



CORNELL
TECH

Deep Learning Clinic (DLC)

Lecture 9

Case Study: Object Detection and Segmentation

Jin Sun

11/19/2019

Today - Object Detection and Segmentation

- **Overview**
- Object Detection
- Segmentation
- Detection + Segmentation: Instance Segmentation
- What's Next

It's All About Objects

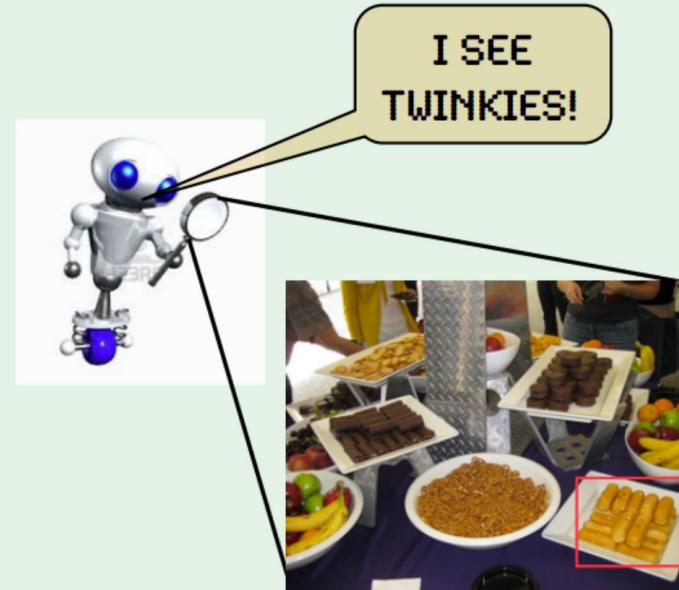
Snack time in the lab



What objects are where?



•
•



robot: "I see a table with twinkies,
pretzels, fruit, and some mysterious
chocolate things..."

Computer Vision Tasks

Classification



CAT

No spatial extent

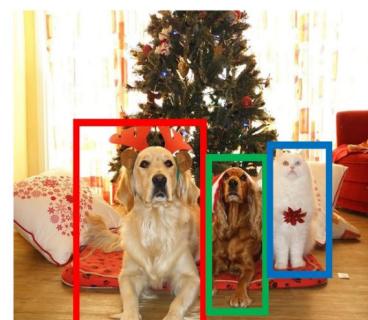
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

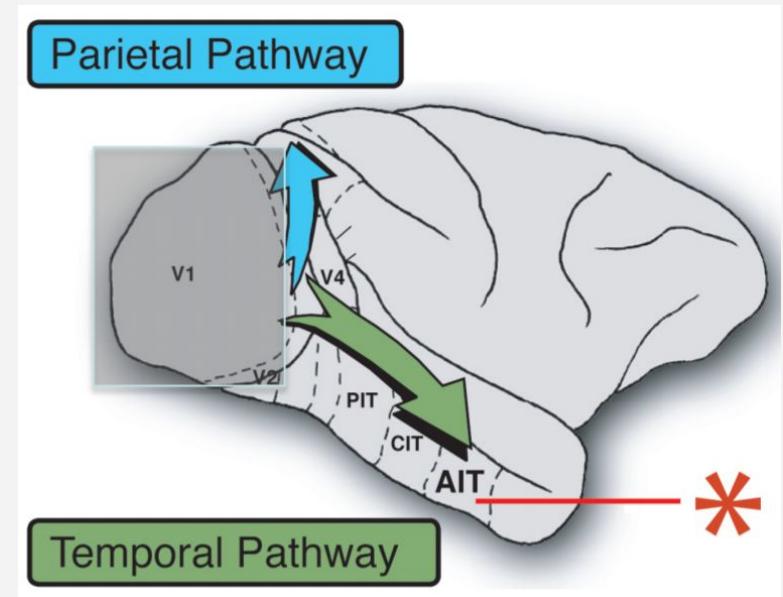
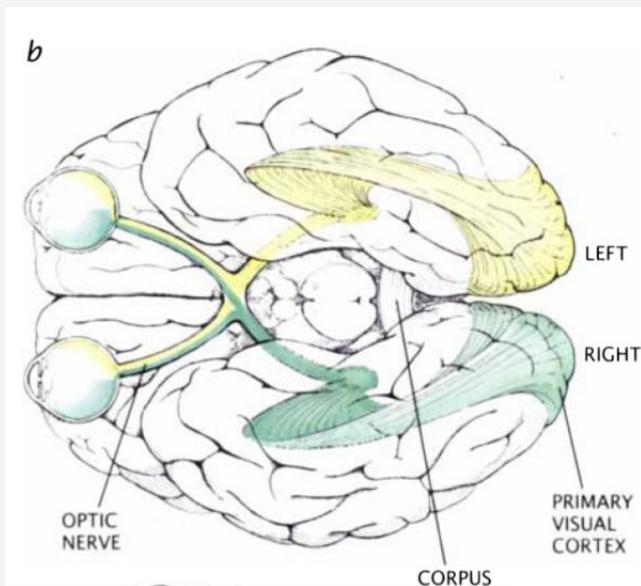
Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

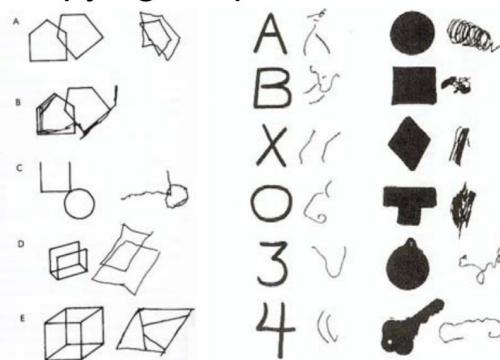
Object Recognition in Brains



http://klab.tch.harvard.edu/academia/classes/Neuro230/2016/slides/hms230_2016_Lecture3.pdf

Lesion Study on Object Recognition

Copying shapes



- Patient cannot name, copy or match simple shapes

- Acuity, color recognition and motion perception are preserved

- Bilateral damage to extrastriate visual areas

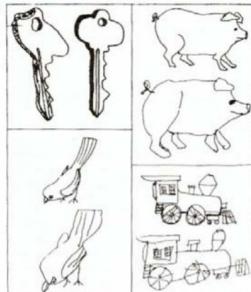
Matching shapes

△		●	▲	■	■
■		■	■	▲	●
X		L	∅	X	N
▲		■	▲	■	●

Warrington 1985

Lesion Study on Object Recognition

Copying from templates



- Subject can copy complex drawings, match complex shapes and use the objects correctly

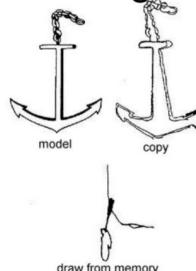
- Subject cannot identify (name) those shapes

- Subject cannot draw from memory

- Acuity, color recognition and motion perception are preserved

- Bilateral lesion of the anterior inferior temporal lobe

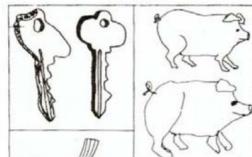
Drawing from memory



Warrington 1985

Lesion Study on Object Recognition

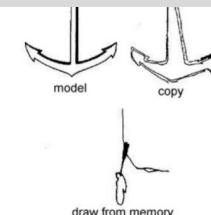
Copying from templates



- Subject can copy complex drawings, match complex shapes and use the objects correctly
- Subject cannot identify (name) those shapes

Object Recognition seems to be effortless for human

But actually it is quite complicated!



