

Object Recognition

Seminar

Rita Osadchy

So what does object recognition involve?



Verification: is that a bus?



Detection: locate the cars in the image



Verification: is that a picture of Mao?



Object categorization



sky

building

flag

face

banner

wall

street lamp

bus

bus

cars

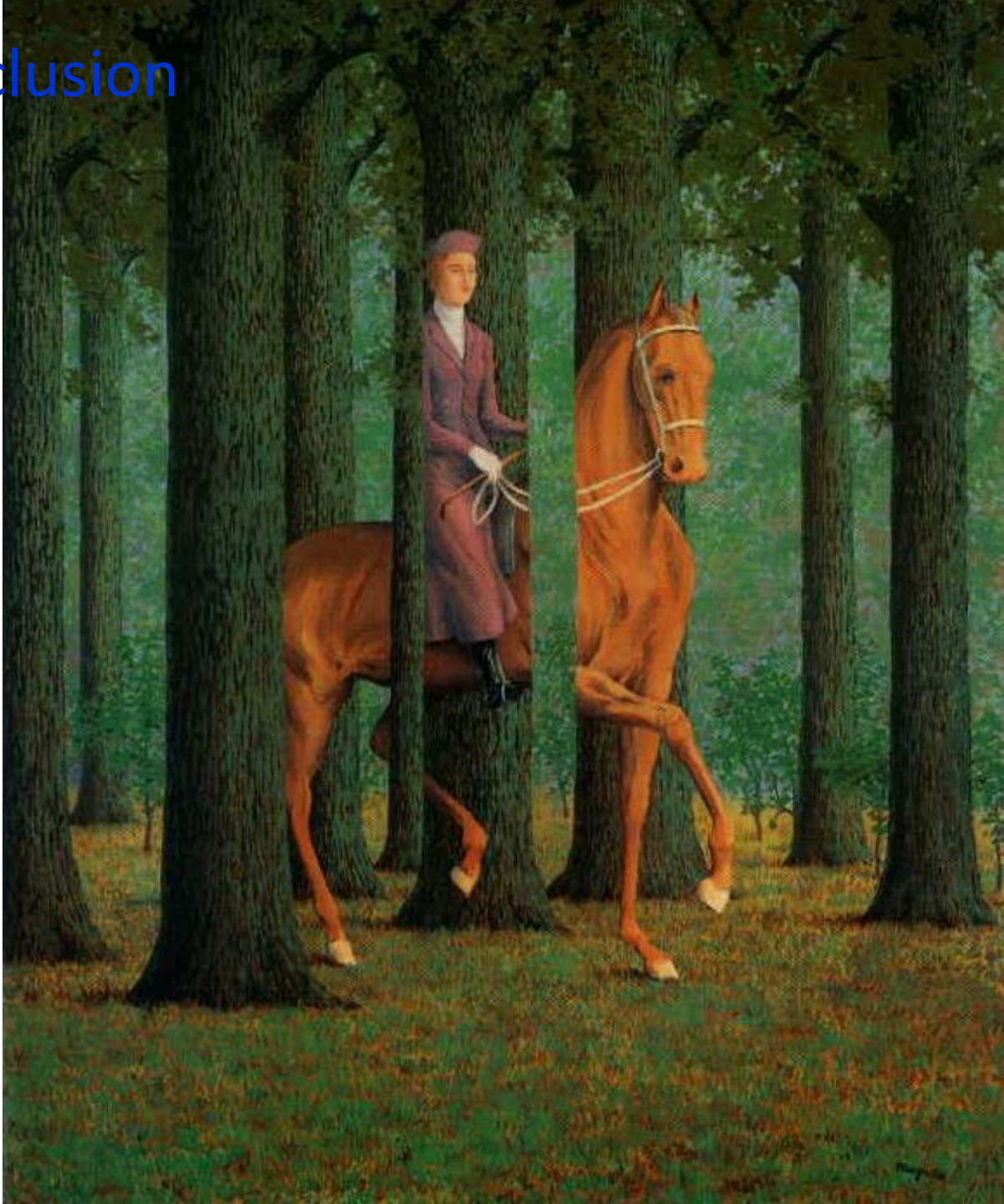
Challenges 1: view point variation



Challenges 2: illumination



Challenges 3: occlusion

