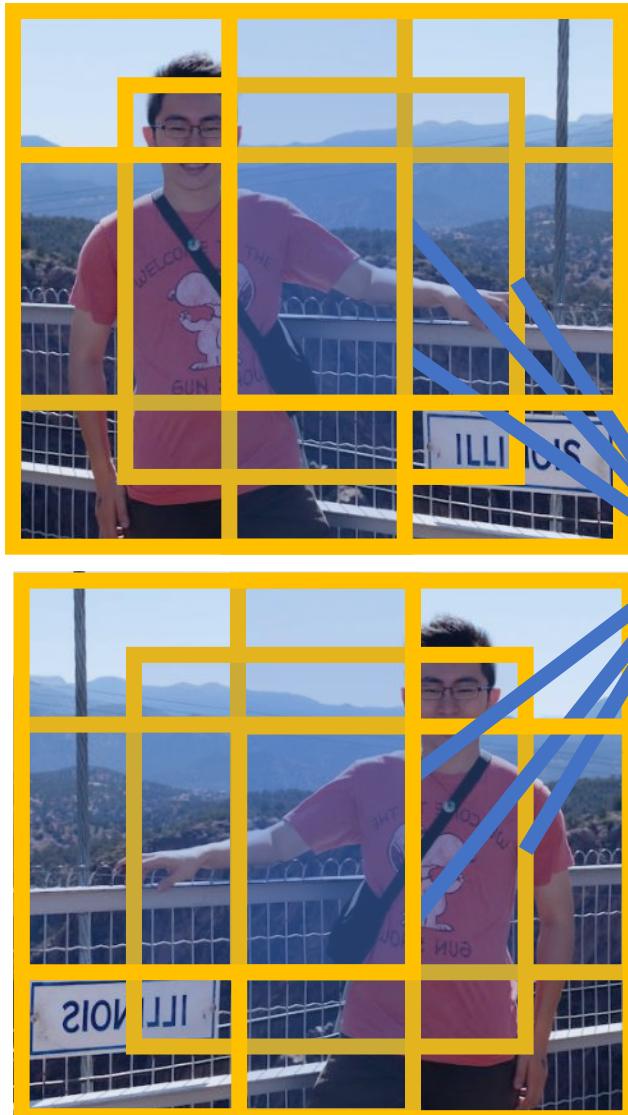
- Similar framework to LeCun'98 but:
  - Bigger model (7 hidden layers, 650,000 units, 60,000,000 params)
  - More data ( $10^6$  vs.  $10^3$  images)
  - GPU implementation (50x speedup over CPU)
    - Trained on two GPUs for a week



A. Krizhevsky, I. Sutskever, and G. Hinton,  
[ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012  
Dr. Sander Ali Khowaja



# Using CNN for Image Classification

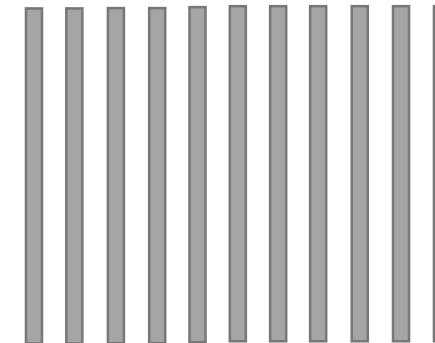


Fixed input size:  
 $224 \times 224 \times 3$

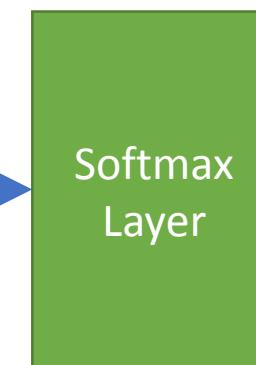
$d = 4096$

Dr. Sander Ali Khowaja

Fully connected layer Fc7  
 $d = 4096$



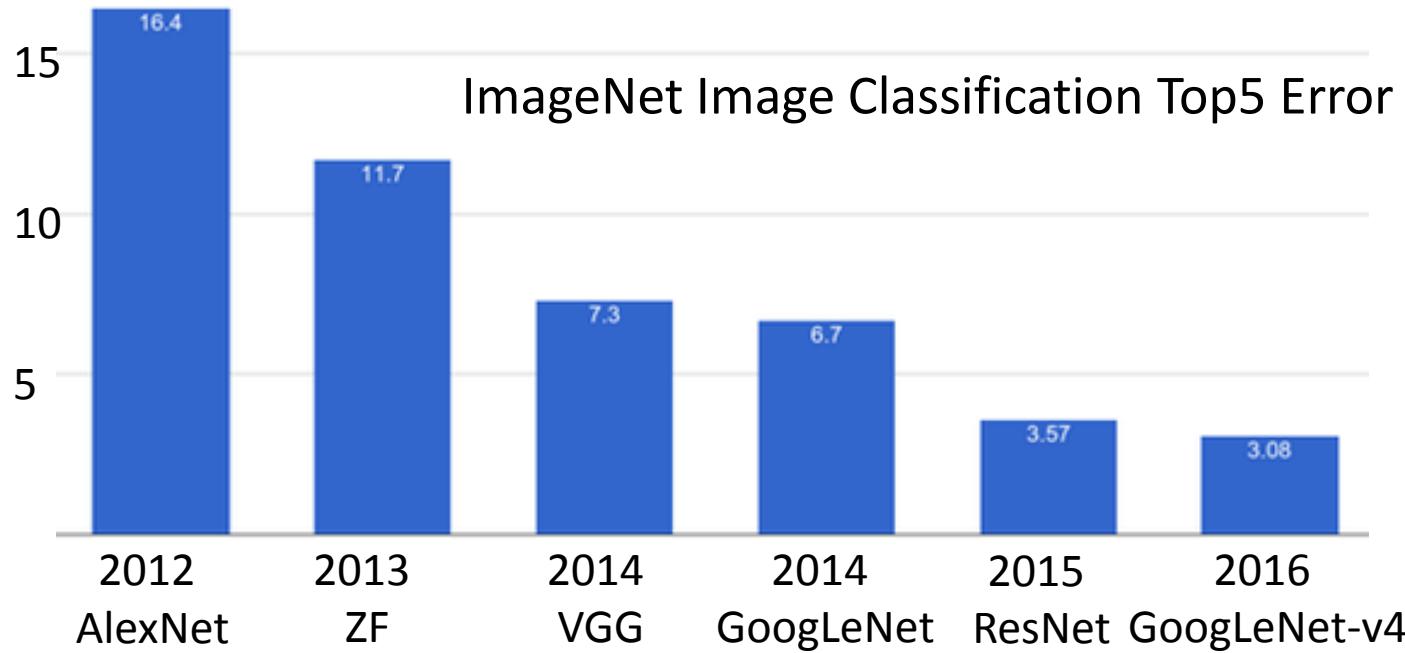
Averaging



"Jia-Bin"



# Progress on ImageNet



# VGG-Net

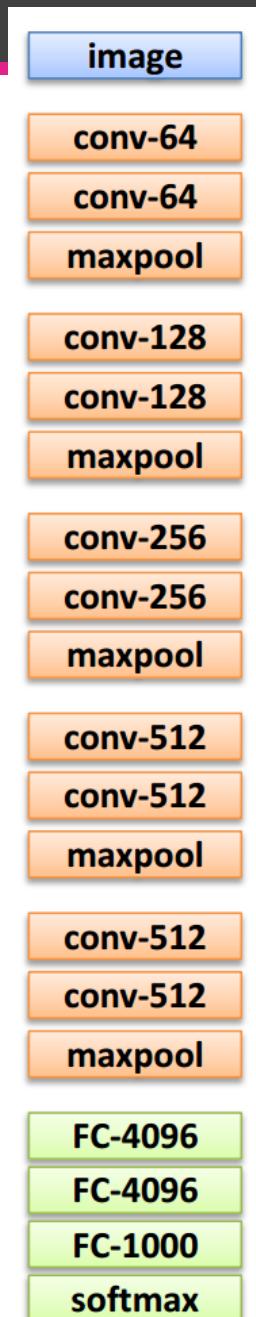
- The deeper, the better

- Key design choices:

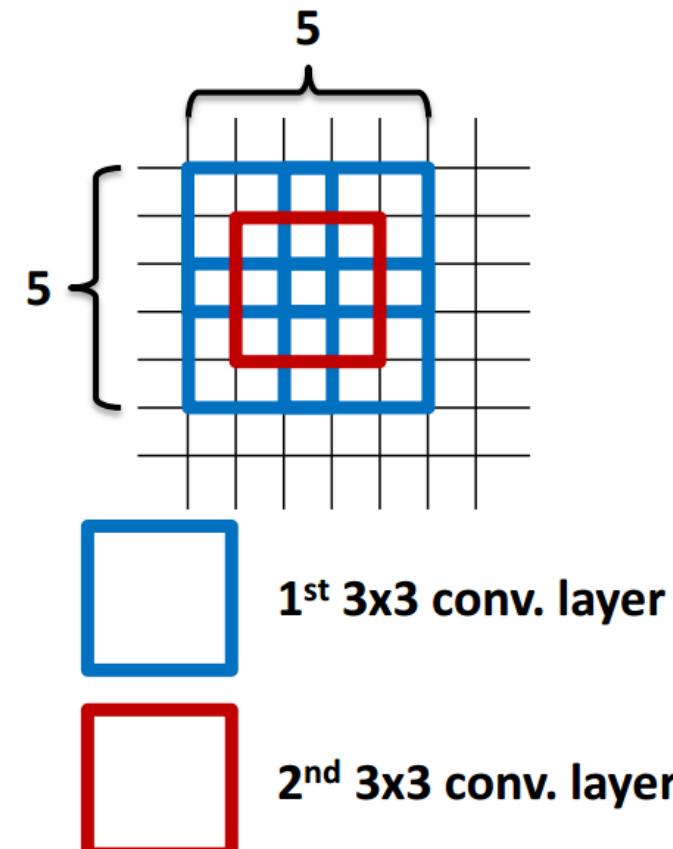
- 3x3 conv. Kernels
  - very small
- conv. stride 1
  - no loss of information

- Other details:

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

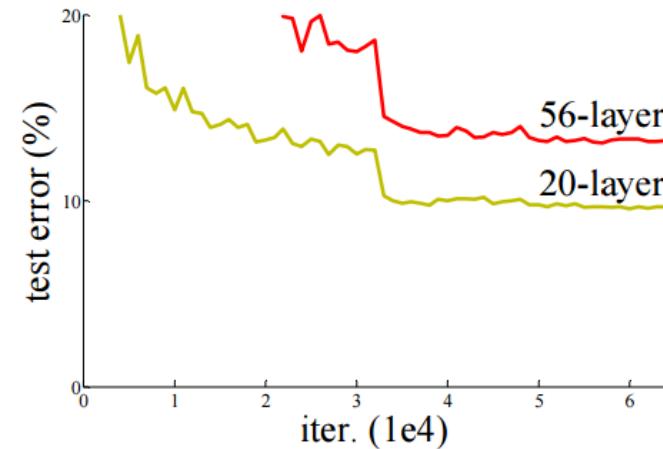
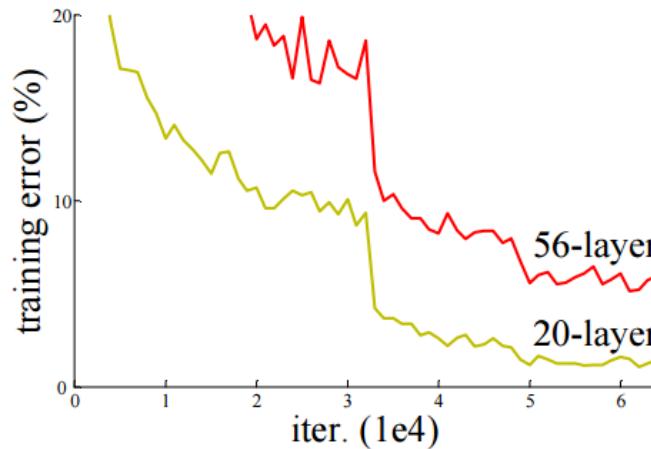


- Why 3x3 layers?
  - Stacked conv. layers have a large receptive field
  - two 3x3 layers – 5x5 receptive field
  - three 3x3 layers – 7x7 receptive field
- More non-linearity
  - Less parameters to learn
  - ~140M per net

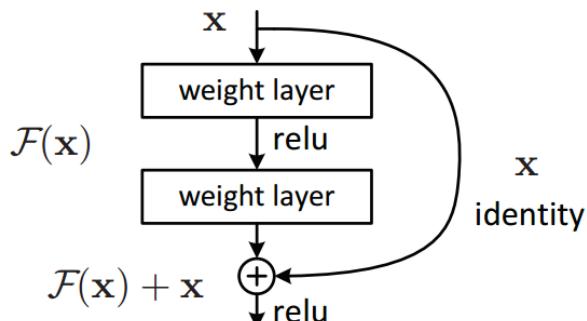


# ResNet

- Can we just increase the #layer?



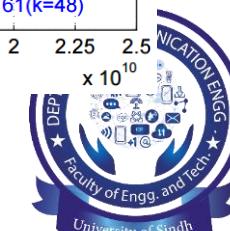
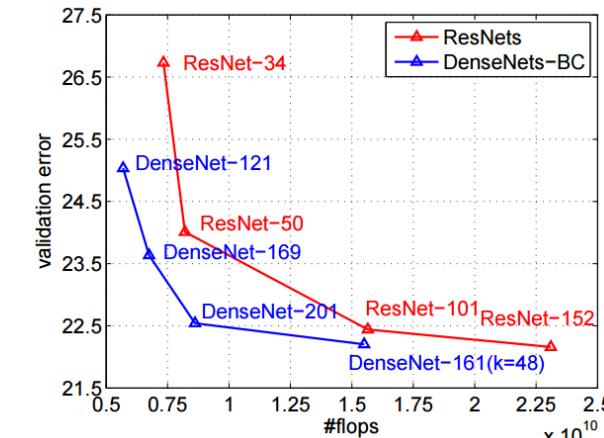
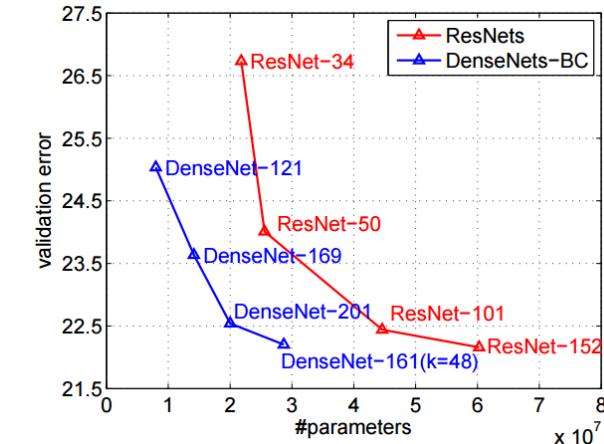
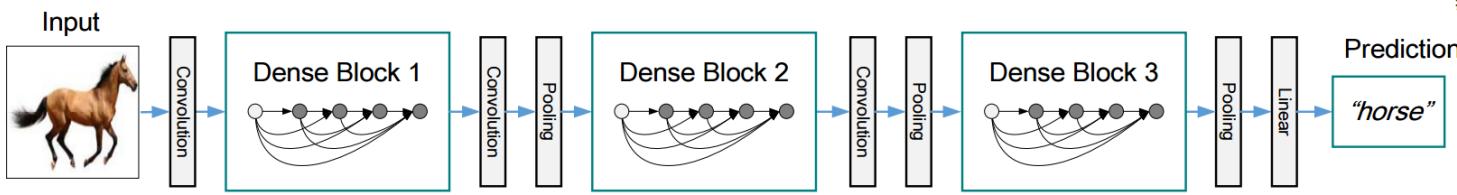
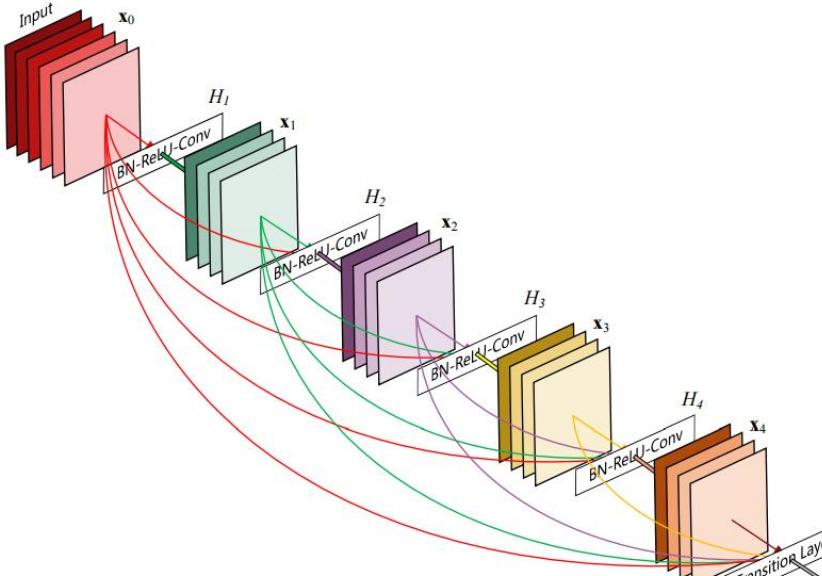
- How can we train  
- Residual learning



method	top-5 err. (test)
VGG [41] (ILSVRC'14)	7.32
GoogLeNet [44] (ILSVRC'14)	6.66
VGG [41] (v5)	6.8
PReLU-net [13]	4.94
BN-inception [16]	4.82
<b>ResNet (ILSVRC'15)</b>	<b>3.57</b>

# DenseNet

- Shorter connections (like ResNet) help
- Why not just connect them all?



# Training Convolutional Neural Networks

- Backpropagation + stochastic gradient descent with momentum
  - [Neural Networks: Tricks of the Trade](#)
- Dropout
- Data augmentation
- Batch normalization
- Initialization
  - Transfer learning

# Training CNN with gradient descent

- A CNN as composition of functions

$$f_{\mathbf{w}}(\mathbf{x}) = f_L(\dots (f_2(f_1(\mathbf{x}; \mathbf{w}_1); \mathbf{w}_2) \dots; \mathbf{w}_L)$$

- Parameters

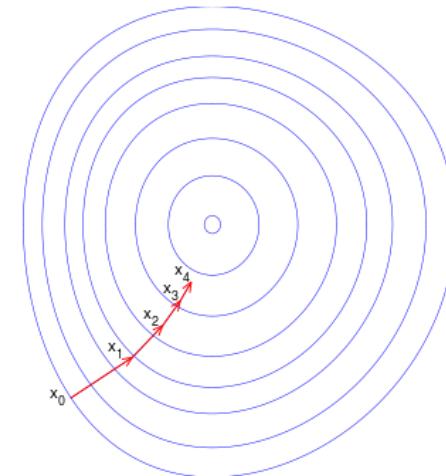
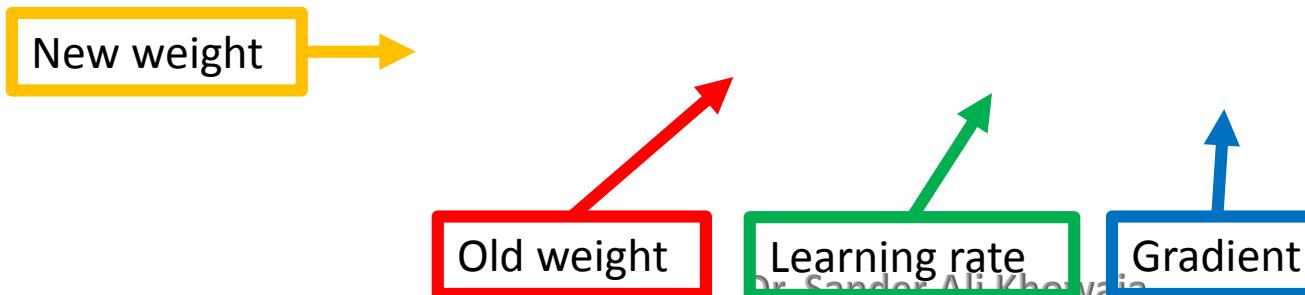
$$\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_L)$$

- Empirical loss function

$$L(\mathbf{w}) = \frac{1}{n} \sum_i l(z_i, f_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient descent

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \frac{\partial f}{\partial \mathbf{w}}(\mathbf{w}^t)$$



# An Illustrative example

$$f(x, y) = xy, \quad \frac{\partial f}{\partial x} = y, \frac{\partial f}{\partial y} = x$$

Example:  $x = 4, y = -3 \Rightarrow f(x, y) = -12$

## Partial derivatives

$$\frac{\partial f}{\partial x} = -3, \quad \frac{\partial f}{\partial y} = 4$$

## Gradient

$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$



$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

Goal: compute the gradient

$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right]$$



$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$

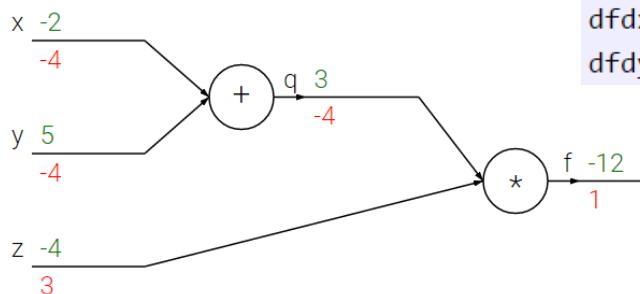
$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

## Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



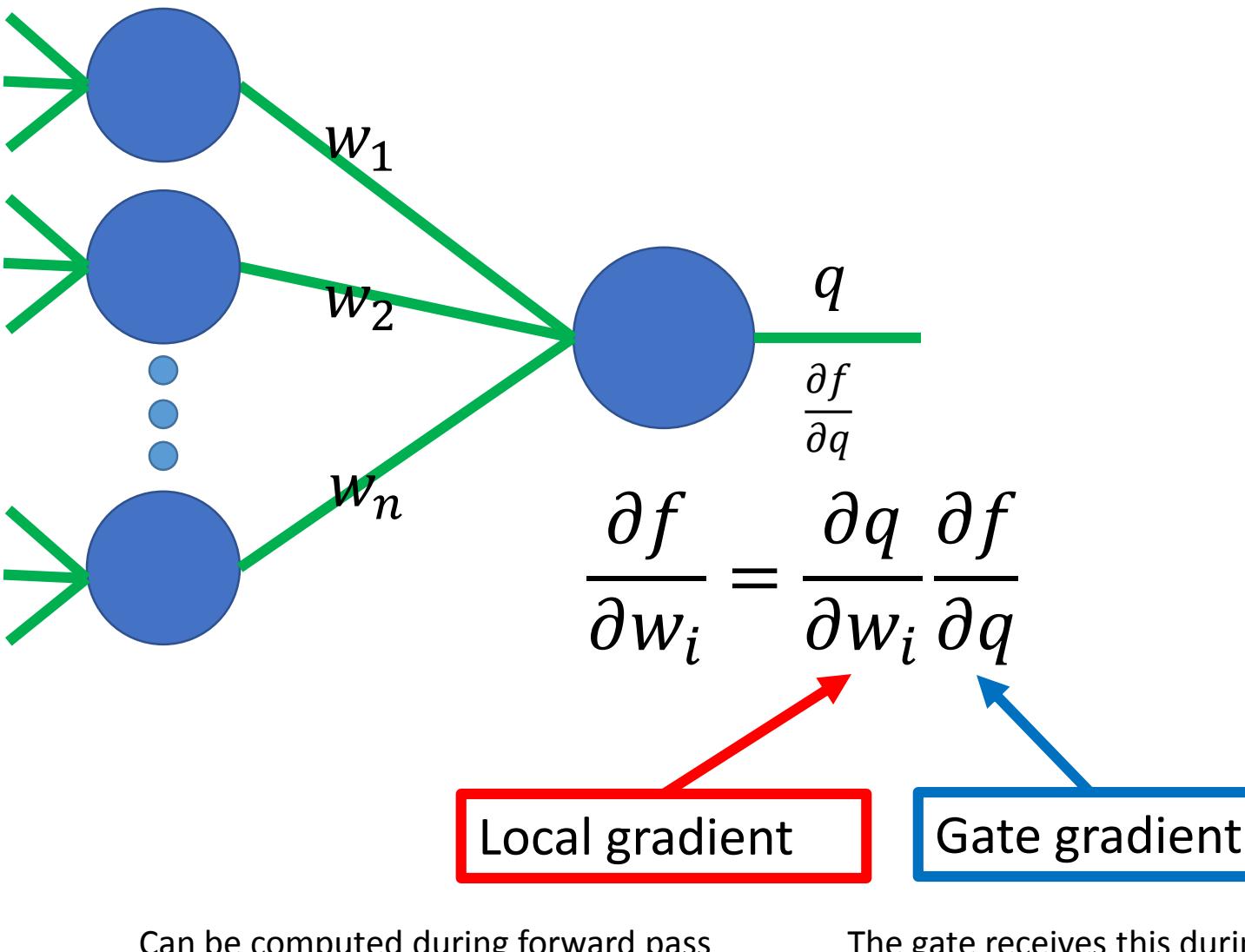
```

# set some inputs
x = -2; y = 5; z = -4

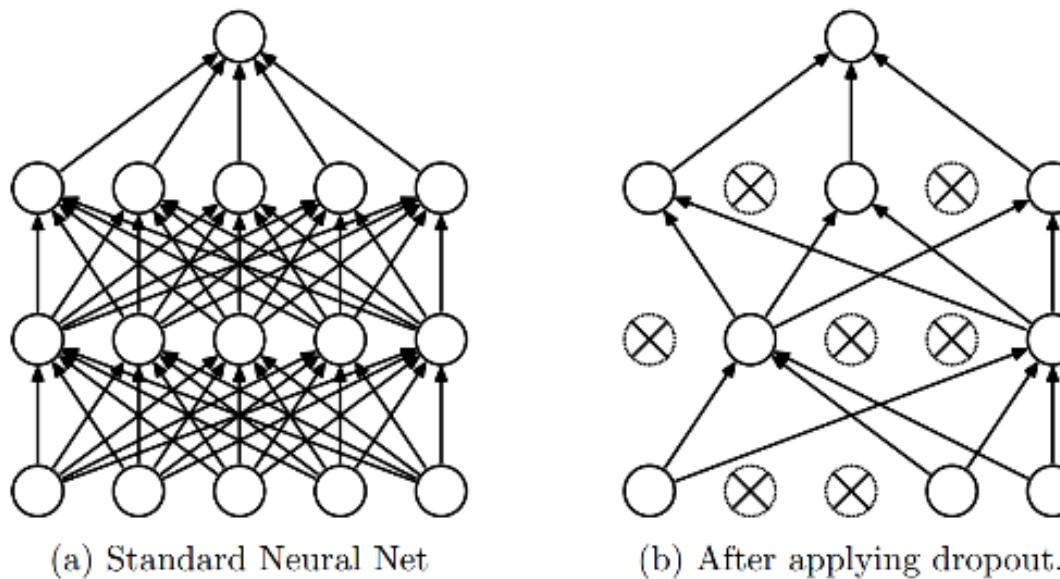
# perform the forward pass
q = x + y # q becomes 3
f = q * z # f becomes -12

# perform the backward pass (backpropagation) in reverse order:
# first backprop through f = q * z
dfdq = q # df/dz = q, so gradient on z becomes 3
dfdz = z # df/dq = z, so gradient on q becomes -4
# now backprop through q = x + y
dfdx = 1.0 * dfdq # dq/dx = 1. And the multiplication here is the chain rule!
dfdy = 1.0 * dfdq # dq/dy = 1
  
```

# Backpropagation (recursive chain rule)



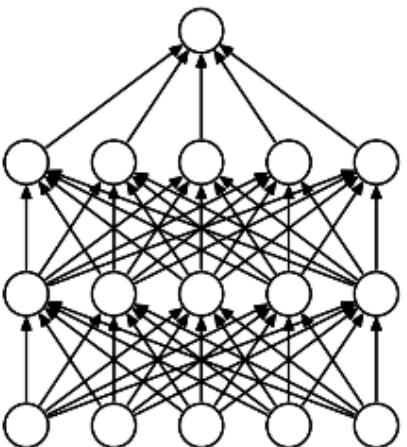
# Dropout



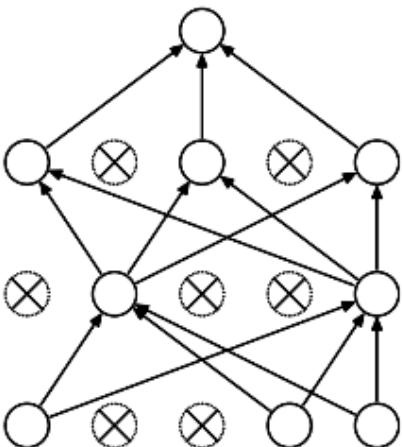
Intuition: successful conspiracies

- 50 people planning a conspiracy
- Strategy A: plan a big conspiracy involving 50 people
  - Likely to fail. 50 people need to play their parts correctly.
- Strategy B: plan 10 conspiracies each involving 5 people
  - Likely to succeed!

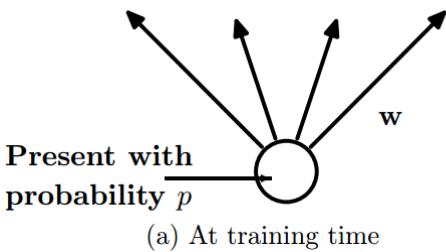
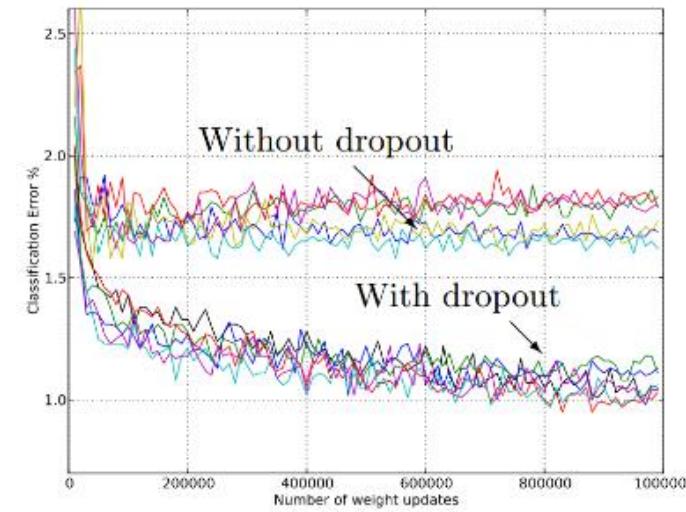
# Dropout



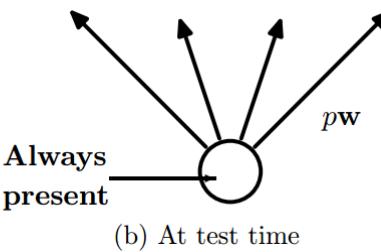
(a) Standard Neural Net



(b) After applying dropout.



(a) At training time



(b) At test time

**Main Idea:** approximately combining exponentially many different neural network architectures efficiently

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
SVM on Fisher Vectors of Dense SIFT and Color Statistics	-	-	27.3
Avg of classifiers over FVs of SIFT, LBP, GIST and CSIFT	-	-	26.2
Conv Net + dropout (Krizhevsky et al., 2012)	40.7	18.2	-
Avg of 5 Conv Nets + dropout (Krizhevsky et al., 2012)	38.1	16.4	16.4

Table 6: Results on the ILSVRC-2012 validation/test set.

