# Segmentation and Detection

CSCI 5299, 4/13/2020
Guest Lecture, Michelina Pallone (Google)
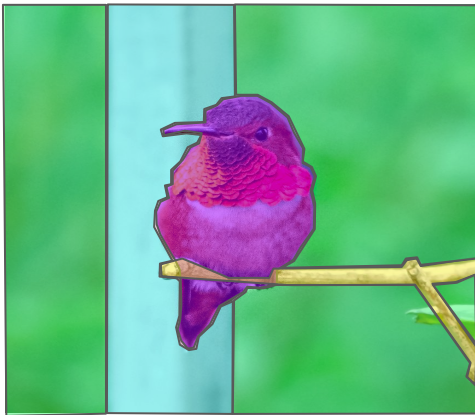
# Classification

- AlexNet
- GoogLeNet
- VGG
- ResNet

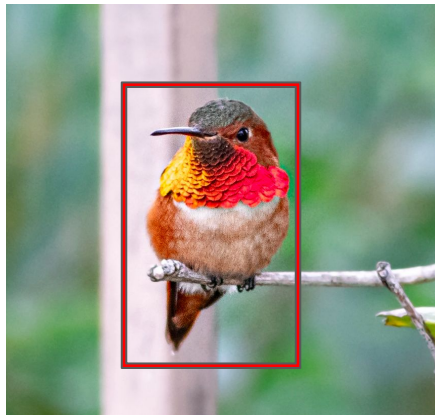## Scores

- Hummingbird 0.8