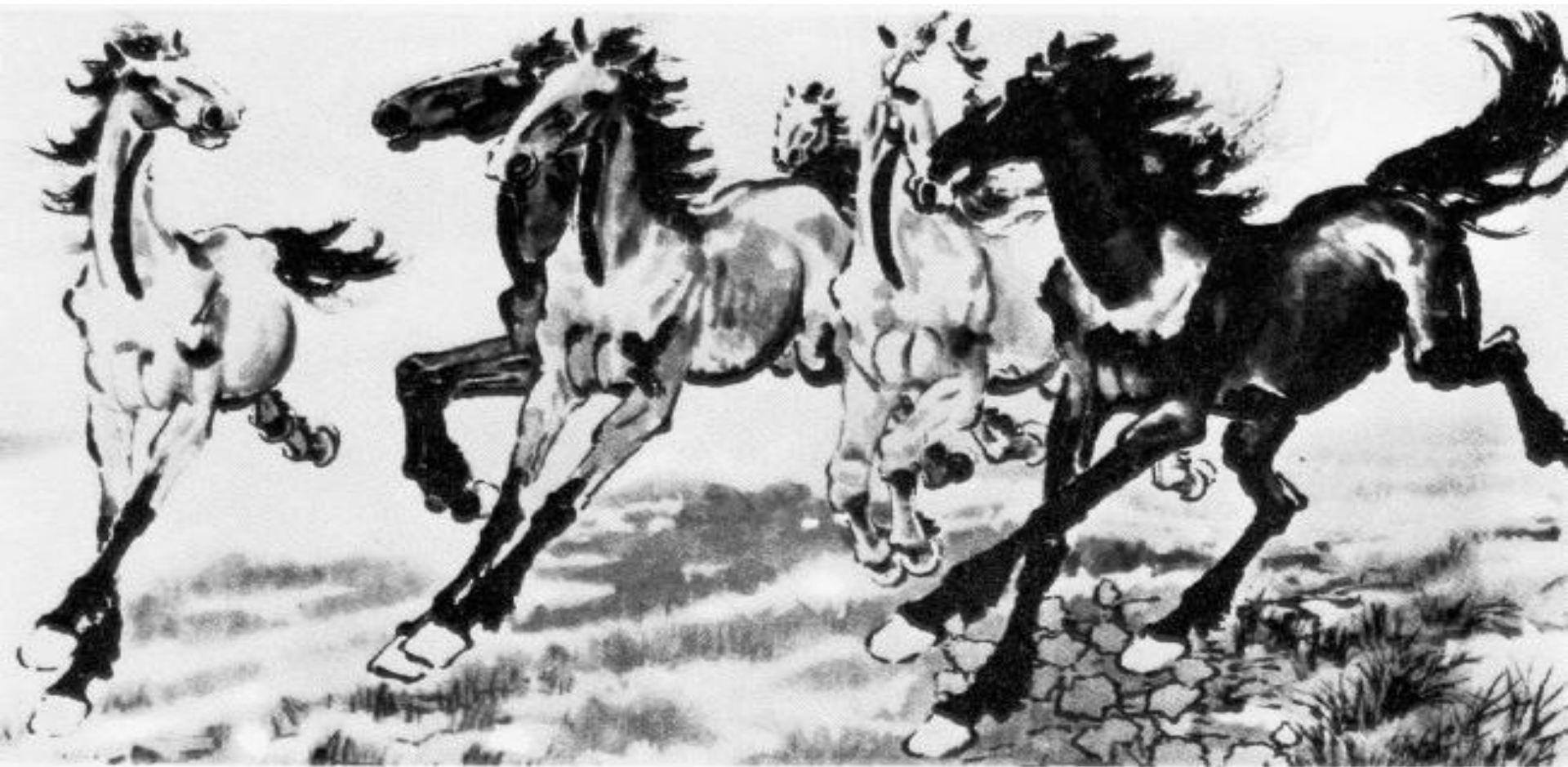


Magritte, 1957

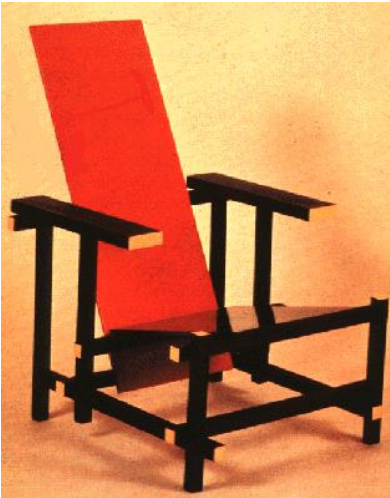
Challenges 4: scale



Challenges 5: deformation

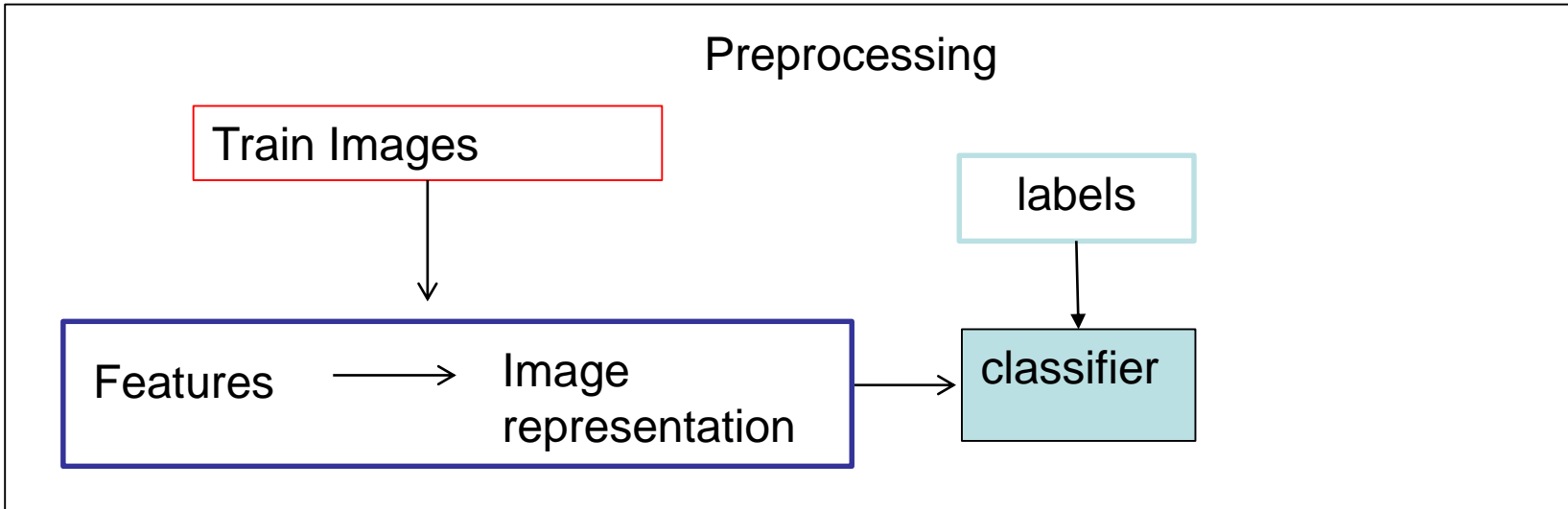


Challenges 7: intra-class variation

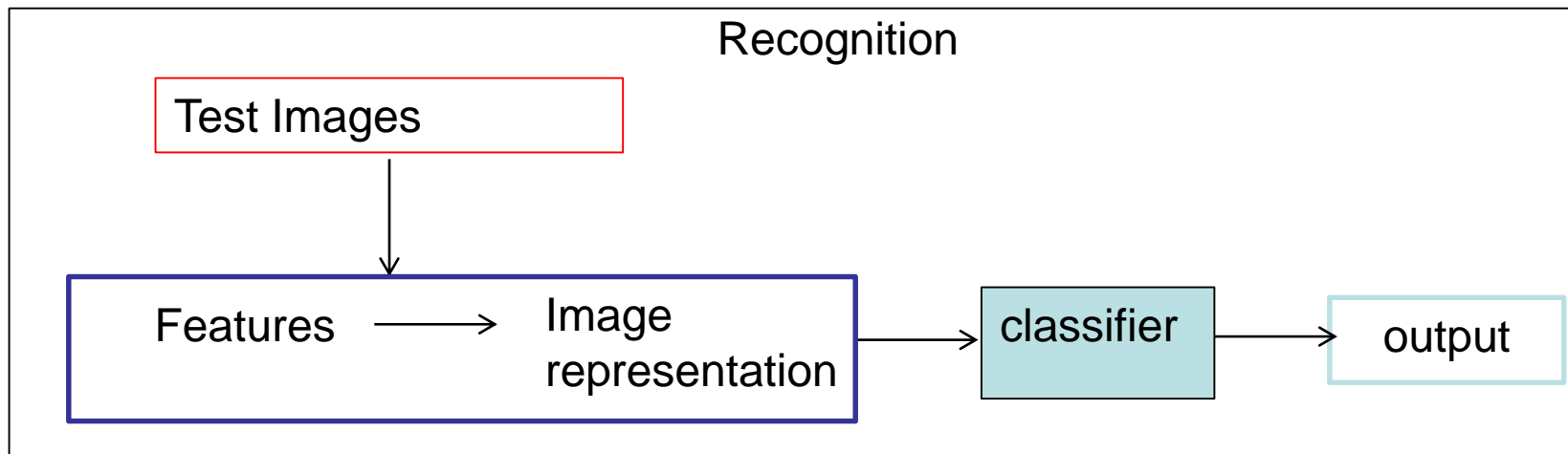


Recognition Steps

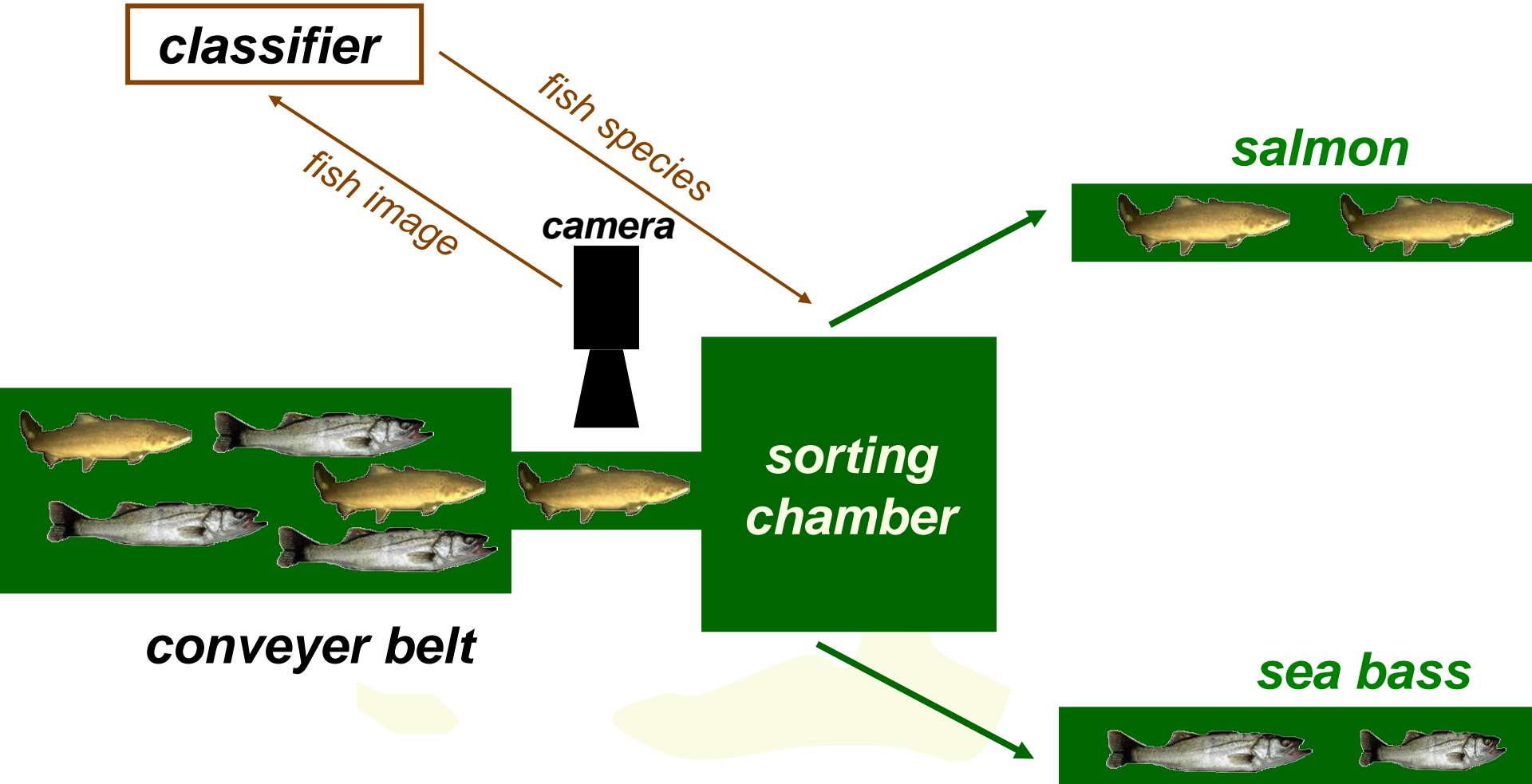
Preprocessing



Recognition

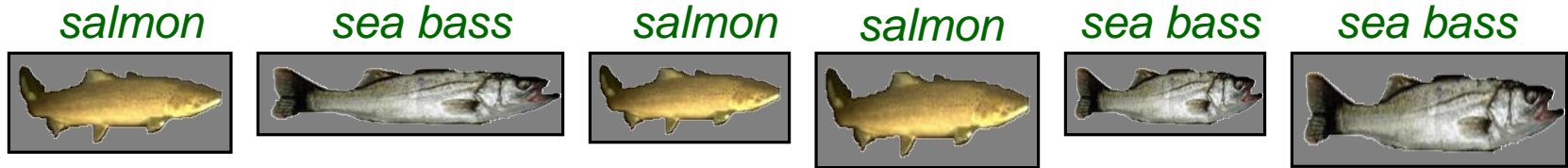


Object Recognition System



How to design a PR system?