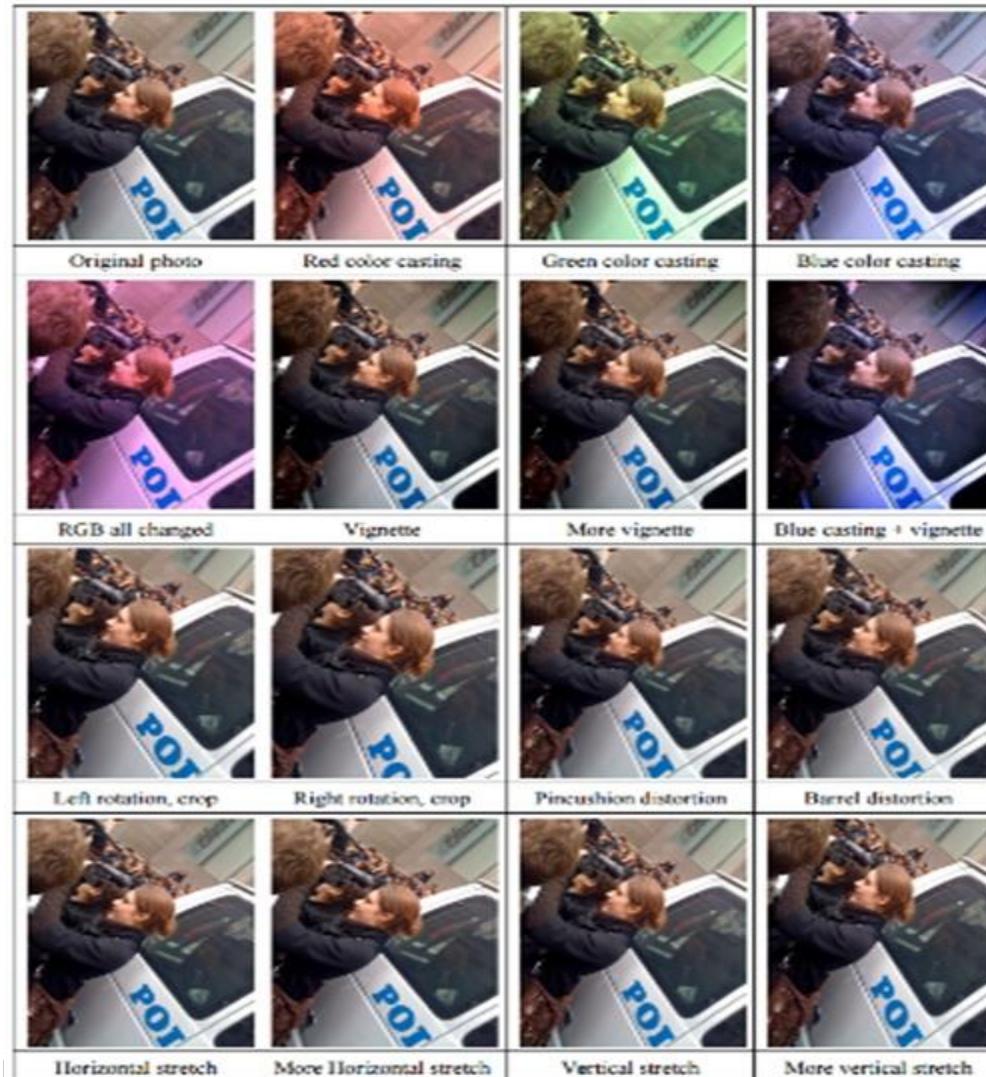
Dropout: A simple way to prevent neural networks from overfitting [[Srivastava JMLR 2014](#)]

Dr. Sander Ali Khawaja

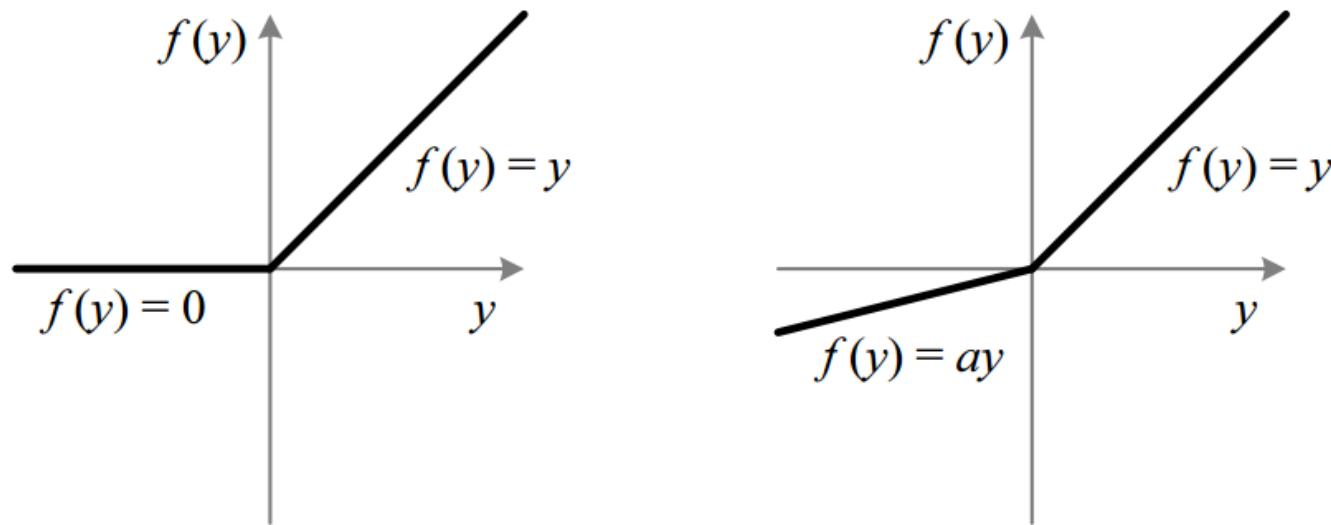


# Data Augmentation (Jittering)

- Create *virtual* training samples
  - Horizontal flip
  - Random crop
  - Color casting
  - Geometric distortion



# Parametric Rectified Linear Unit



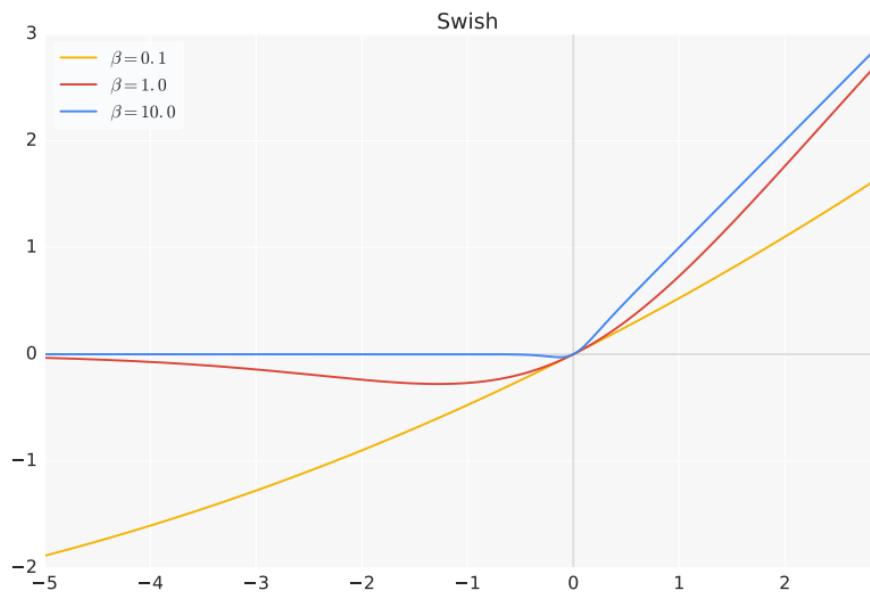
	team	top-5 (test)
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	<b>MSRA, PReLU-nets</b>	<b>4.94</b>

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification [[He et al. 2015](#)]

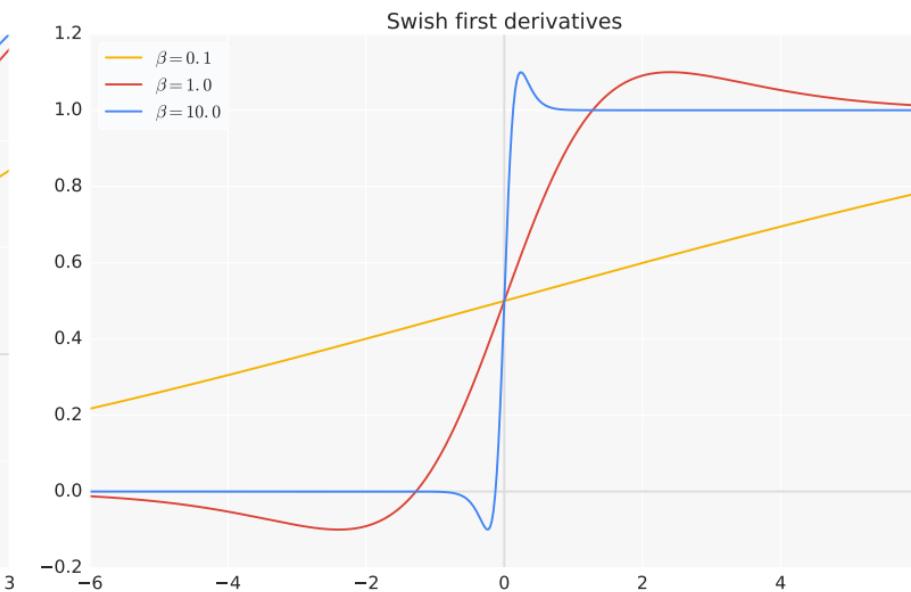
Dr. Sander Ali Khawaja



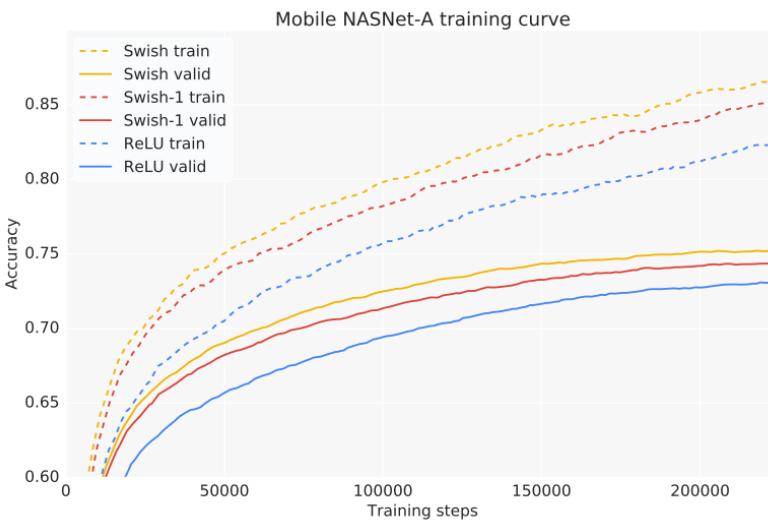
# Swish



The Swish activation function



First derivatives of Swish

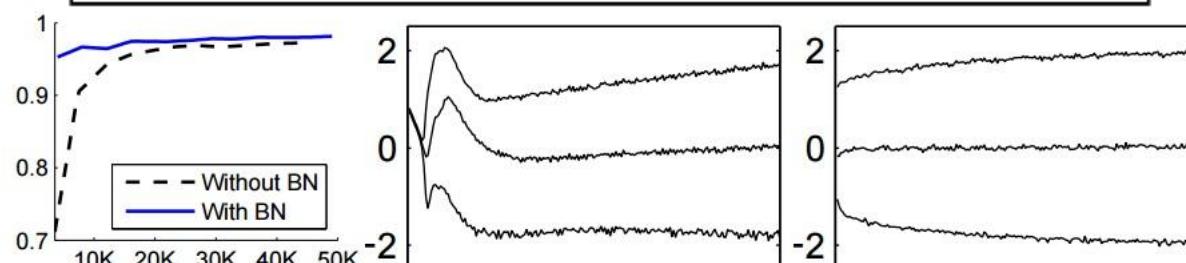


Model	Top-1 Acc. (%)			Top-5 Acc. (%)		
LReLU	73.8	73.9	74.2	91.6	91.9	91.9
PReLU	74.6	74.7	74.7	92.4	92.3	92.3
Softplus	74.0	74.2	74.2	91.6	91.8	91.9
ELU	74.1	74.2	74.2	91.8	91.8	91.8
SELU	73.6	73.7	73.7	91.6	91.7	91.7
GELU	74.6	-	-	92.0	-	-
ReLU	73.5	73.6	73.8	91.4	91.5	91.6
Swish-1	74.6	74.7	74.7	92.1	92.0	92.0
Swish	74.9	74.9	75.2	92.3	92.4	92.4

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;  
Parameters to be learned:  $\gamma, \beta$   
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\begin{aligned}\mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i && // \text{mini-batch mean} \\ \sigma_{\mathcal{B}}^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 && // \text{mini-batch variance} \\ \hat{x}_i &\leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} && // \text{normalize} \\ y_i &\leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) && // \text{scale and shift}\end{aligned}$$



(a)

(b) Without BN

(c) With BN

Batch Normalization: Accelerating Deep Network Training by  
Reducing Internal Covariate Shift [[Ioffe and Szegedy 2015](#)]

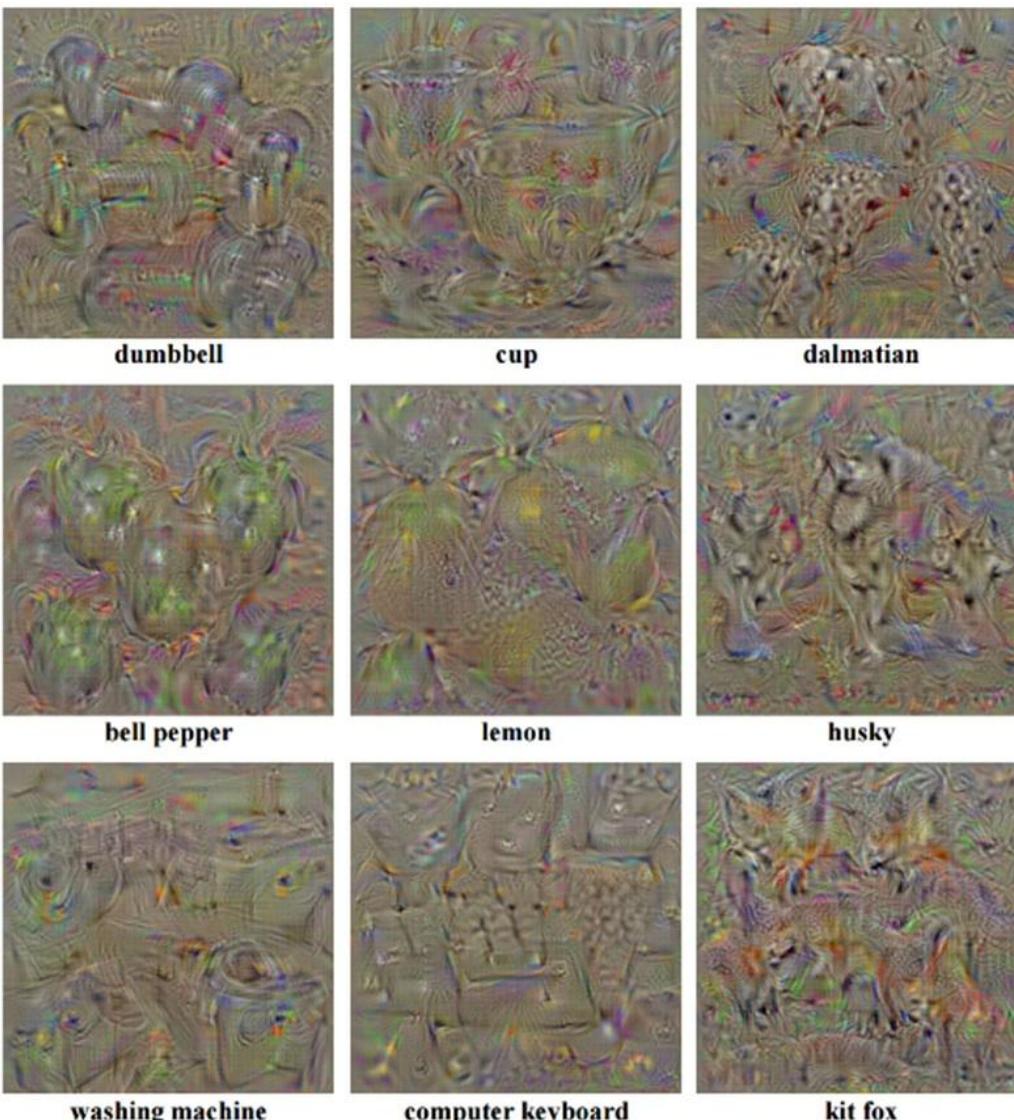


# Understanding and Visualizing CNN

- Find images that maximize some class scores
- Individual neuron activation
- Breaking CNNs

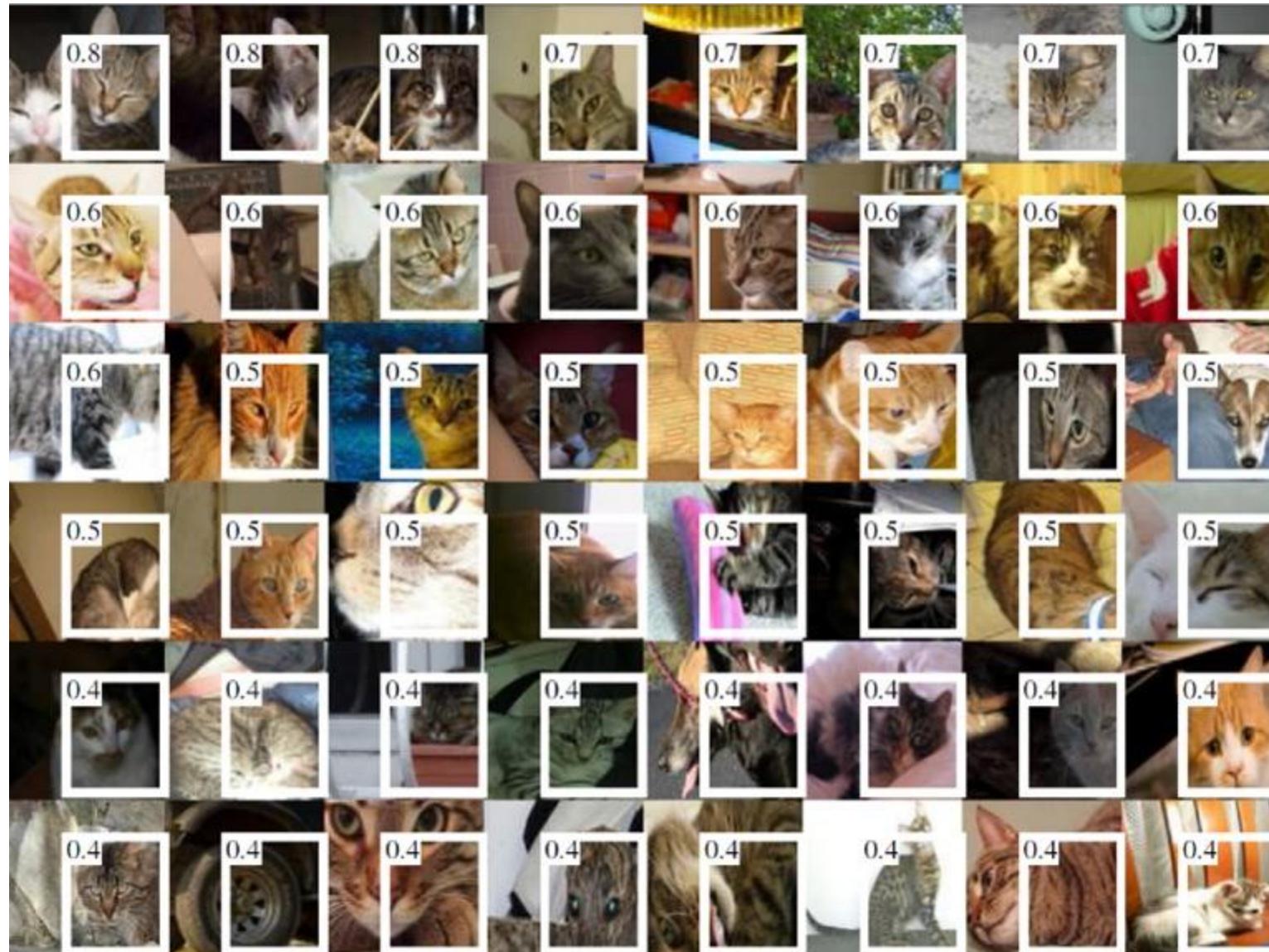


# Find images that maximize some class scores



person: HOG template

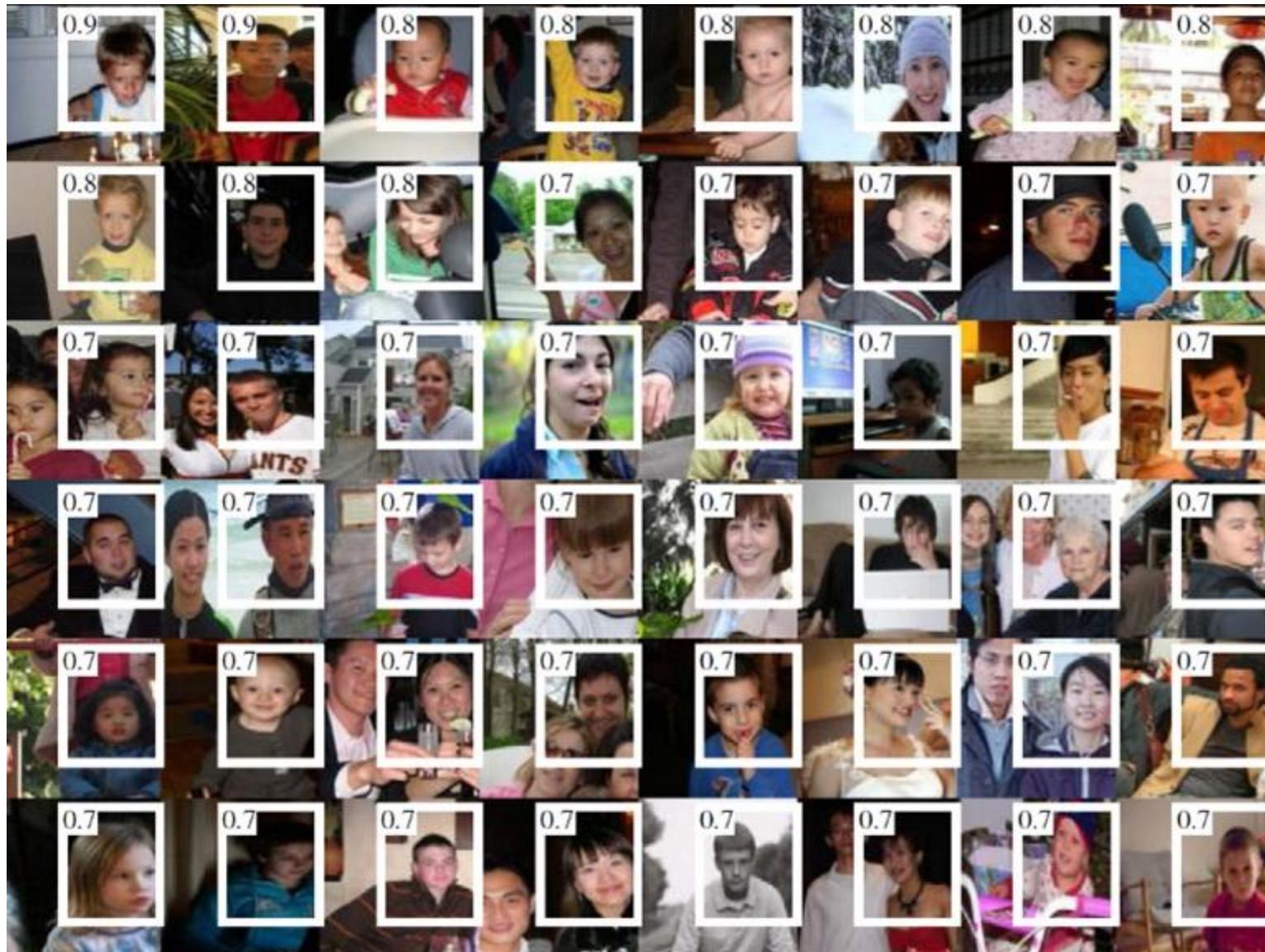
# Individual Neuron Activation



RCNN [Girshick et al. CVPR 2014]  
Prashant Arora



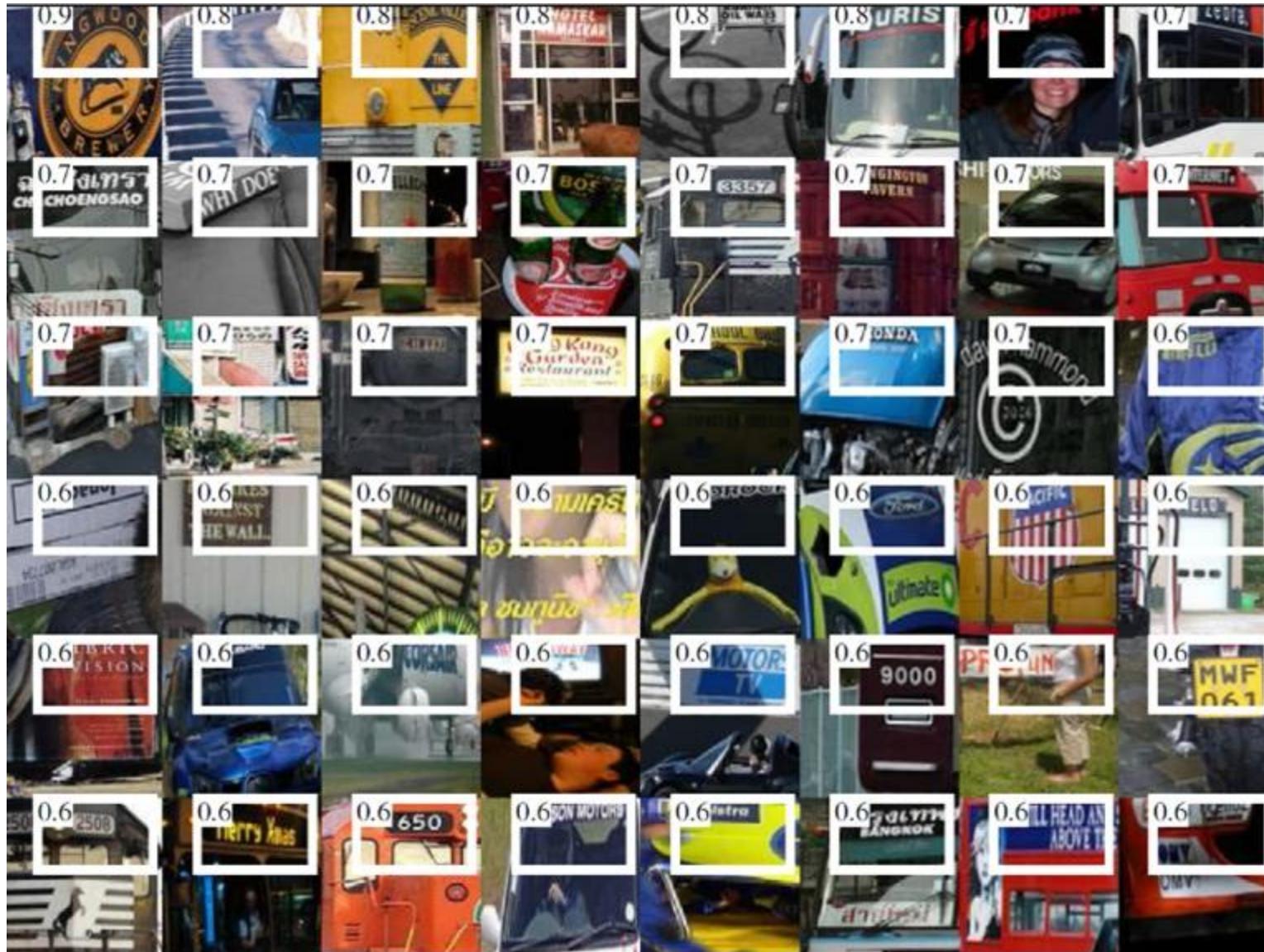
# Individual Neuron Activation



RCNN [[Girshick et al. CVPR 2014](#)]  
Processor: A. Khanwala



# Individual Neuron Activation

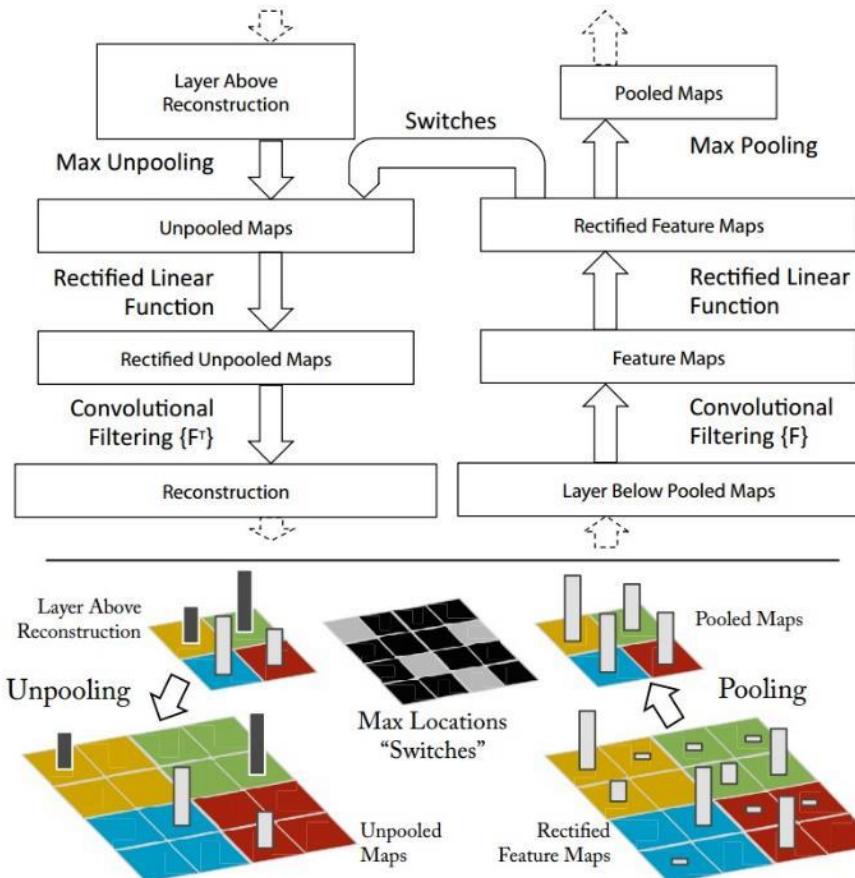


RCNN [Girshick et al. CVPR 2014]  
Processor: AM Khawaja

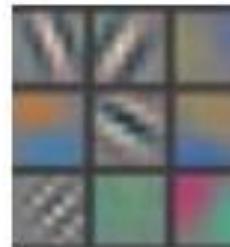


# Map activation back to the input pixel space

- What input pattern originally caused a given activation in the feature maps?