Damage to the anterior part of the temporal lobe

Warrington 1985

Why This Is A Hard Problem

Viewpoint



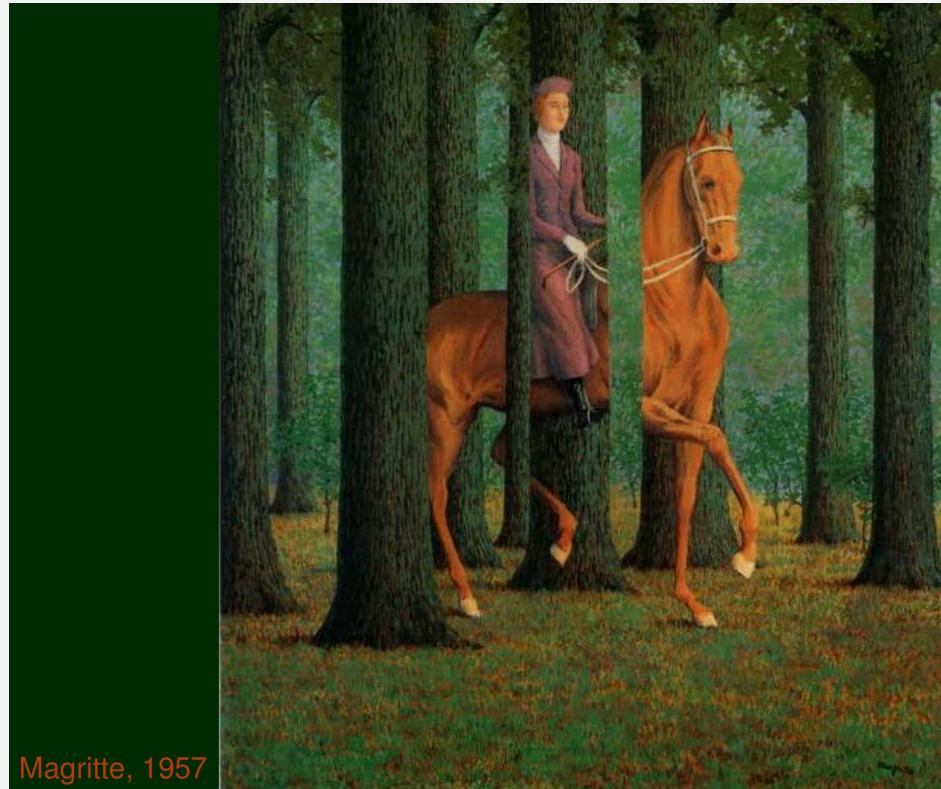
Why This Is A Hard Problem

Illumination



Why This Is A Hard Problem

Occlusion



Why This Is A Hard Problem

Deformation



Why This Is A Hard Problem

Background Clutter



Why This Is A Hard Problem

In-class variations



Slides from Rob Fergus

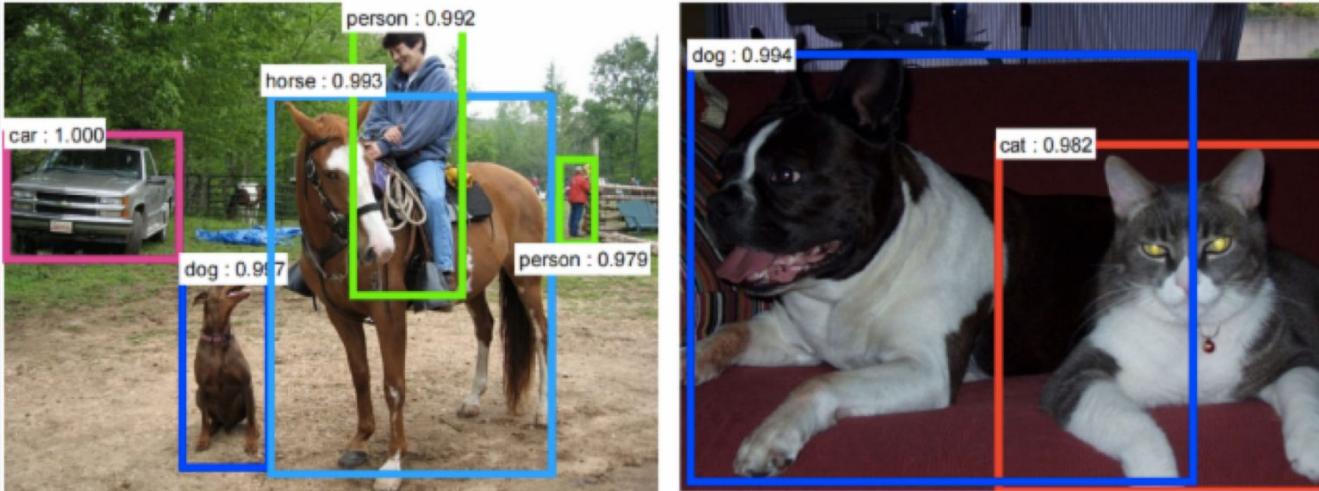
Applications of Object Recognition

- Face Detection
- Face Verification
- Assistive Driving
- Cashier-free Shopping
- Drones
- And much more!

Today - Object Detection and Segmentation

- Overview
- **Object Detection**
- Segmentation
- Detection + Segmentation: Instance Segmentation
- What's Next

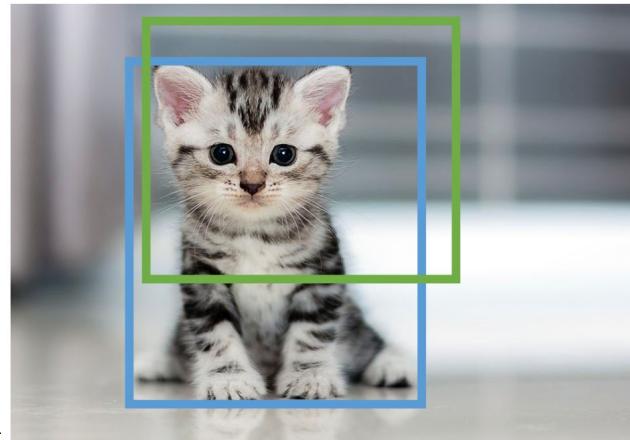
Task Formulation



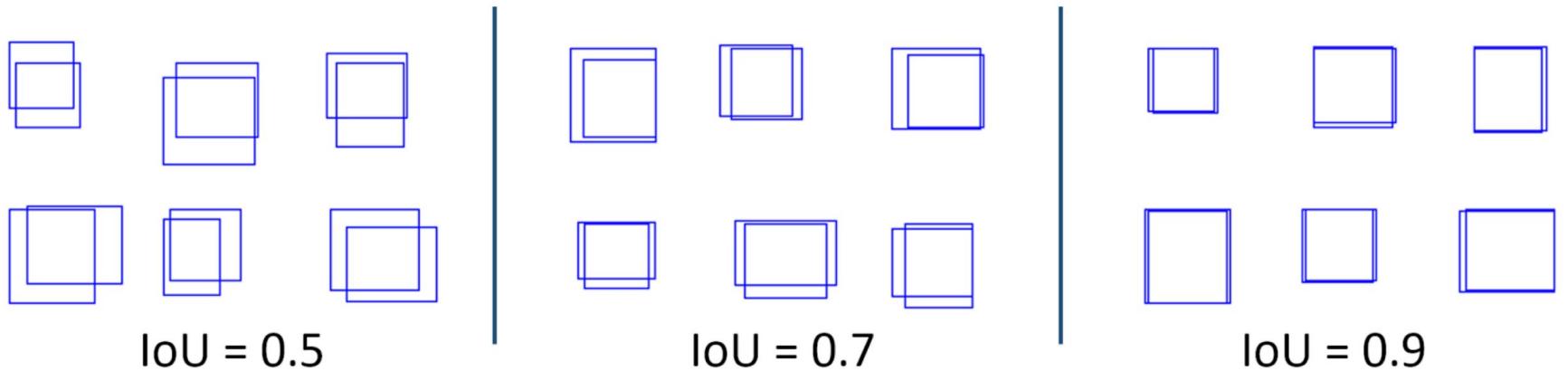
Evaluation Metric



$$\text{Overlap} = \frac{\text{True box} \cap \text{Predicted}}{\text{True box} \cup \text{Predicted}} > \text{threshold} \quad \checkmark$$

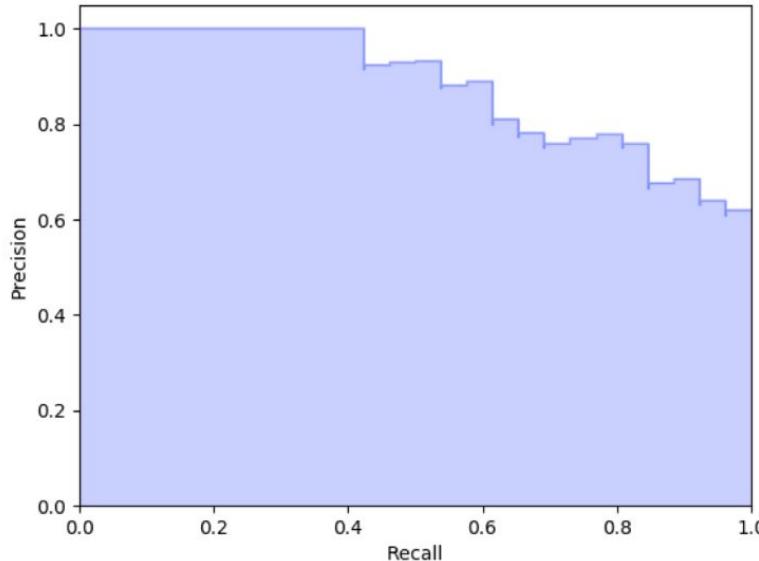


Evaluation Metric



Edge boxes, Zitnick and Dollar

Evaluation Metric



Precision:
true positive detections /
total detections

Recall:
true positive detections /
total positive test instances

Average Precision: Area under curve

Simple Detector: Sliding Window with Correlation

Example 1: Find Waldo!

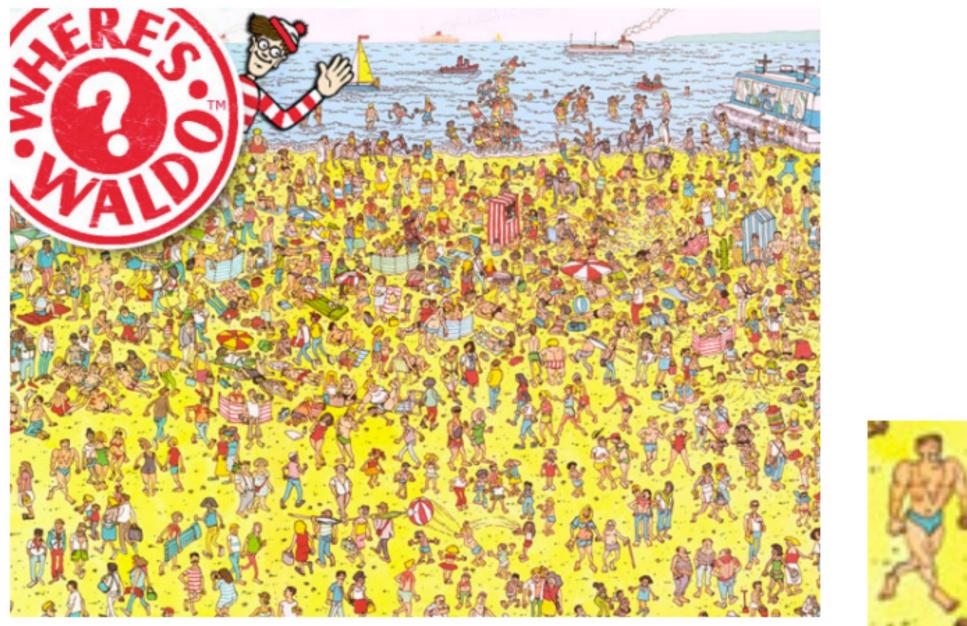
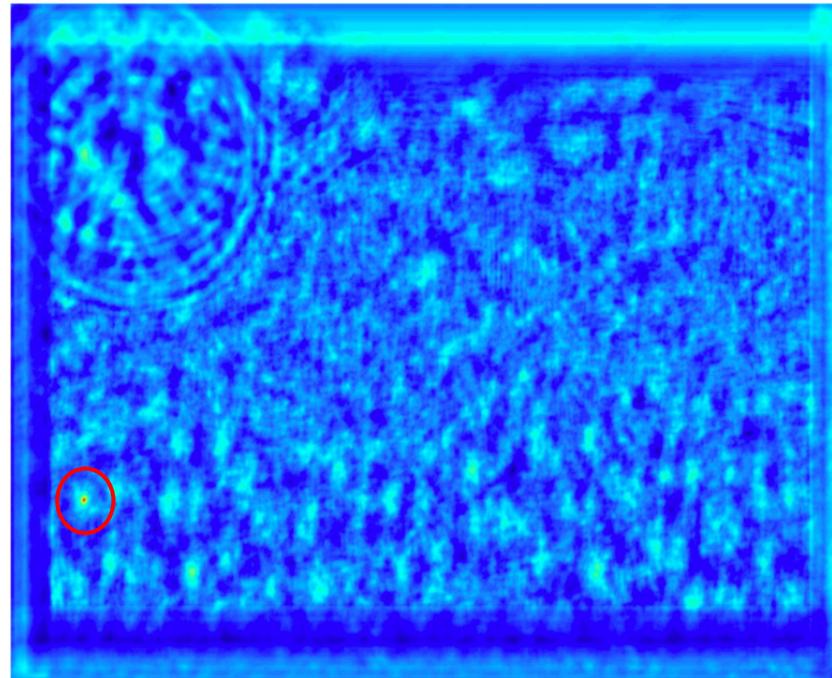


image 1

slide credit: Chris McIntosh

Simple Detector: Sliding Window with Correlation

2. Result of normalized cross-correlation



slide credit: Chris McIntosh

Simple Detector: Sliding Window with Correlation

Example 2: Find all persons?