- **Collect data** and classify by hand



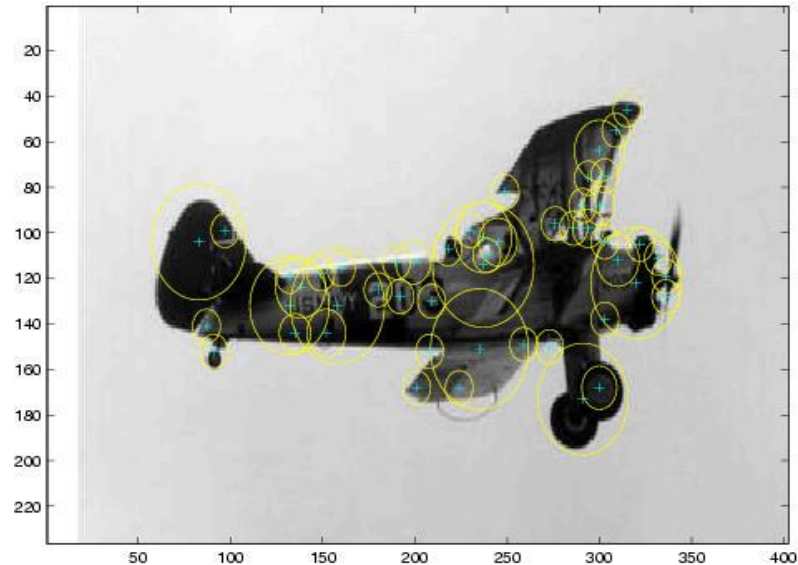
- **Preprocess** by segmenting fish from background



- **Extract** possibly discriminating **features**
 - length, lightness, width, number of fins, etc.
- **Classifier design**
 - **Choose model**
 - **Train classifier** on part of collected data (**training** data)
- **Test classifier** on the rest of collected data (**test** data)
i.e. the data not used for training
 - Should classify **new** data (new fish images) well

Interest Point Detectors

- Basic requirements:
 - Sparse
 - Informative
 - Repeatable
- Invariance
 - Rotation
 - Scale (Similarity)
 - Affine

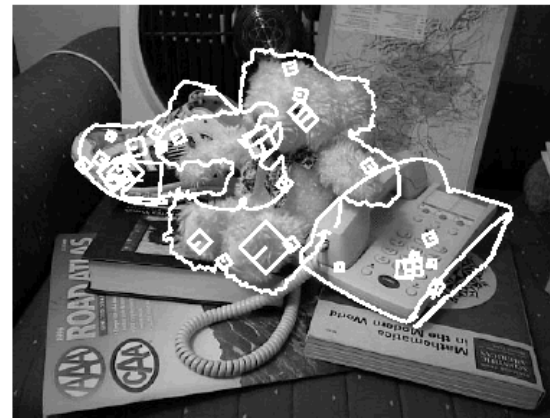
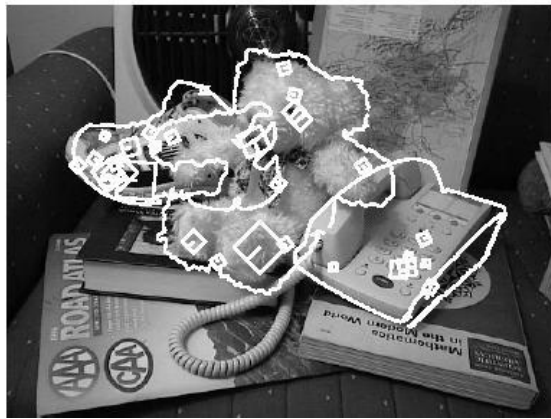


Recognizing Specific Objects

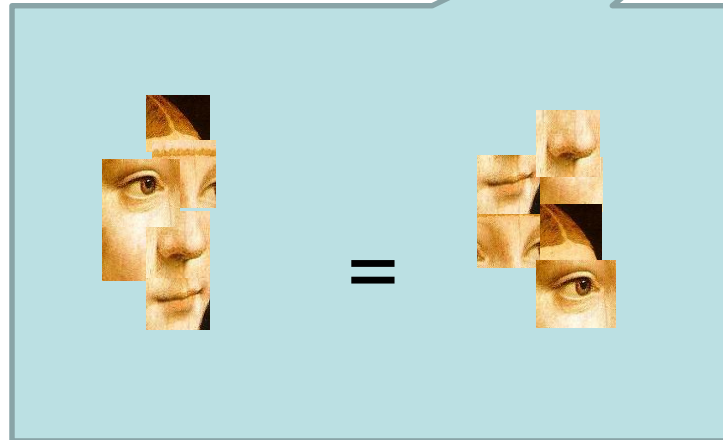
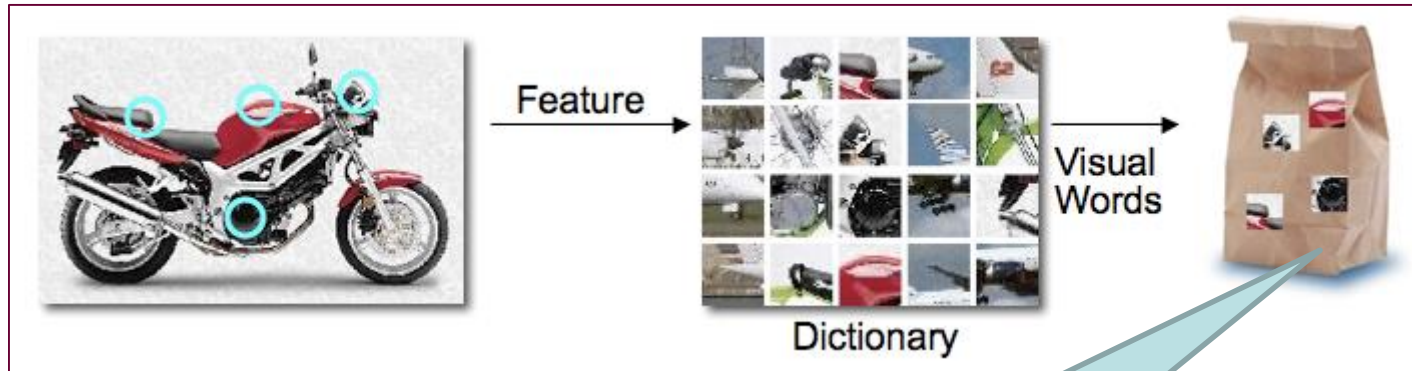
Learned models of local features, and got object outline from



Objects may then be found under occlusion and 3D rotation



Bag of Features



Bag of Features



Pros: fairly flexible and computationally efficient

Cons: problems with large clutter



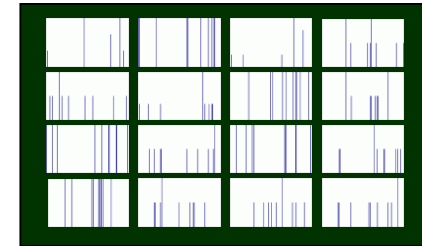
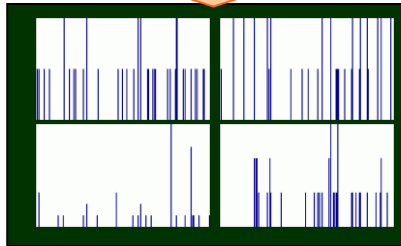
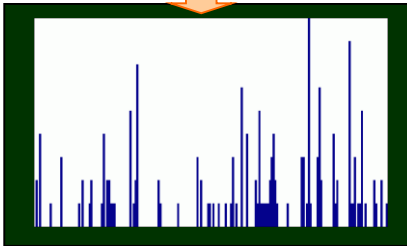
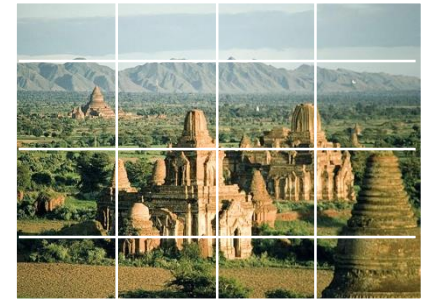
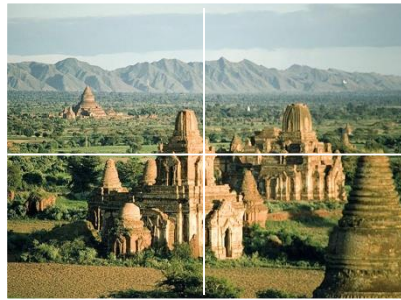
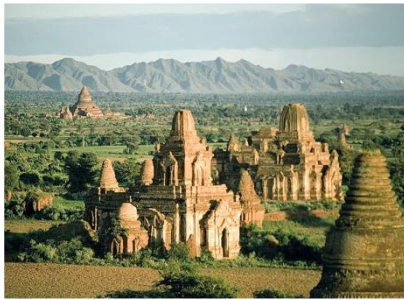
Different objects, but Similar representations;



Similar objects, different representations;