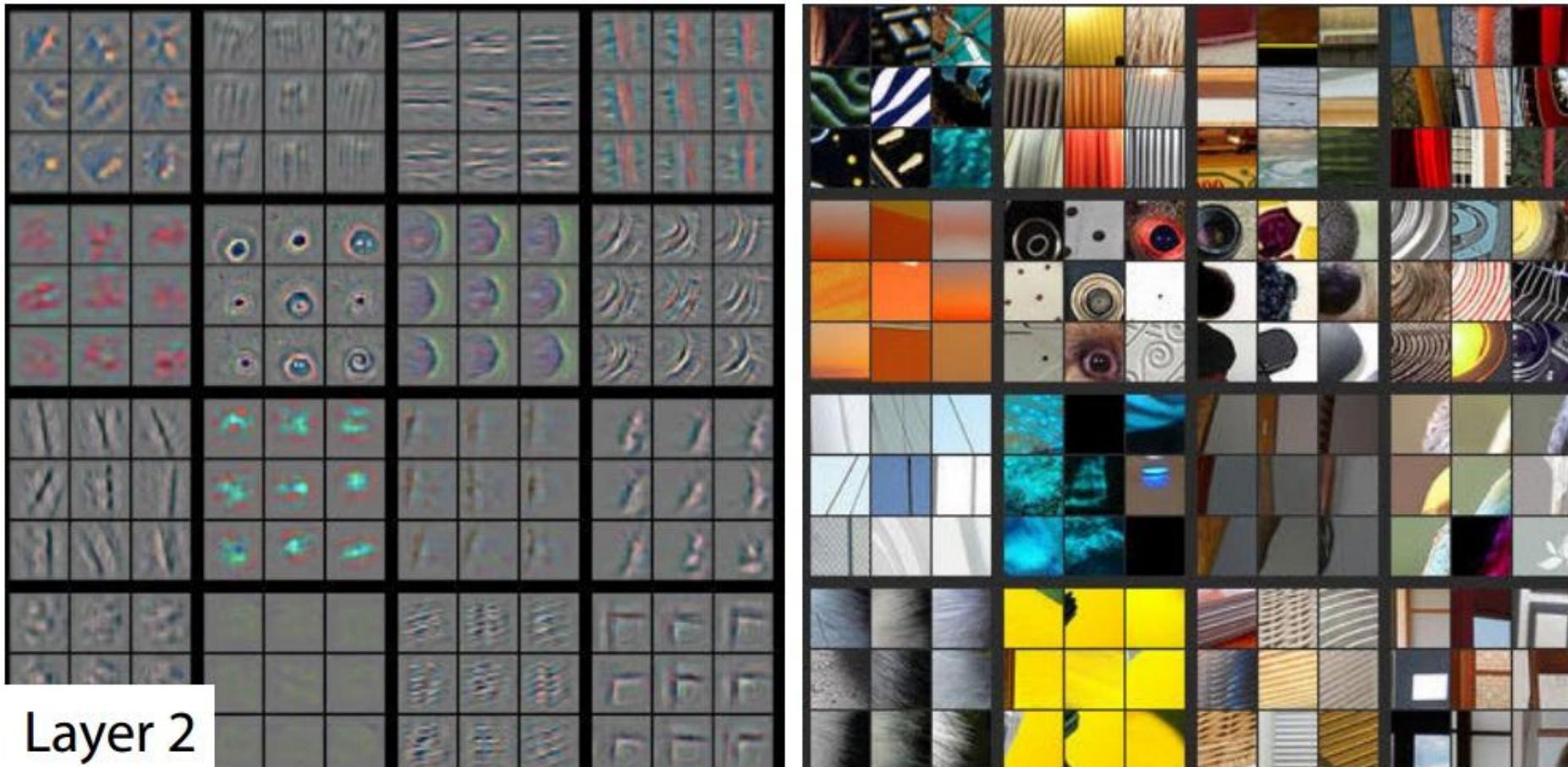
# Layer 1



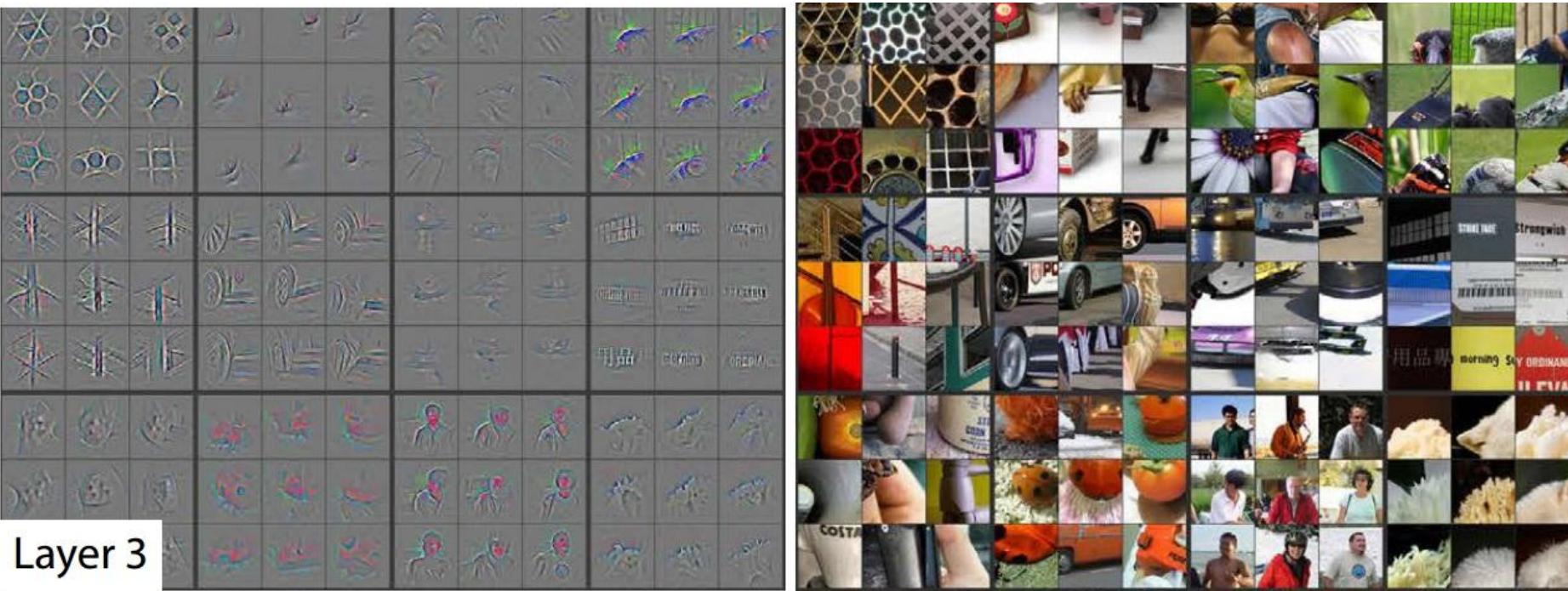
Layer 1



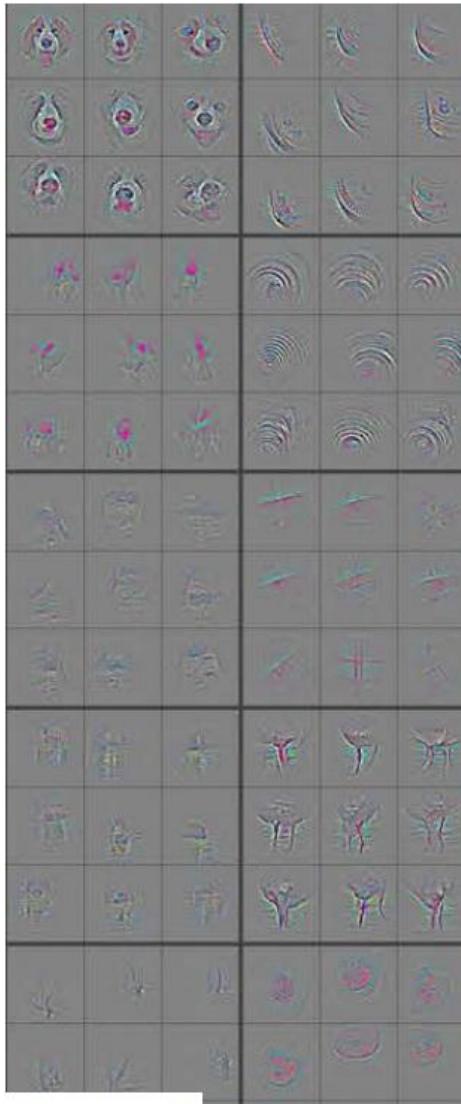
# Layer 2



# Layer 3

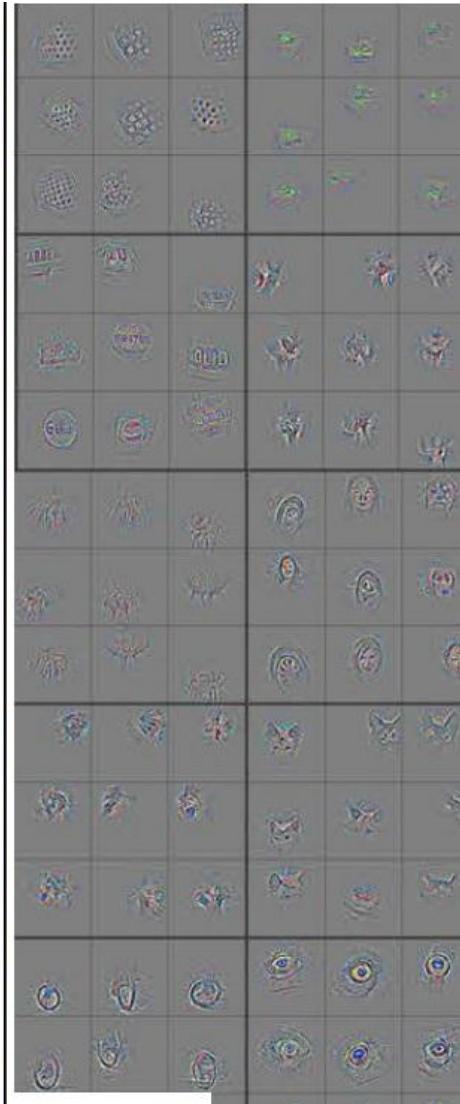


# Layer 4 and 5



Layer 4

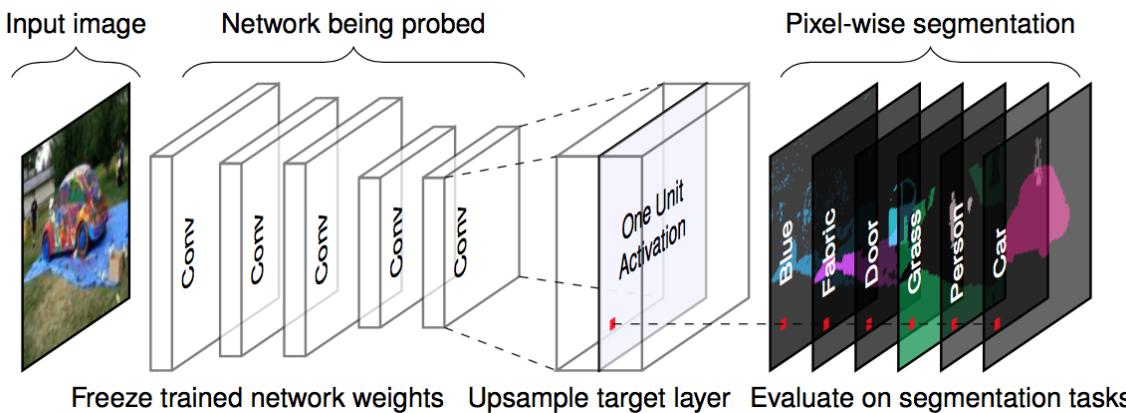
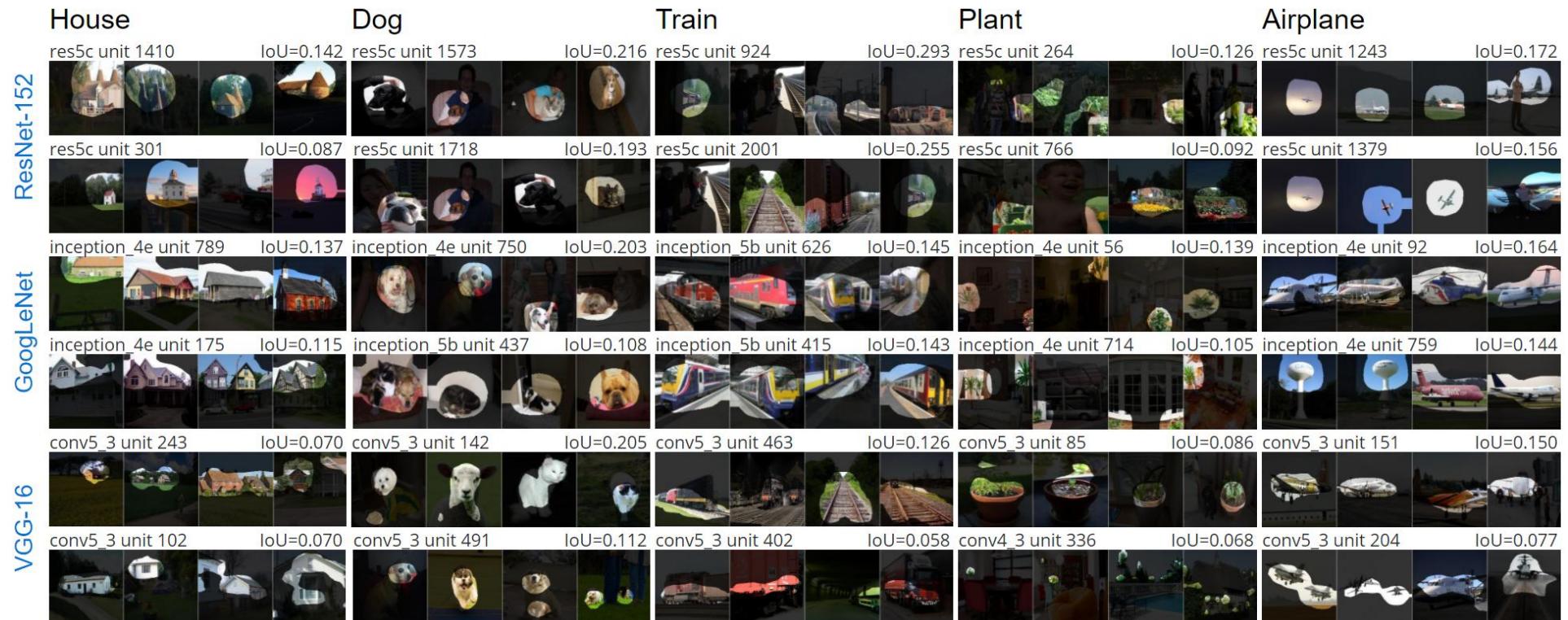
Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]  
Dr. Sander Ali Khowaja



Layer 5



# Network Dissection



# Deep learning library

- TensorFlow
  - Research + Production
- PyTorch
  - Research
- Caffe2
  - Production



# Things to remember

- Convolutional neural networks
  - A cascade of conv + ReLU + pool
  - Representation learning
  - Advanced architectures
  - Tricks for training CNN
- Visualizing CNN
  - Activation
  - Dissection



# Resources

- <http://deeplearning.net/>
  - Hub to many other deep learning resources
- <https://github.com/ChristosChristofidis/awesome-deep-learning>
  - A resource collection deep learning
- <https://github.com/kjw0612/awesome-deep-vision>
  - A resource collection deep learning for computer vision
- <http://cs231n.stanford.edu/syllabus.html>
  - Nice course on CNN for visual recognition



# Things to remember

- Overview
  - Neuroscience, Perceptron, multi-layer neural networks
- Convolutional neural network (CNN)
  - Convolution, nonlinearity, max pooling
  - CNN for classification and beyond
- Understanding and visualizing CNN
  - Find images that maximize some class scores; visualize individual neuron activation, input pattern and images; breaking CNNs
- Training CNN
  - Dropout; data augmentation; batch normalization; transfer learning



Geoffery Hinton is a legend (its just fo



## 'Godfather of AI' Geoffrey Hinton quits Google and warns over dangers of misinformation

The neural network pioneer says dangers of chatbots were 'quite scary' and warns they could be exploited by 'bad actors'



Dr Geoffrey Hinton, the 'godfather of AI', has left Google. Photograph: Linda Nylind/The Guardian

The man often touted as the godfather of AI has quit [Google](#), citing concerns over the flood of misinformation, the possibility for AI to upend the job market, and the "existential risk" posed by the creation of a true digital intelligence.

# Object Detection



**DOG, DOG, CAT**

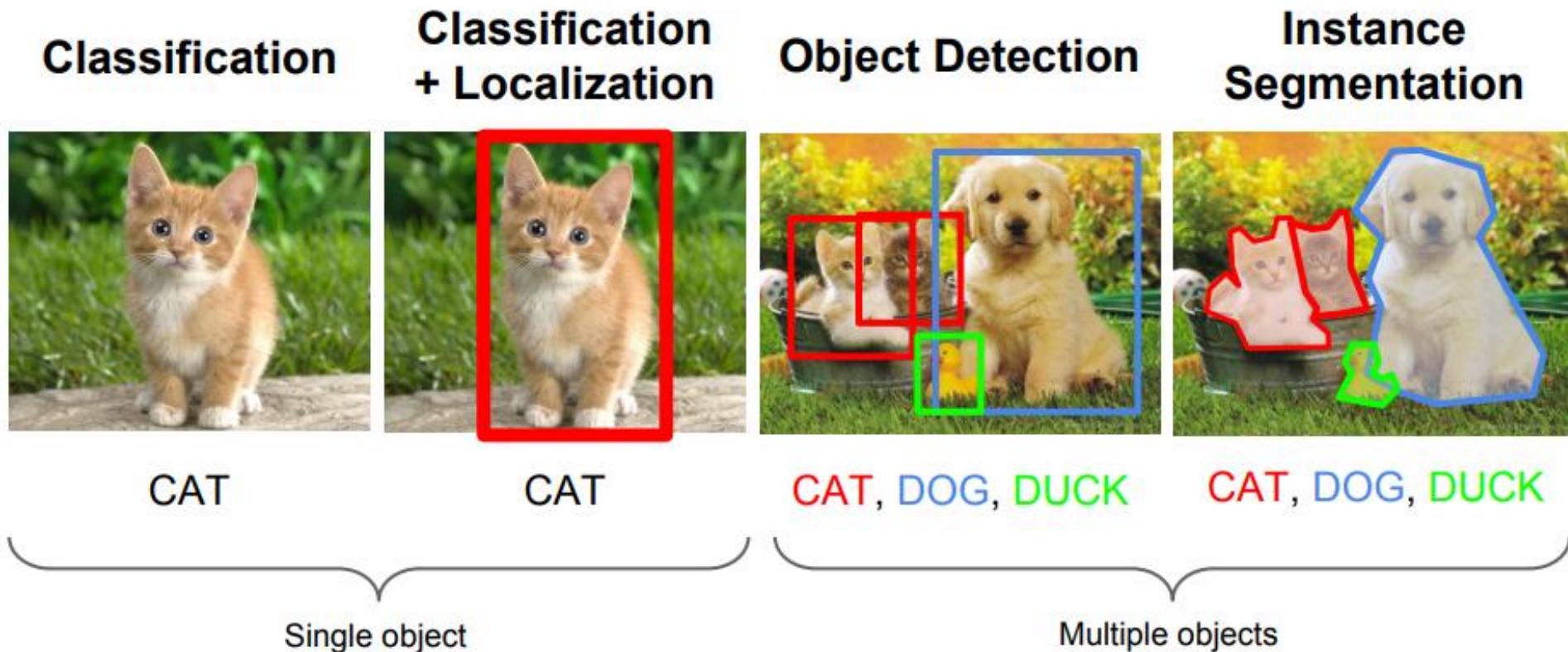


Many slides from D. Hoiem, J. Hays, J. Johnson, R. Girshick

Dr. Sander Ali Khowaja



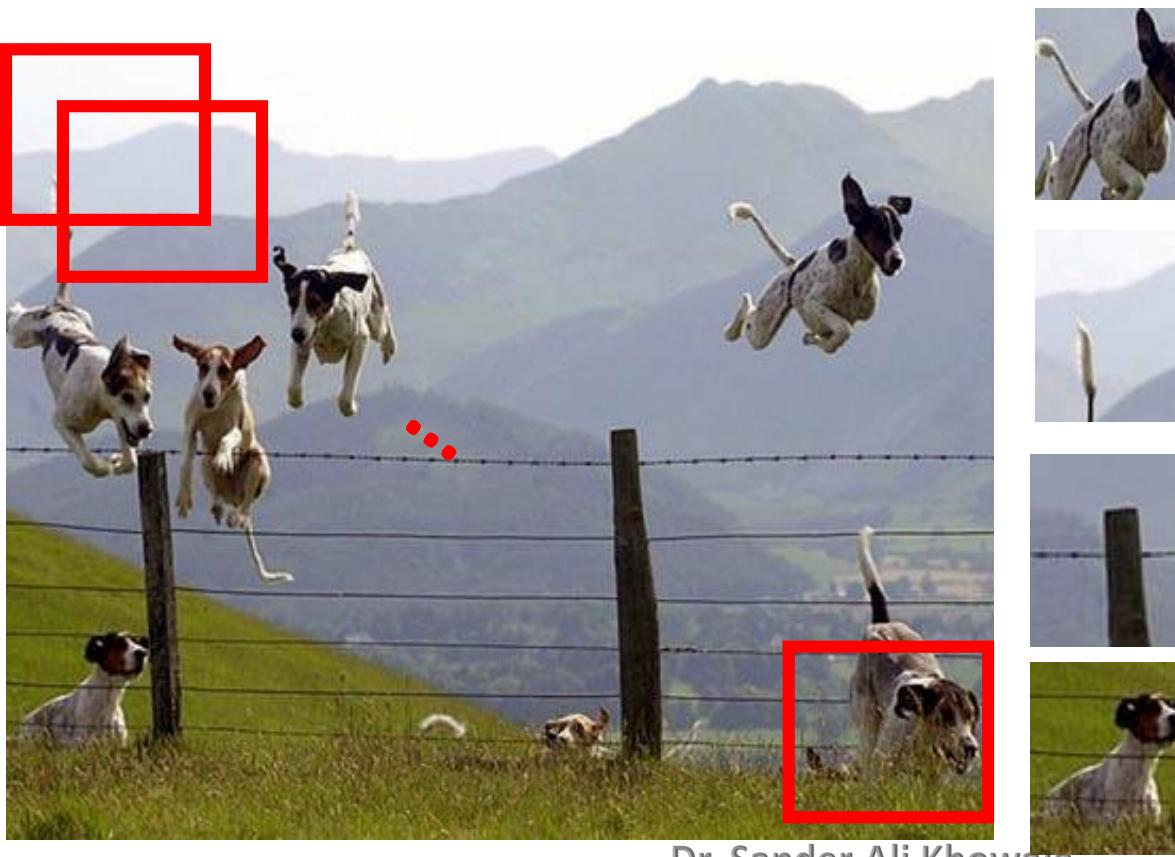
# Roadmap



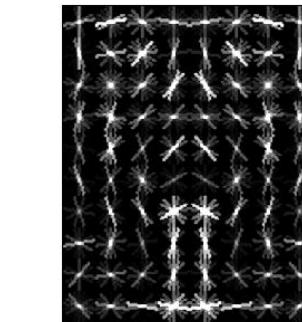


# Object Category Detection

- Focus on object search: “Where is it?”
- Build templates that quickly differentiate object patch from background patch



Dog Model



Object or  
Non-Object?



# Challenges in modeling the object class



Illumination



Object pose



Clutter



Occlusions



Intra-class  
appearance



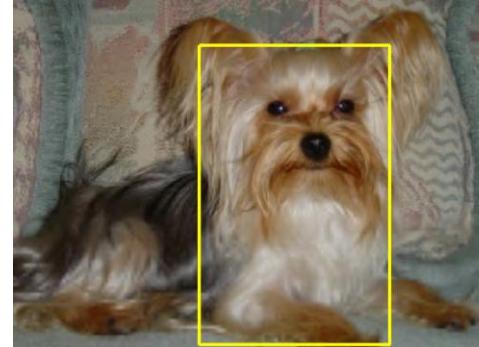
Viewpoint

# Challenges in modeling the non-object class

True  
Detections



Bad  
Localization



Confused with  
Similar Object



Misc. Background



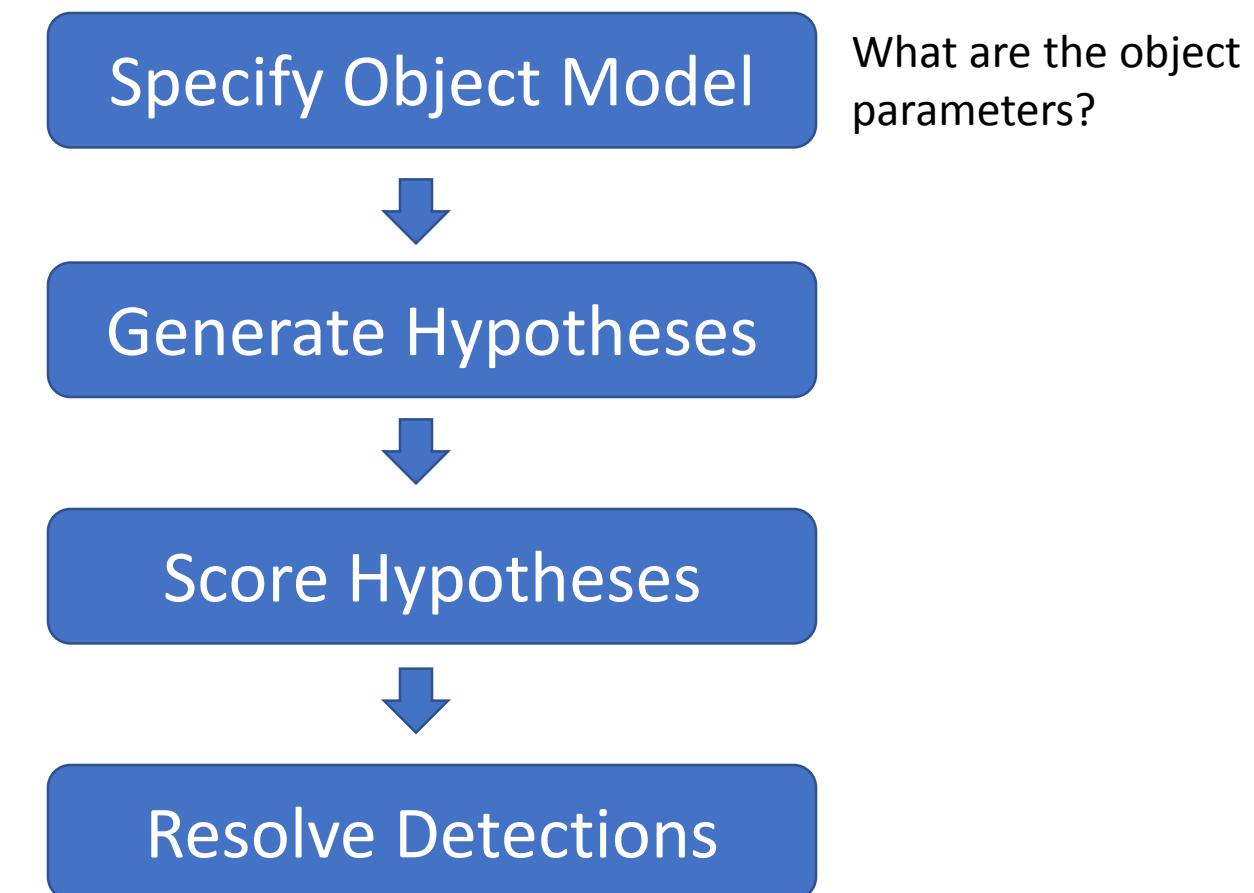
Confused with  
Dissimilar Objects



Dr. Sander Ali Khowaja



# General Process of Object Recognition



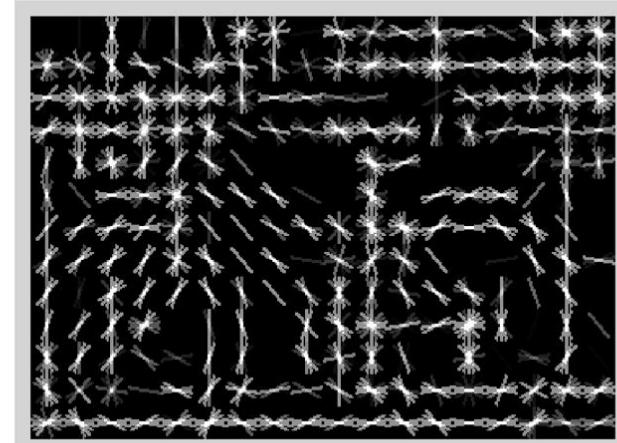
# Specifying an object model

## 1. Statistical Template in Bounding Box

- Object is some  $(x,y,w,h)$  in image
- Features defined wrt bounding box coordinates



Image

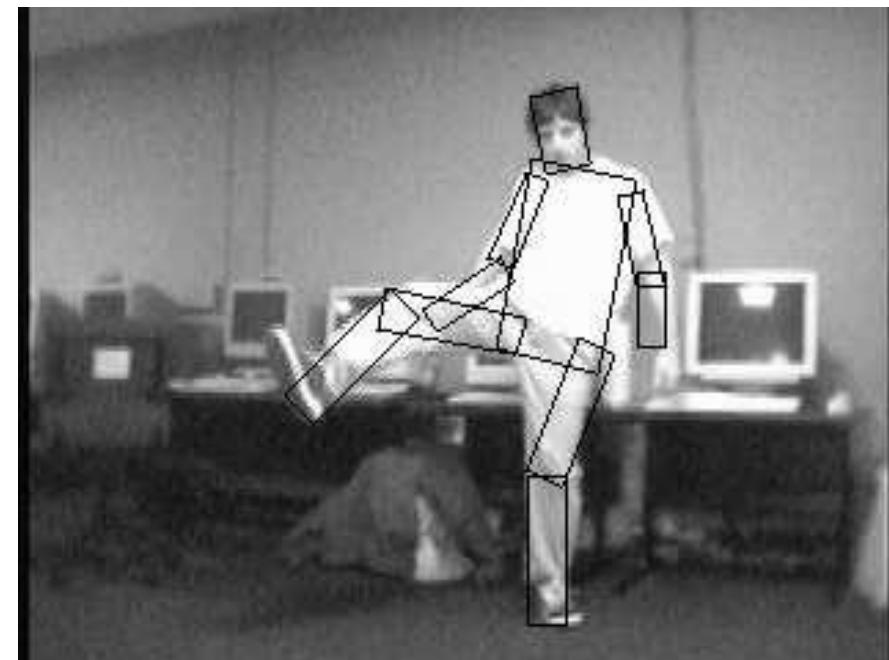
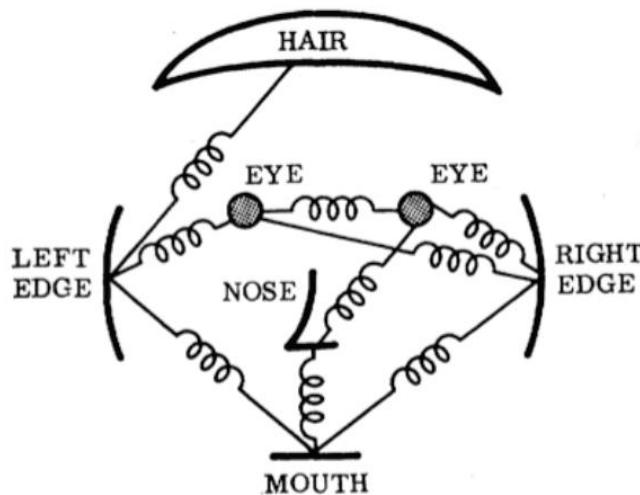


Template Visualization

# Specifying an object model

## 2. Articulated parts model

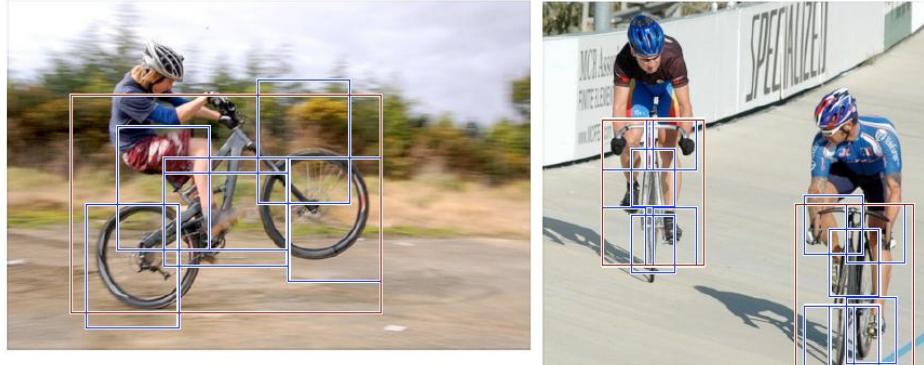
- Object is configuration of parts
- Each part is detectable