slide credit: Chris McIntosh



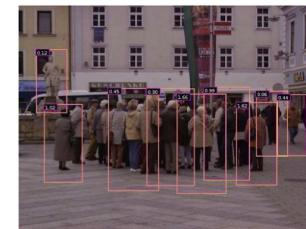
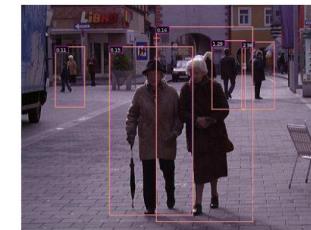
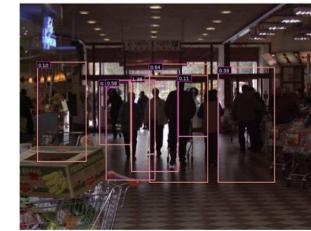
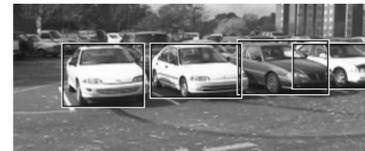
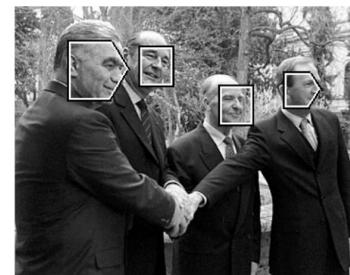
A template for all instances?

We need features!

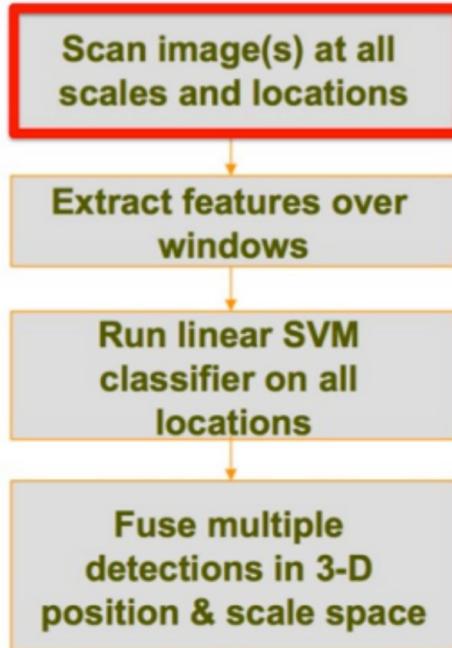
Pre-Deep Learning Era



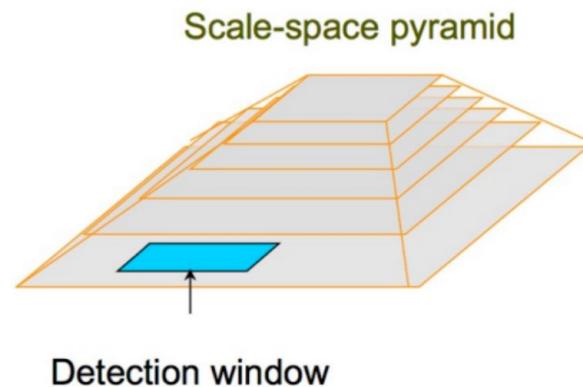
- Wide variety of poses
- Variable appearance
- Complex backgrounds
- Unconstrained illumination
- Occlusions, different scales
- ...



Histogram of Oriented Gradients (HOG)

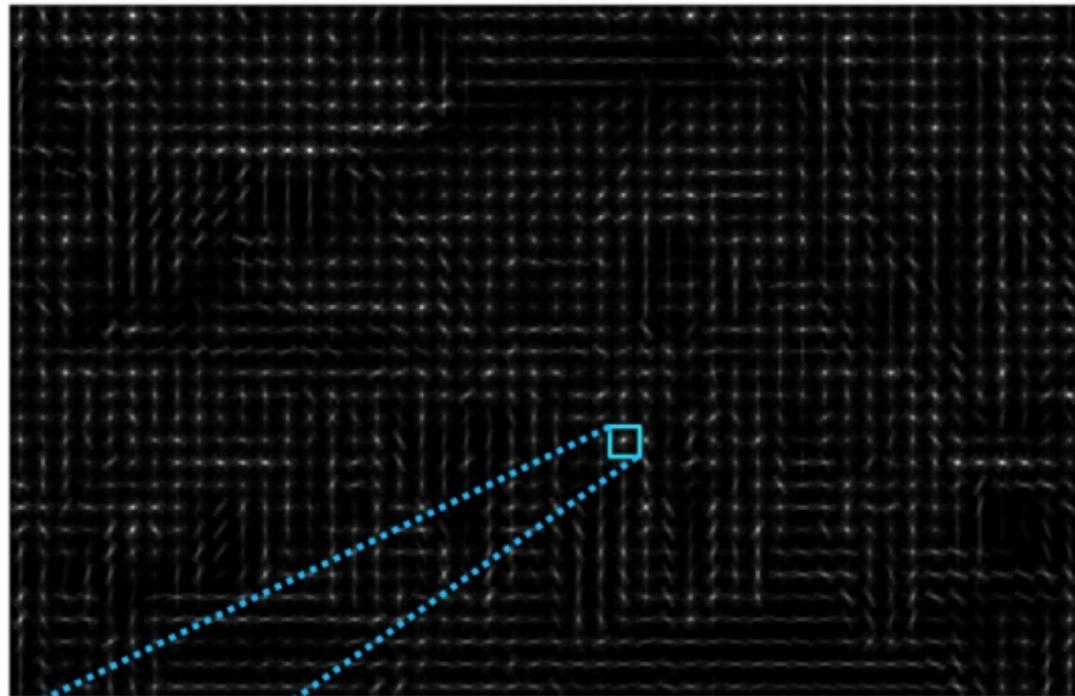
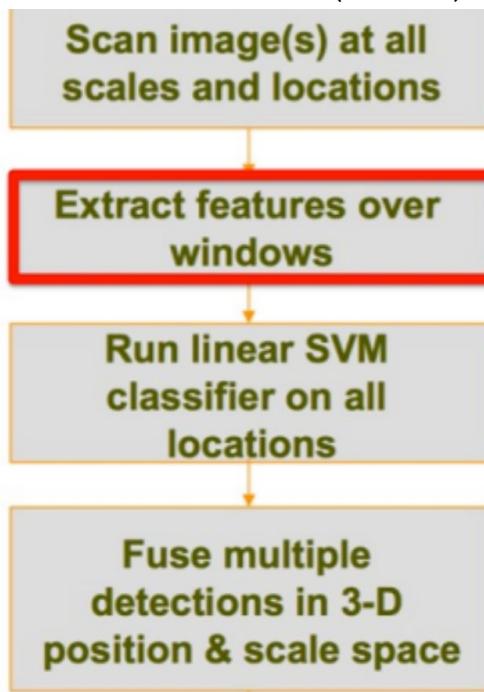


locations

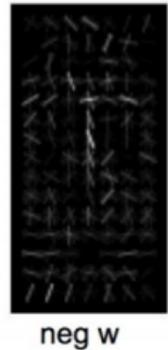
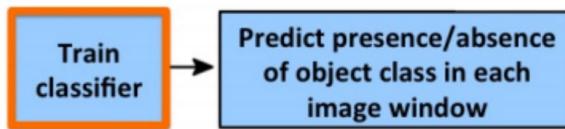


scales

Histogram of Oriented Gradients (HOG)



9-dim feature vector

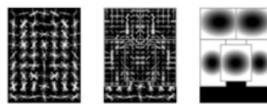
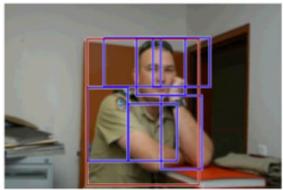
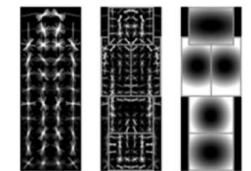


$$w^T \cdot x + b = 0$$

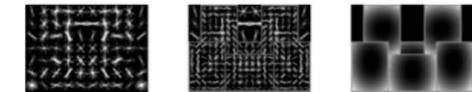
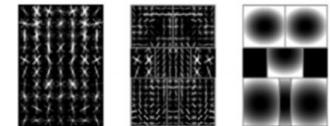
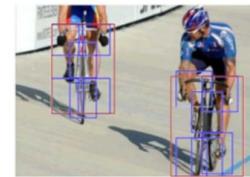
Train classifier. SVM (Support Vector Machines) is typically used.