Beyond Bags of Features

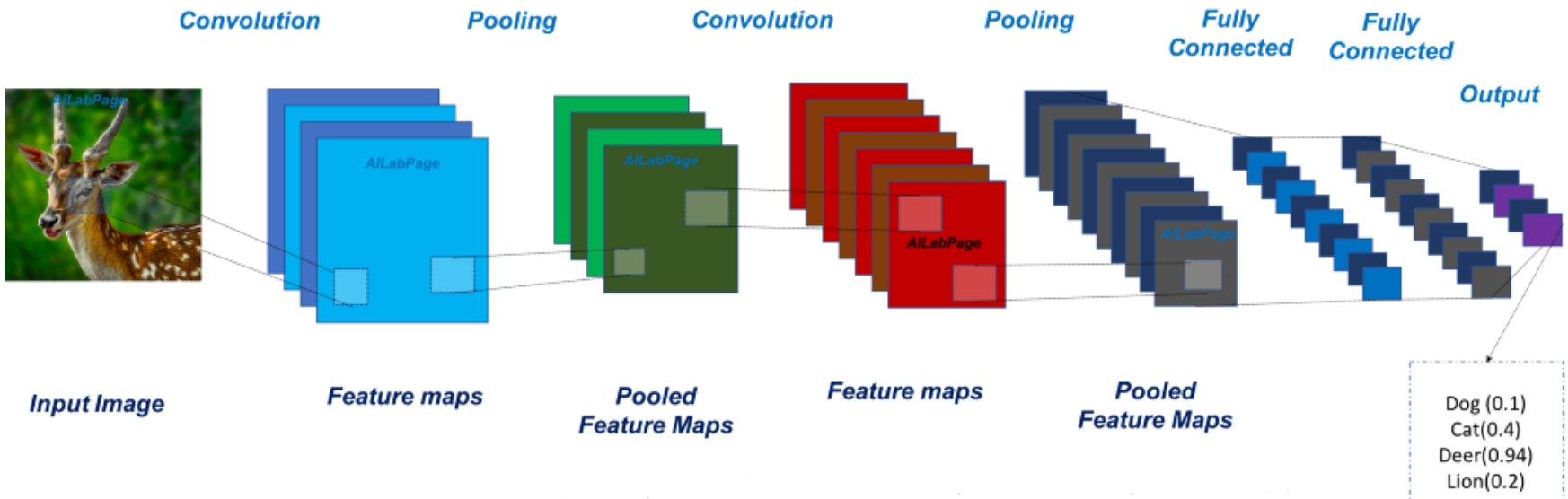
- Computing bags of features on sub-windows of the whole image.



Convolutional Neural Networks

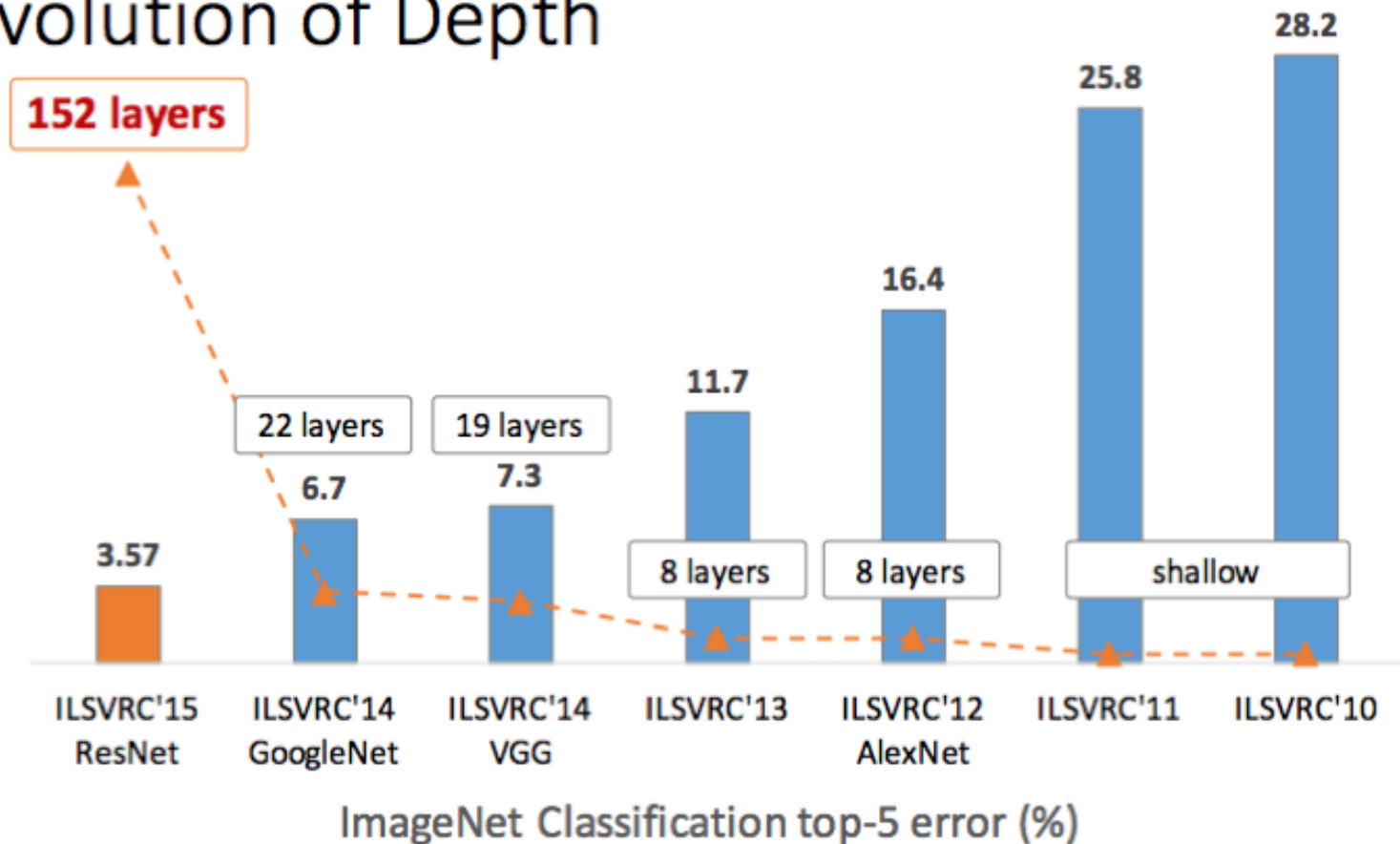
- Learn all in one deep architecture:
 - low level features
 - high level representations
 - context
 - classifiers
- Efficient Classification
- Efficient Detection
- Scalable to very large sets and large number of categories