- Robin 0.1
- Apple 0.05

# Semantic Segmentation



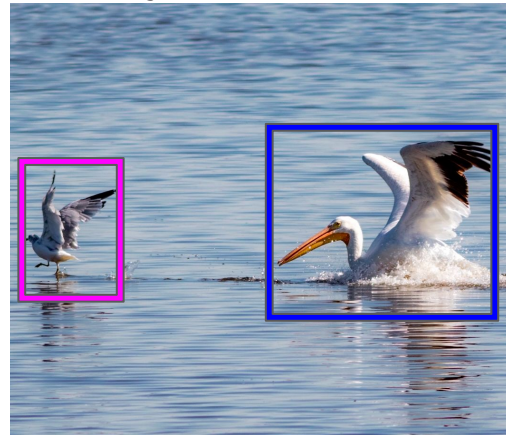hummingbird, post, branch, leaves

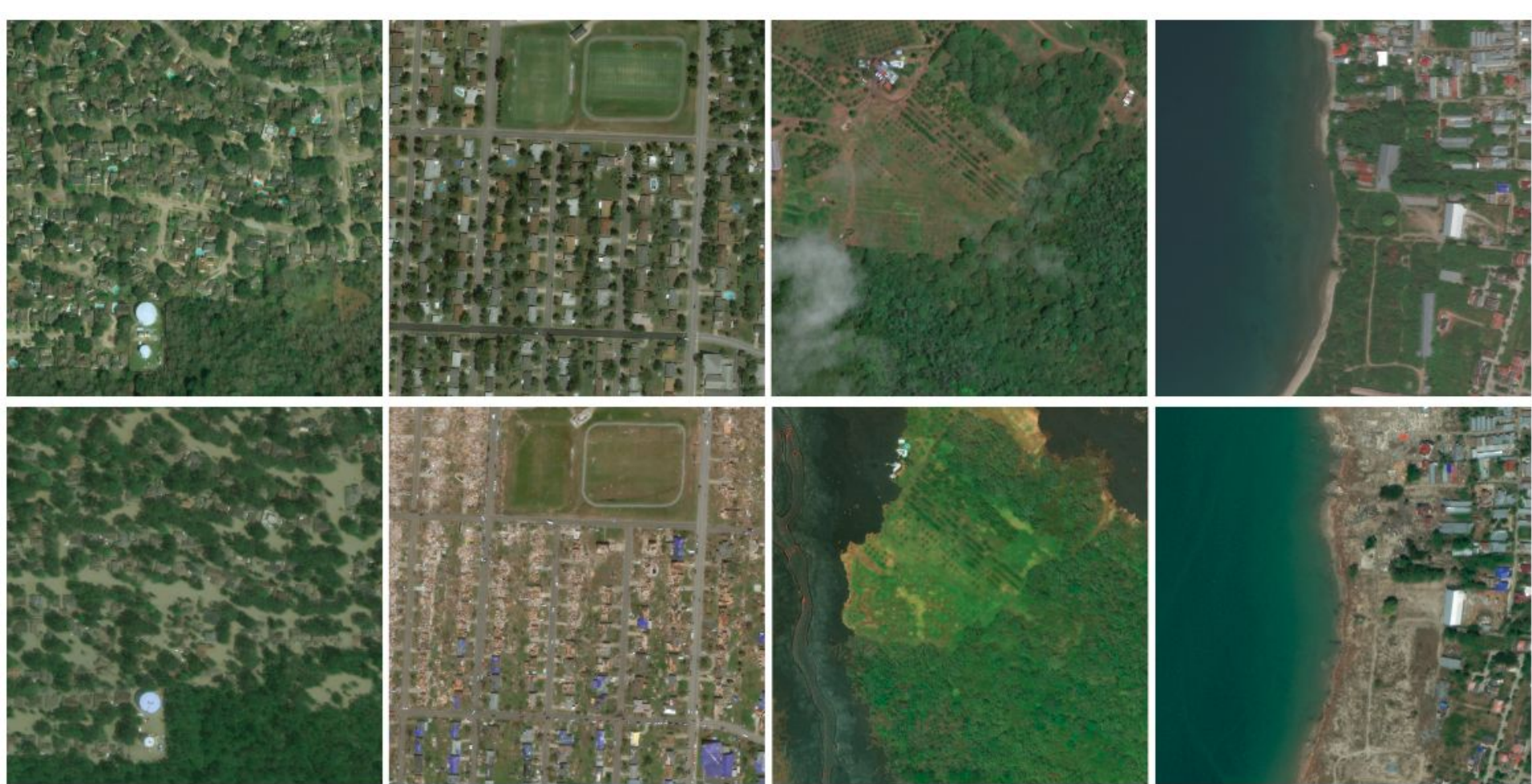Pixelwise labels, no objects

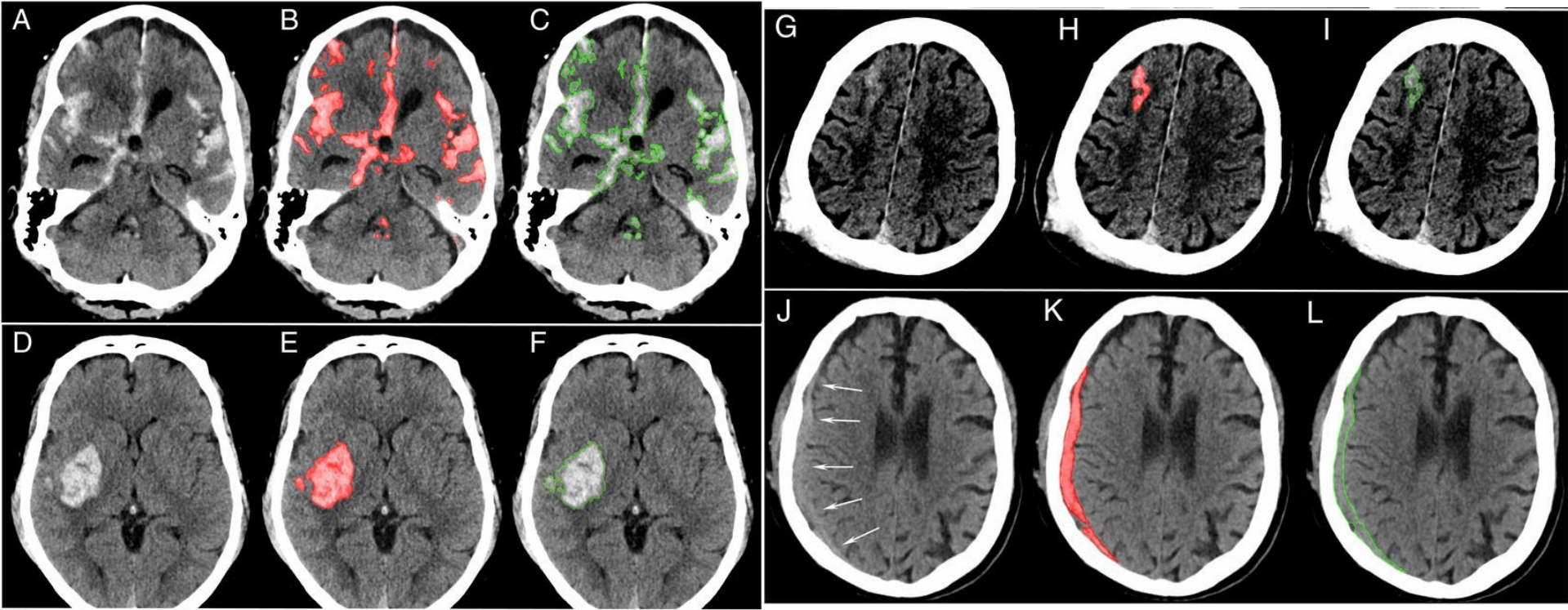# Localization



Hummingbird

Single object

# Object Detection



gull, pelican

Multi-object

[Building Damage Assessment](#)

Cityscapes data set
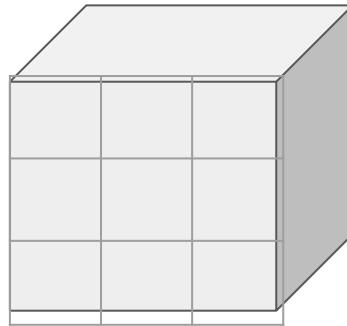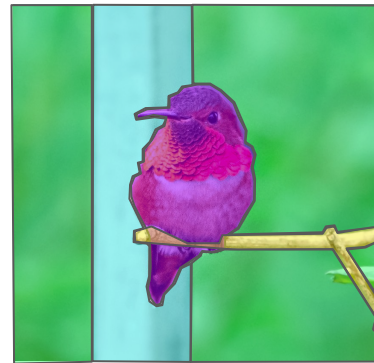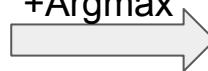
# Segmentation