Deformable Parts Model

Mixture Model Example - Person

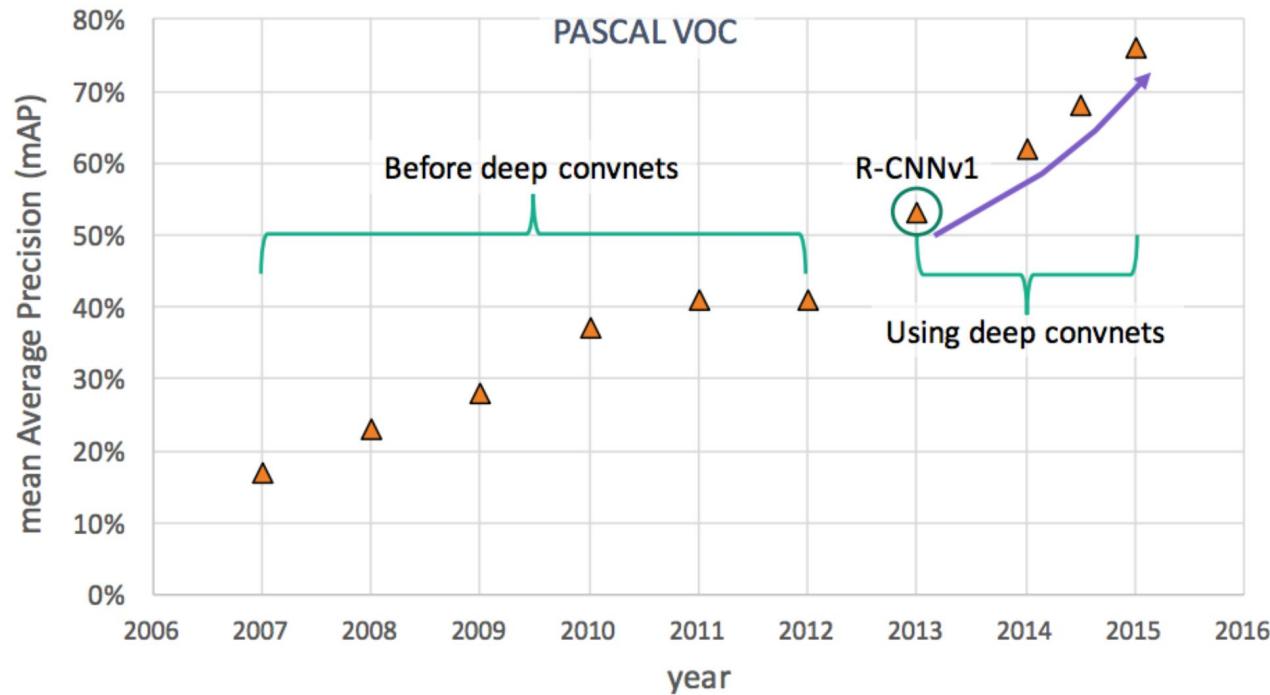


Mixture Model Example - Bicycle



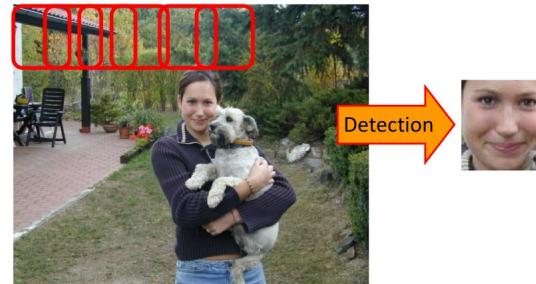
Slide credit: Mubarak Shah

A Significant Boost from Deep Learning

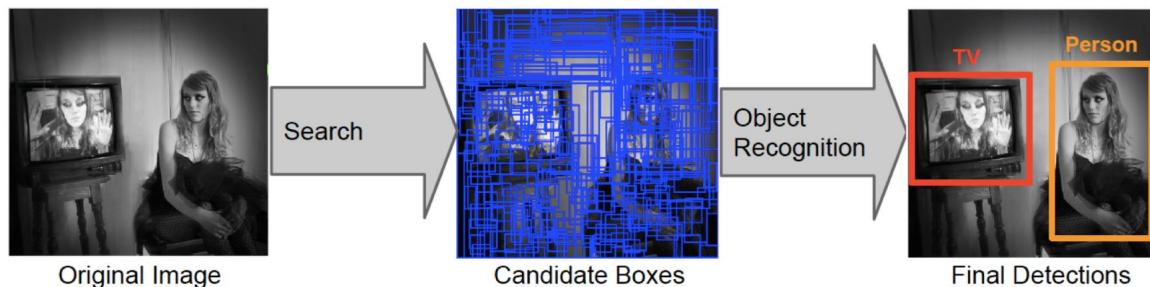


Two Types of Detection Paradigm

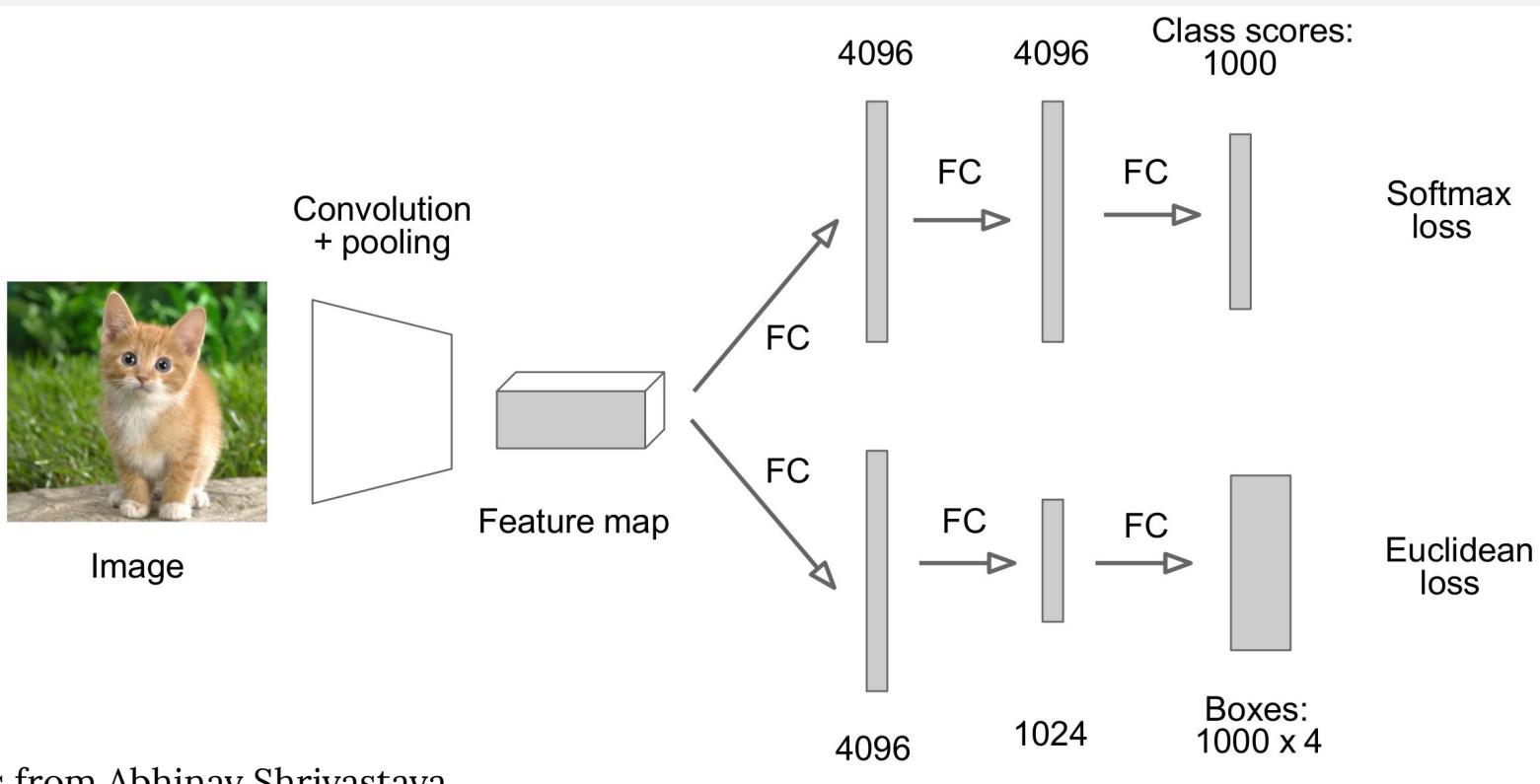
Single-stage Detectors (I-Stage)



Two-stage Detectors (II-Stage):



Single Stage - Overfeat



Single Stage - YOLO

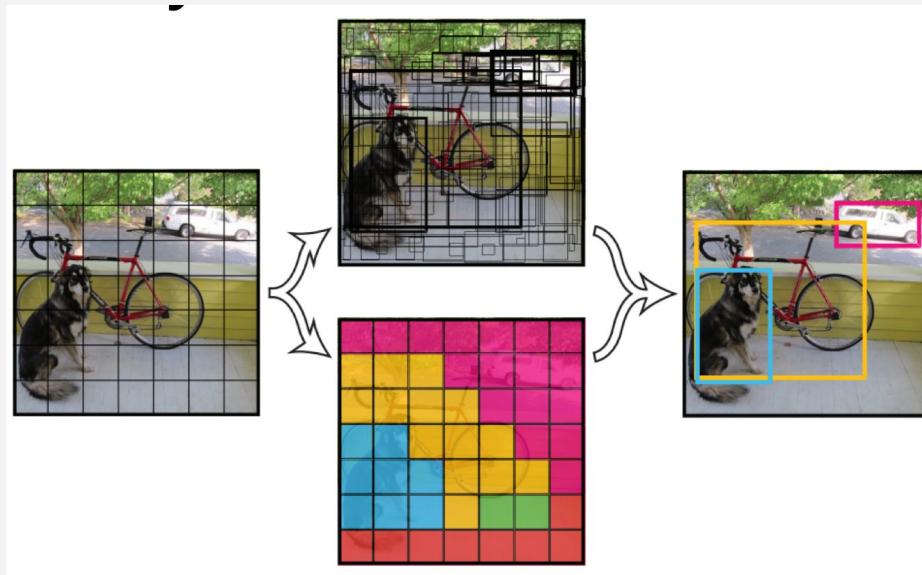
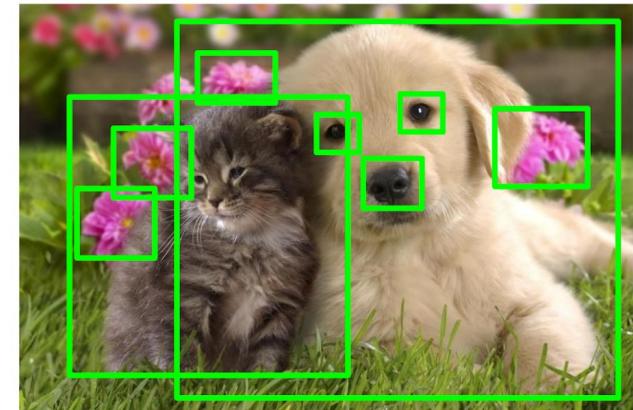


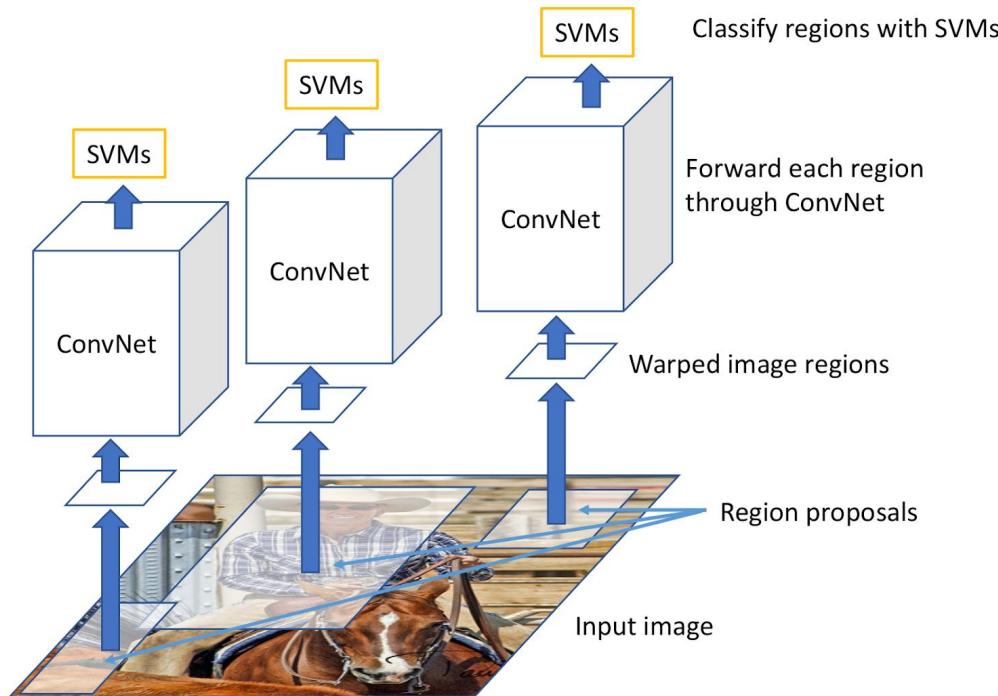
Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an even grid and simultaneously predicts bounding boxes, confidence in those boxes, and class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