Convolutional Neural Networks



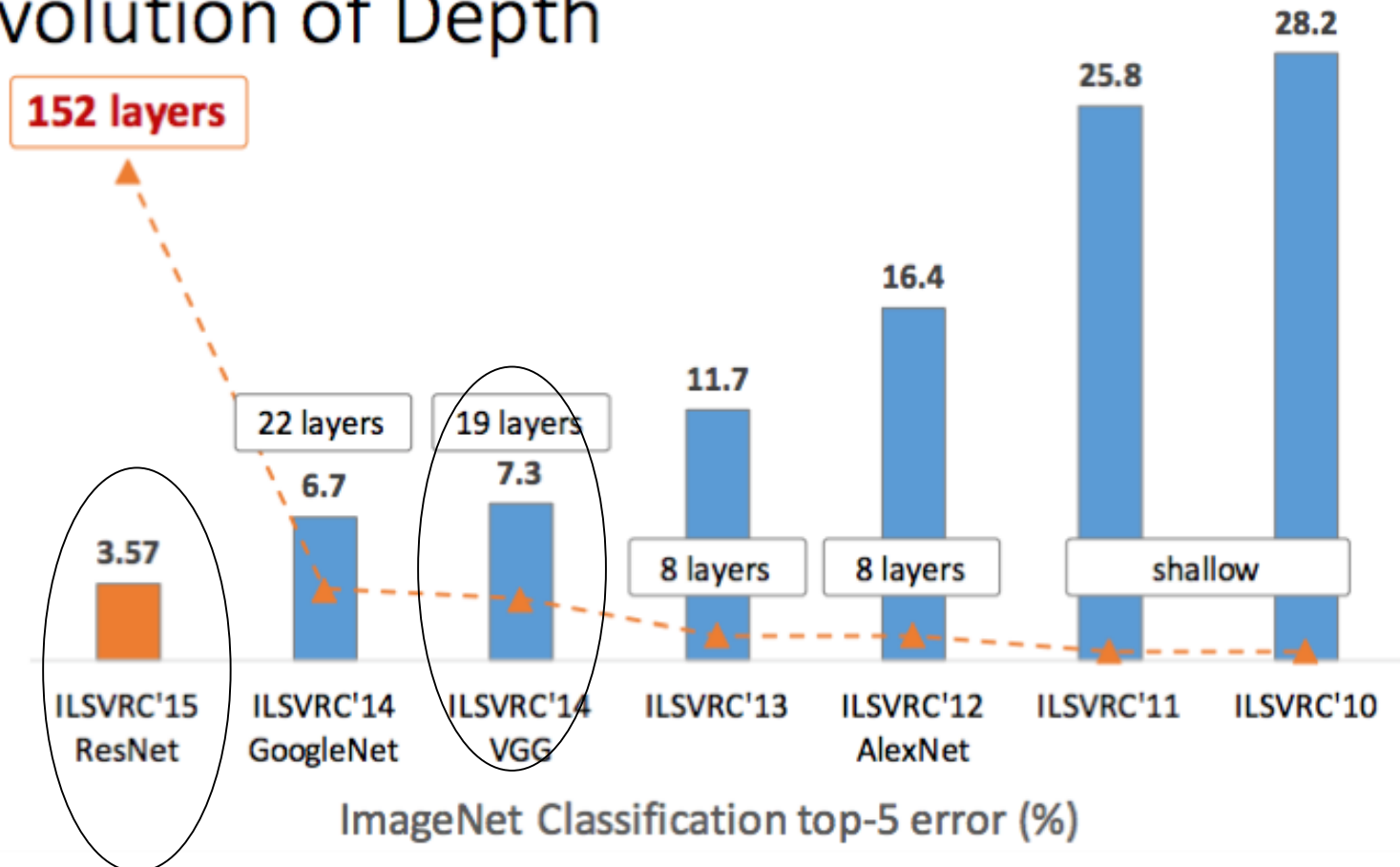
Very Deep Networks

Revolution of Depth



Very Deep Networks

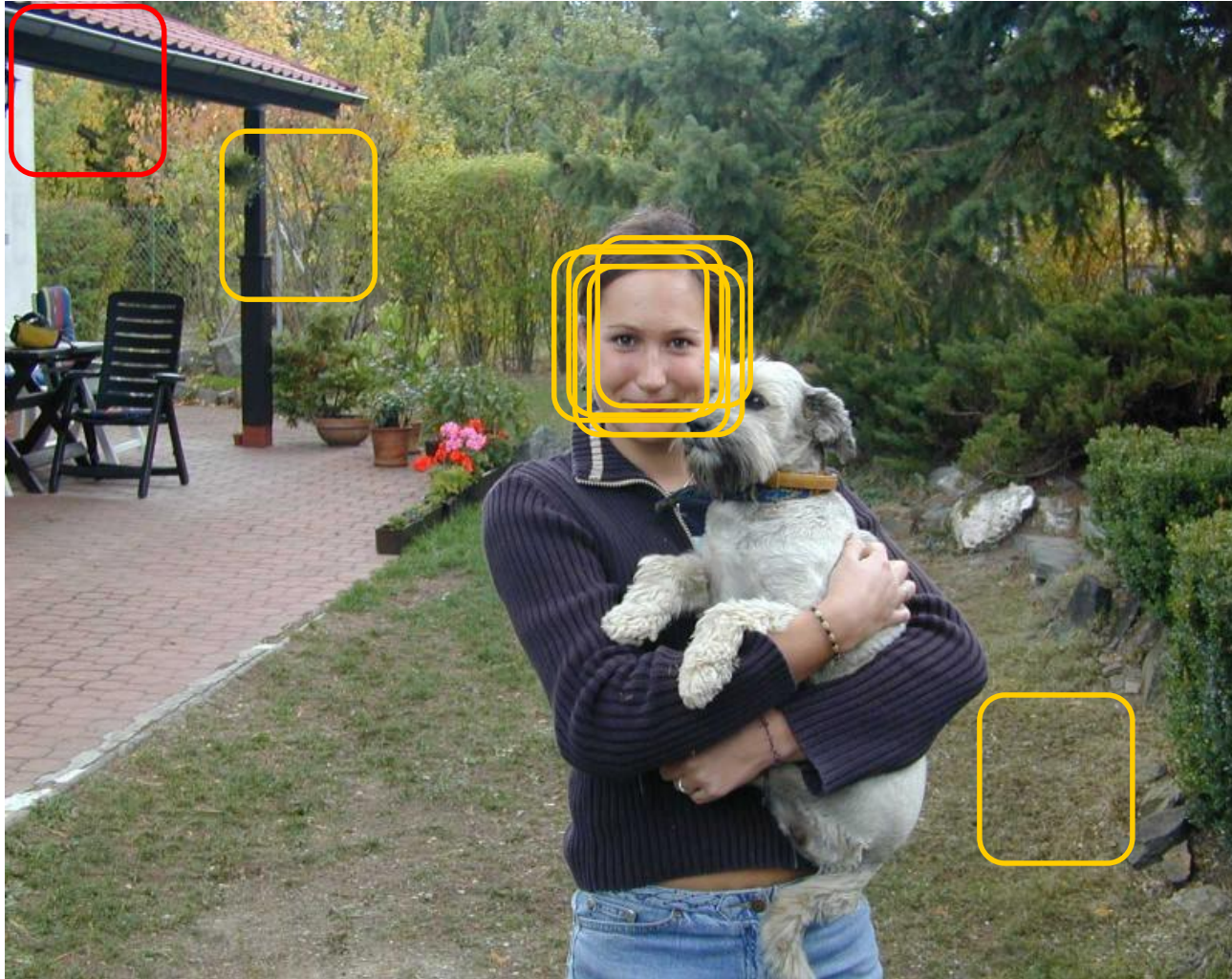
Revolution of Depth



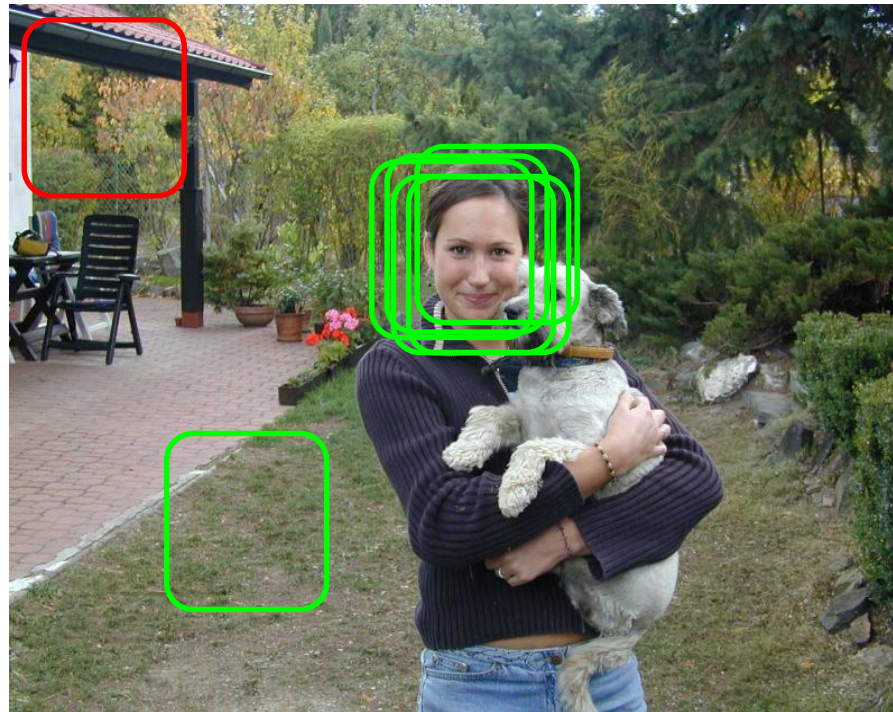
- K. Simonyan and A. Zisserman Very Deep Convolutional Networks for Large-Scale Image Recognition
- K. He, X. Zhang, S. Ren, and J. Sun: Deep Residual Learning for Image Recognition