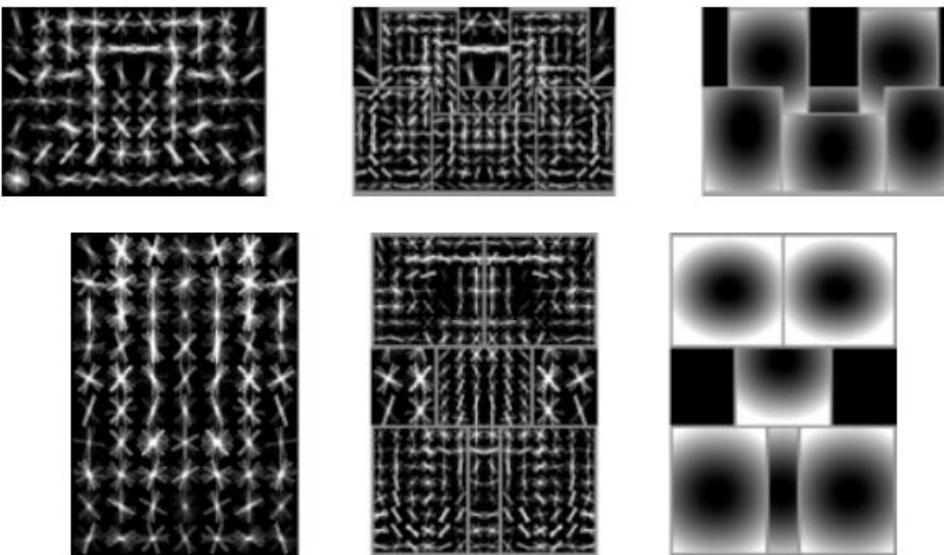
# Specifying an object model

## 3. Hybrid template/parts model

Detections



Template Visualization



root filters  
coarse resolution

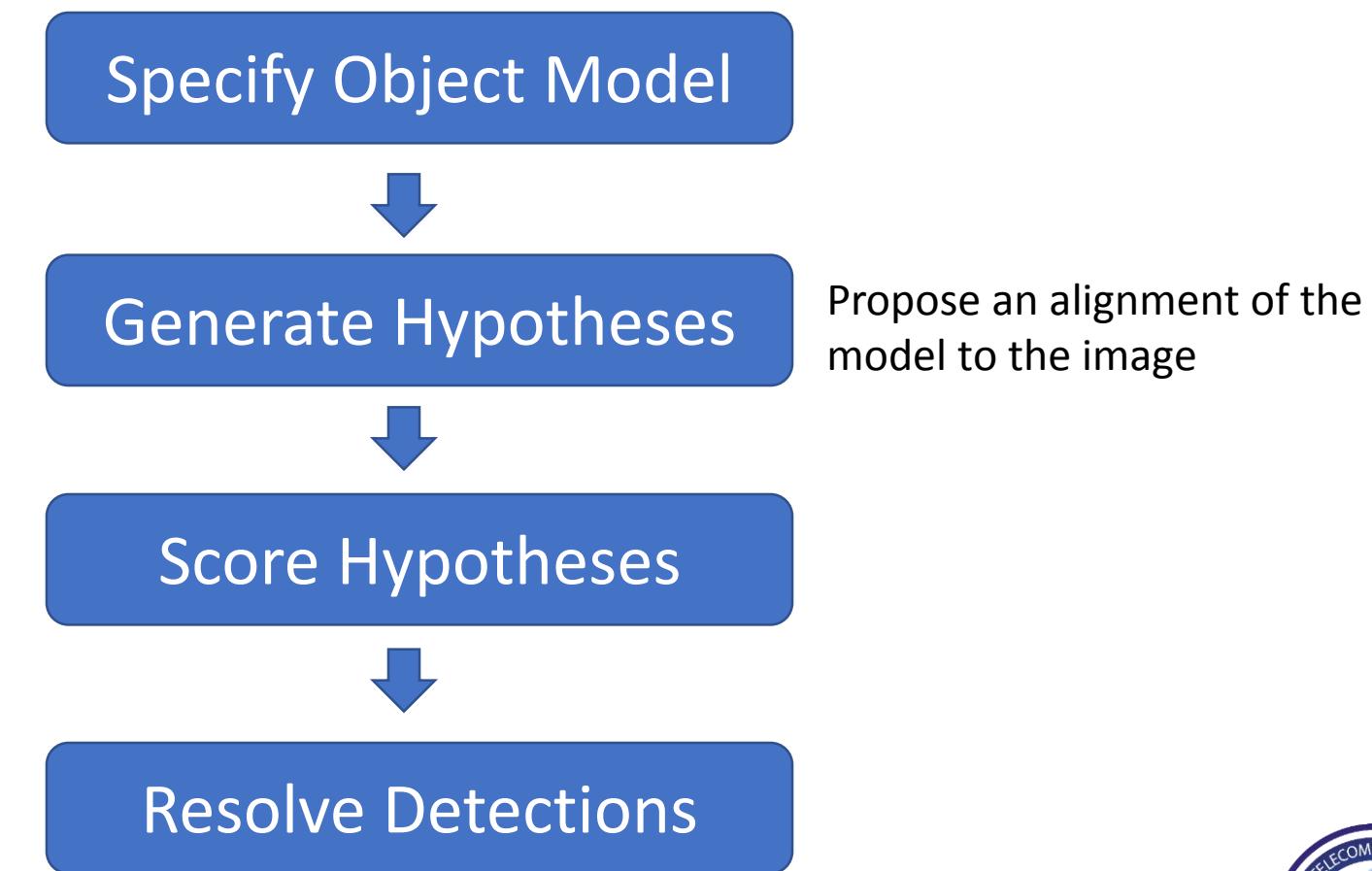
part filters  
finer resolution

deformation  
models

Felzenszwalb et al. 2008



# General Process of Object Recognition



# Generating hypotheses

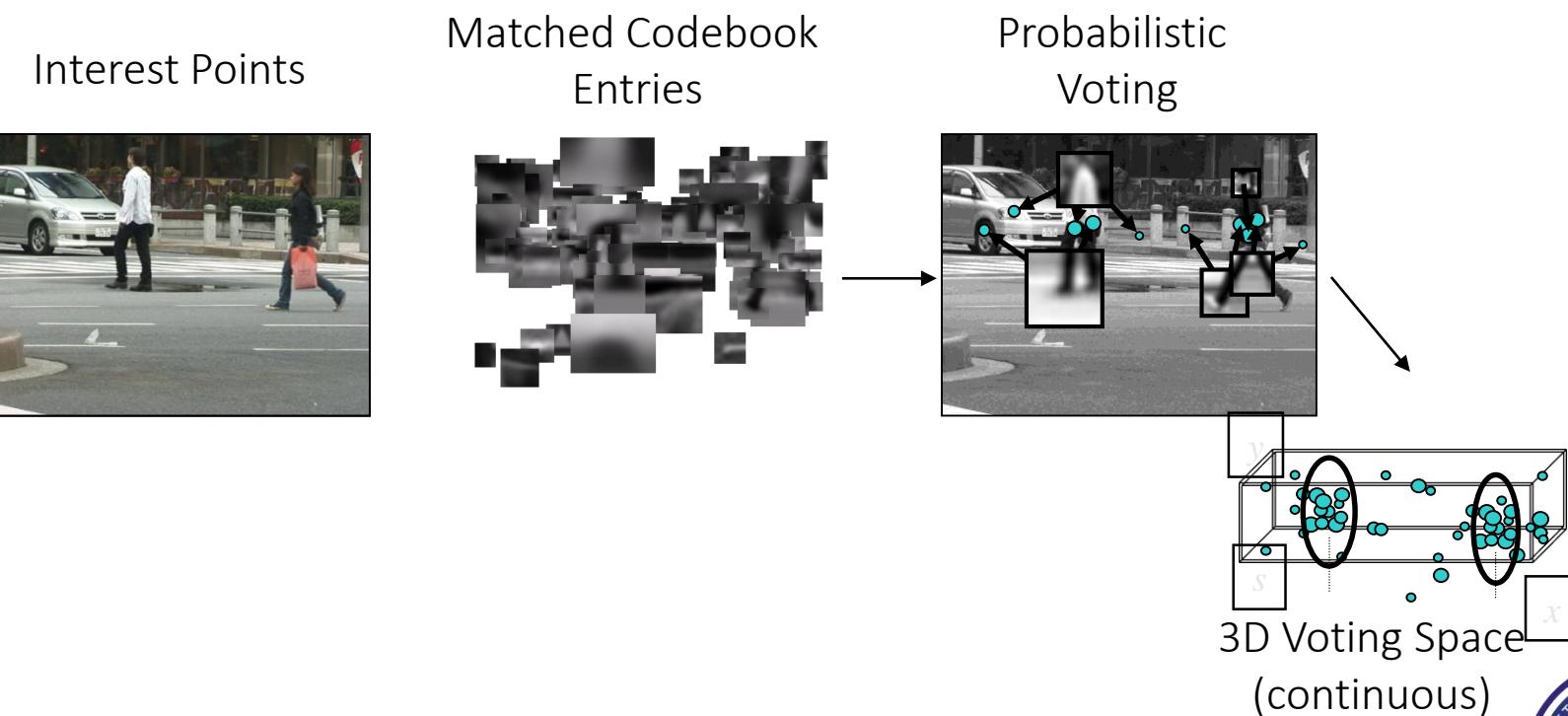
## 1. Sliding window

- Test patch at each location and scale



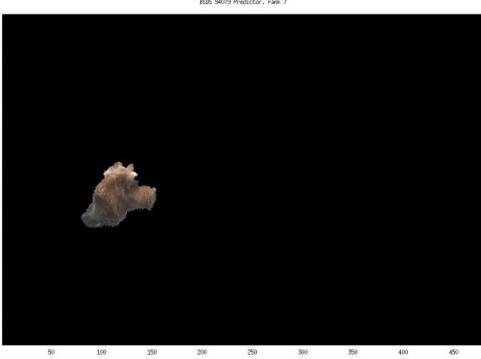
# Generating hypotheses

## 2. Voting from patches/keypoints

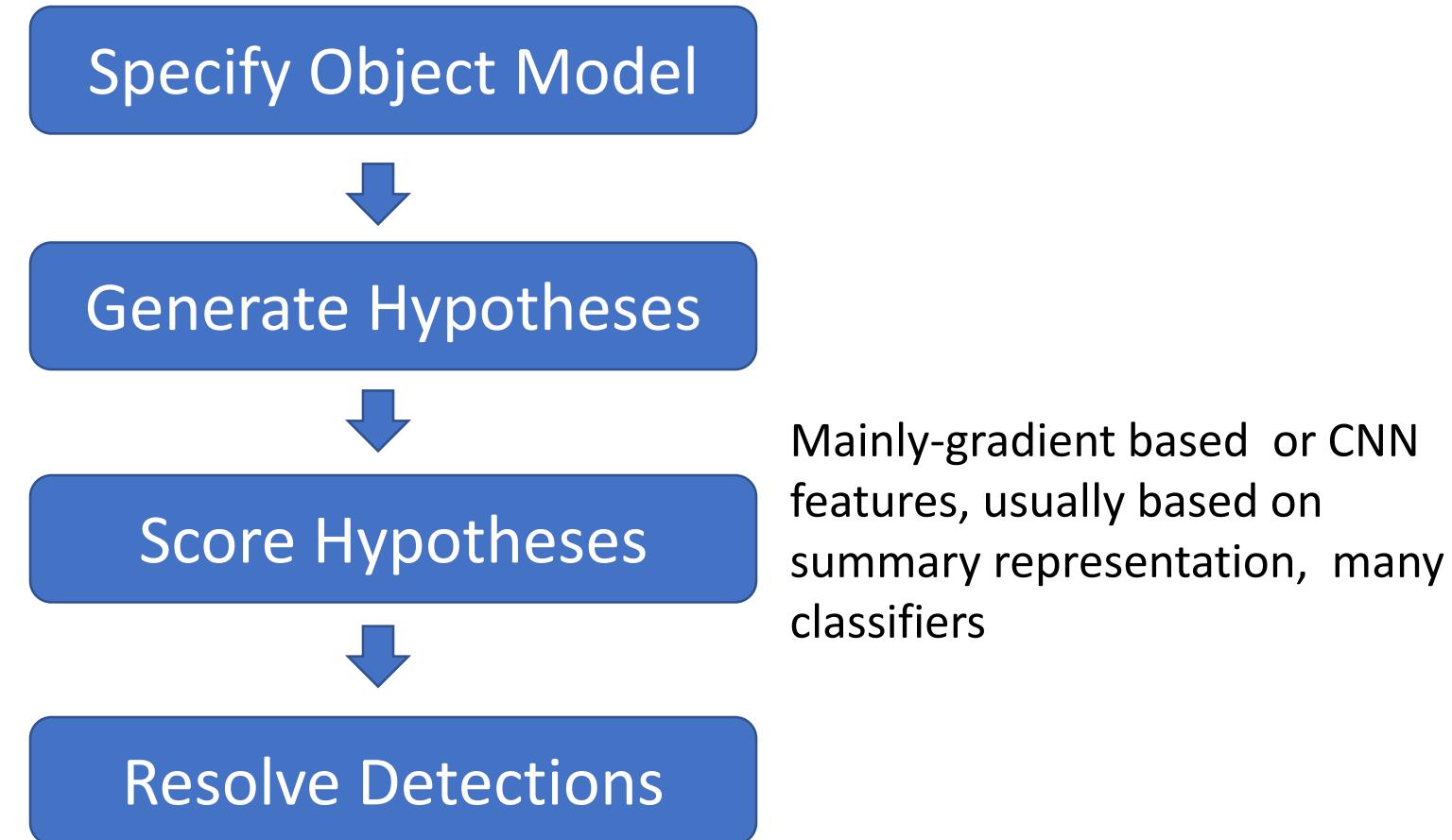


# Generating hypotheses

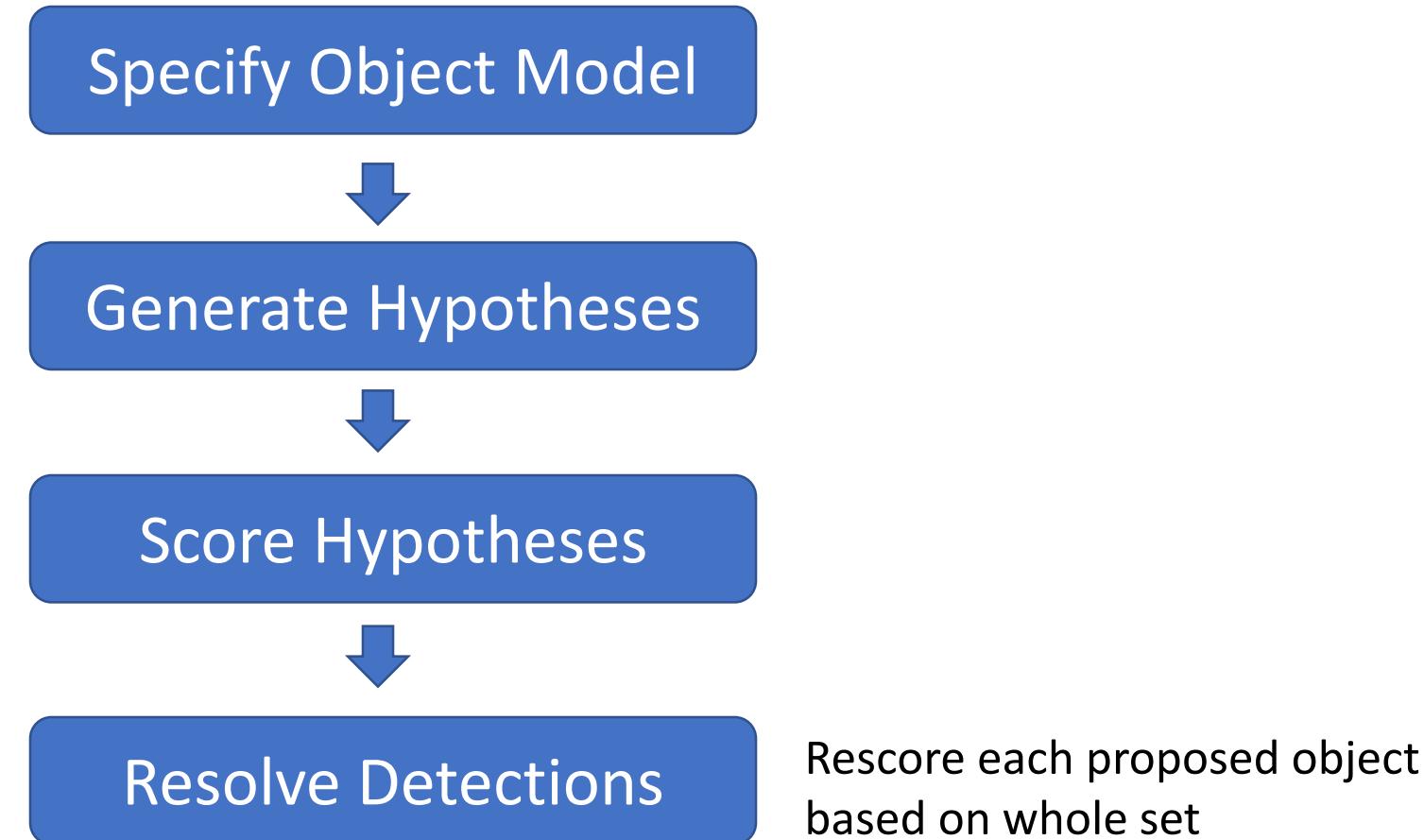
## 3. Region-based proposal



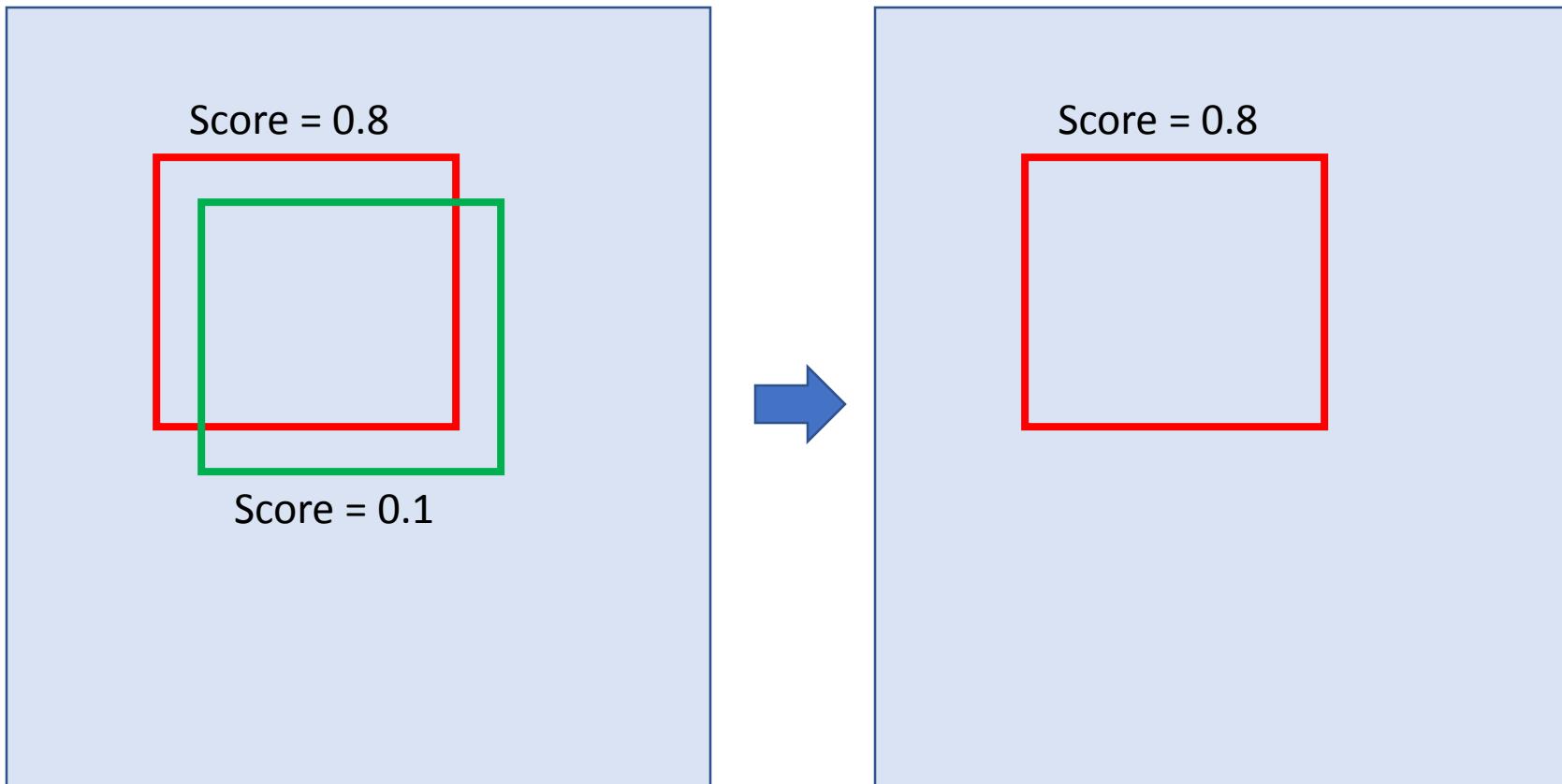
# General Process of Object Recognition



# General Process of Object Recognition

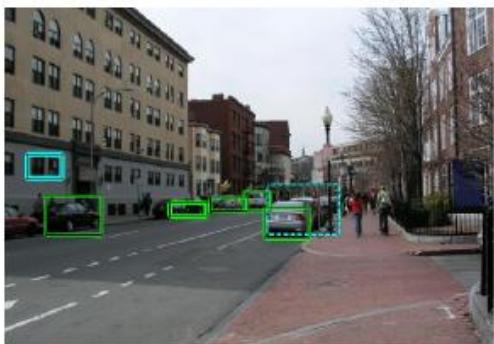


## 1. Non-max suppression



# Resolving detection scores

## 2. Context/reasoning



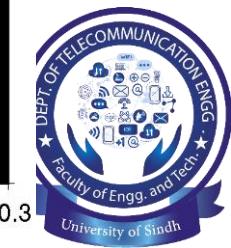
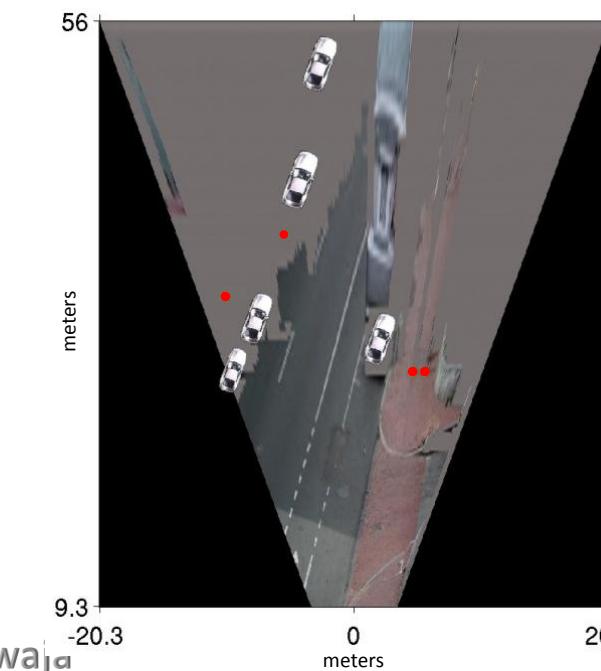
(g) Car Detections: Local



(h) Ped Detections: Local

Hoiem et al. 2006

Dr. Sander Ali Khawaja



# Object category detection in computer vision

Goal: detect all pedestrians, cars, monkeys, etc in image



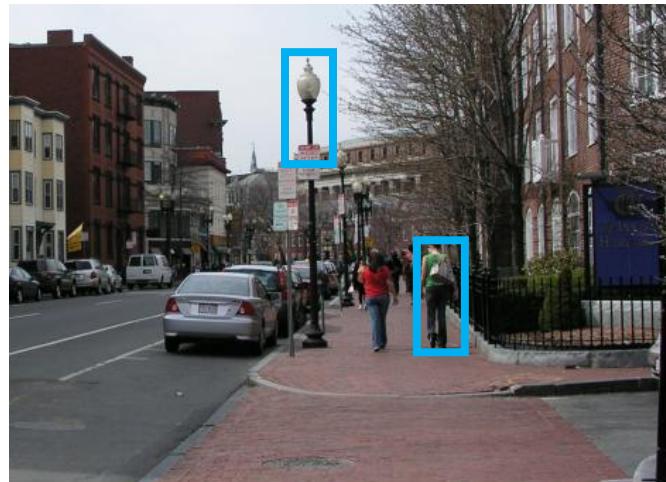
Dr. Sander Ali Khowaja



# Basic Steps of Category Detection

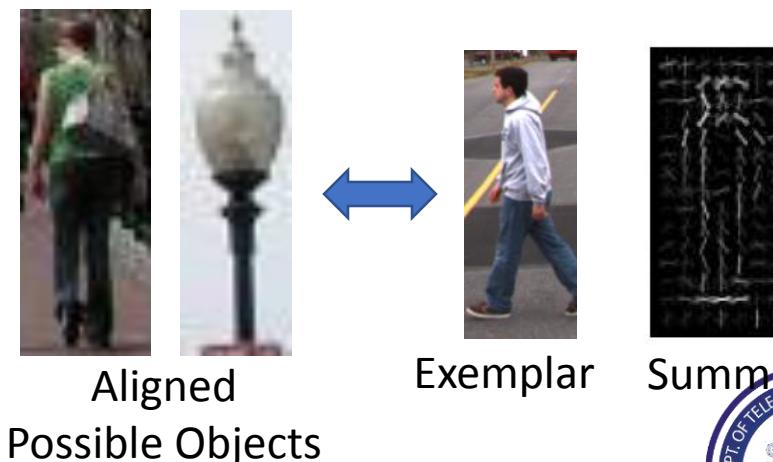
## 1. Align

- E.g., choose position, scale orientation
- How to make this tractable?



## 2. Compare

- Compute similarity to an example object or to a summary representation
- Which differences in appearance are important?





Dr. Sander Ali Khowaja



# Each window is separately classified



Dr. Sander Ali Khawaja



# Statistical Template

- Object model = sum of scores of features at fixed positions



$$+3 \text{ } +2 \text{ } -2 \text{ } -1 \text{ } -2.5 = -0.5 > 7.5$$

?

Non-object



$$+4 \text{ } +1 \text{ } +0.5 \text{ } +3 \text{ } +0.5 = 10.5 > 7.5$$

?

Object

# CNN as feature extractor



Image credit: Justin Johnson

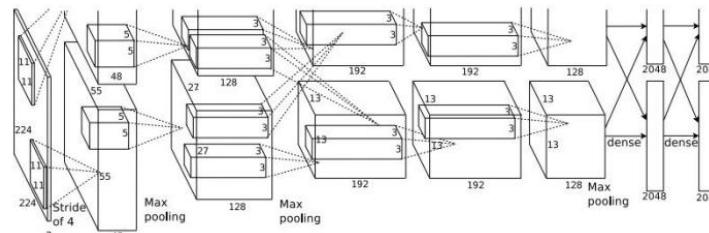
Dr. Sander Ali Khowaja



# CNN as feature extractor



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

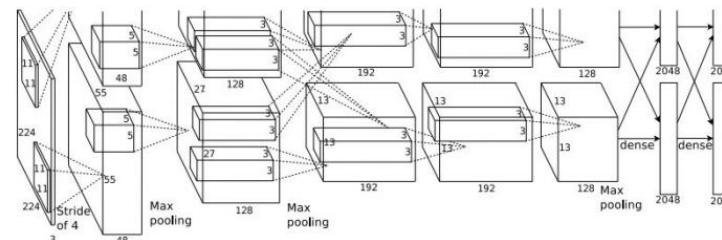


Dog? NO  
Cat? NO  
Background? YES

# CNN as feature extractor



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

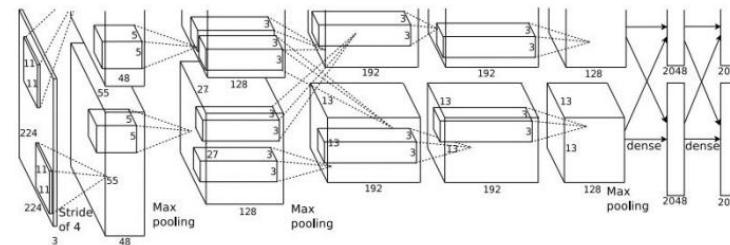


Dog? YES  
Cat? NO  
Background? NO

# CNN as feature extractor



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

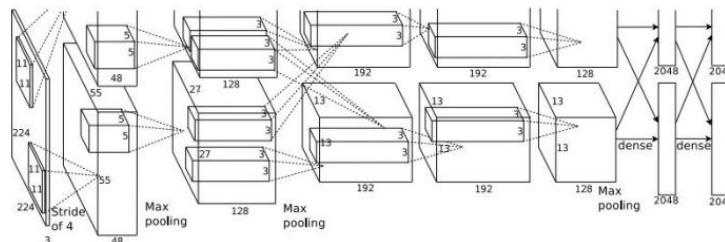


Dog? YES  
Cat? NO  
Background? NO

# CNN as feature extractor



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? YES  
Background? NO

# CNN as feature extractor

- What could be the problems?



# CNN as feature extractor

- What could be the problems?
  - Suppose we have a  $600 \times 600$  image, if sliding window size is  $20 \times 20$ , then have  $(600-20+1) \times (600-20+1) = \sim 330,000$  windows

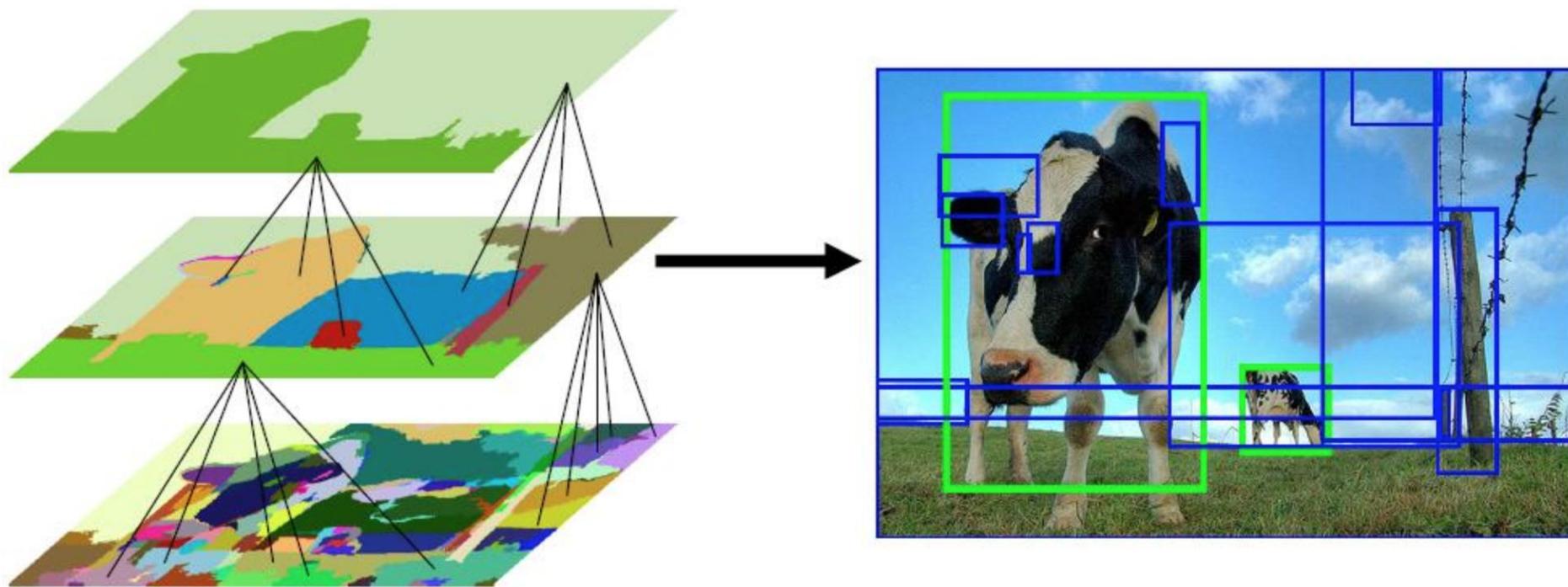
- What could be the problems?
  - Suppose we have a  $600 \times 600$  image, if sliding window size is  $20 \times 20$ , then have  $(600-20+1) \times (600-20+1) = \sim 330,000$  windows
  - Sometimes we want to have more accurate results -> multi-scale detection
    - Resize image
    - Multi-scale sliding window

- What could be the problems?
  - Suppose we have a  $600 \times 600$  image, if sliding window size is  $20 \times 20$ , then have  $(600-20+1) \times (600-20+1) = \sim 330,000$  windows
  - Sometimes we want to have more accurate results -> multi-scale detection
    - Resize image
    - Multi-scale sliding window
  - For each image, we need to do the forward pass in the CNN for  $\sim 330,000$  times. -> Slow!!!