Deep CNN

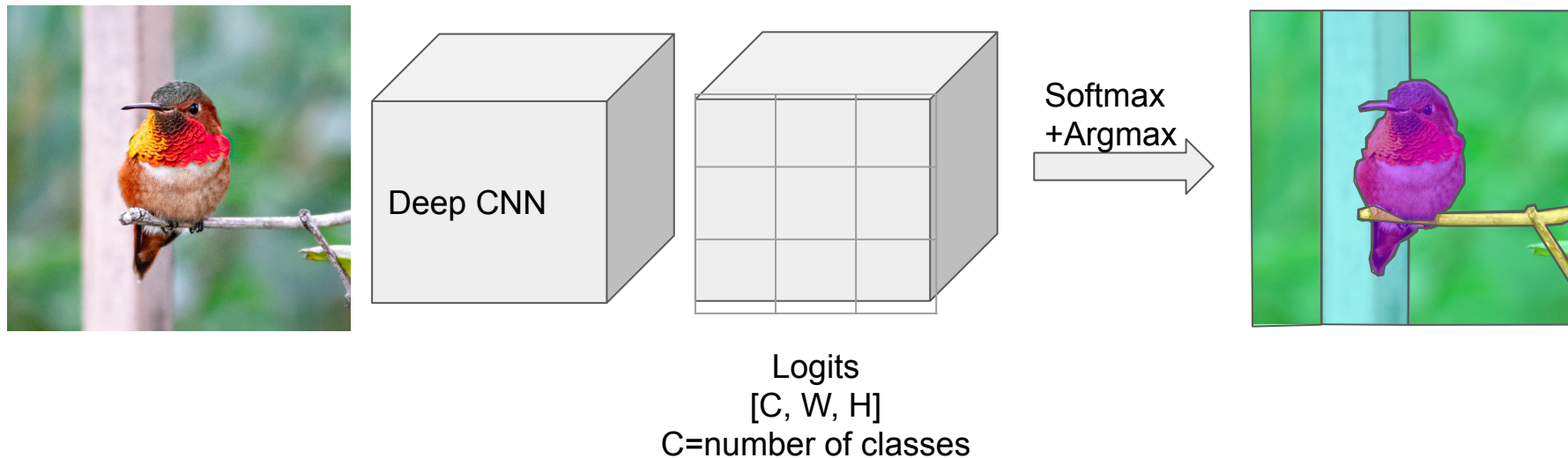Logits
[C, W, H]
C=number of classes

Softmax
+Argmax

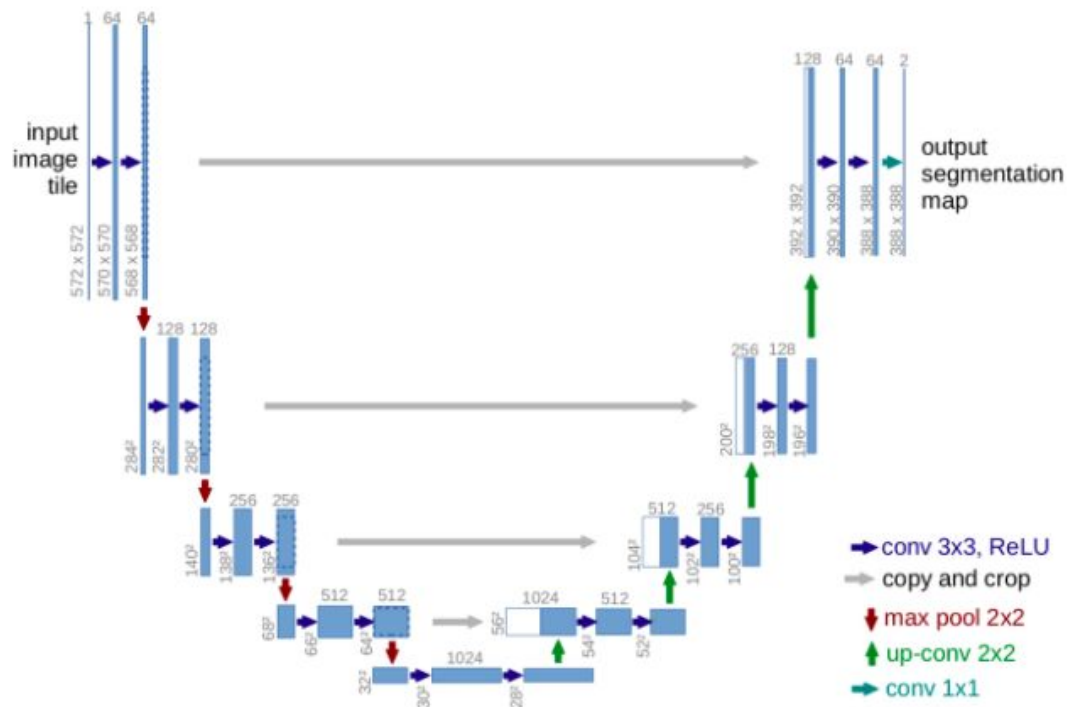# Segmentation



Logits
[C, W, H]
C=number of classes

Softmax +Argmax

- Deep CNN includes downsampling to reduce computational costs.
  Then upsample back to the original size.