Detection

Apply classifier at Scale / position range to search over



Detection

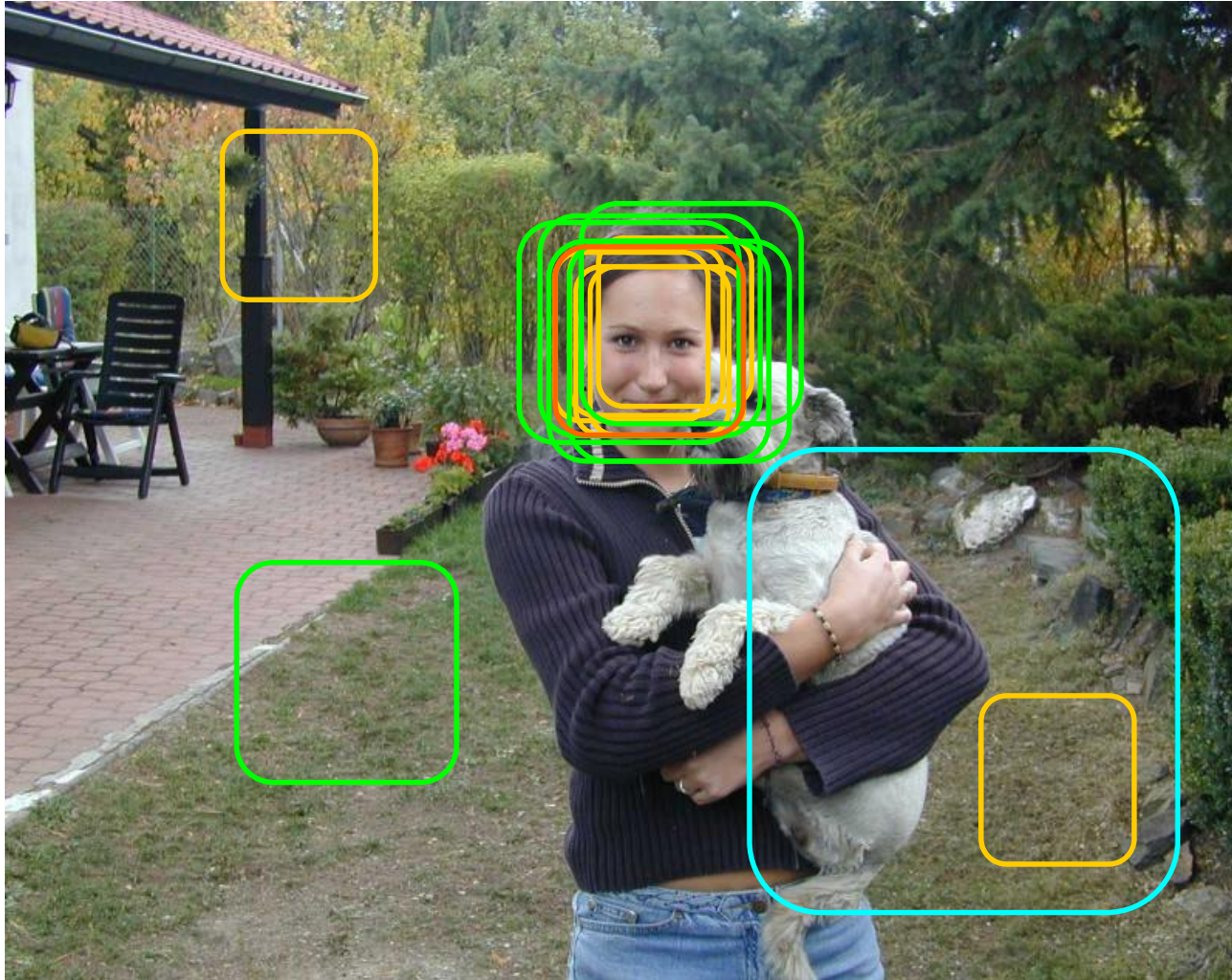


Detection



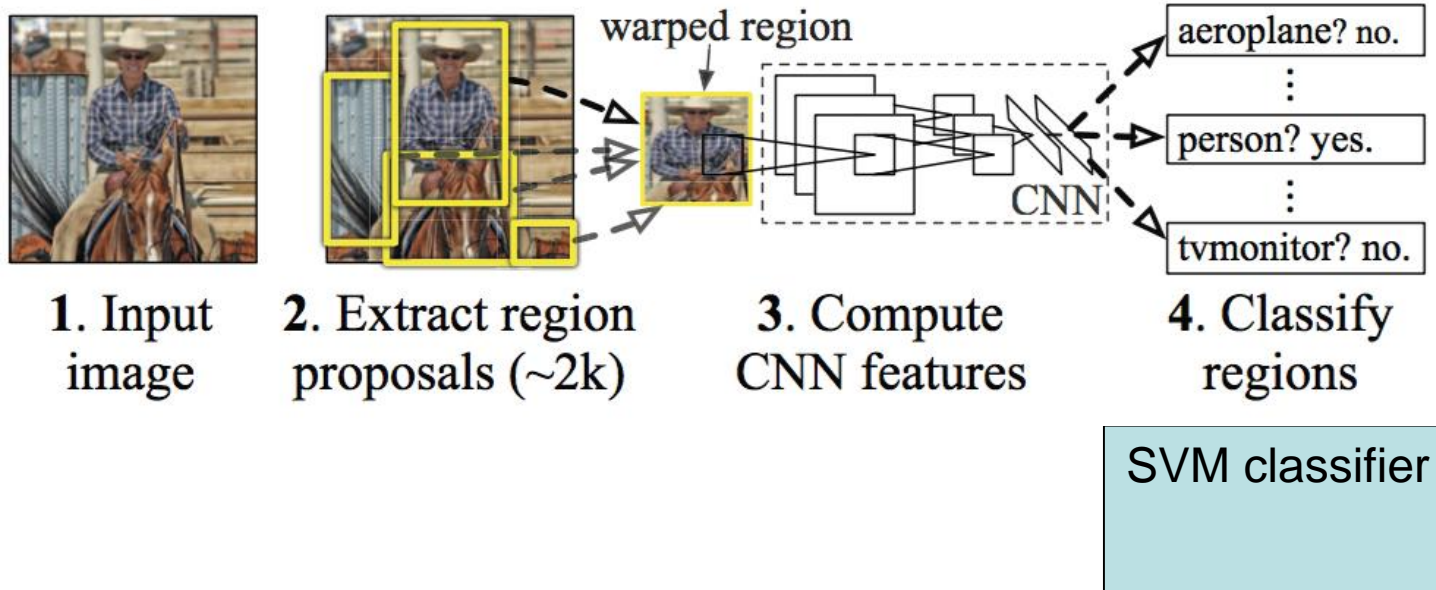
Detection

- Combine detection over space and scale.

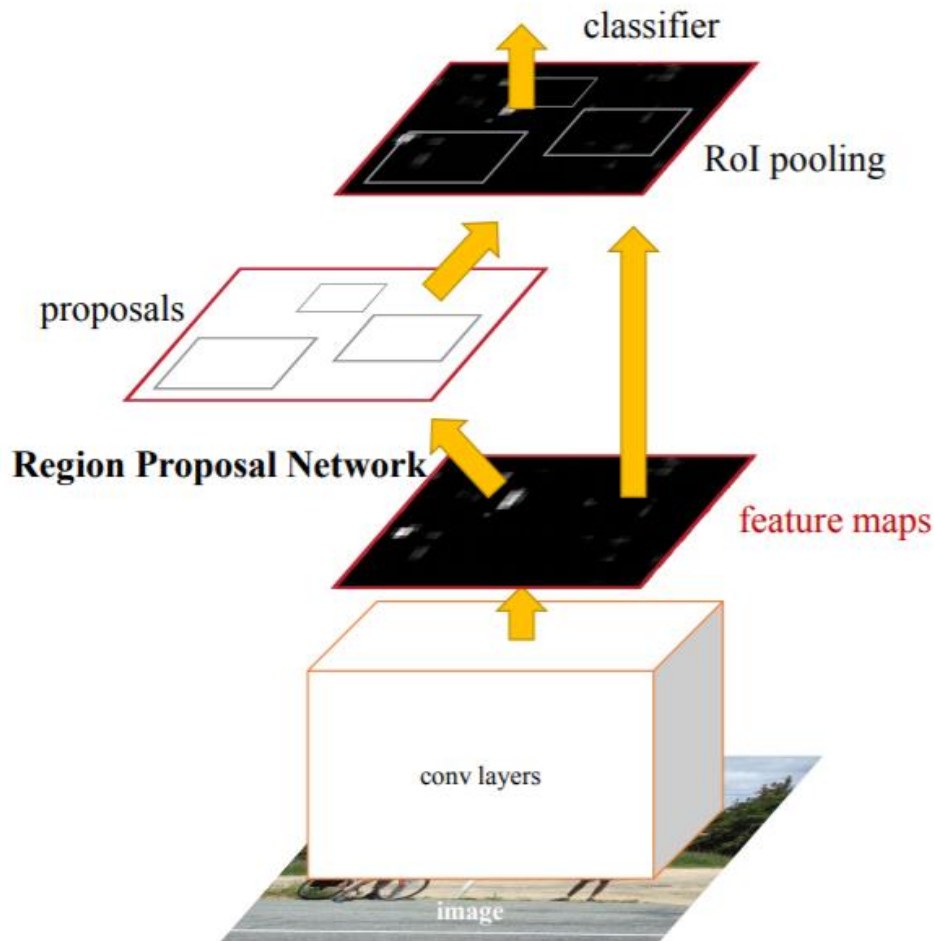


Deep Learning in Object Detection

R-CNN: *Regions with CNN features*



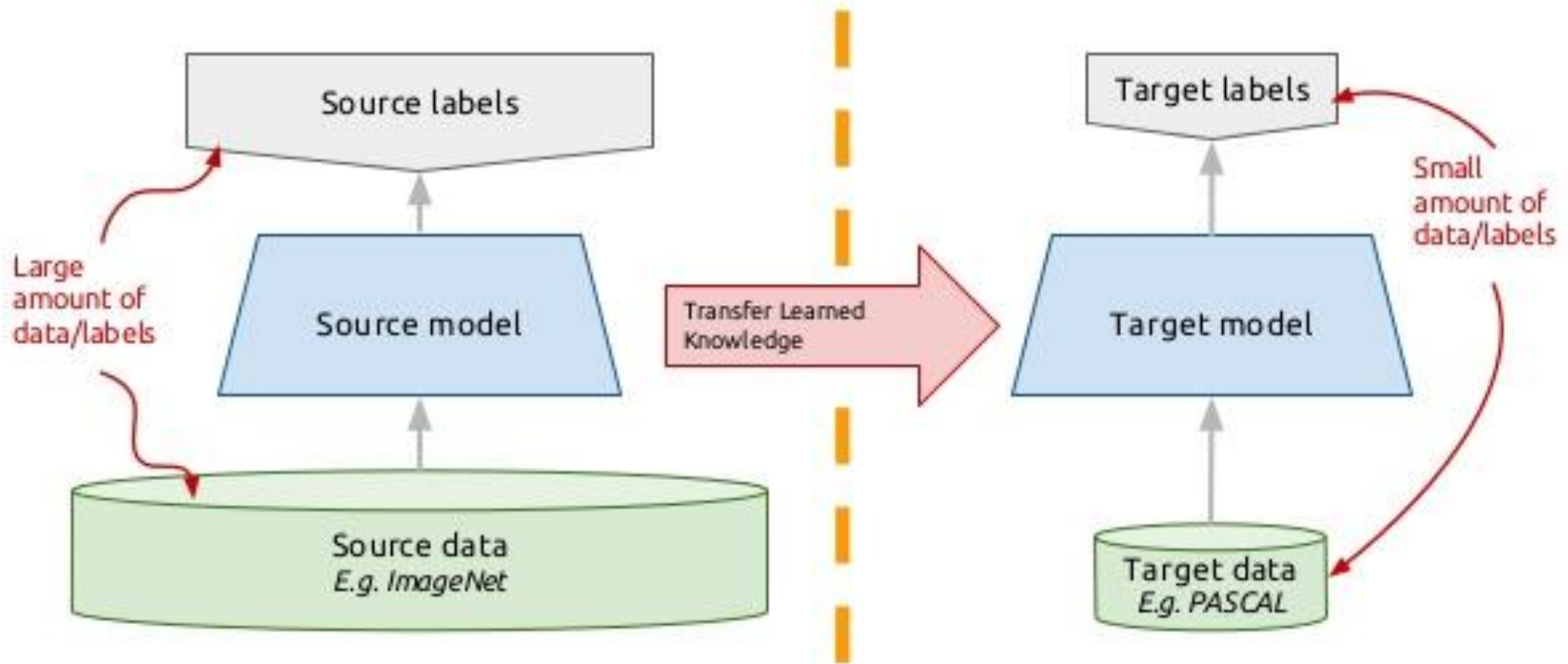
Faster R-CNN







Faster R-CNN is a single, unified network for object detection. The RPN module serves as the 'attention' of this unified network.