Region Proposals in Two Stage Detectors

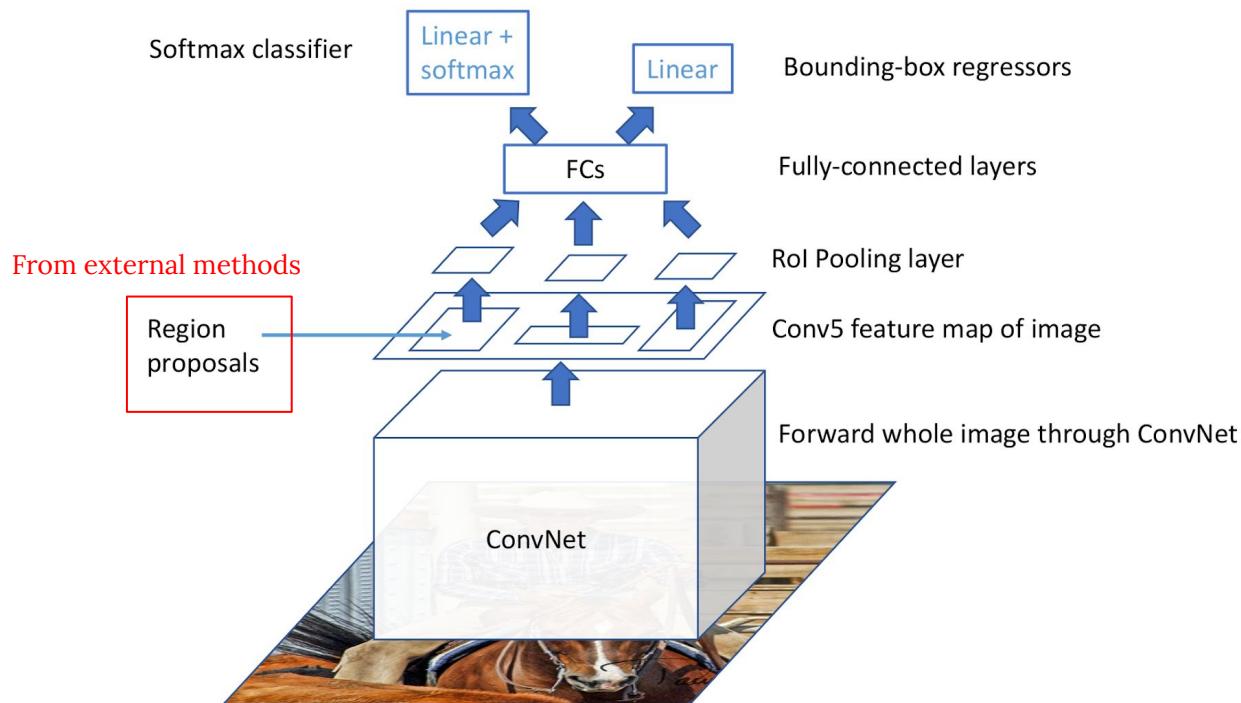
- Find “blobby” image regions that are likely to contain objects
- “Class-agnostic” object detector
- Look for “blob-like” regions



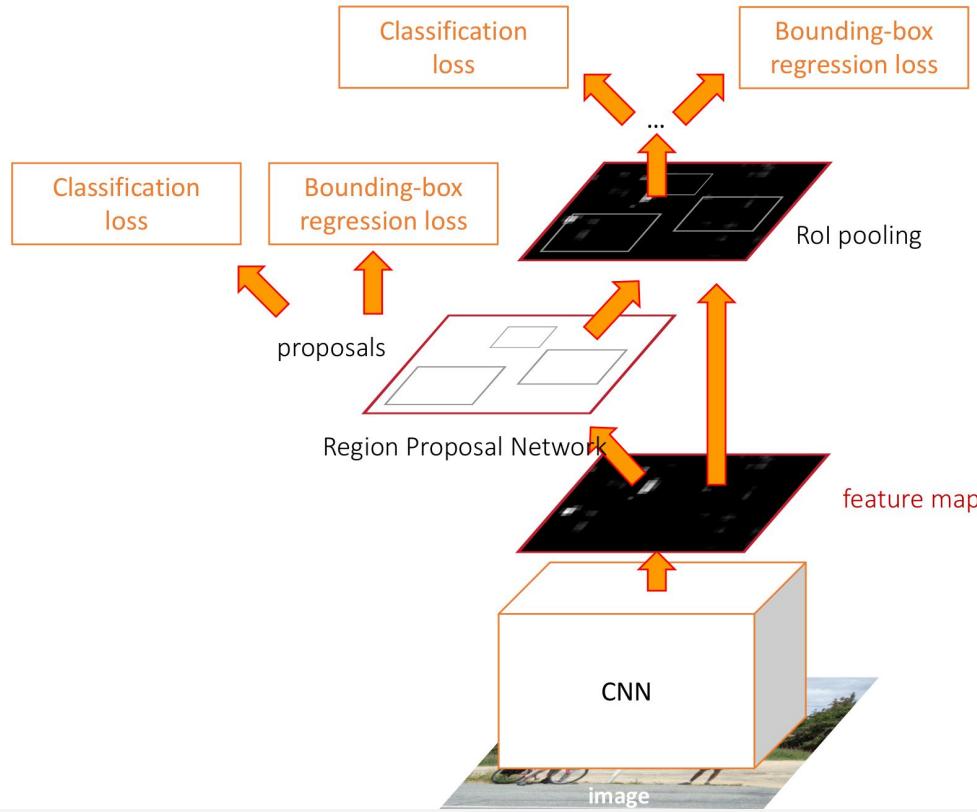
RCNN: Region proposals + CNN features



Fast RCNN



Faster RCNN



Credit: Girshick, He, Lazebnik

All the Details

Better Training

Online hard example mining

Network Variants

Feature Pyramid Network, RetinaNet

Ensemble

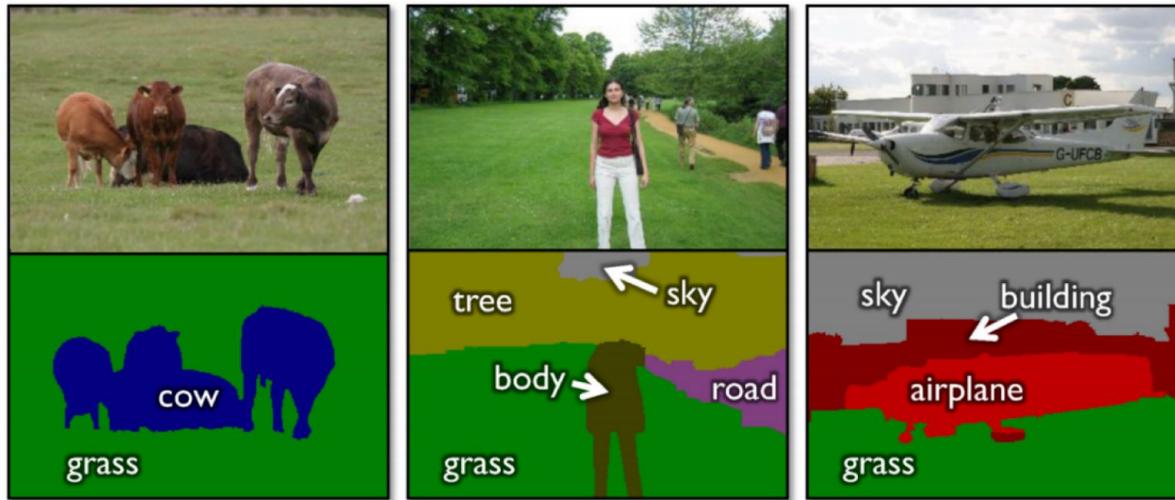
NMS

Remove duplicate detections

Today - Object Detection and Segmentation

- Overview
- Object Detection
- **Segmentation**
- Detection + Segmentation: Instance Segmentation
- What's Next

Task Formulation



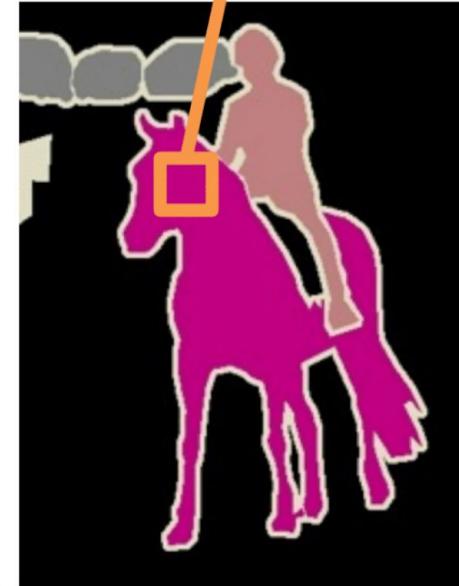
object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Evaluation Metric

Ground Truth: Horse



Prediction:
Horse



Pixel accuracy, Class accuracy, IOU

41
Slide credit: Abhinav Gupta

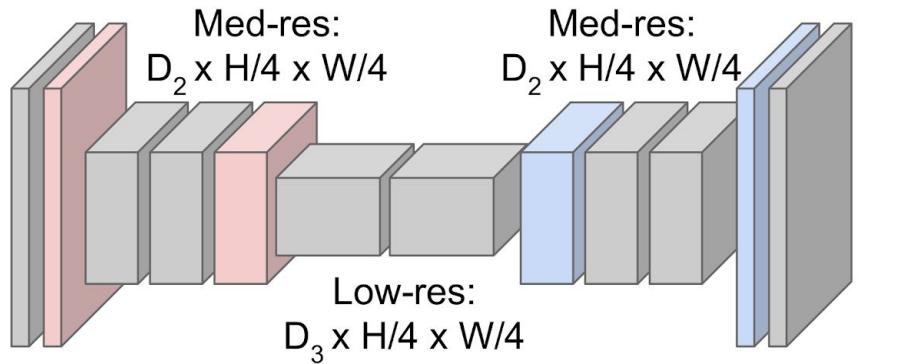
Network Structure for Dense Predictions

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Input:
 $3 \times H \times W$