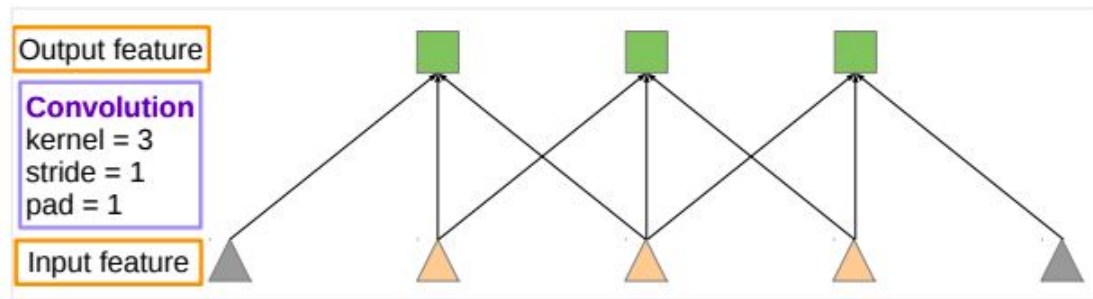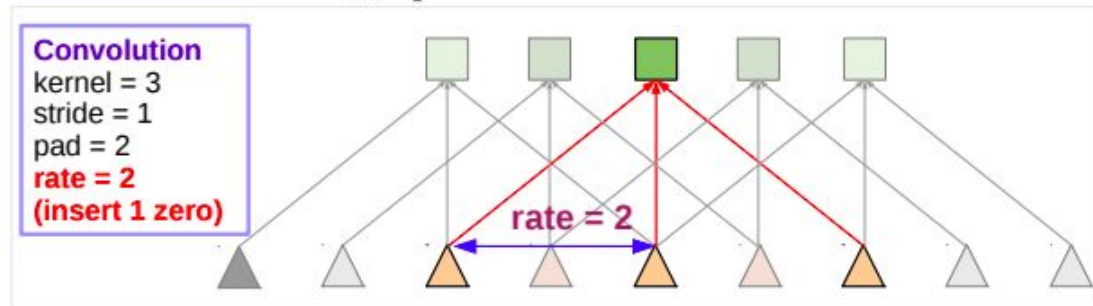- Sum (or average) cross entropy loss on each pixel in the output

# (2015) [U-Net](U-Net)



- Influential in medical imaging

- Skip connections all low level features to propagate up to higher level features

# (2016) DeepLab



Output feature

**Convolution**
kernel = 3
stride = 1
pad = 1

Input feature

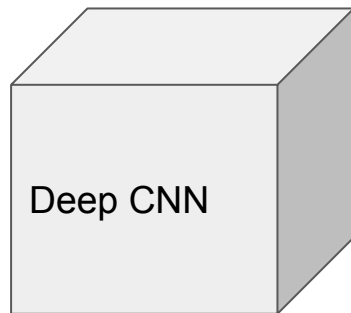(a) Sparse feature extraction

**Convolution**
kernel = 3
stride = 1
pad = 2
rate = 2
(insert 1 zero)

rate = 2

(b) Dense feature extraction

- DeepLab, DeepLab V3
- Atrous (Dilated) Convolutions
- Conditional Random Field (CRF)