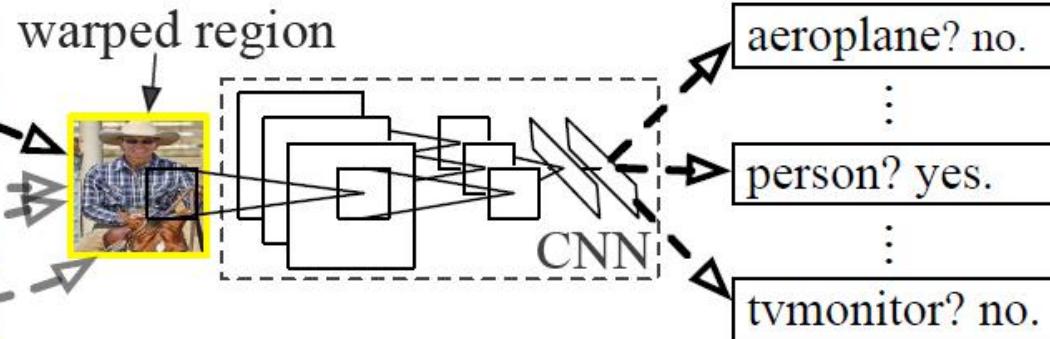
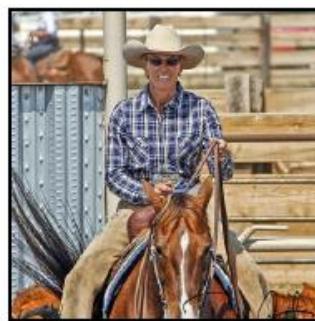


# Region Proposal

- Solution
  - Use some fast algorithms to filter out some regions first, only feed the potential region (region proposals) into CNN
  - E.g. selective search



# R-CNN (Girshick et al. CVPR 2014)



1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

- Replace sliding windows with “selective search” region proposals (Uijlings et al. IJCV 2013)
- Extract rectangles around regions and resize to 227x227
- Extract features with fine-tuned CNN (that was initialized with network trained on ImageNet before training)
- Classify last layer of network features with SVM, refine bounding box localization (bbox regression) simultaneously

<http://arxiv.org/pdf/1311.2524.pdf>

Dr. Sander Ali Khowaja



# Bounding Box Regression

- Intuition
  - If you observe part of the object, according to the seen examples, you should be able to refine the localization
  - E.g. given the red box below, since you've seen many airplanes, you know this is not a good localization, you will adjust it to the green one



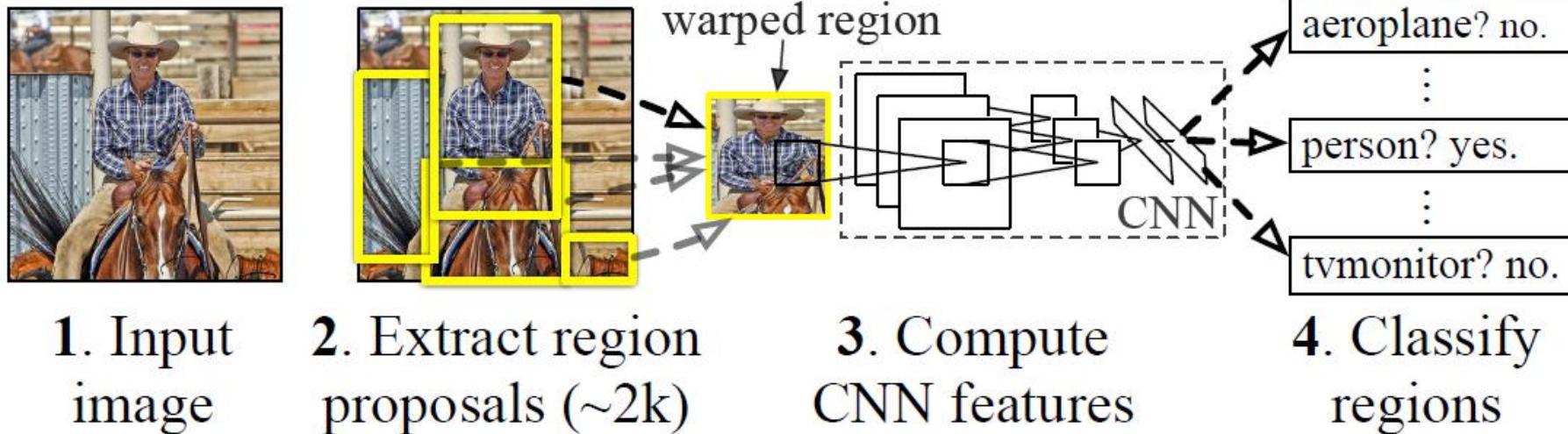
# Bounding Box Regression

- Intuition

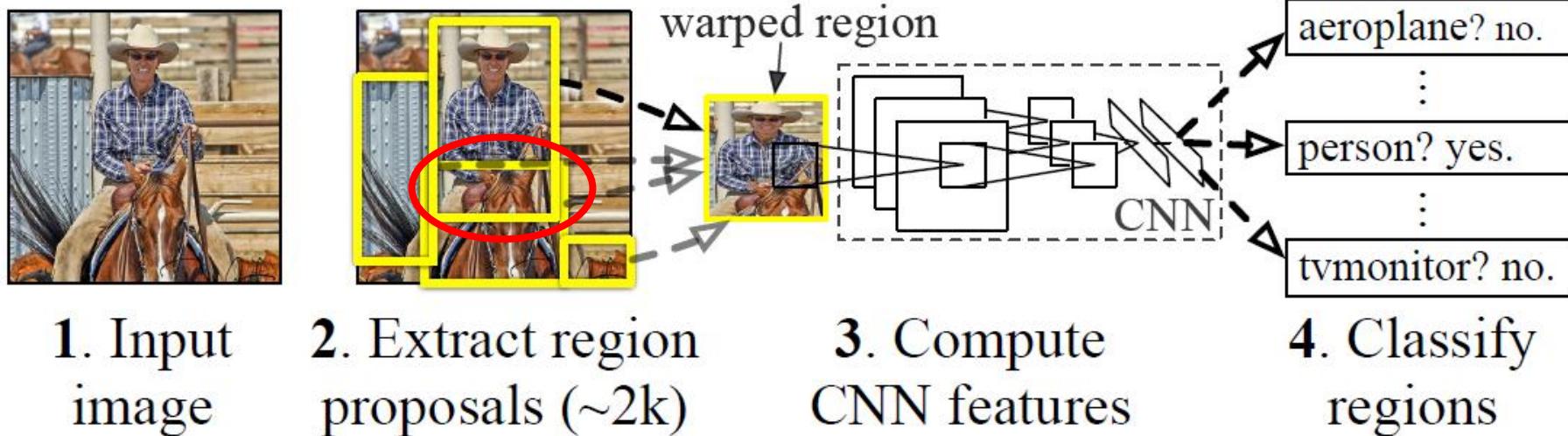
- If you observe part of the object, according to the seen examples, you should be able to refine the localization
- E.g. given the red box below, since you've seen many airplanes, you know this is not a good localization, you will adjust it to the green one



- What could be the problems?

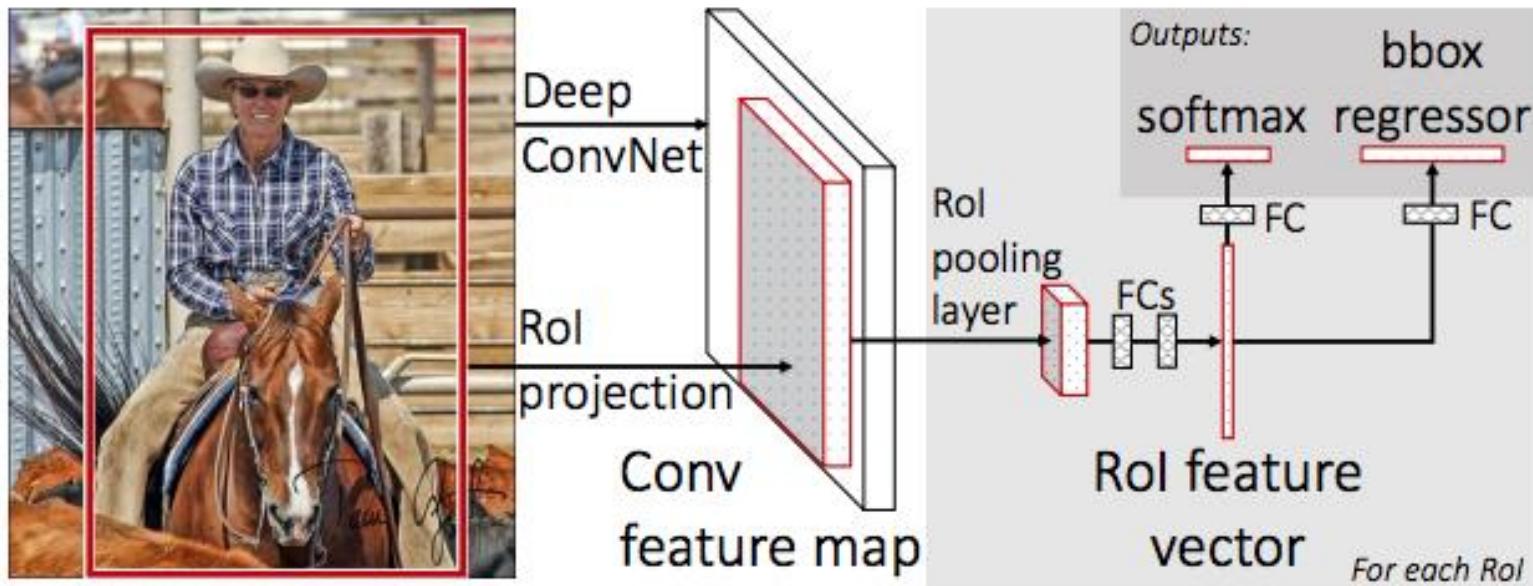


- What could be the problems?
  - Repetitive computation! For overlapping regions, we feed it multiple times into CNN



- Solution

- Why not feed the whole image into CNN only once! Then crop features instead of image itself

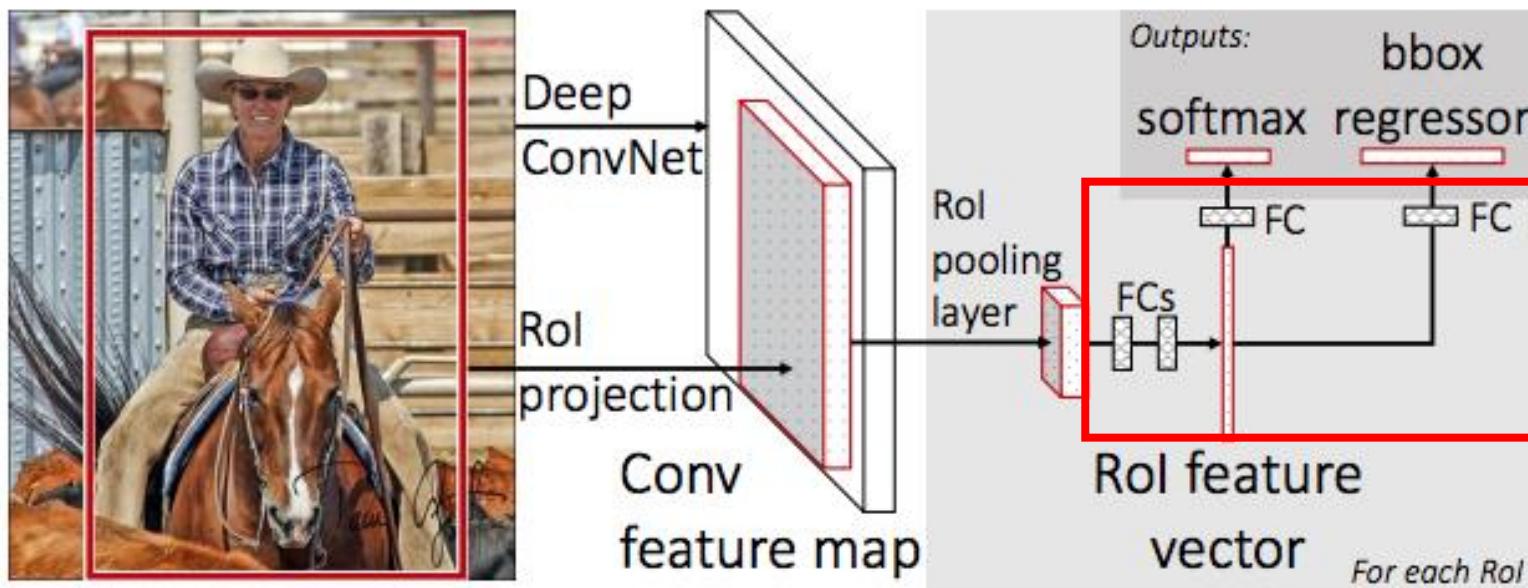


<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



- How to crop features?
  - Since we have fully-connected layers, the size of feature map for each bounding box should be a fixed number

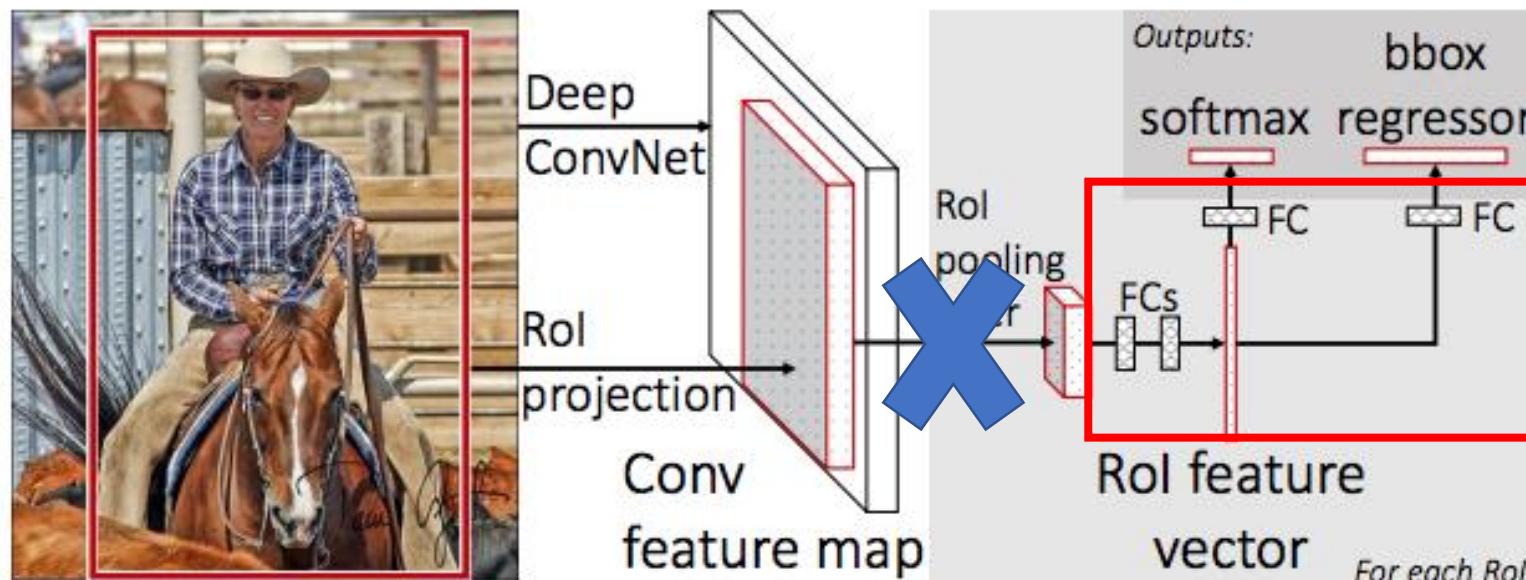


<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



- How to crop features?
  - Since we have fully-connected layers, the size of feature map for each bounding box should be a fixed number
  - Resize/Interpolate the feature map as fixed size?
    - Not optimal. This operation is hard to backprop -> we cannot train the conv layers for this problem



<https://arxiv.org/pdf/1504.08083.pdf>

- How to crop features?
  - Since we have fully-connected layers, the size of feature map for each bounding box should be a fixed number
  - Resize/Interpolate the feature map as fixed size?
    - Not optimal. This operation is hard to backprop -> we cannot train the conv layers for this problem
  - ROI (Region of Interest) Pooling

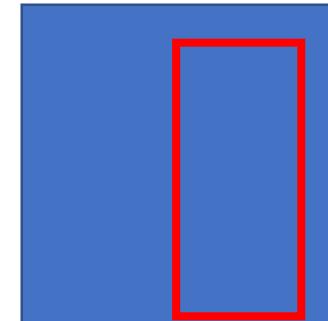
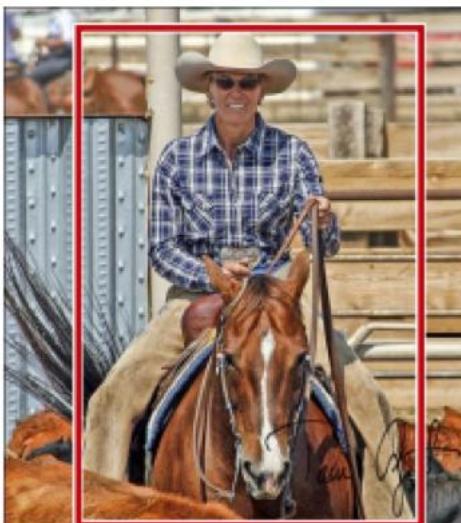
<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



# Rol Pooling

- Step 1: Get bounding box for feature map from bounding box for image
  - Due to the (down)convolution / pooling operations, feature map would have a smaller size than the original image



Feature map

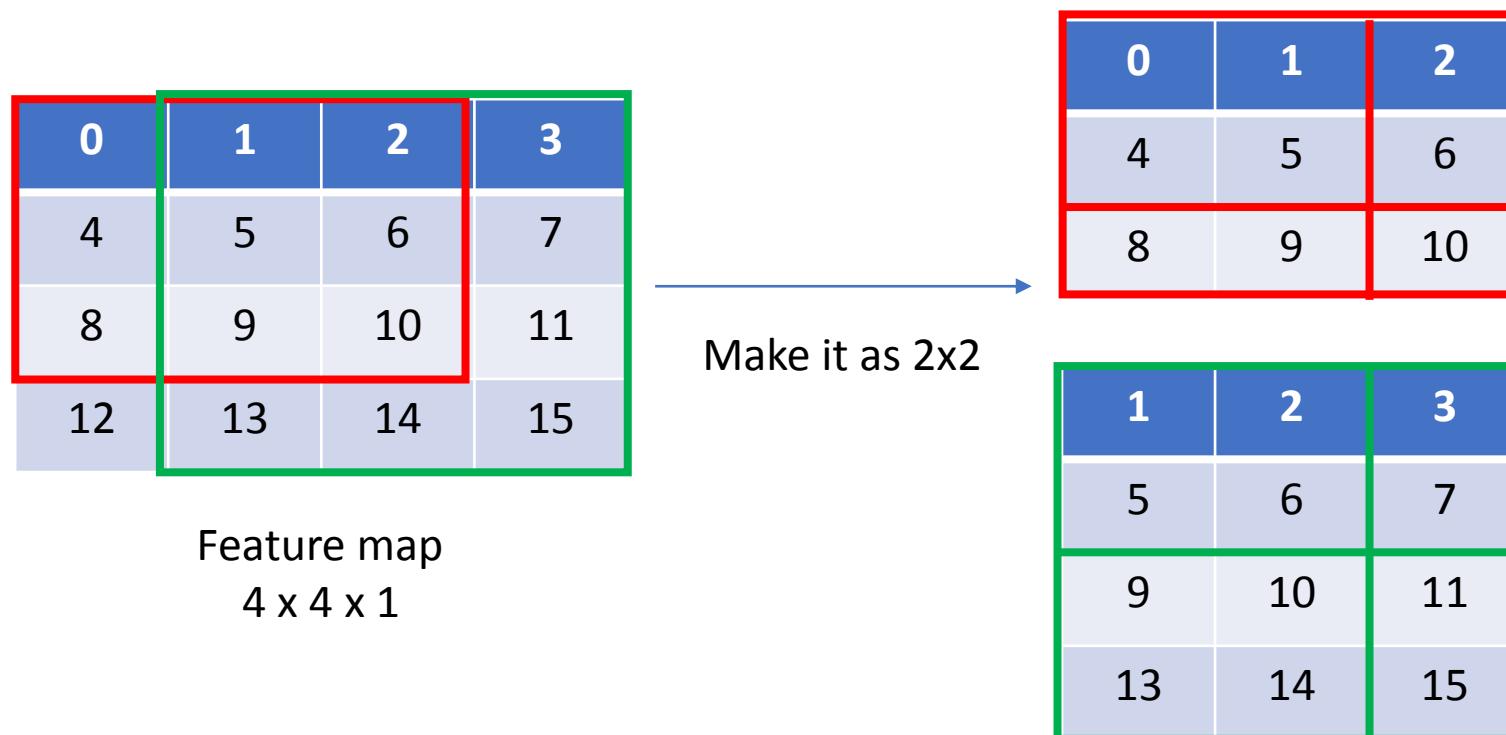
<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



# RoI Pooling

- Step 2: Divide cropped feature map into fixed number of sub-regions
  - The last column and last row might be smaller



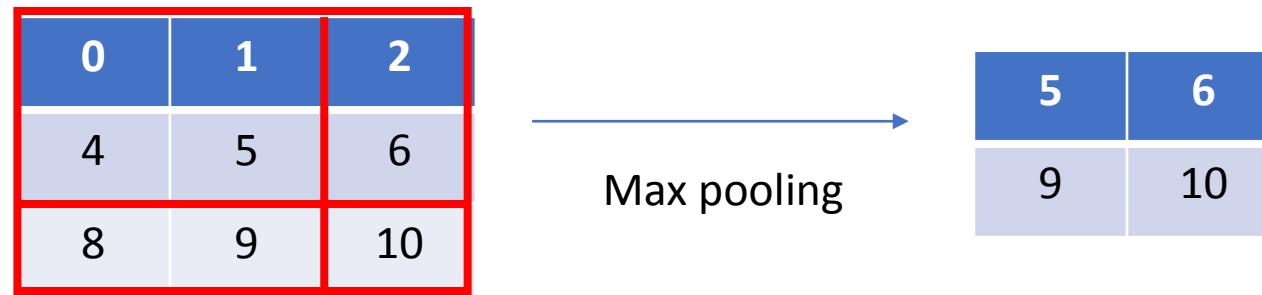
<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



# Rol Pooling

- Step 3: For each sub-region, perform max pooling (pick the max one)



<https://arxiv.org/pdf/1504.08083.pdf>

Dr. Sander Ali Khowaja



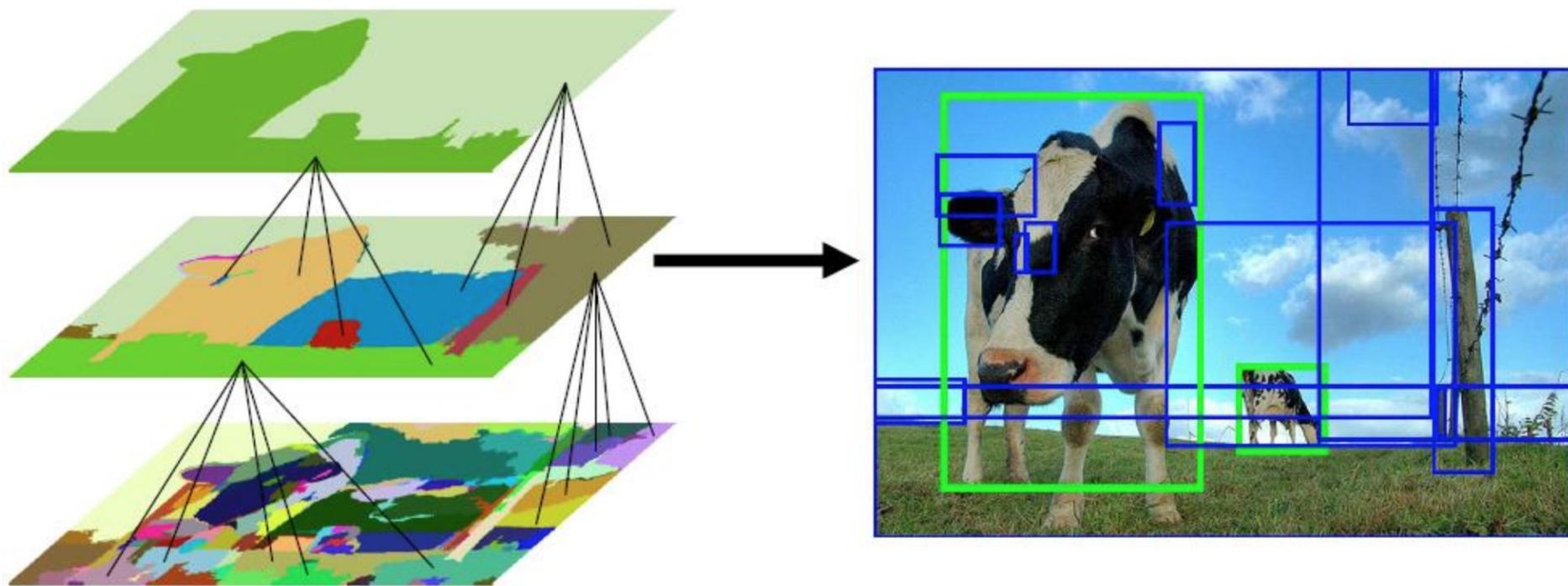
# Fast R-CNN (Girshick et al. ICCV 2015)

- What could be the problems?



# Fast R-CNN (Girshick et al. ICCV 2015)

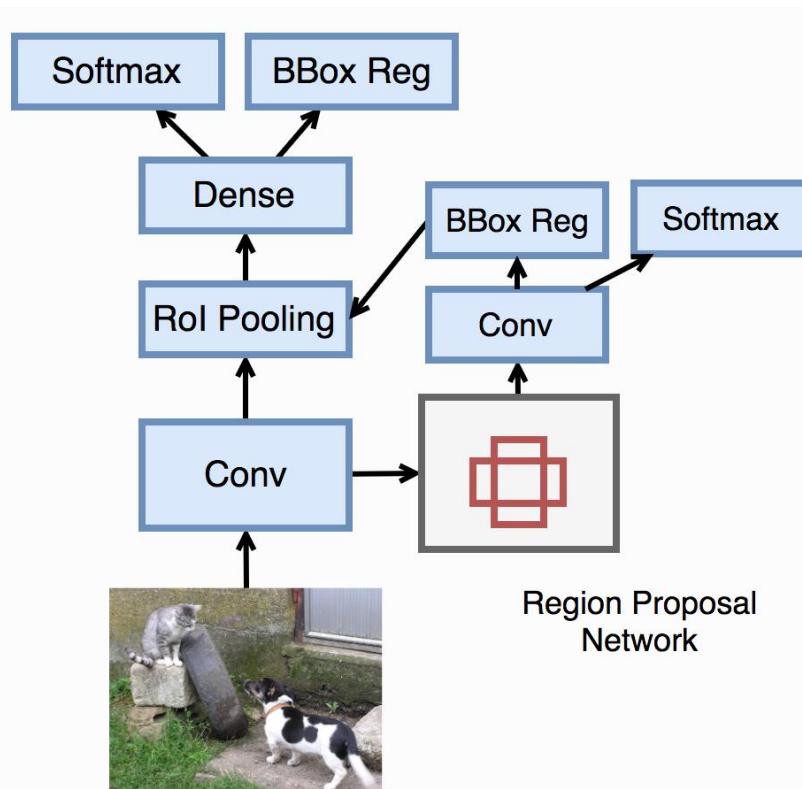
- What could be the problems?
  - Why we need the region proposal pre-processing step?  
That's not “deep learning” at all. Not cool!