Transfer Learning

Transfer learning: idea



One Shot Learning

		same	"cow" (speaker #1)	"cow" (speaker #2)	same
		different	"cow" (speaker #1)	"cat" (speaker #2)	different
		same	"can" (speaker #1)	"can" (speaker #2)	same
		different	"can" (speaker #1)	"cab" (speaker #2)	different

Verification tasks (training)



One-shot tasks (test)

Few-Shot Learning

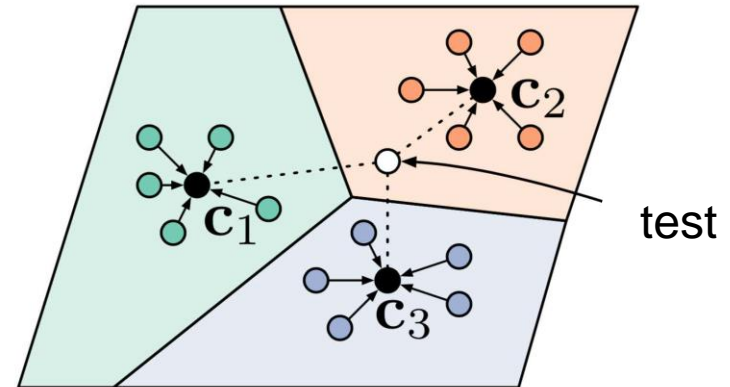
• Prototypical Networks

$$p_{\phi}(y = k | \mathbf{x}) = \frac{\exp(-d(f_{\phi}(\mathbf{x}), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_{\phi}(\mathbf{x}), \mathbf{c}_{k'}))}$$

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_{\phi}(\mathbf{x}_i)$$

$$S_k = \{(\mathbf{x}_i, y_i) | y_i = k, (\mathbf{x}_i, y_i) \in D_{train}\}$$

$$\phi \equiv \Theta$$



- Maps examples to embedding such that examples of a given class are close together
- Calculates a prototype (mean vector) for every class
- Maps test instances to the same embedding
- Uses softmax over distance to prototype for label prediction

Describing Objects with Attributes

- Shift the goal of recognition from naming to describing

otter

black:	yes
white:	no
brown:	yes
stripes:	no
water:	yes
eats fish:	yes



Discover/detect new categories

polar bear

black:	no
white:	yes
brown:	no
stripes:	no
water:	yes
eats fish:	yes

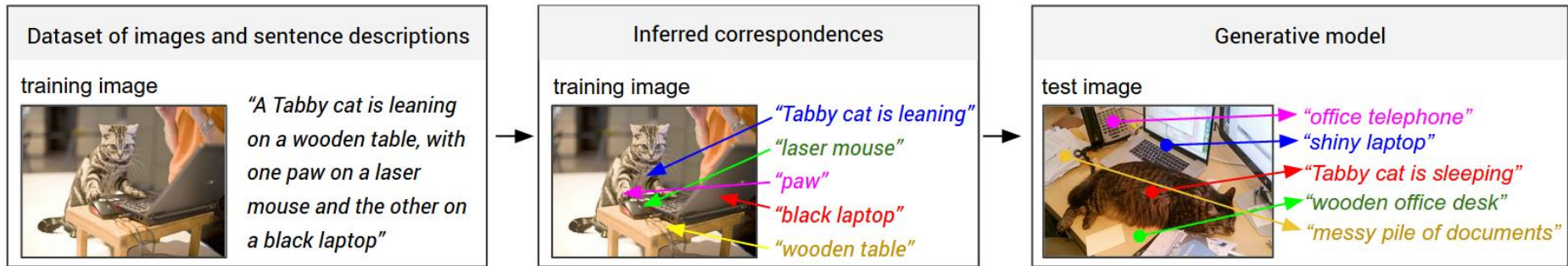


Improvement



Attributes	Presence		Rating	
	walrus	polar bear	walrus	polar bear
Spot	no	no	less relevant	irrelevant
Blue	no	no	irrelevant	less relevant
Swim	yes	yes	highly relevant	relevant
Coastal	yes	yes	relevant	highly relevant

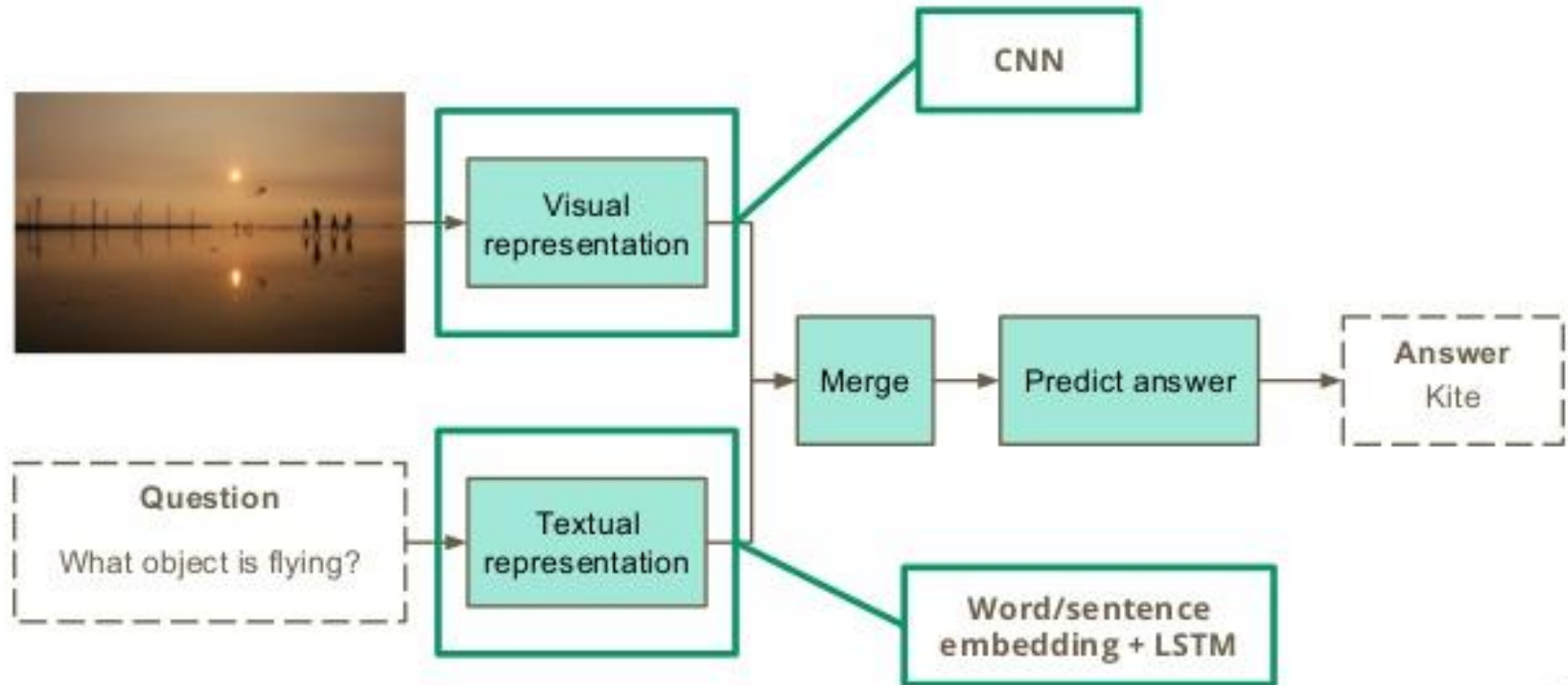
Image Descriptions (Captioning)



Uses CNNs and RNNs

VQA: Visual Question Answering

VQA: Common approach



Results

Where is the child sitting?