# Localization



**Class Scores**
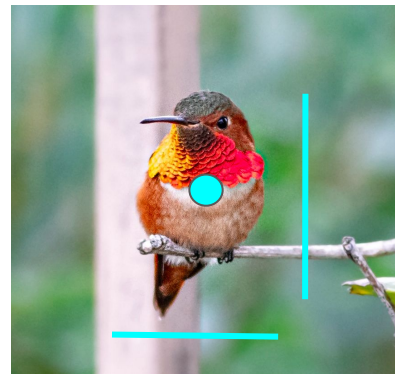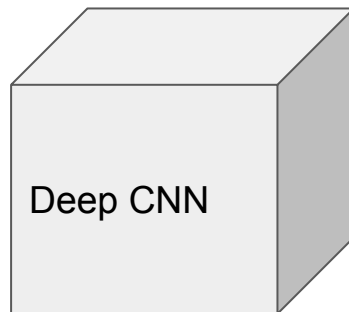Class
Confidence

**Box Coordinates**
X
Y
W
H

Humminging Bird
0.7

# Localization


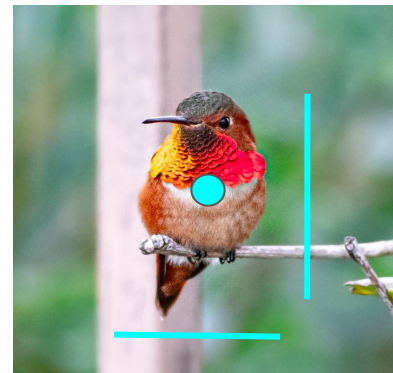
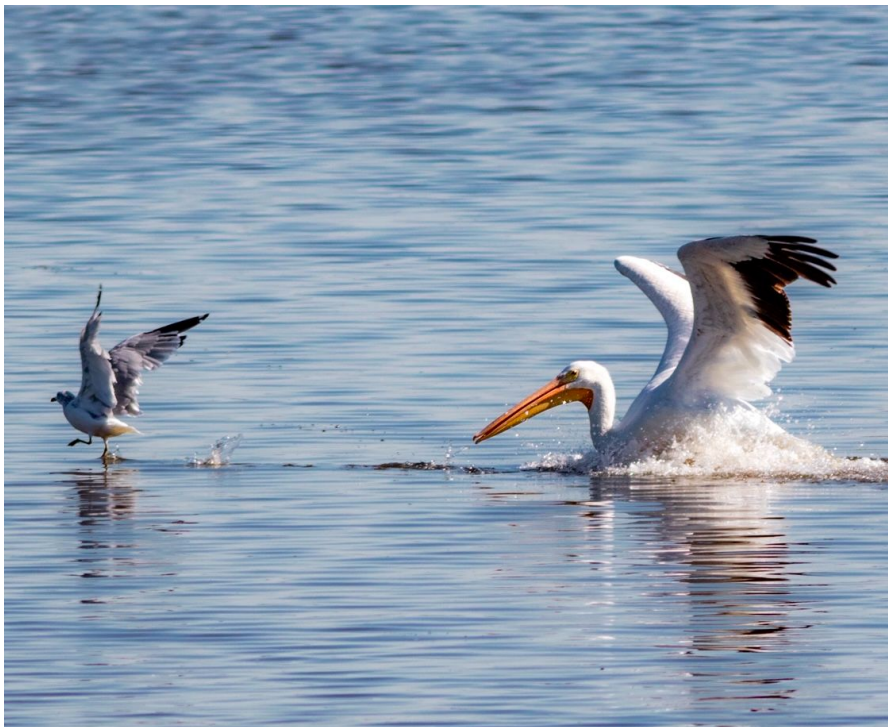**Class Scores**
Class
Confidence