High-res:
 $D_1 \times H/2 \times W/2$

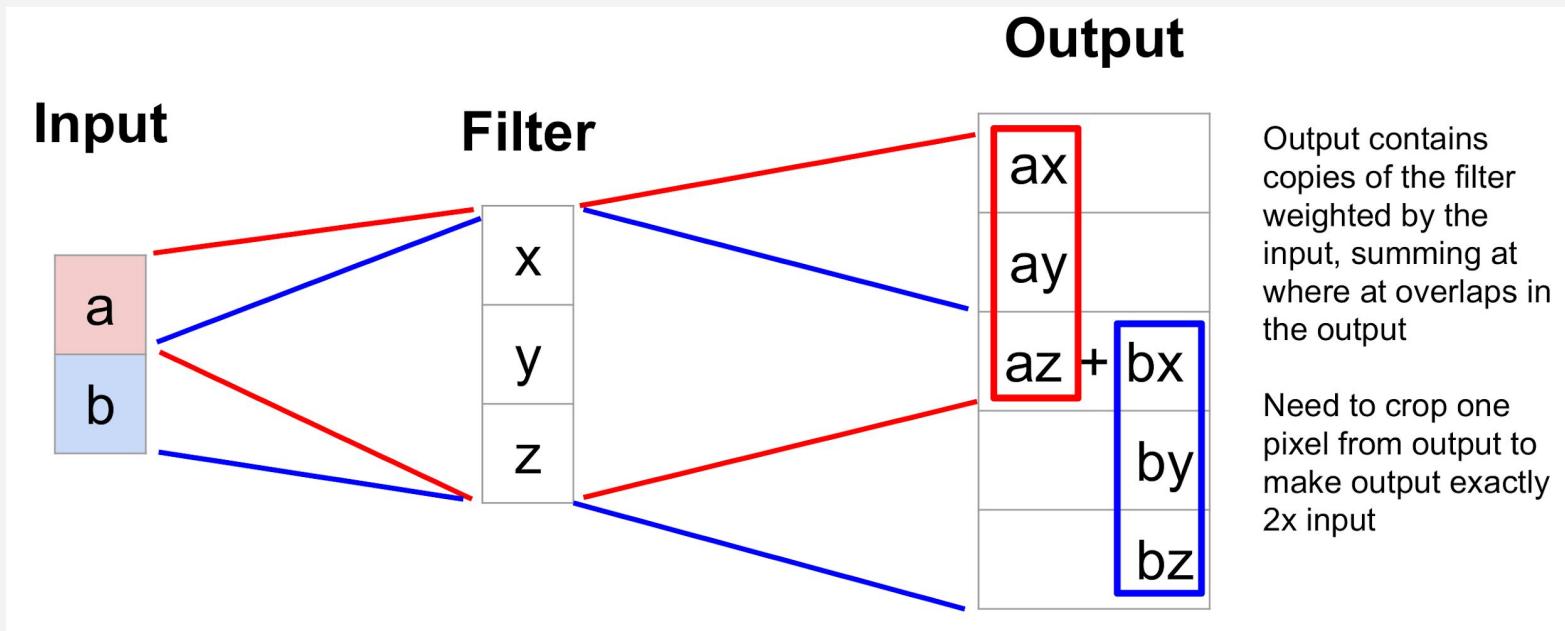


Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

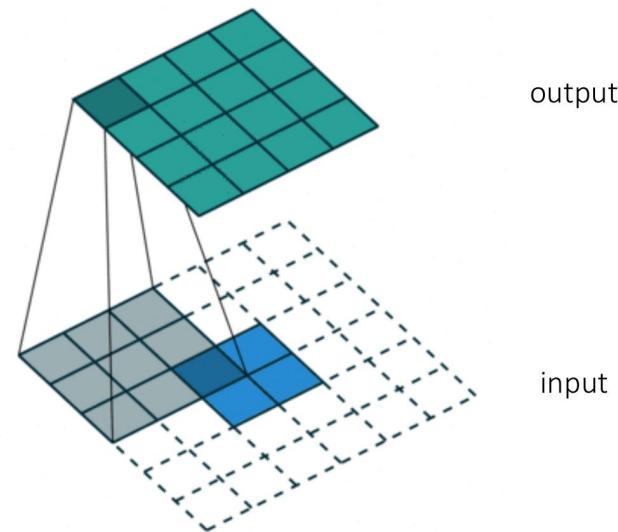
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Upconvolution



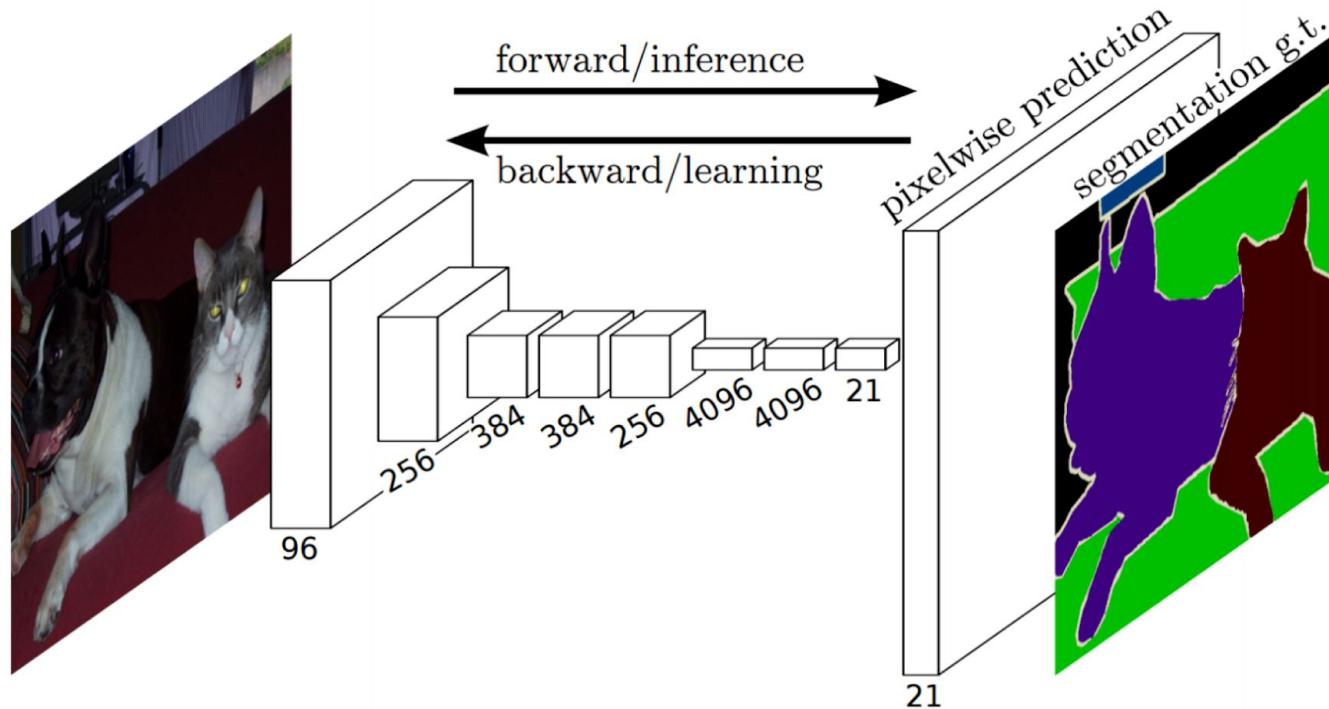
Upconvolution

Transposed convolution

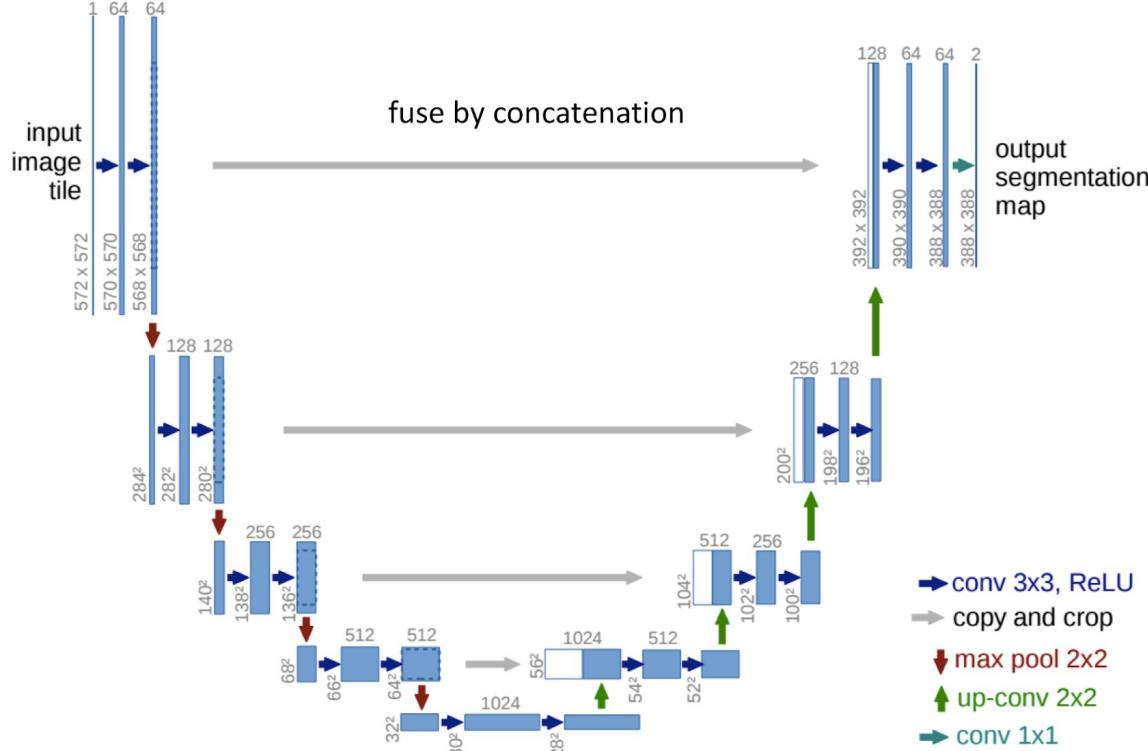


Credit: Lazebnik; V. Dumoulin and F. Visin, [A guide to convolution arithmetic for deep learning](#), arXiv 2018