fridge

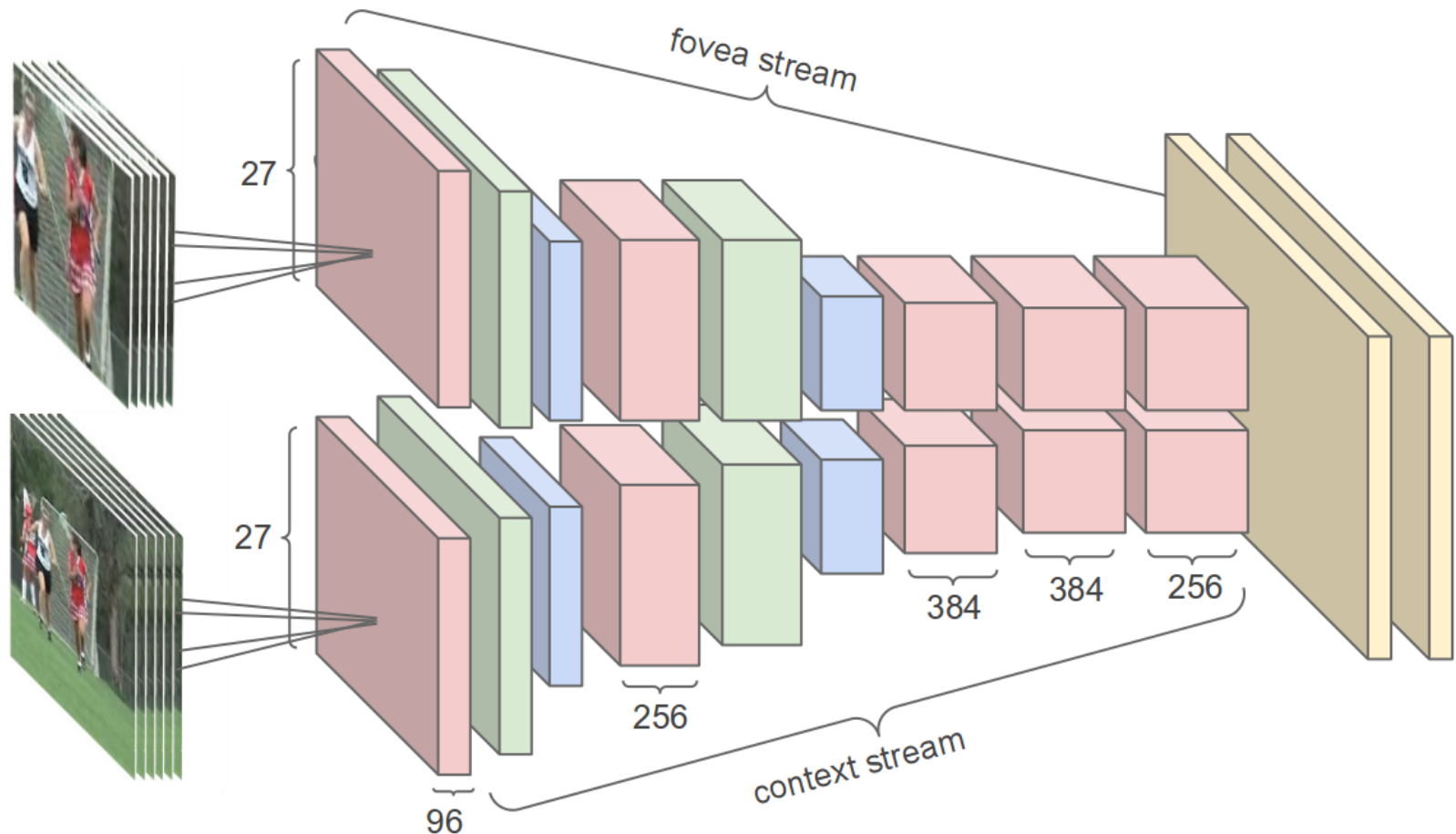


arms



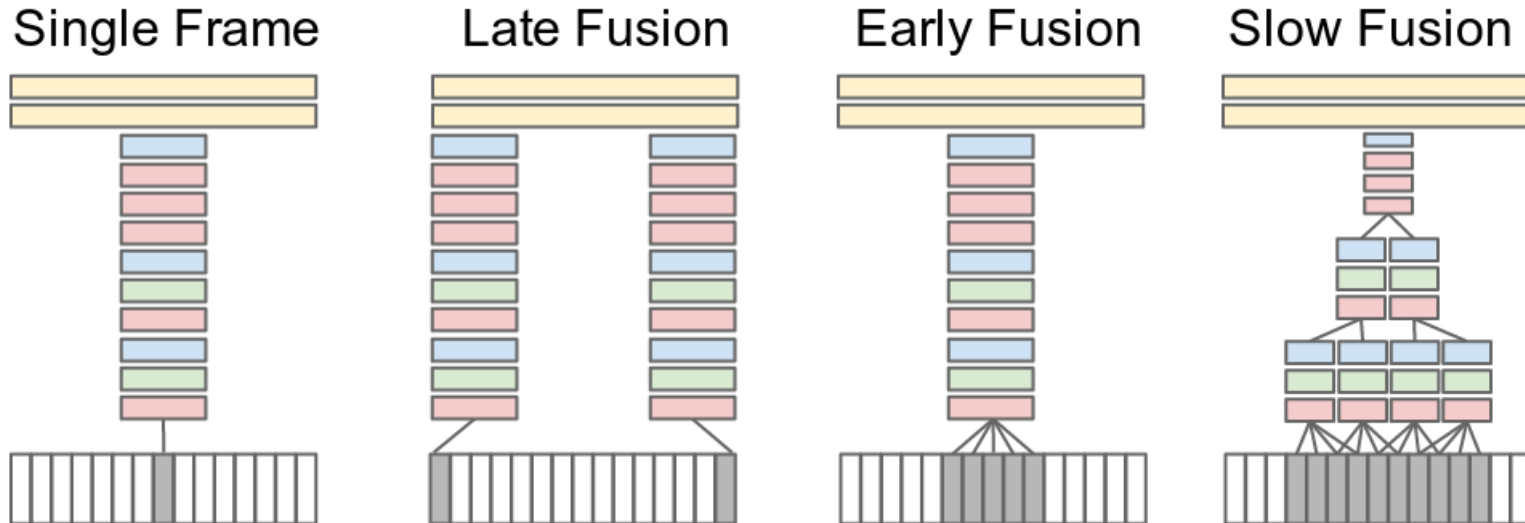
Video Classification

- Using CNN – Naïve Approach



Video Classification

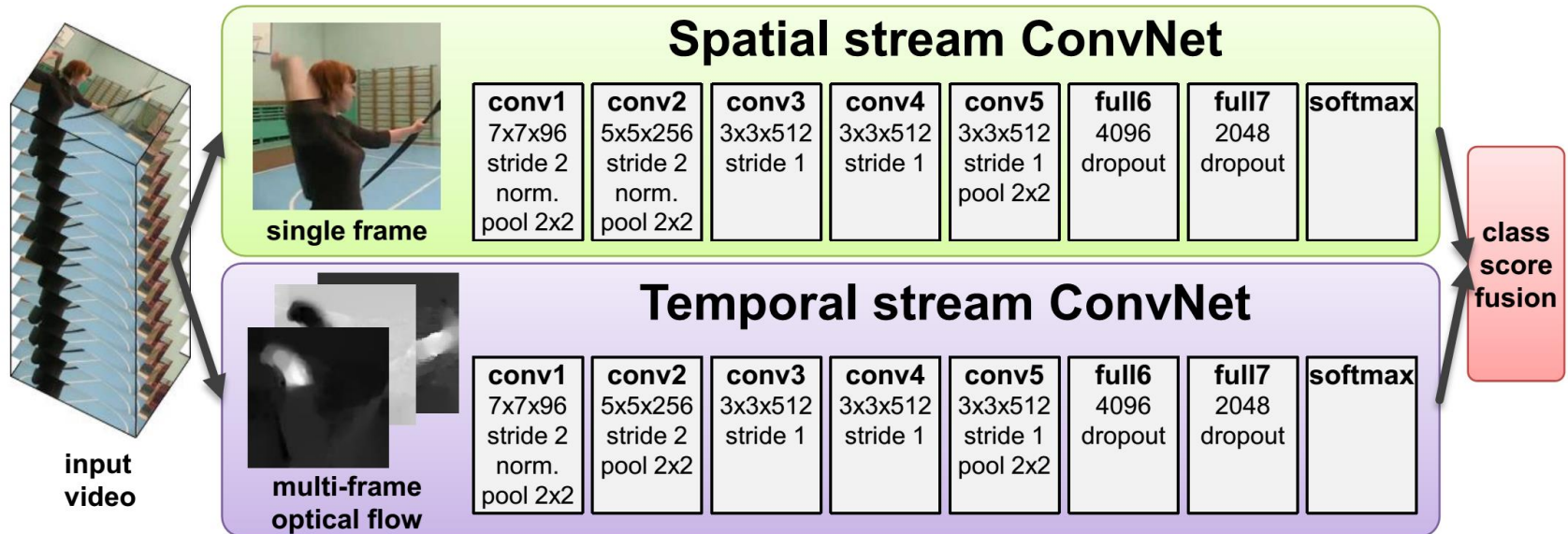
- Using CNN – Naïve Approach



Temporal fusion

Video Classification

- Modern Approaches



Multi-task training

Training Set



Random Batch



Continual Learning

Training Set
at Time T_1



Training Set at
Time T_2



Training Set
at Time T_3



...

Continual Learning

- Tasks are learned sequentially over time.
- At time T_i , the data of the tasks $T_1 \dots T_{i-1}$ is no longer available.
- Leads to forgetting of previously learned tasks.
- Termed “Catastrophic Forgetting”.

Possible Solution:

- Measure the importance of learnable parameters in NN for the task, constrain their change in learning a new task.
- This could work because deep networks are highly overparametrized

Adversarial Examples


 x

“panda”

57.7% confidence

 $+ .007 \times$

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

“nematode”

8.2% confidence

 $=$

 $x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Adversarial Examples: Imperceptible Noise