# Faster R-CNN (Ren et al. NIPS 2015)

- Solution

- Why not generate region proposals using CNN??! -> RPN



<https://arxiv.org/pdf/1506.01497.pdf>

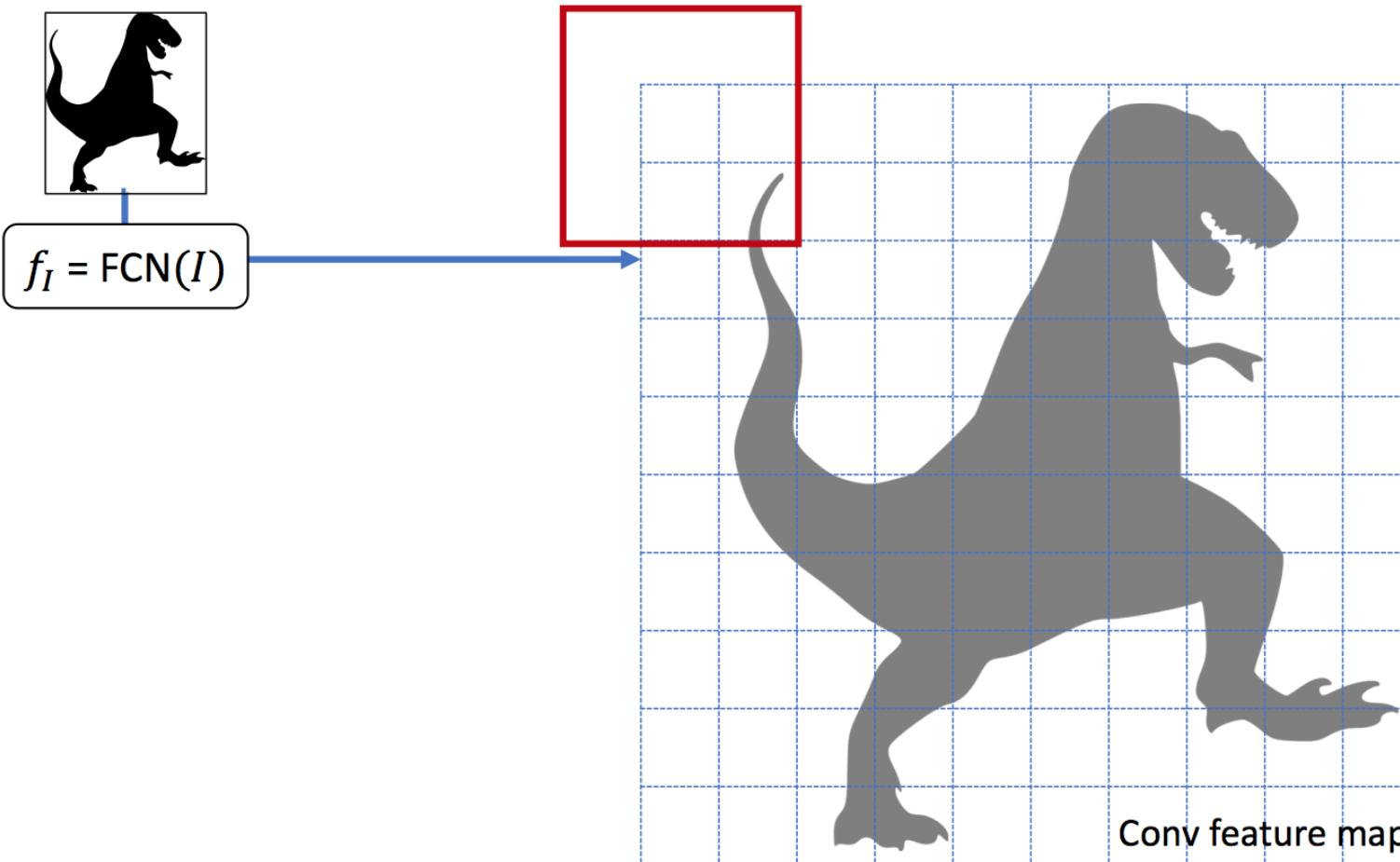
Dr. Sander Ali Khawaja

Image credit:

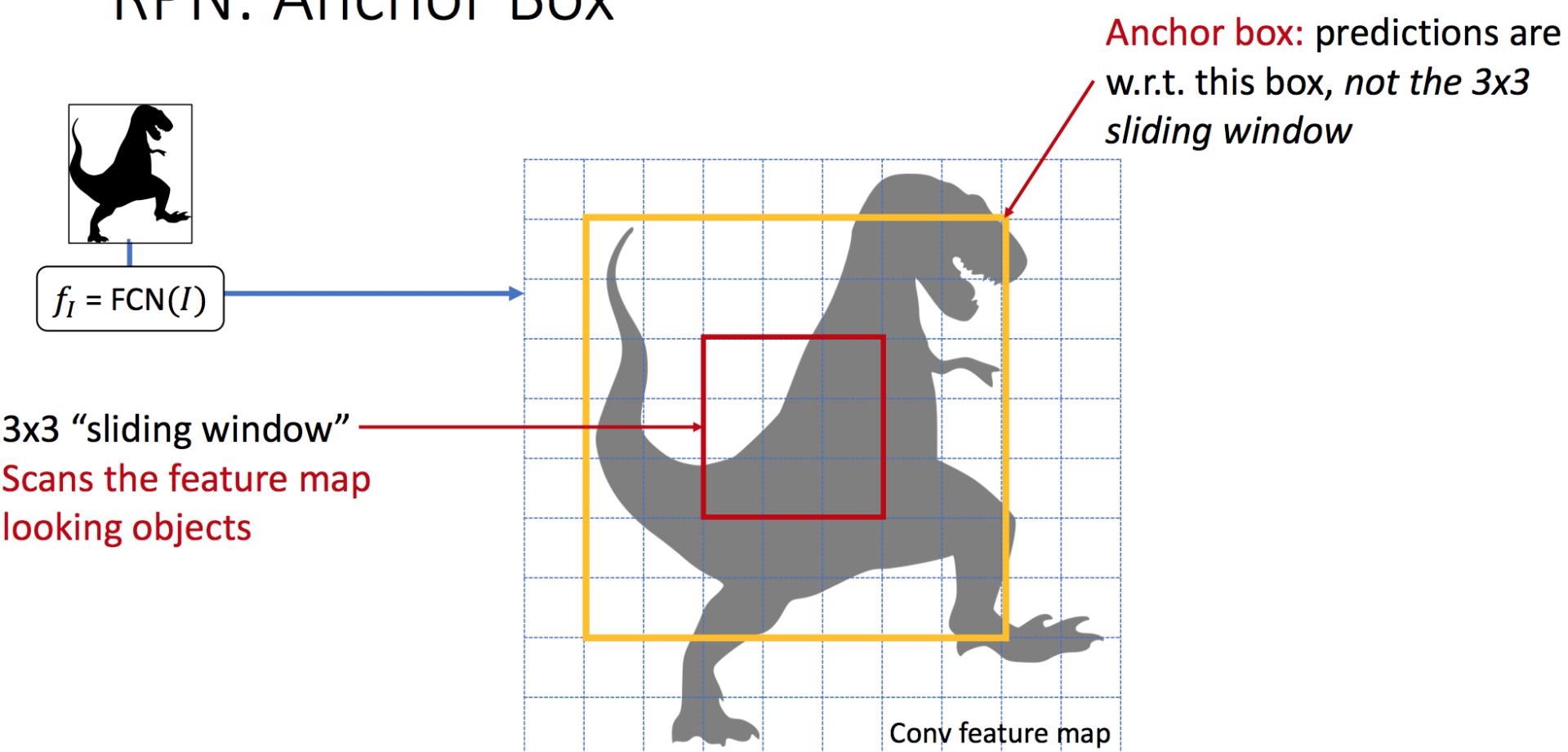
[http://zh.gluon.ai/chapter\\_computer-vision/object-detection.html](http://zh.gluon.ai/chapter_computer-vision/object-detection.html)



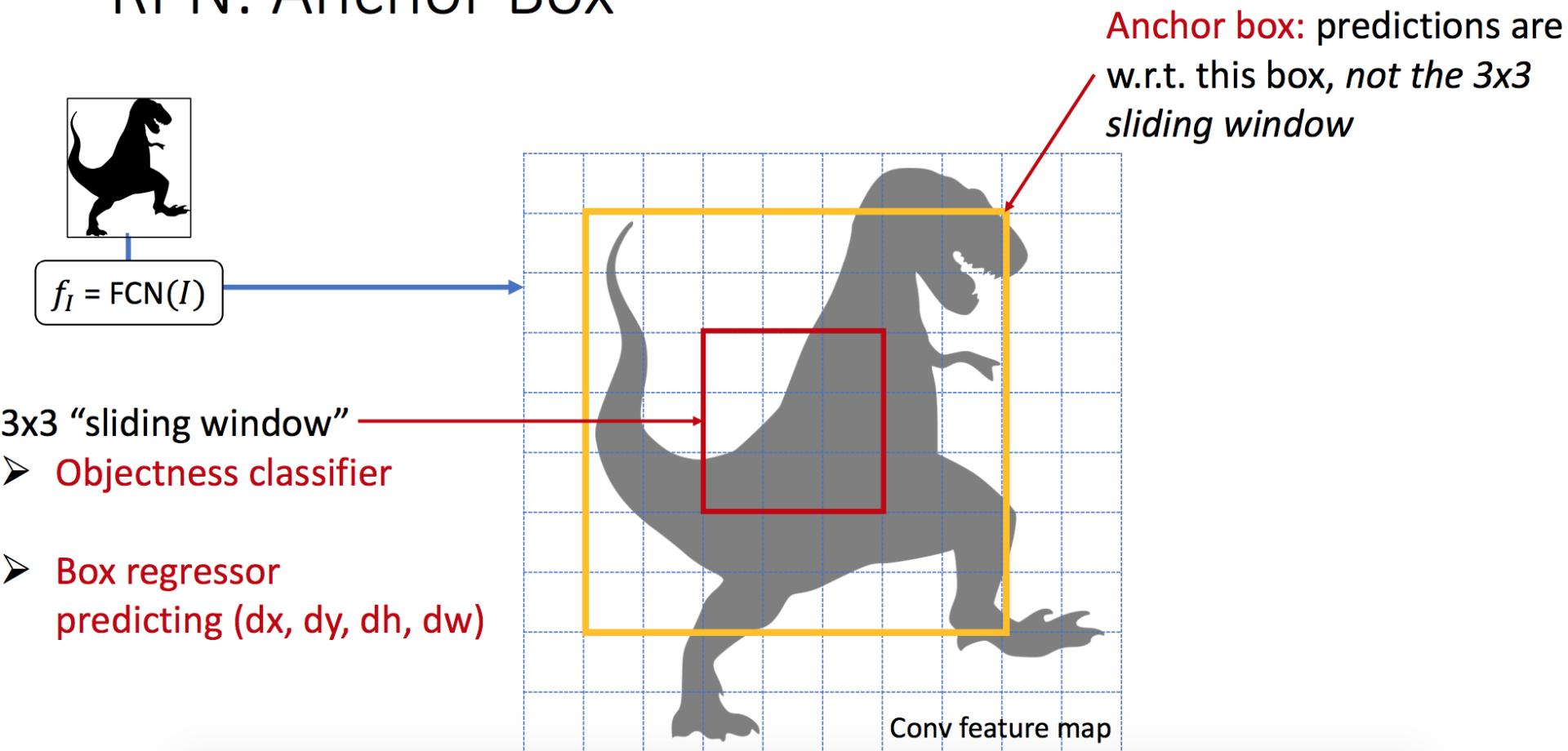
## RPN: Region Proposal Network



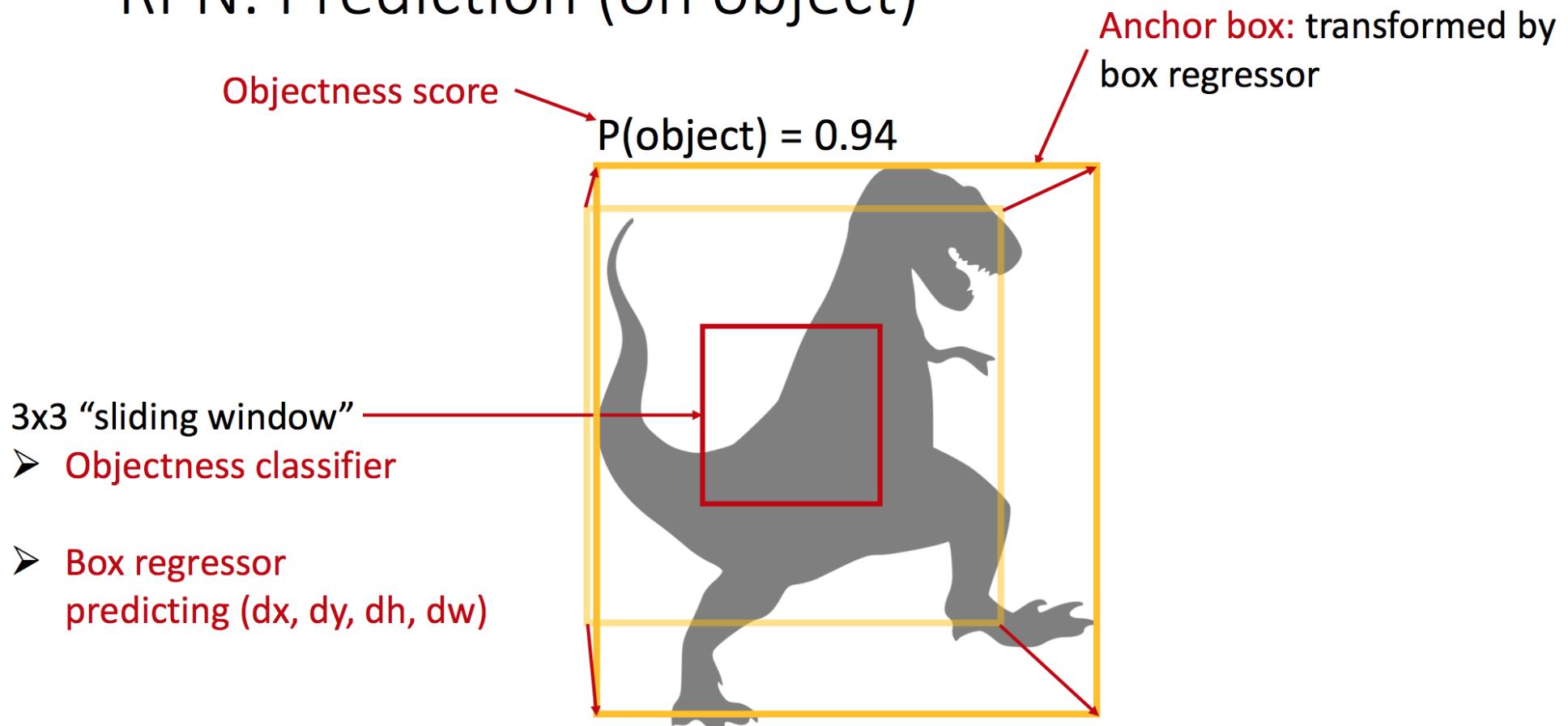
## RPN: Anchor Box



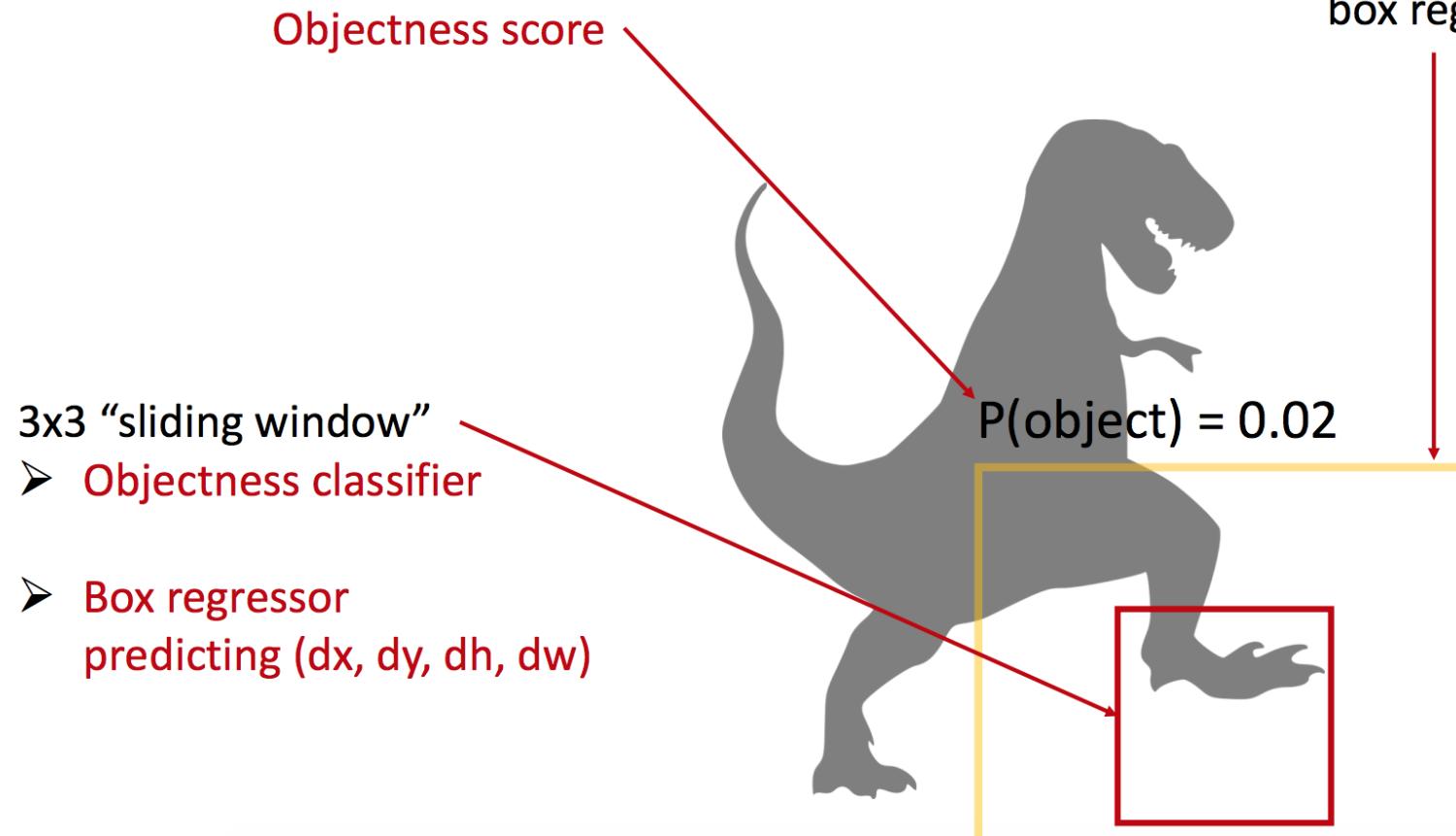
## RPN: Anchor Box



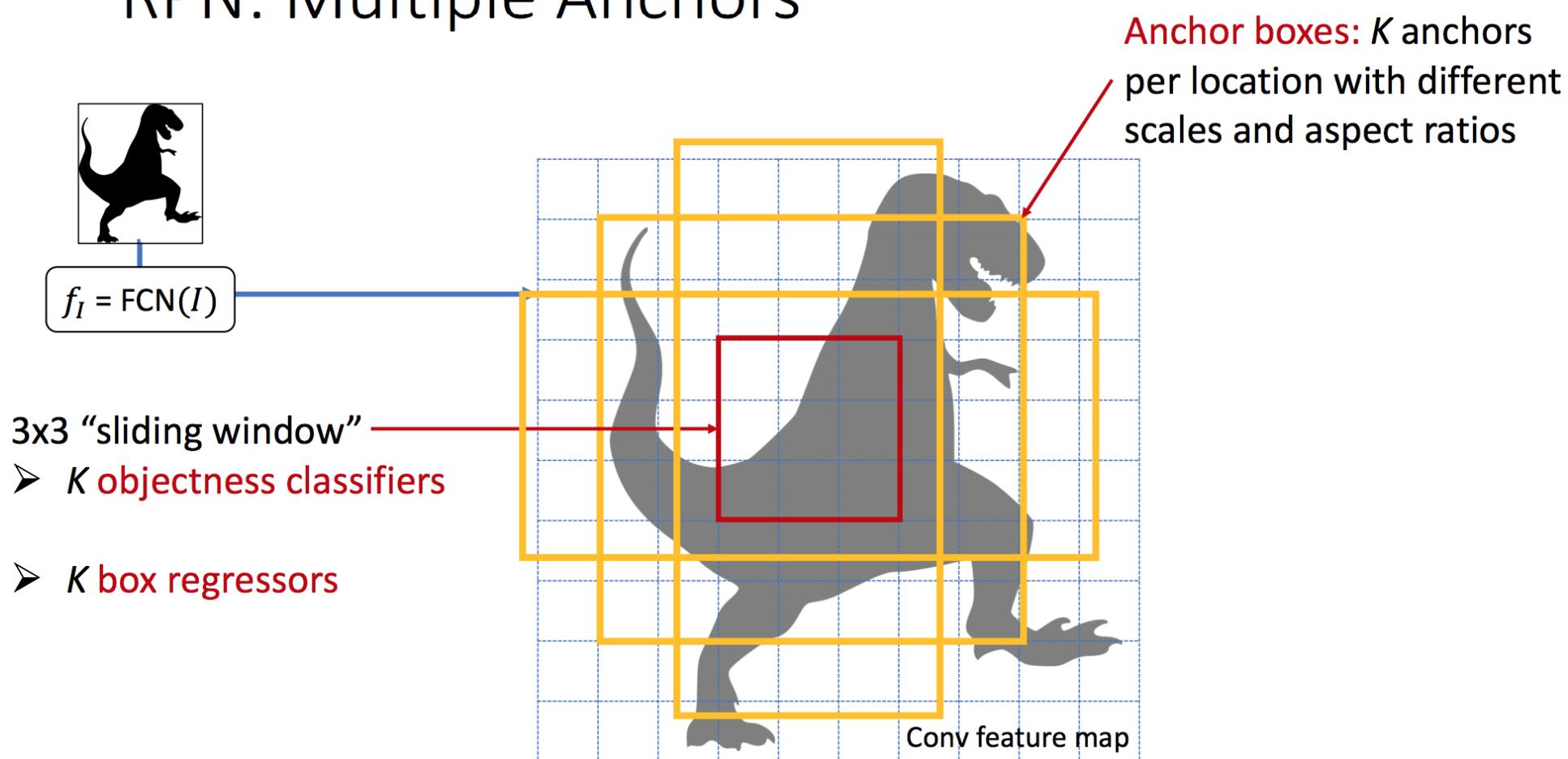
## RPN: Prediction (on object)



## RPN: Prediction (off object)

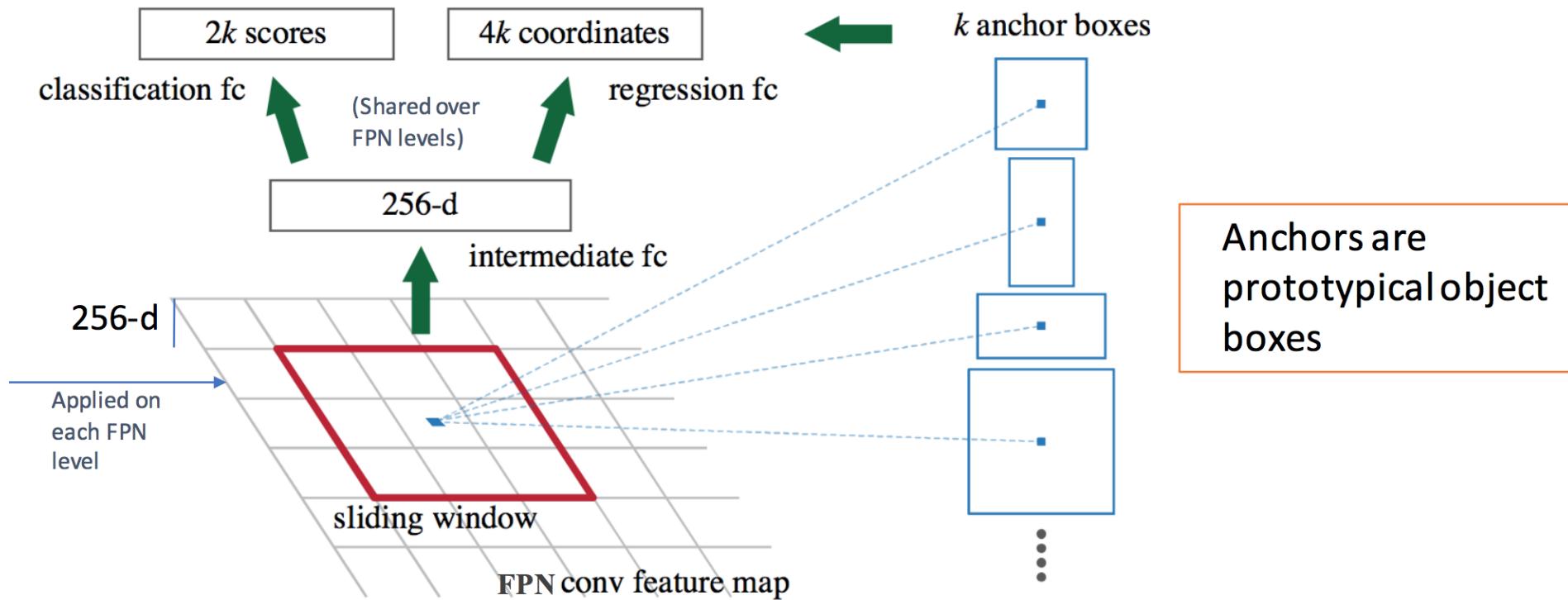


## RPN: Multiple Anchors



- Solution

- Why not generate region proposals using CNN??!



# Faster R-CNN (Ren et al. NIPS 2015)

- What could be the problems

<https://arxiv.org/pdf/1506.01497.pdf>

Dr. Sander Ali Khowaja



# Faster R-CNN (Ren et al. NIPS 2015)

- What could be the problems
  - Two-stage detection pipeline is still too slow to apply on real-time videos

<https://arxiv.org/pdf/1506.01497.pdf>

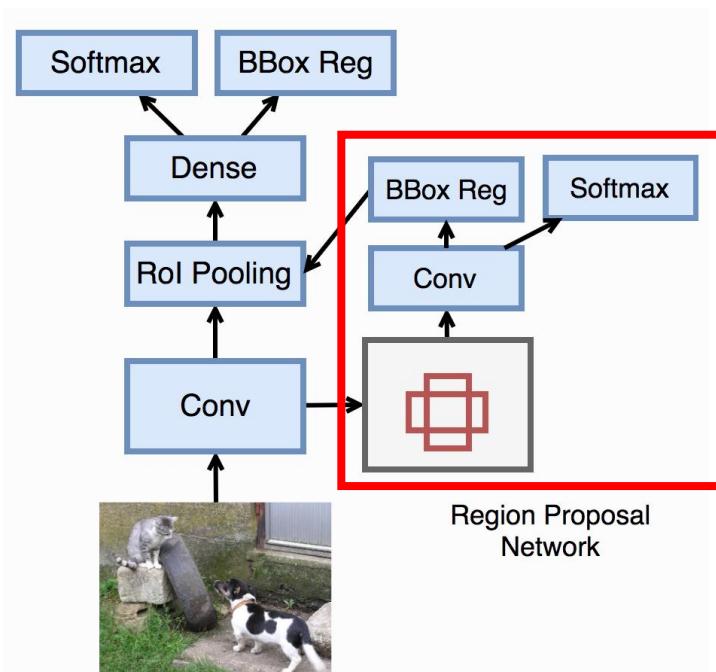
Dr. Sander Ali Khowaja



# One-stage detection

- Solution

- Don't generate object proposals!
- Consider a tiny subset of the output space by design; directly classify this small set of boxes



Dr. Sander Ali Khawaja

Image credit:

[http://zh.gluon.ai/chapter\\_computer-vision/object-detection.html](http://zh.gluon.ai/chapter_computer-vision/object-detection.html)



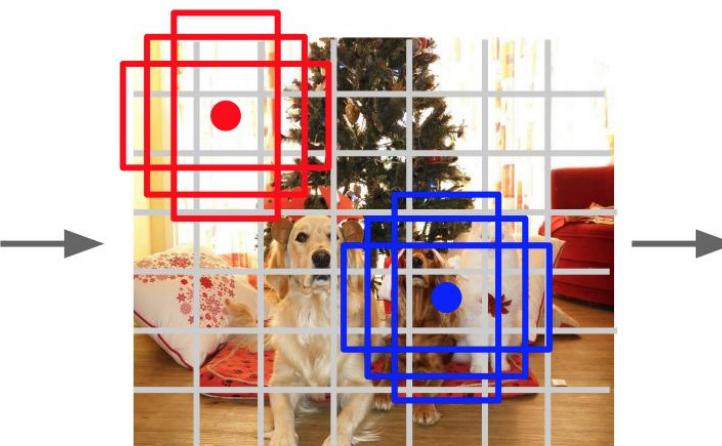
# One-stage detection

- Solution

Go from input image to tensor of scores with one big convolutional network!



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
( $dx$ ,  $dy$ ,  $dh$ ,  $dw$ , confidence)
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

# One-stage detection

- What could be the problems?



# One-stage detection

- What could be the problems?
  - The extreme foreground-background class imbalance -> we have a lot more negative examples.



# One-stage detection

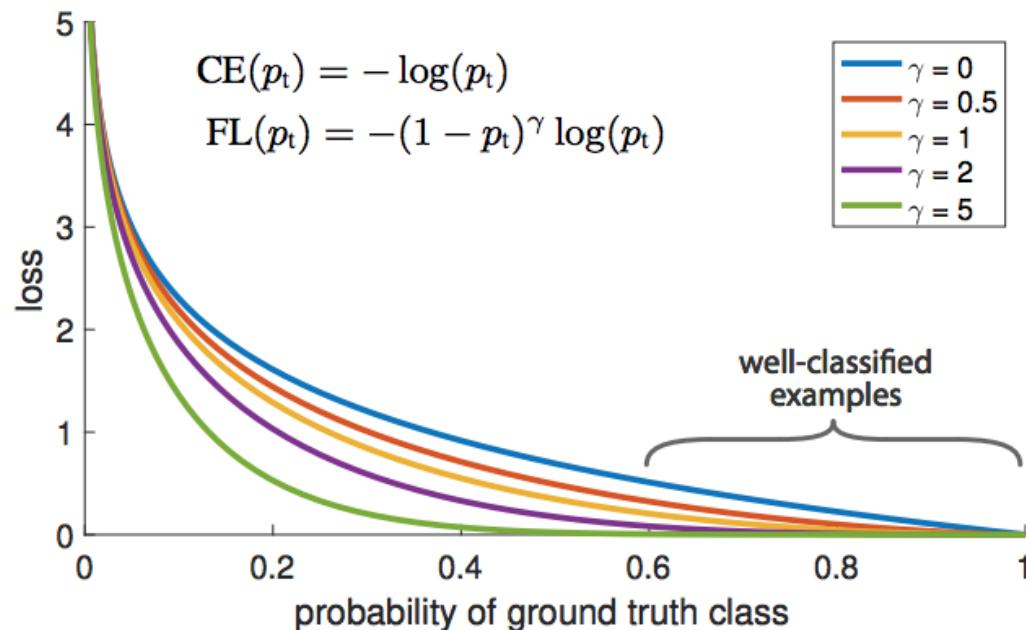
- What could be the problems?
  - The extreme foreground-background class imbalance -> we have a lot more negative examples.
  - Even though they have small loss values, the gradients overwhelm the model



# Focal Loss for Dense Object Detection (Lin et al. ICCV 2017)

- Solution

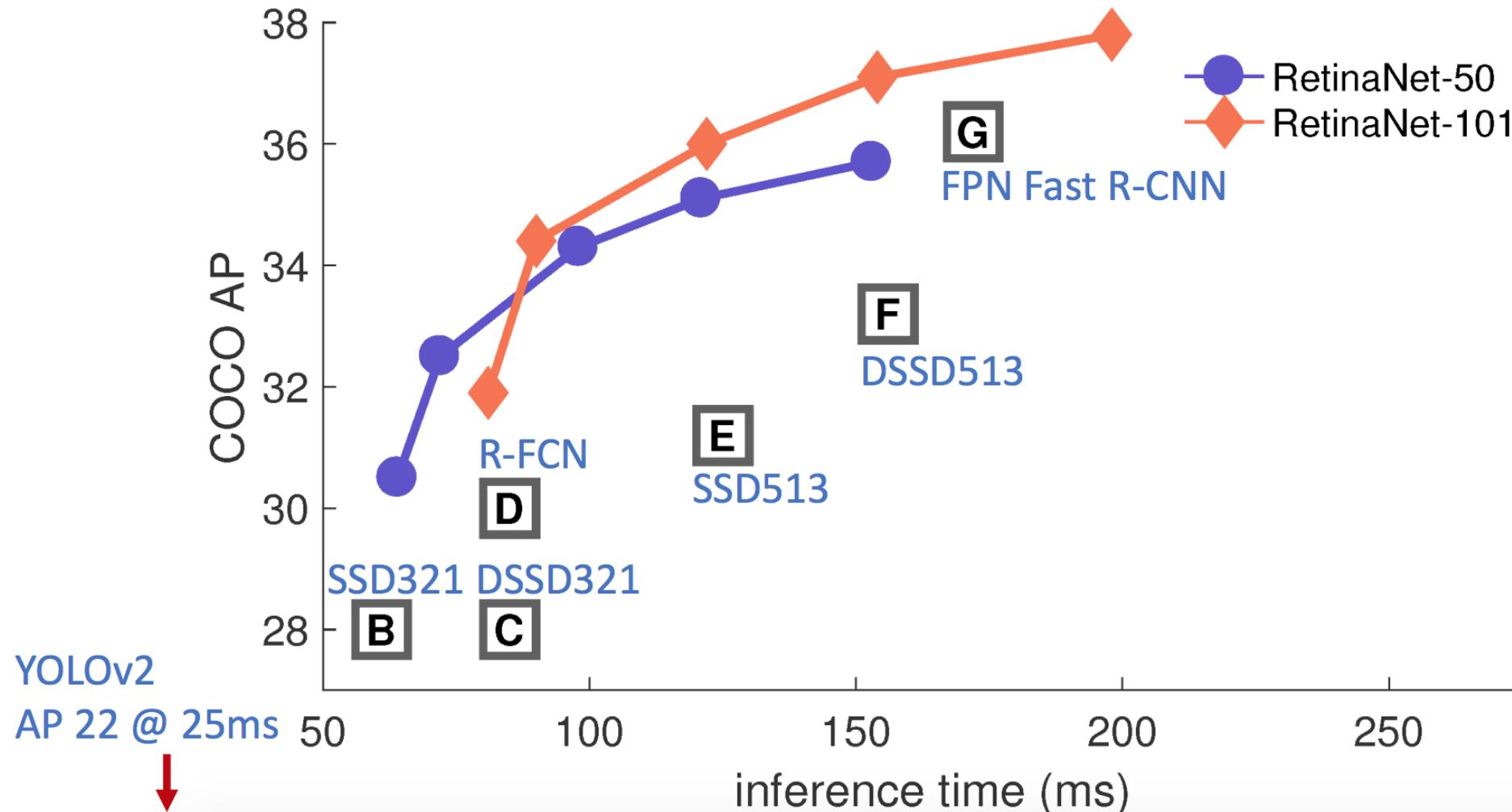
- For easy examples, we down-weight it loss, so that the gradients from these example have smaller impact to the model