Fully Convolutional Networks



UNet

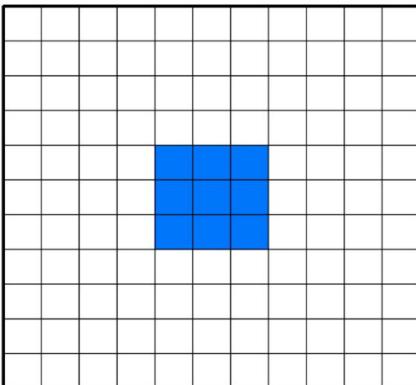


Credit: Lazebnik, Badrinarayanan et al., [SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation](#)

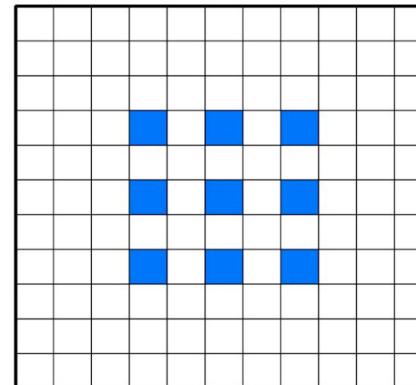
Dilated Convolutions

Instead of reducing spatial resolution of feature maps, use a large sparse filter

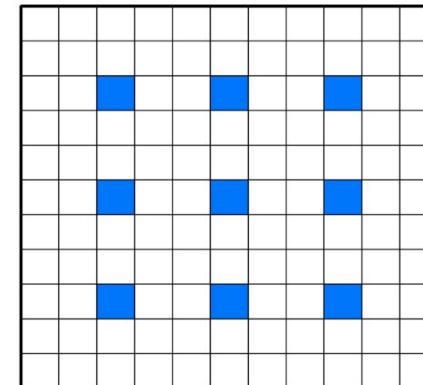
Dilation factor 1



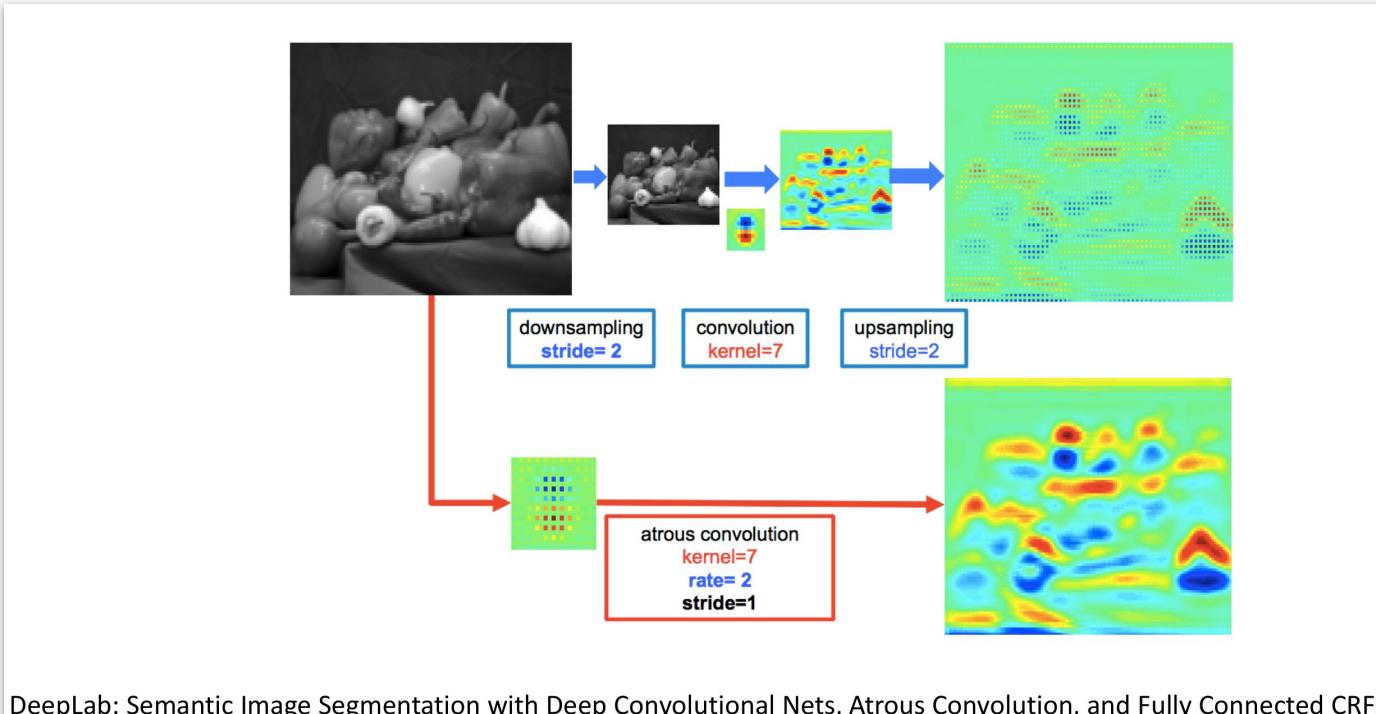
Dilation factor 2



Dilation factor 3



Deeplab



DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs

Today - Object Detection and Segmentation

- Overview
- Object Detection
- Segmentation
- **Detection + Segmentation: Instance Segmentation**
- What's Next

Task Formulation

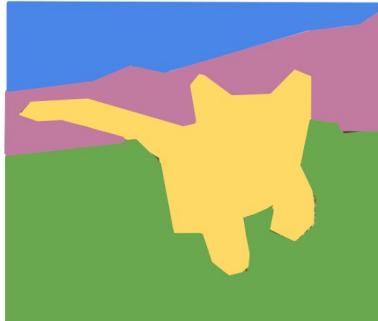
Classification



CAT

No spatial extent

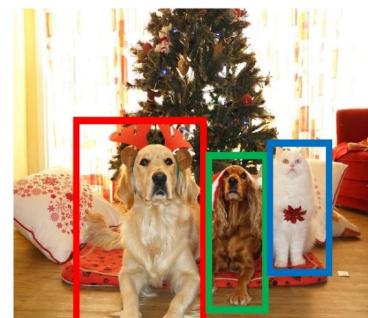
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

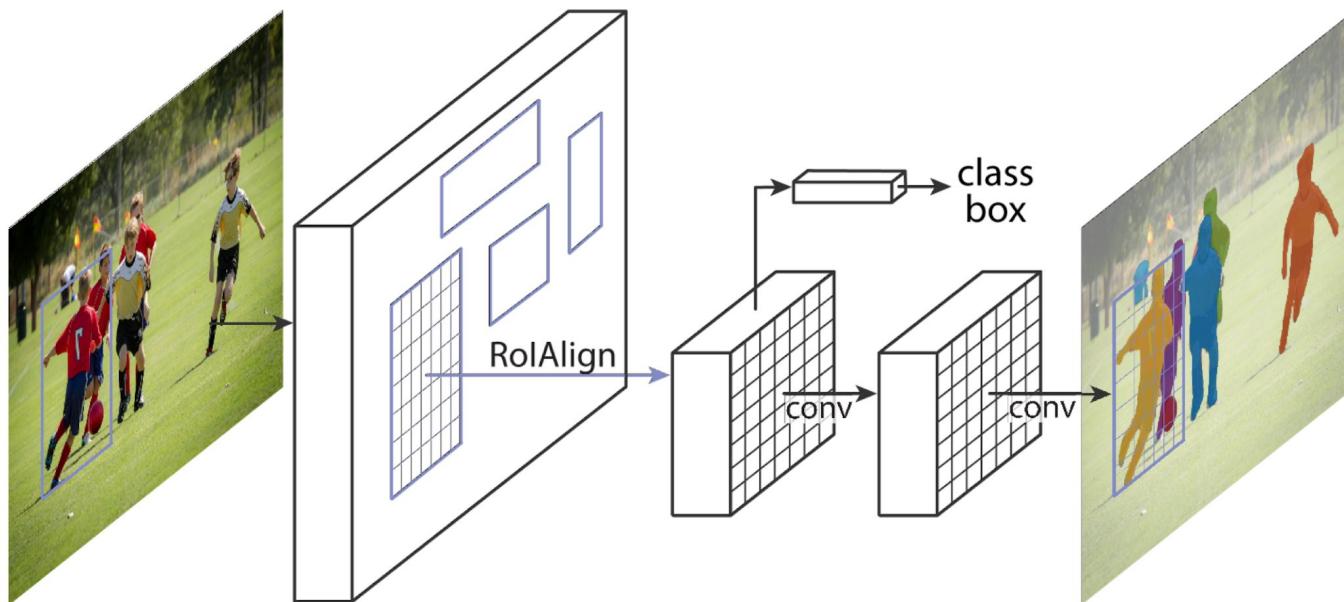
Instance Segmentation



DOG, DOG, CAT

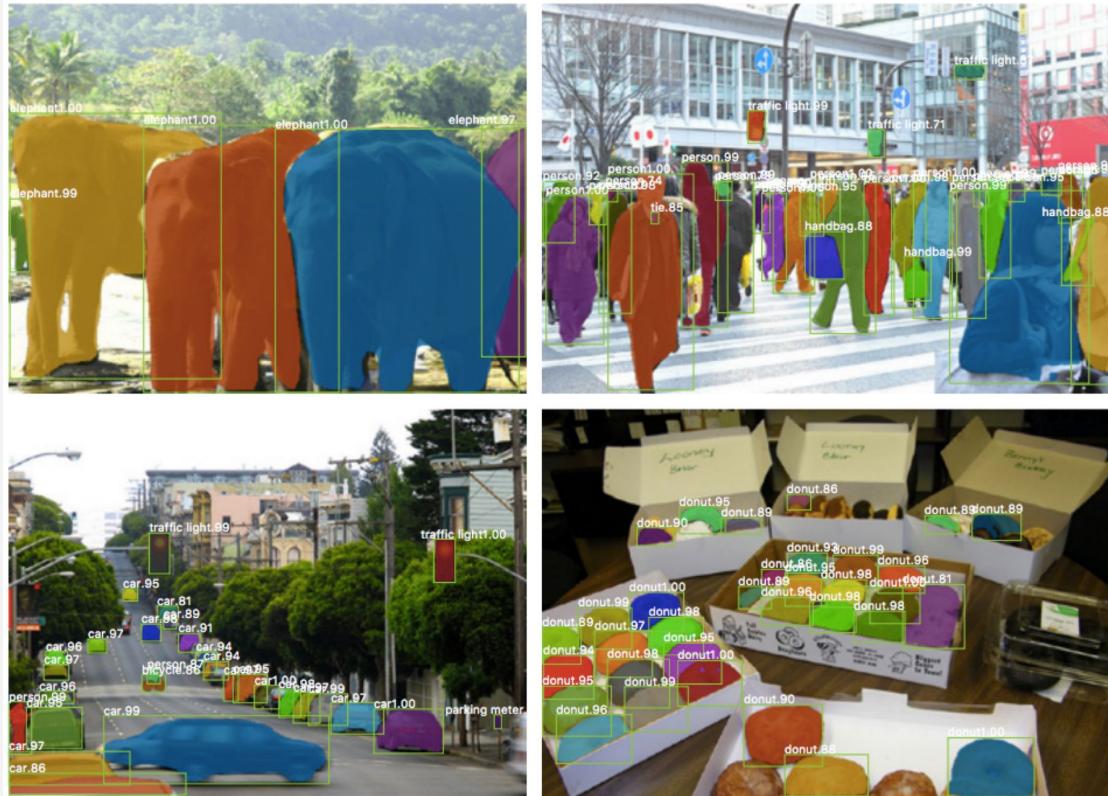
This image is CC0 public domain

Mask RCNN



Mask branch: separately predict segmentation for each possible class

Mask RCNN Results



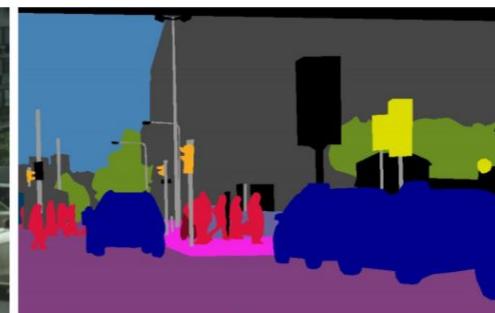
Today - Object Detection and Segmentation

- Overview
- Object Detection
- Segmentation
- Detection + Segmentation: Instance Segmentation
- **What's Next**

Panoptic Segmentation



(a) image



(b) semantic segmentation

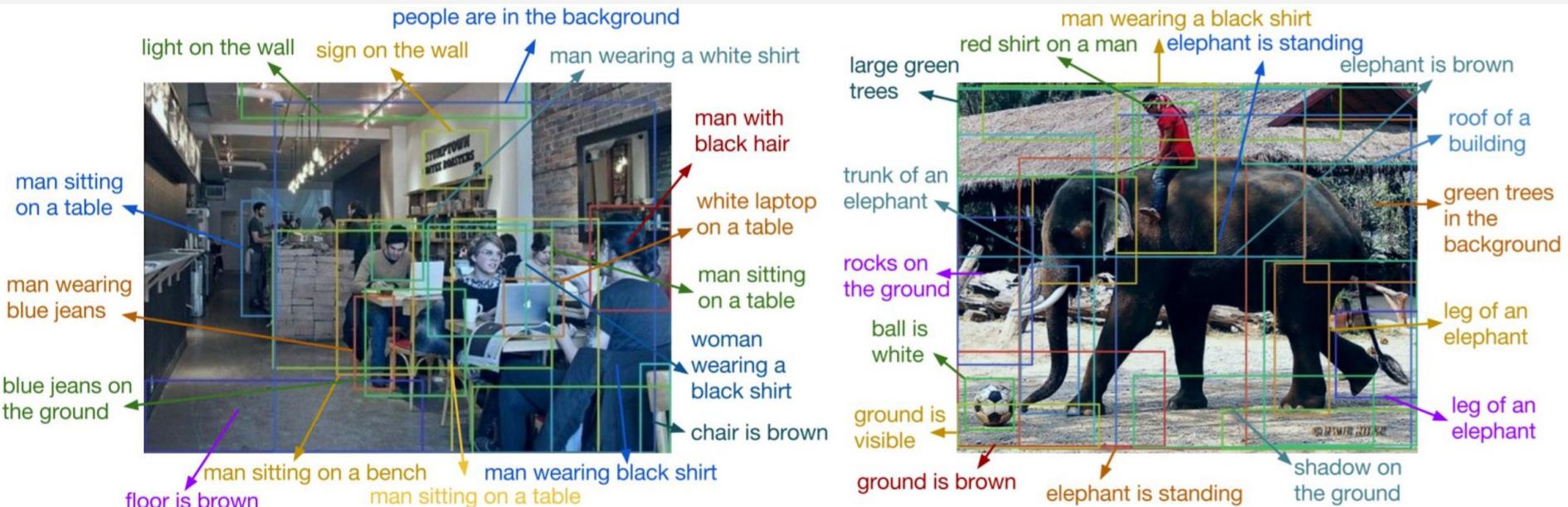


(c) instance segmentation



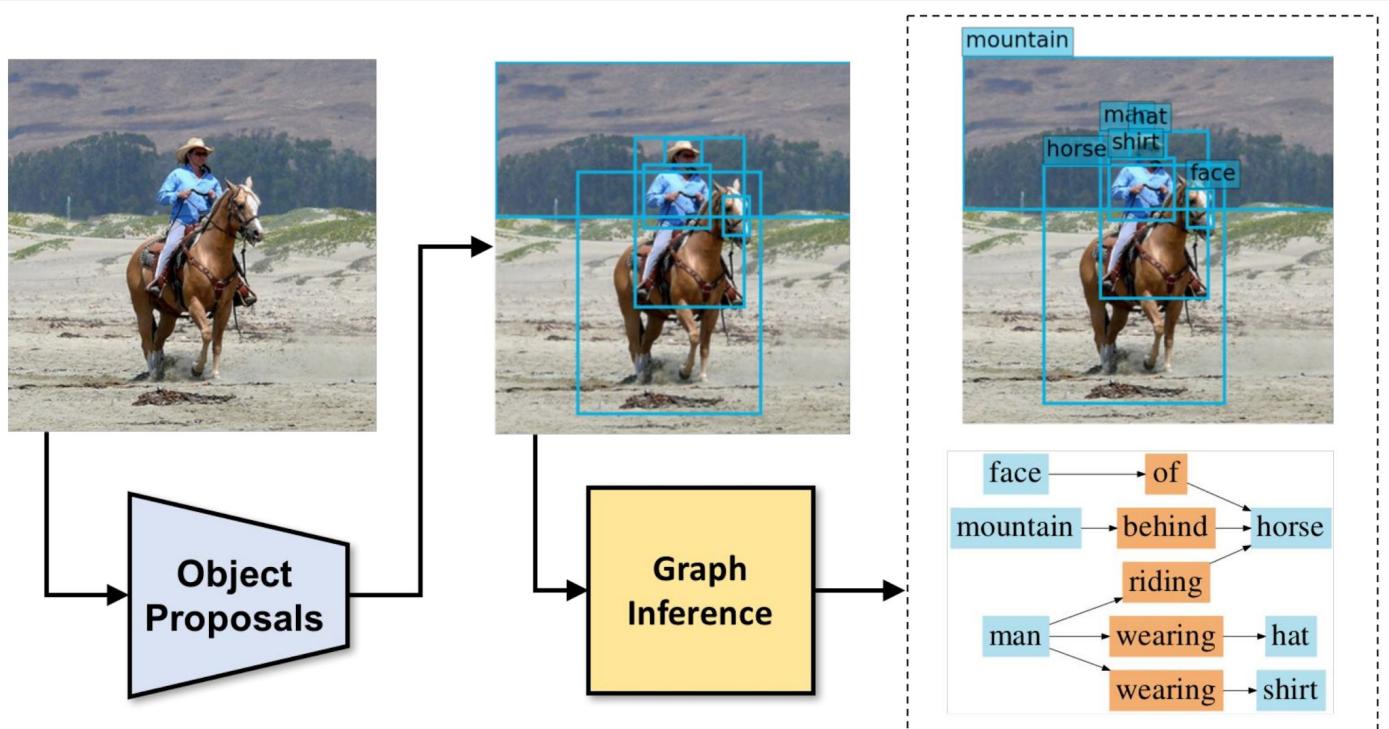
(d) panoptic segmentation

Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016
Figure copyright IEEE, 2016. Reproduced for educational purposes.

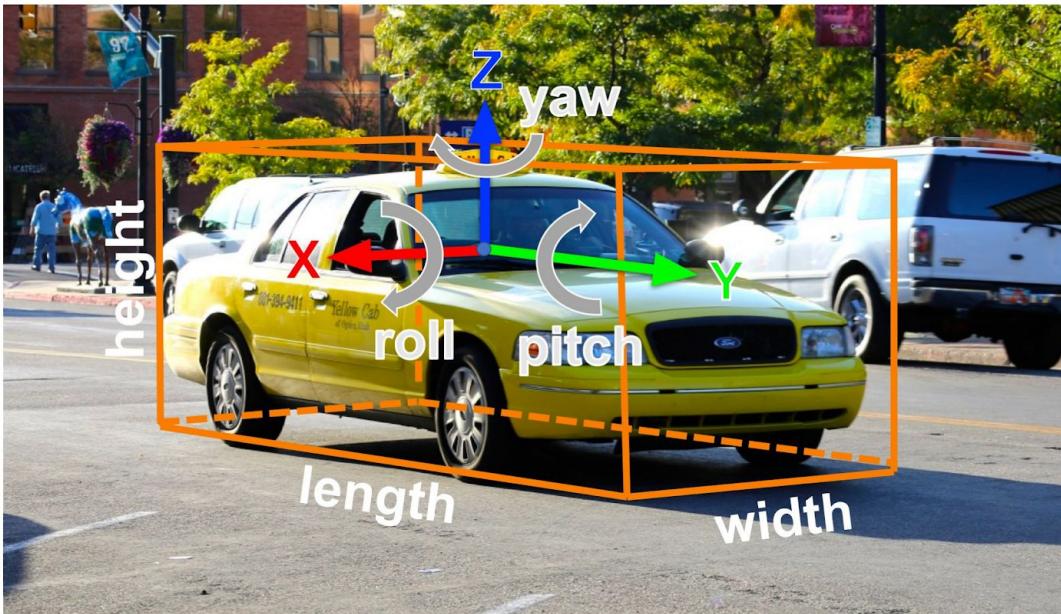
Scene Graph



Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017
Figure copyright IEEE, 2018. Reproduced for educational purposes.

Slides from Stanford cs231n

3D Detection



2D Object Detection:

2D bounding box

(x, y, w, h)

3D Object Detection:

3D oriented bounding box

$(x, y, z, w, h, l, r, p, y)$

Simplified bbox: no roll & pitch

Much harder problem than 2D object detection!

This image is CC0 public domain

Summary - Object Detection and Segmentation

- Overview
- Object Detection
- Segmentation
- Detection + Segmentation: Instance Segmentation
- What's Next