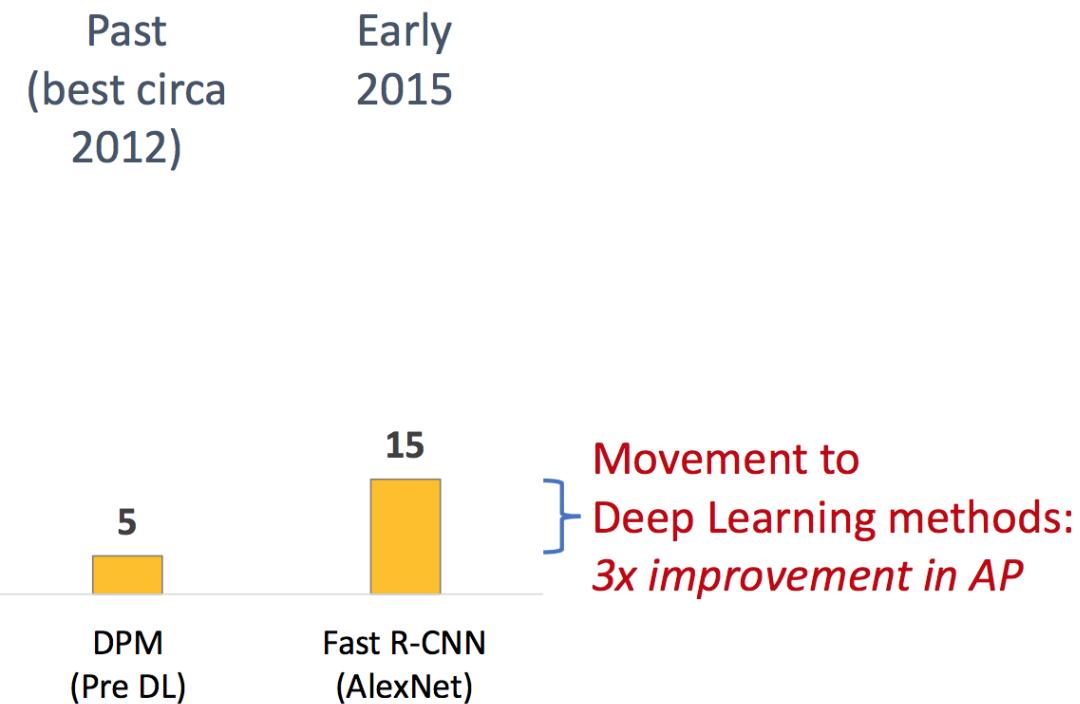
<https://arxiv.org/pdf/1708.02002.pdf>

# One-stage detection

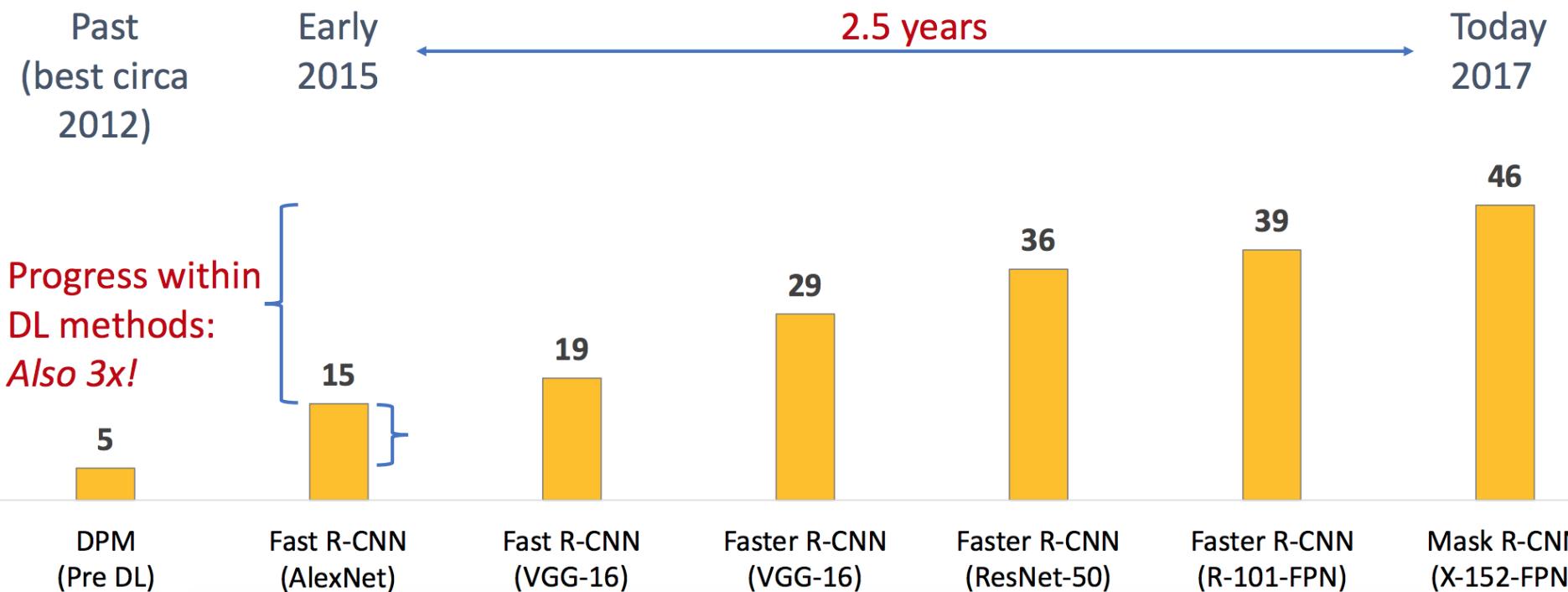
## Speed/Accuracy Tradeoff



# COCO Object Detection Average Precision (%)

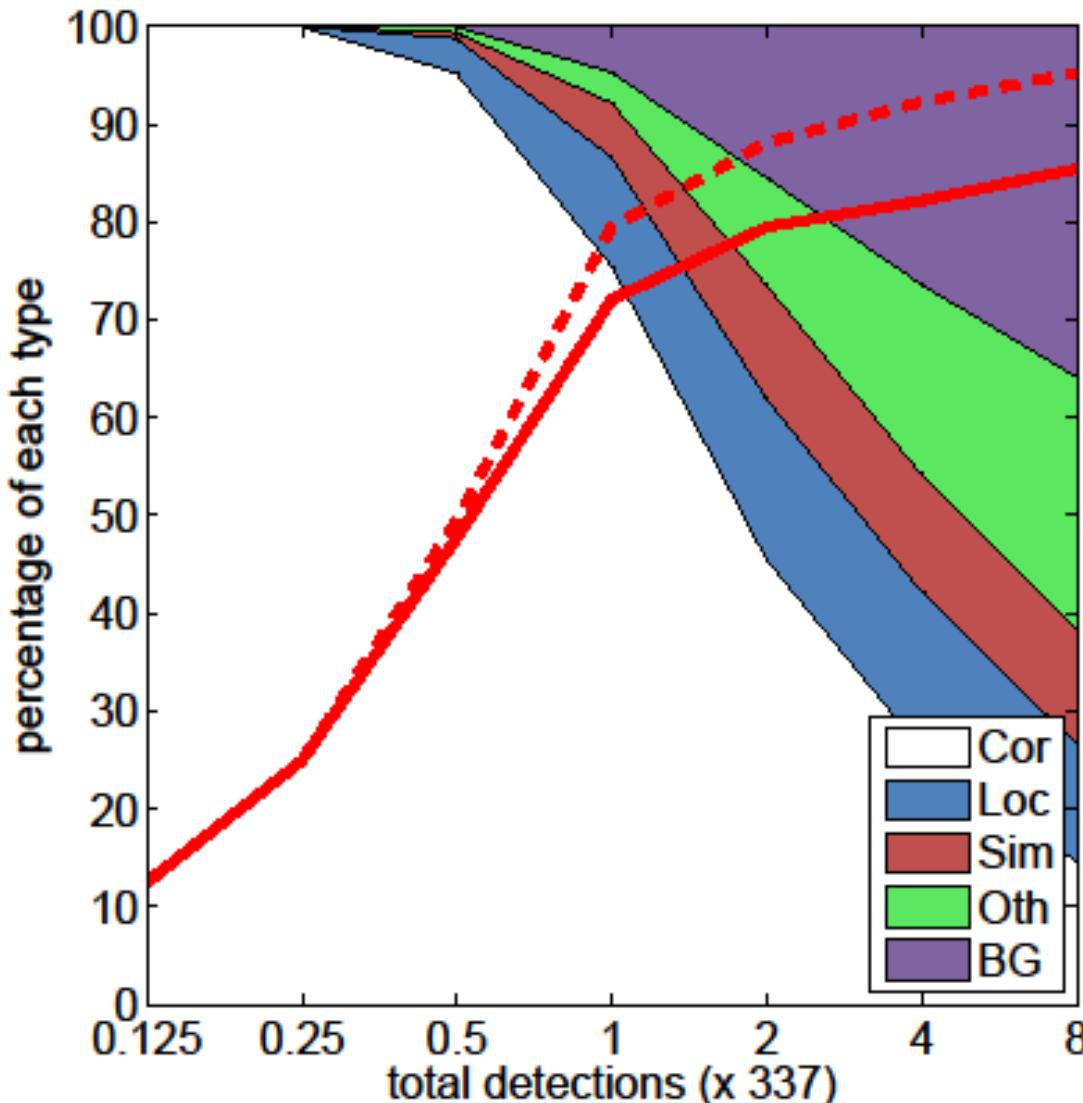


# COCO Object Detection Average Precision (%)



# Mistakes are often reasonable

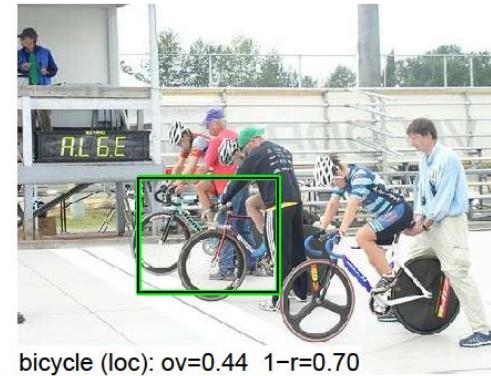
Bicycle: AP = 0.73



R-CNN results

Dr. Sander Ali Khowaja

Confident Mistakes

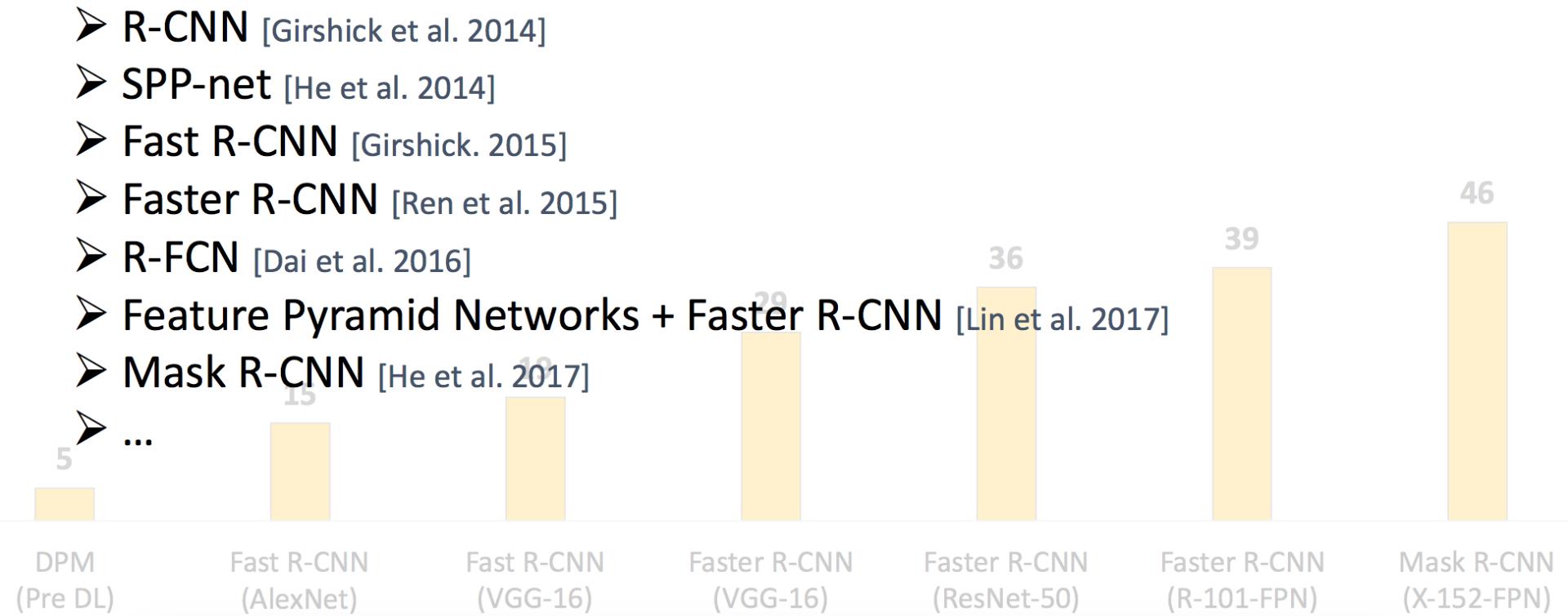


# Influential Works in Detection

- Sung-Poggio (1994, 1998) : ~2100 citations
  - Basic idea of statistical template detection (I think), bootstrapping to get “face-like” negative examples, multiple whole-face prototypes (in 1994)
- Rowley-Baluja-Kanade (1996-1998) : ~4200
  - “Parts” at fixed position, non-maxima suppression, simple cascade, rotation, pretty good accuracy, fast
- Schneiderman-Kanade (1998-2000,2004) : ~2250
  - Careful feature/classifier engineering, excellent results, cascade
- Viola-Jones (2001, 2004) : ~20,000
  - Haar-like features, Adaboost as feature selection, hyper-cascade, very fast, easy to implement
- Dalal-Triggs (2005) : ~11000
  - Careful feature engineering, excellent results, HOG feature, online code
- Felzenszwalb-Huttenlocher (2000): ~1600
  - Efficient way to solve part-based detectors
- Felzenszwalb-McAllester-Ramanan (2008,2010)? ~4000
  - Excellent template/parts-based blend



# Influential Works in Detection



# Fails in commercial face detection

- Things iPhoto thinks are faces

## Who's in These Photos?

The photos you uploaded were grouped automatically so you can quickly label and notify friends in these pictures. (Friends can always untag themselves.)



Who is this?

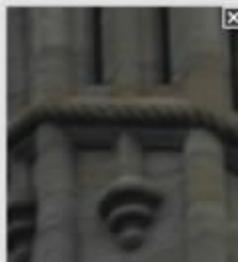
Who is this?



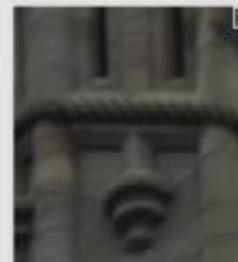
Unnamed people

4 Group(s), 67 Face(s)

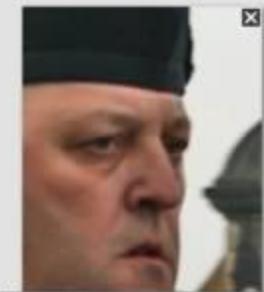
Select someone you know and add a name, or click the "x" to ignore that person.



Add a name



Add a name



Add a name



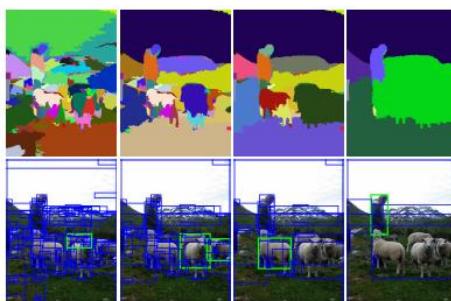
[http://www.oddee.com/item\\_98248.aspx](http://www.oddee.com/item_98248.aspx) Mr. Sander Ali Khawaja



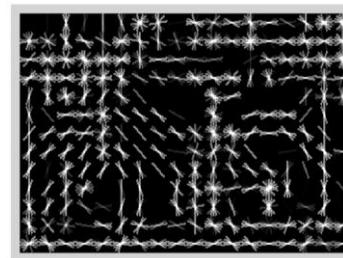
# Summary: statistical templates



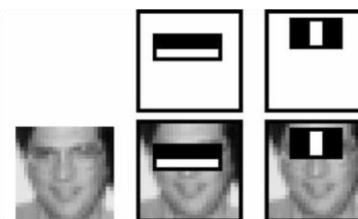
Sliding window: scan image pyramid



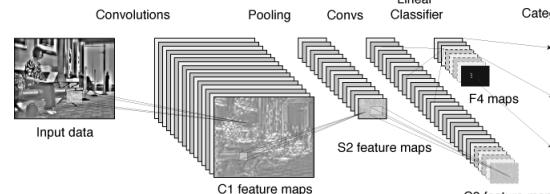
Region proposals: edge/region-based, resize to fixed window



HOG



Fast randomized features



CNN features  
Dr. Sander Ali Khowaja

SVM

Boosted stumps

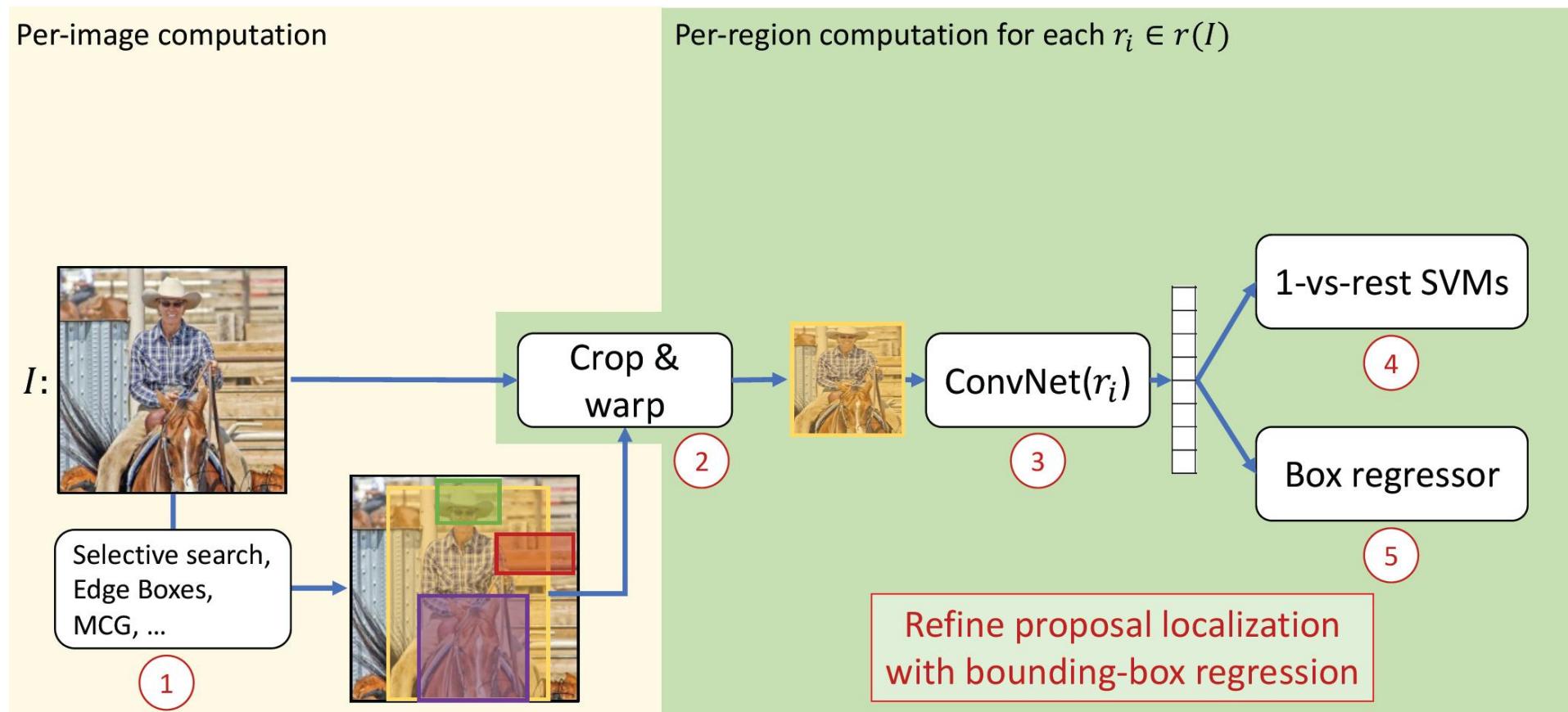
Neural network

Non-max suppression

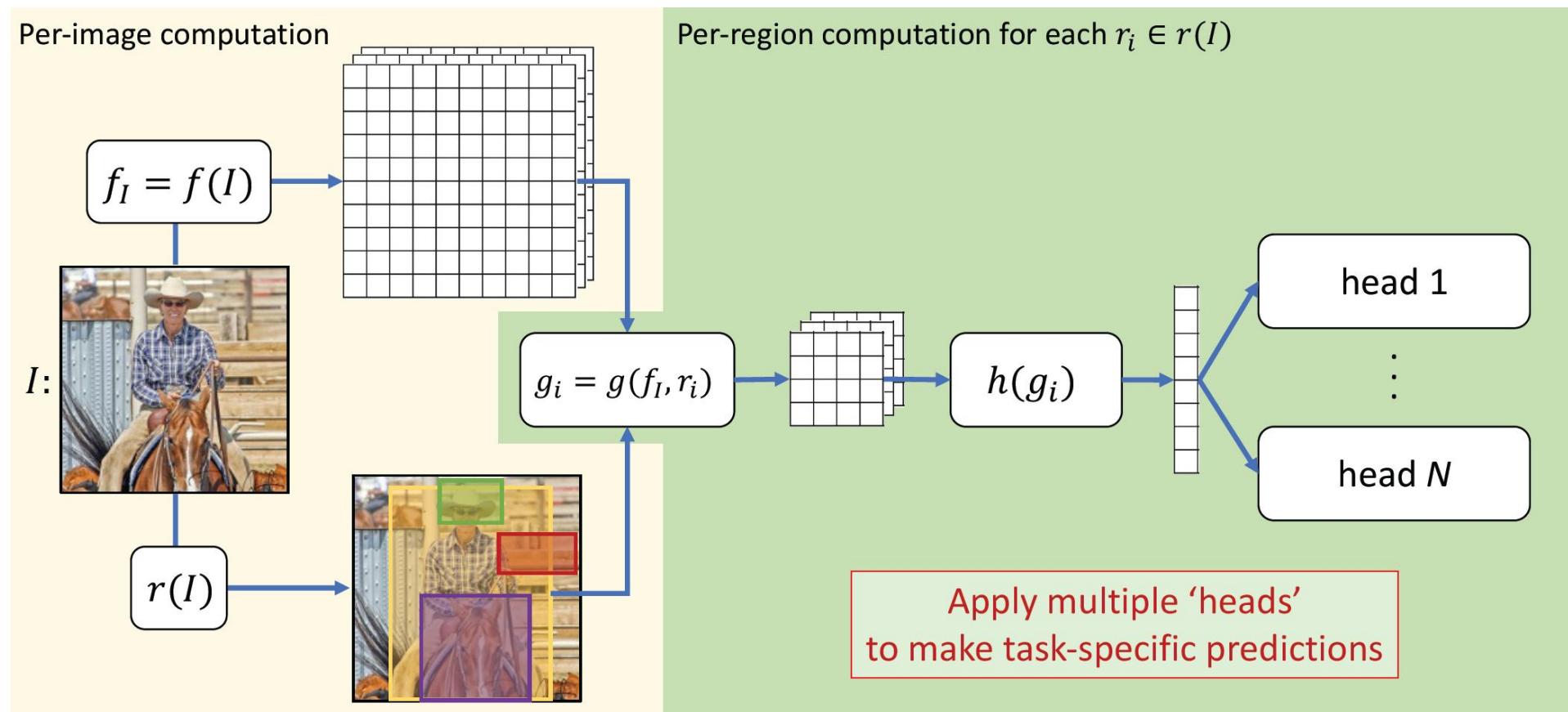
Segment or refine localization



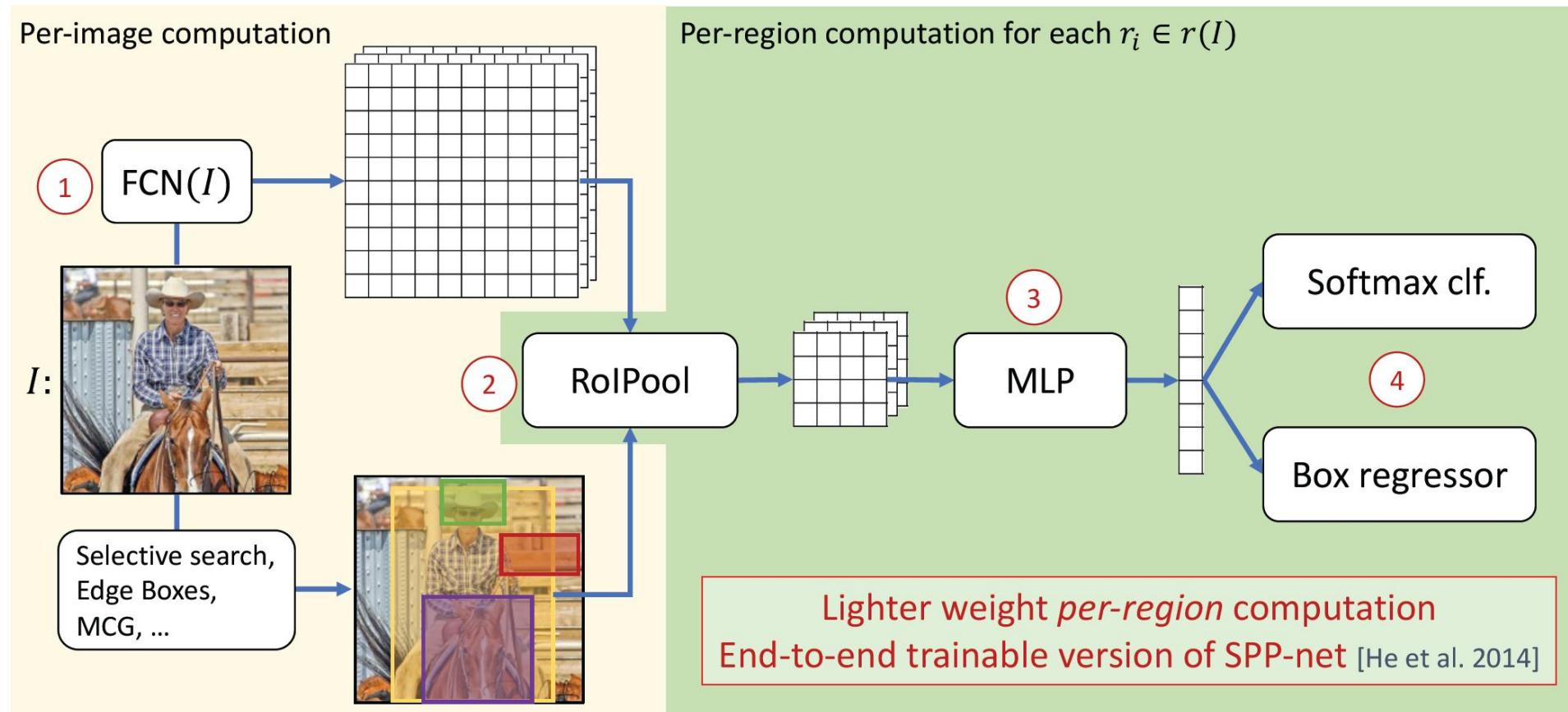
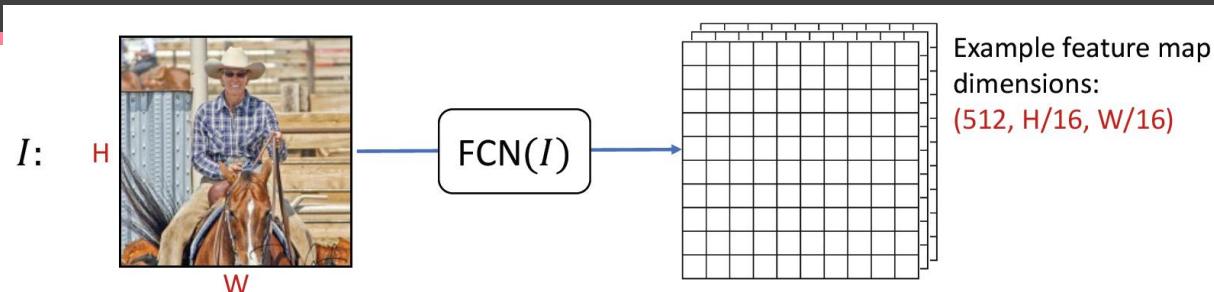
# Region-based CNN



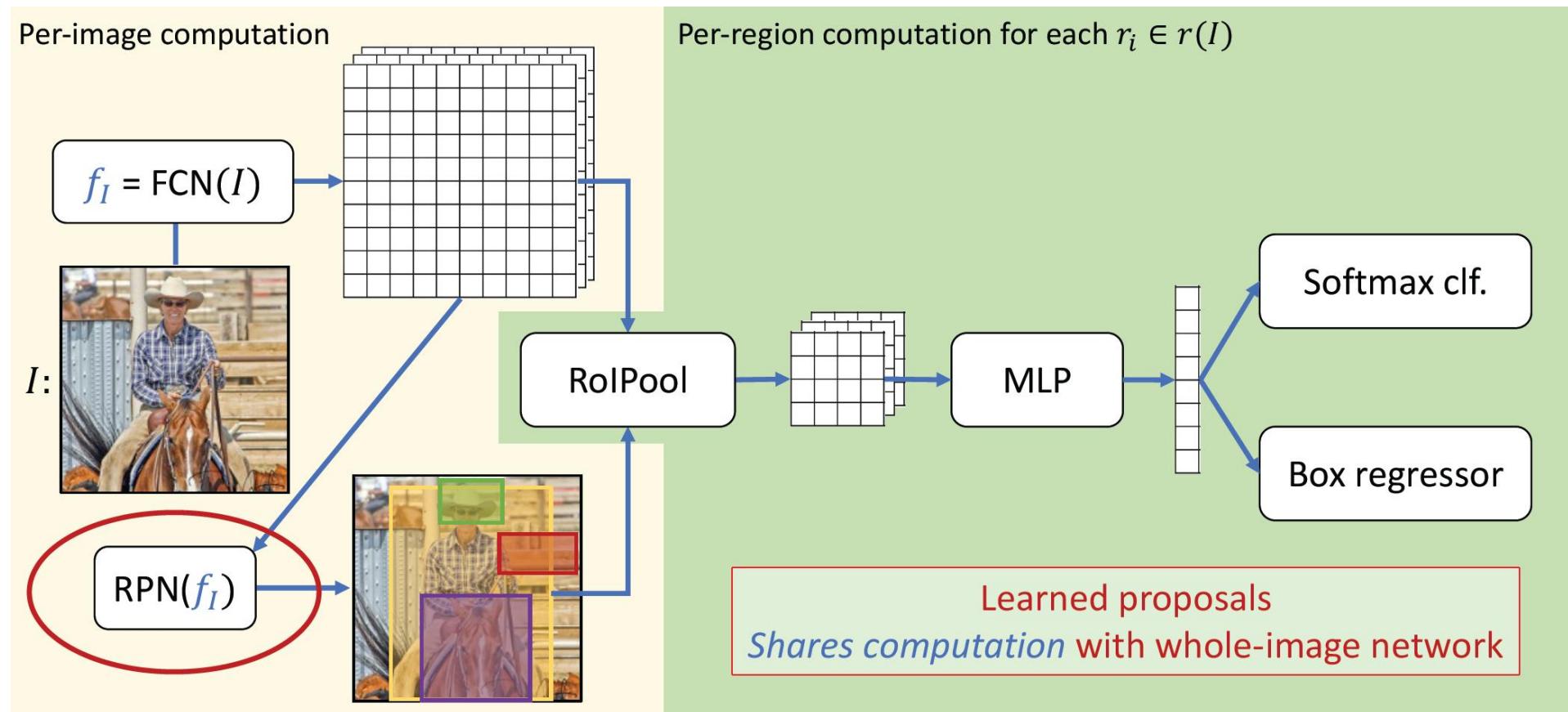
# Generalized R-CNN



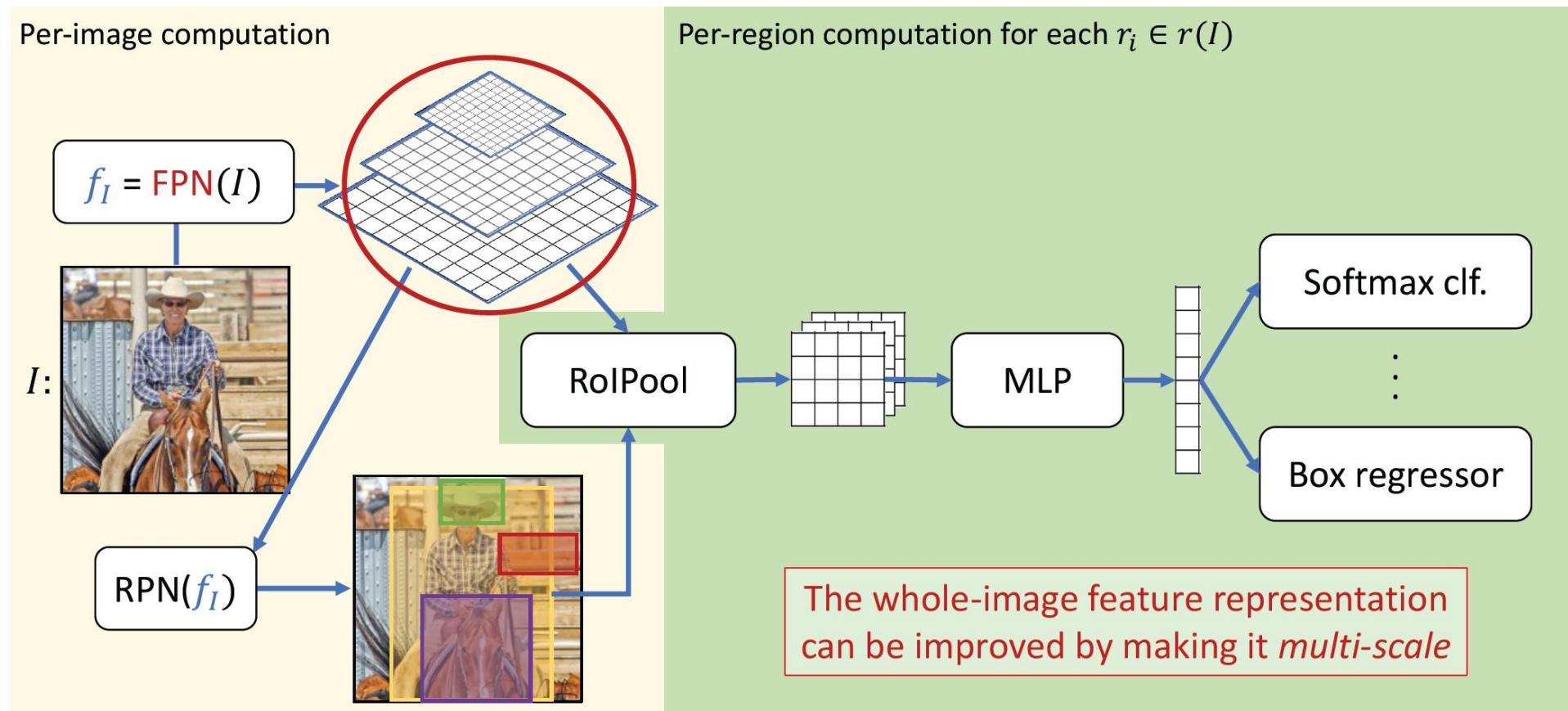
# Fast R-CNN



# Faster RCNN



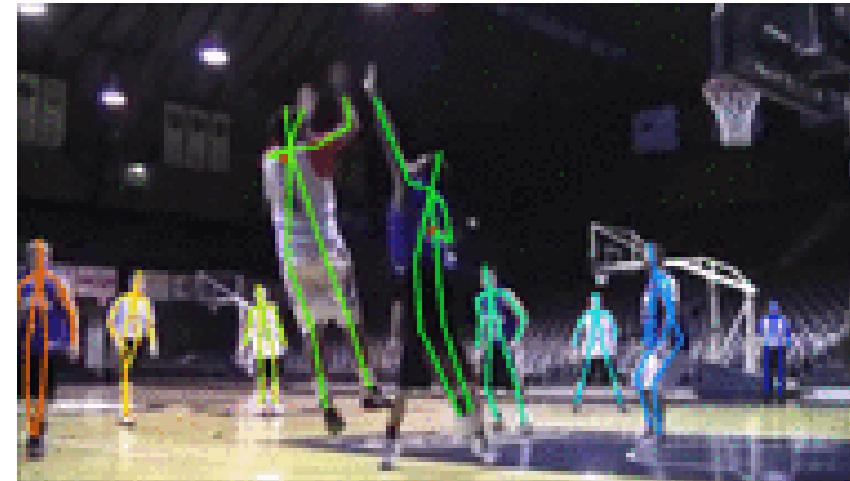
# Faster RCNN + Feature Pyramid Network



# Part/keypoint Prediction



# PoseTrack



# Semantic Segmentation



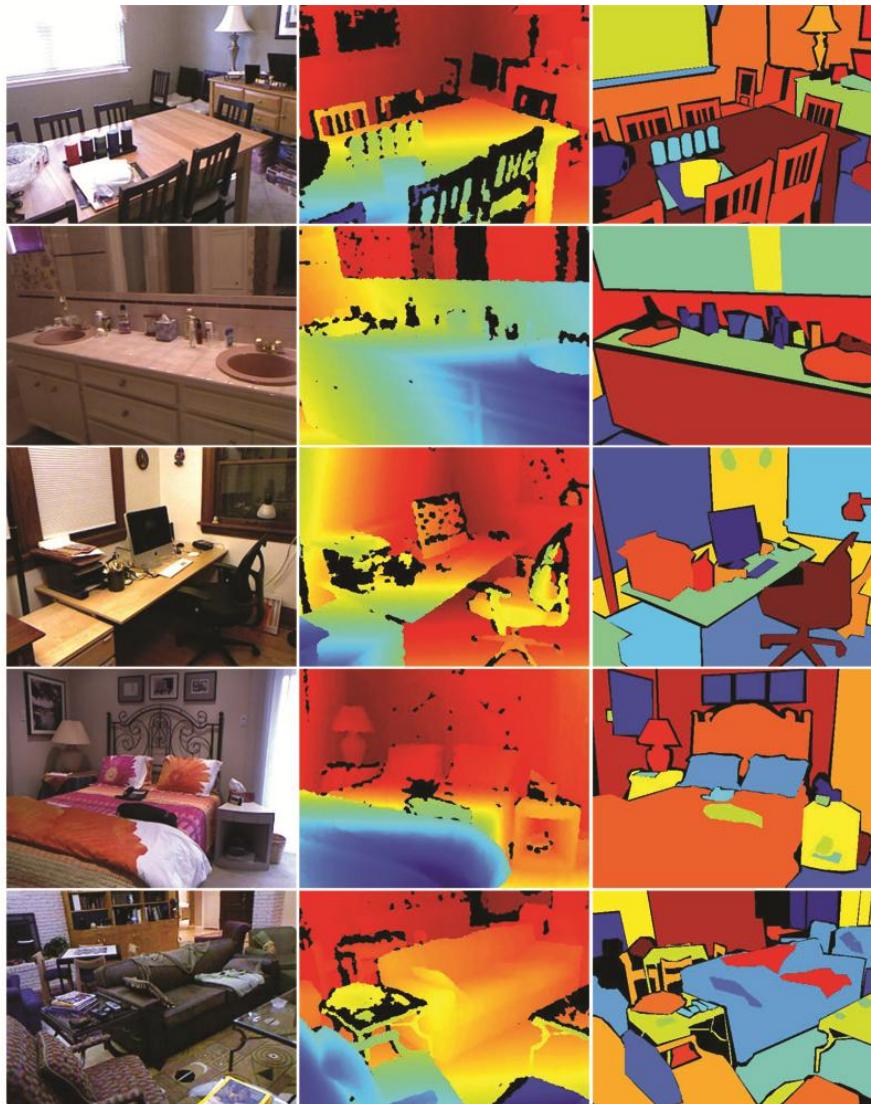
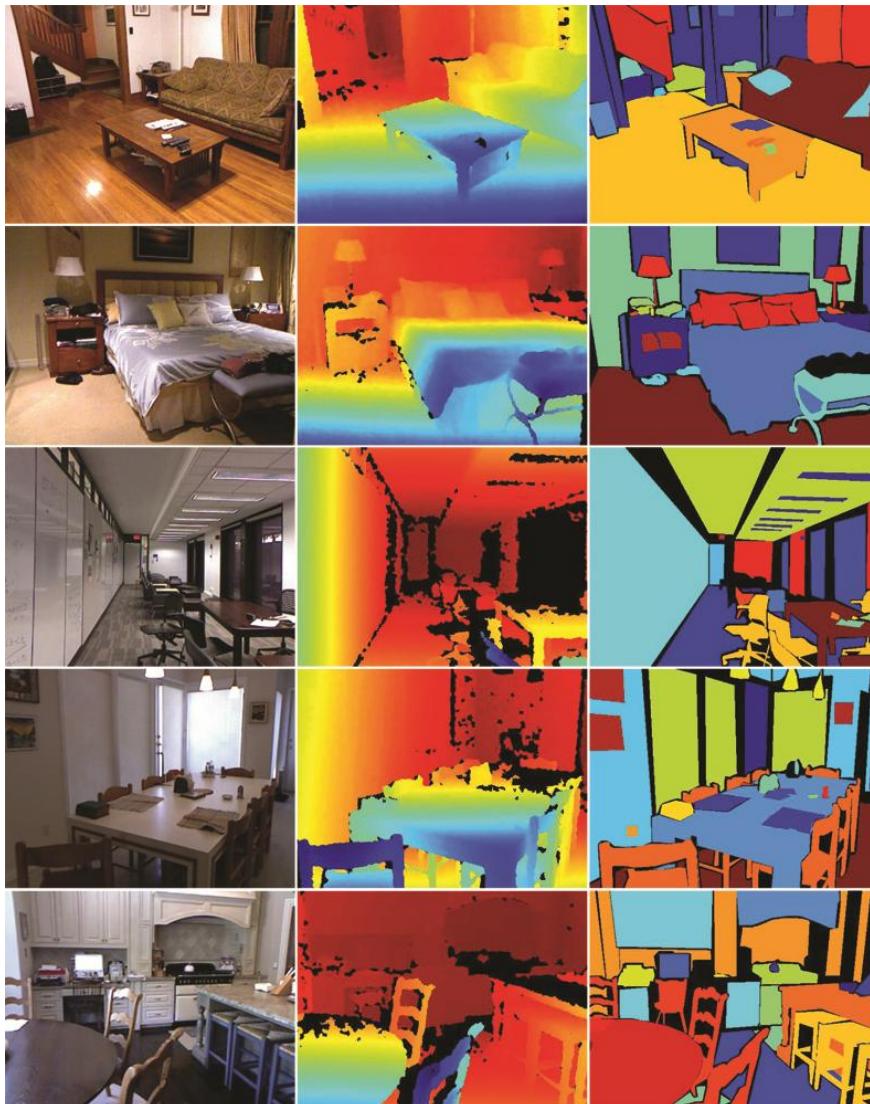
<http://mscoco.org/dataset/#detections-challenge2016>



<https://www.cityscapes-dataset.com/examples/#fine-annotations>



# Semantic Segmentation

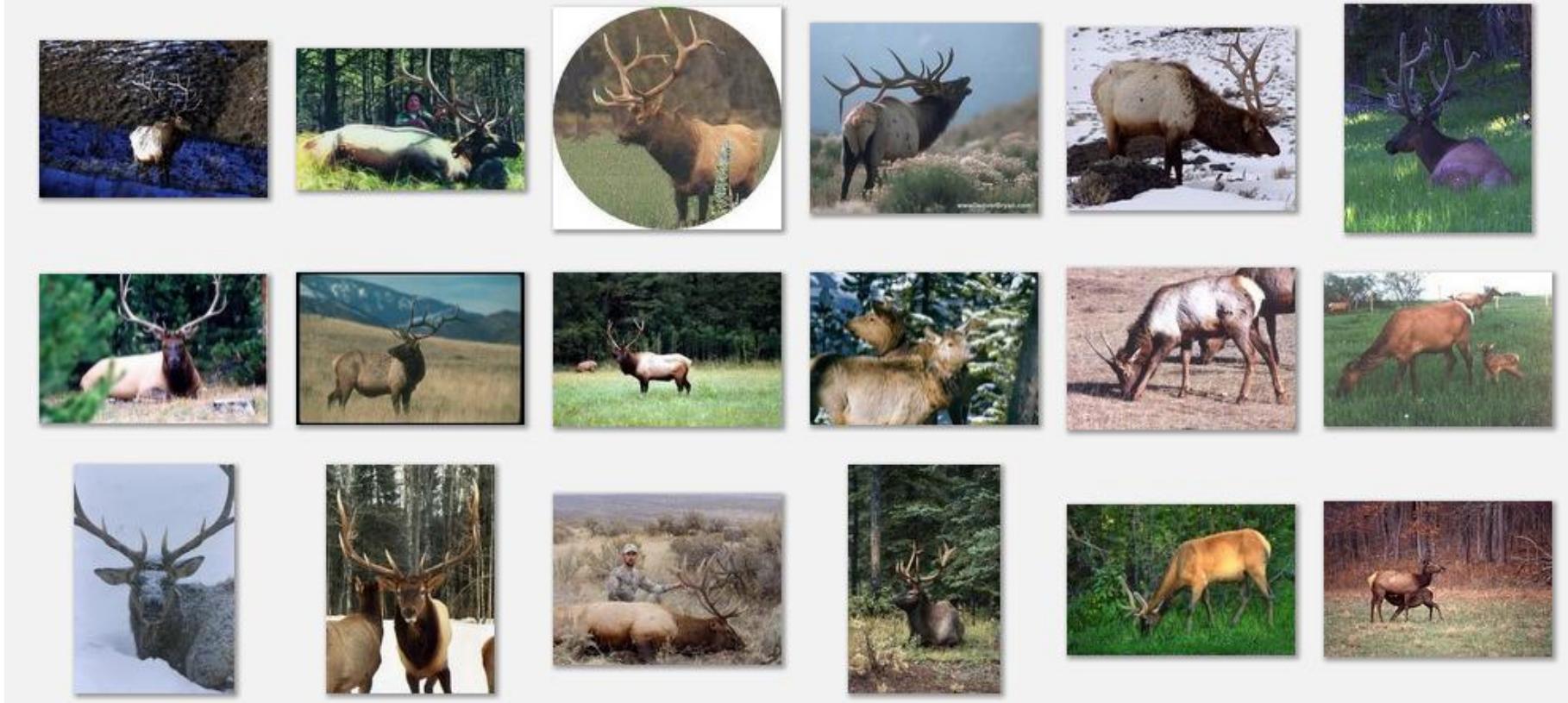


[http://cs.nyu.edu/~silberman/datasets/nyu\\_depth\\_v2.html](http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html)

Dr. Sander A. H. Mowafa



# Deformable objects



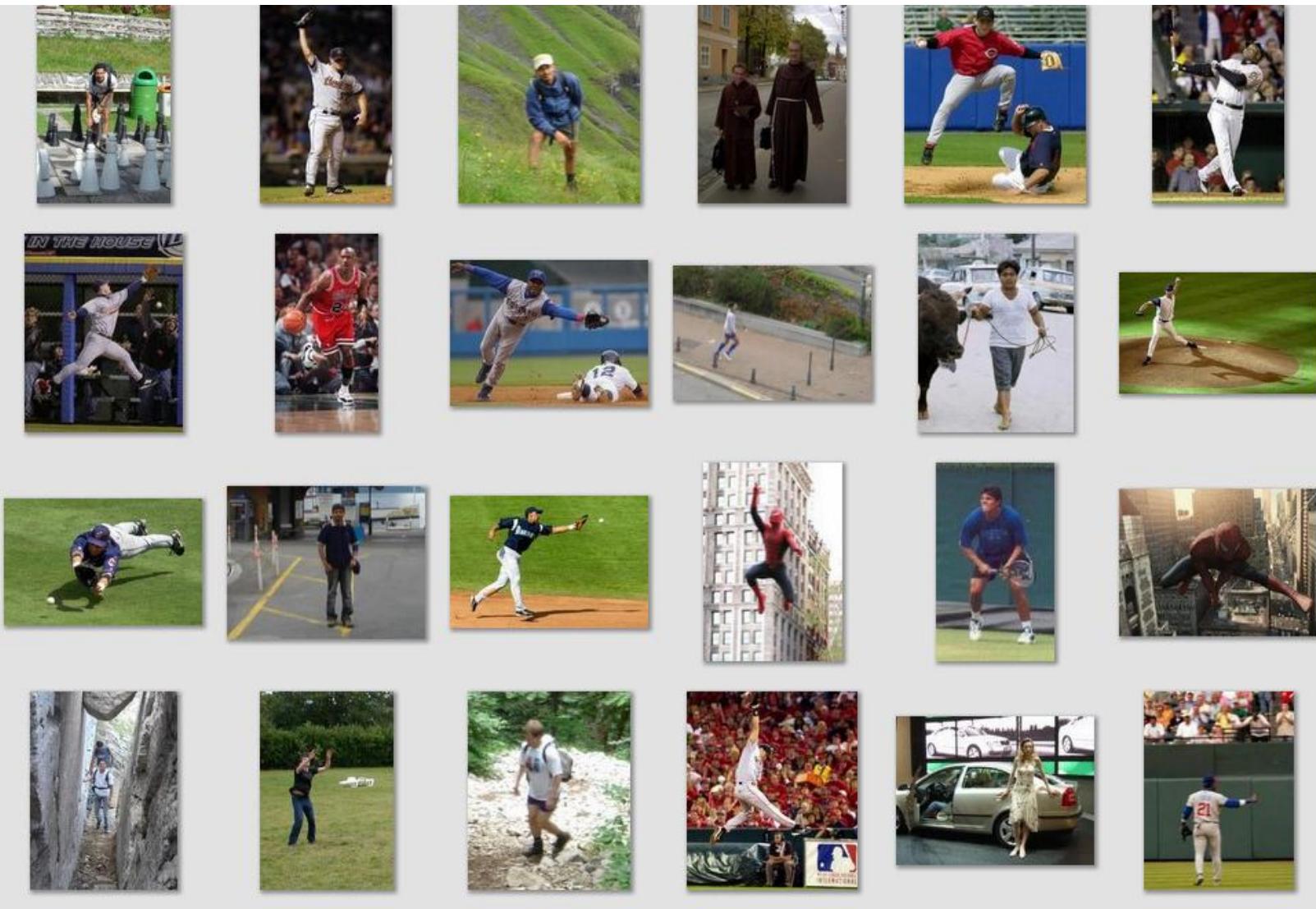
Images from Caltech-256

# Graphical models vs. sequential/parallel prediction

- Advantages of BP/graphcut/etc
  - Elegant
  - Relations are explicitly modeled
  - Exact inference in some cases
- Advantages of sequential/parallel prediction
  - Simple procedures for training and inference
  - Learns how much to rely on each prediction
  - Can model very complex relations



# Deformable objects



Images from D. Ramanan's dataset  
Dr. Sander Ali Khowaja

Slide Credit: Duan Tran



# Compositional objects



Dr. Sander Ali Khowaja



# What if you want to label every pixel?

“Stuff” can be hard to capture with bounding boxes

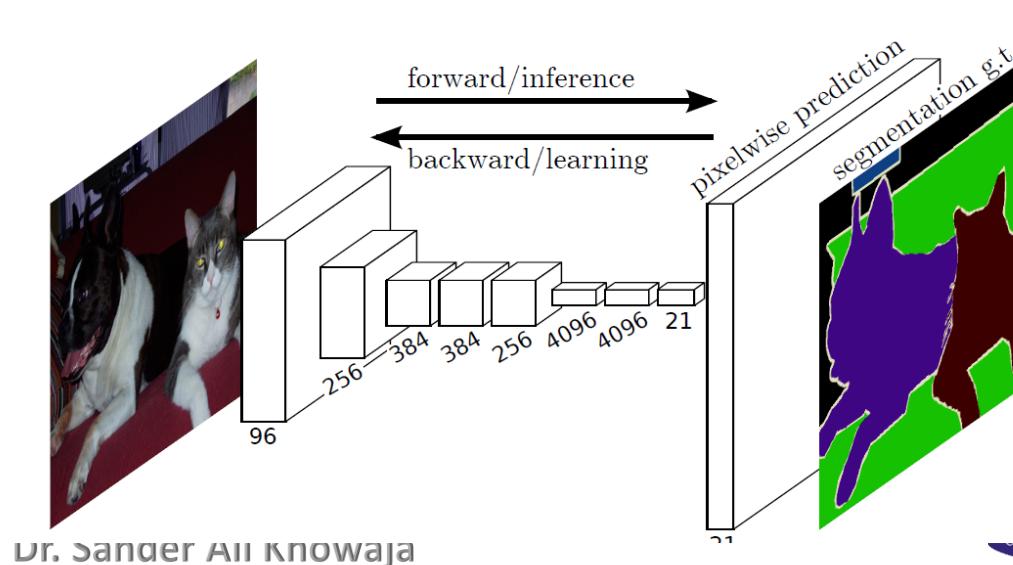
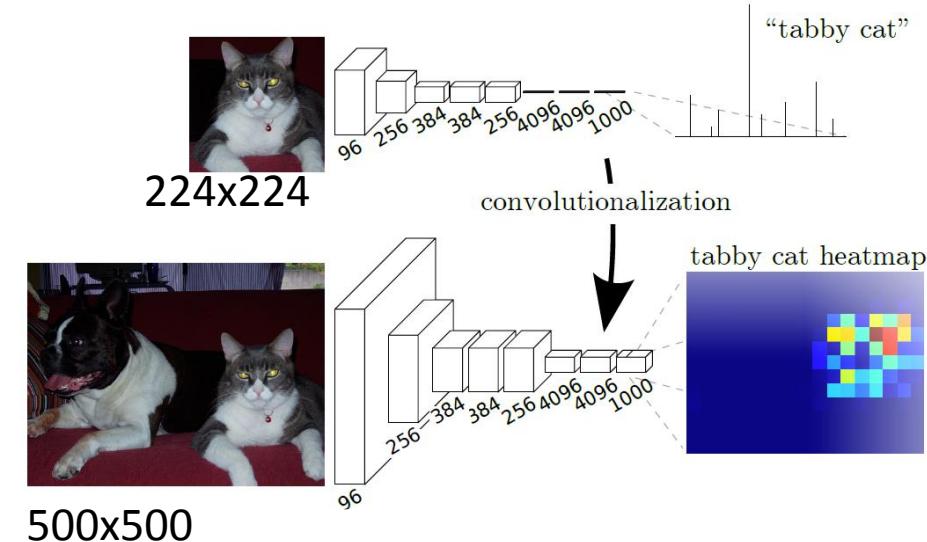


<https://www.cityscapes-dataset.com/examples/#fine-annotations>

Dr. Sander Arik Khowaja

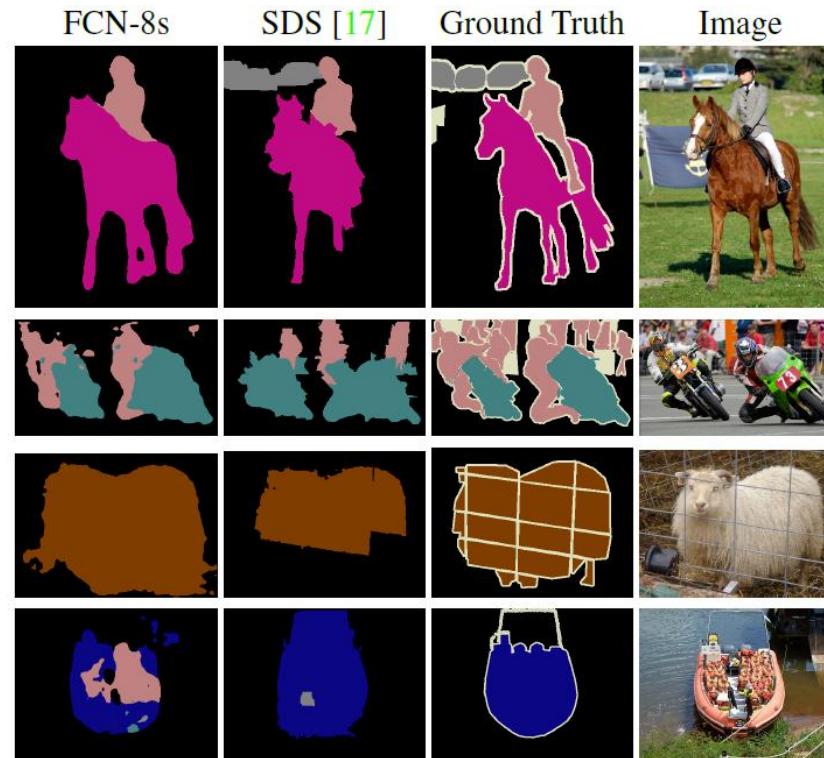


- Use network trained for classification as pre-trained network for pixel labeling
- Convert fully connected layers into convolutions
- Add features from earlier conv layers to improve resolution
- Fine-tune for pixel labeling task



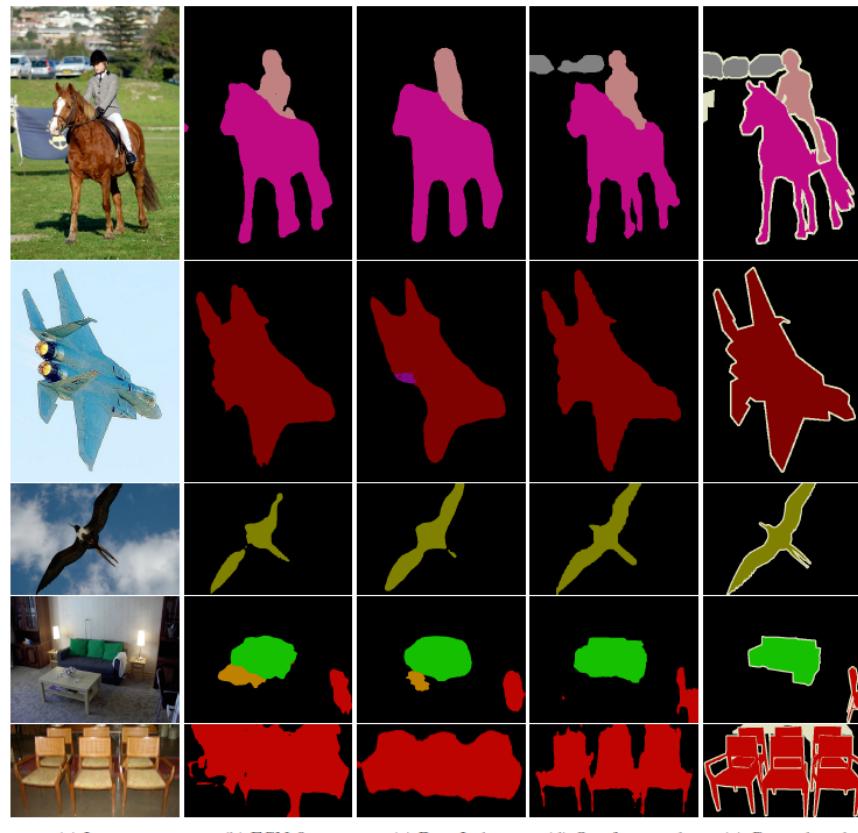
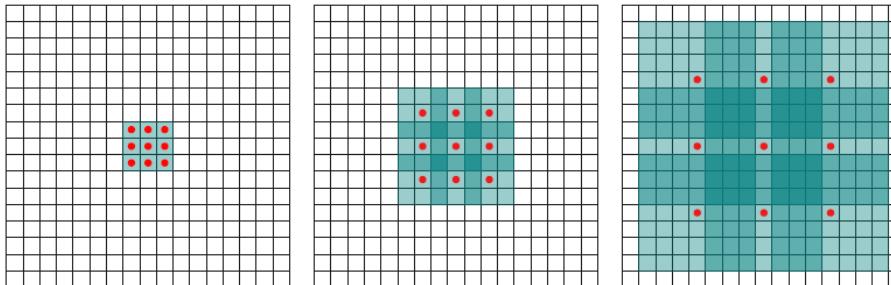
# “Fully convolutional” results

- Takes advantage of pre-training from classification
- Applied to objects and scenes (NYUd v2)
- But feature pooling reduces spatial sensitivity and resolution



# Dilated Convolutions – Yu Kolton 2016

- Replacing last two pooling layers with “dilated convolution” that filters a sparse 3x3 grid of pixels
- Enables large receptive field with few parameters
- Improves resolution



# Dilated Convolutions results



(a) Image

(b) Front end

(c) + Context

(d) + CRF-RNN

(e) Ground truth

	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean IoU
DeepLab++	89.1	38.3	88.1	63.3	69.7	87.1	<b>83.1</b>	85	29.3	76.5	56.5	79.8	77.9	85.8	82.4	57.4	84.3	54.9	80.5	64.1	72.7
DeepLab-MSc++	89.2	46.7	88.5	63.5	68.4	87.0	81.2	86.3	32.6	80.7	62.4	81.0	81.3	84.3	82.1	56.2	84.6	58.3	76.2	67.2	73.9
CRF-RNN	90.4	<b>55.3</b>	88.7	<b>68.4</b>	69.8	88.3	82.4	85.1	32.6	78.5	<b>64.4</b>	79.6	81.9	<b>86.4</b>	81.8	<b>58.6</b>	82.4	53.5	77.4	<b>70.1</b>	74.7
Front end	86.6	37.3	84.9	62.4	67.3	86.2	81.2	82.1	32.6	77.4	58.3	75.9	81	83.6	82.3	54.2	81.5	50.1	77.5	63	71.3
Context	89.1	39.1	86.8	62.6	68.9	88.2	82.6	87.7	33.8	81.2	59.2	81.8	87.2	83.3	83.6	53.6	84.9	53.7	80.5	62.9	73.5
Context + CRF	91.3	39.9	<b>88.9</b>	64.3	69.8	88.9	82.6	89.7	34.7	82.7	59.5	83	88.4	84.2	85	55.3	86.7	54.4	<b>81.9</b>	63.6	74.7
Context + CRF-RNN	<b>91.7</b>	39.6	87.8	63.1	<b>71.8</b>	<b>89.7</b>	82.9	<b>89.8</b>	<b>37.2</b>	<b>84</b>	63	<b>83.3</b>	<b>89</b>	83.8	<b>85.1</b>	56.8	<b>87.6</b>	<b>56</b>	80.2	64.7	<b>75.3</b>