**Box Coordinates**
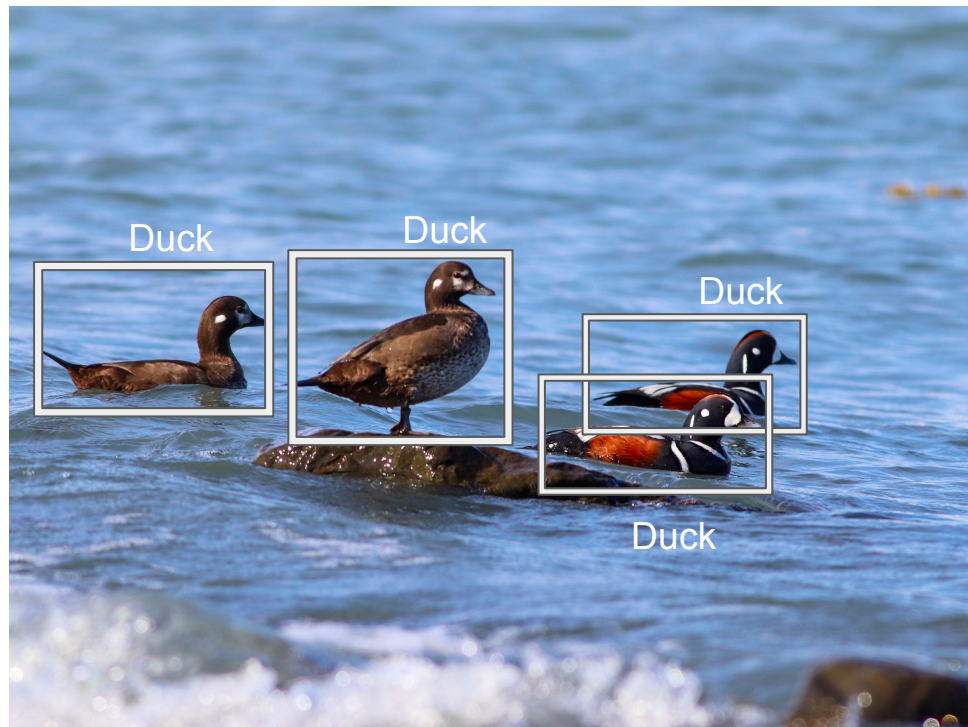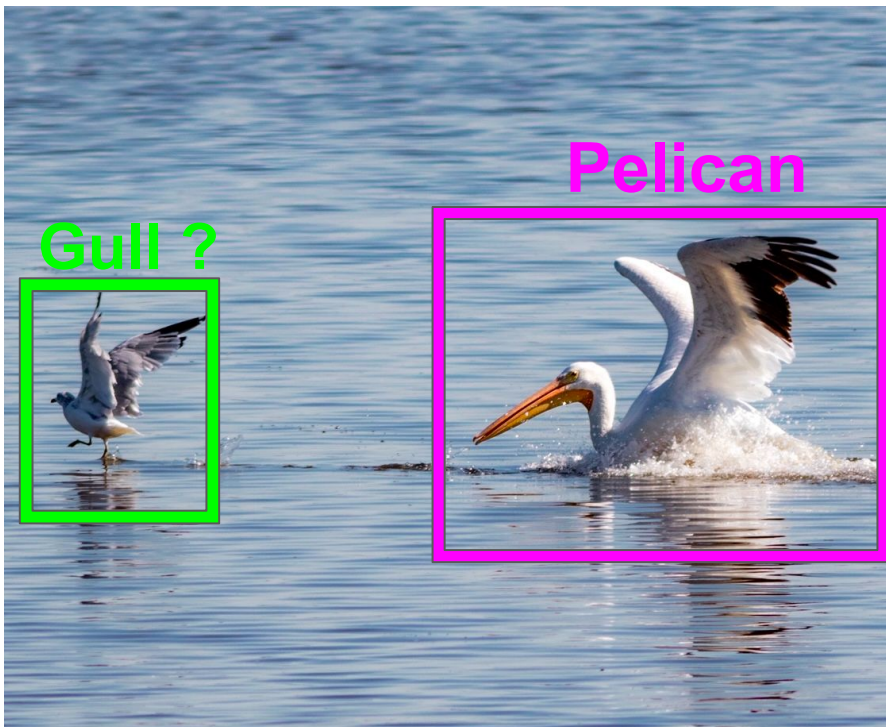X
Y
W
H

Humminging Bird
0.7

- Two losses: same softmax loss for classes and scores that you would use in image classification, some kind of regression loss for box coordinates.
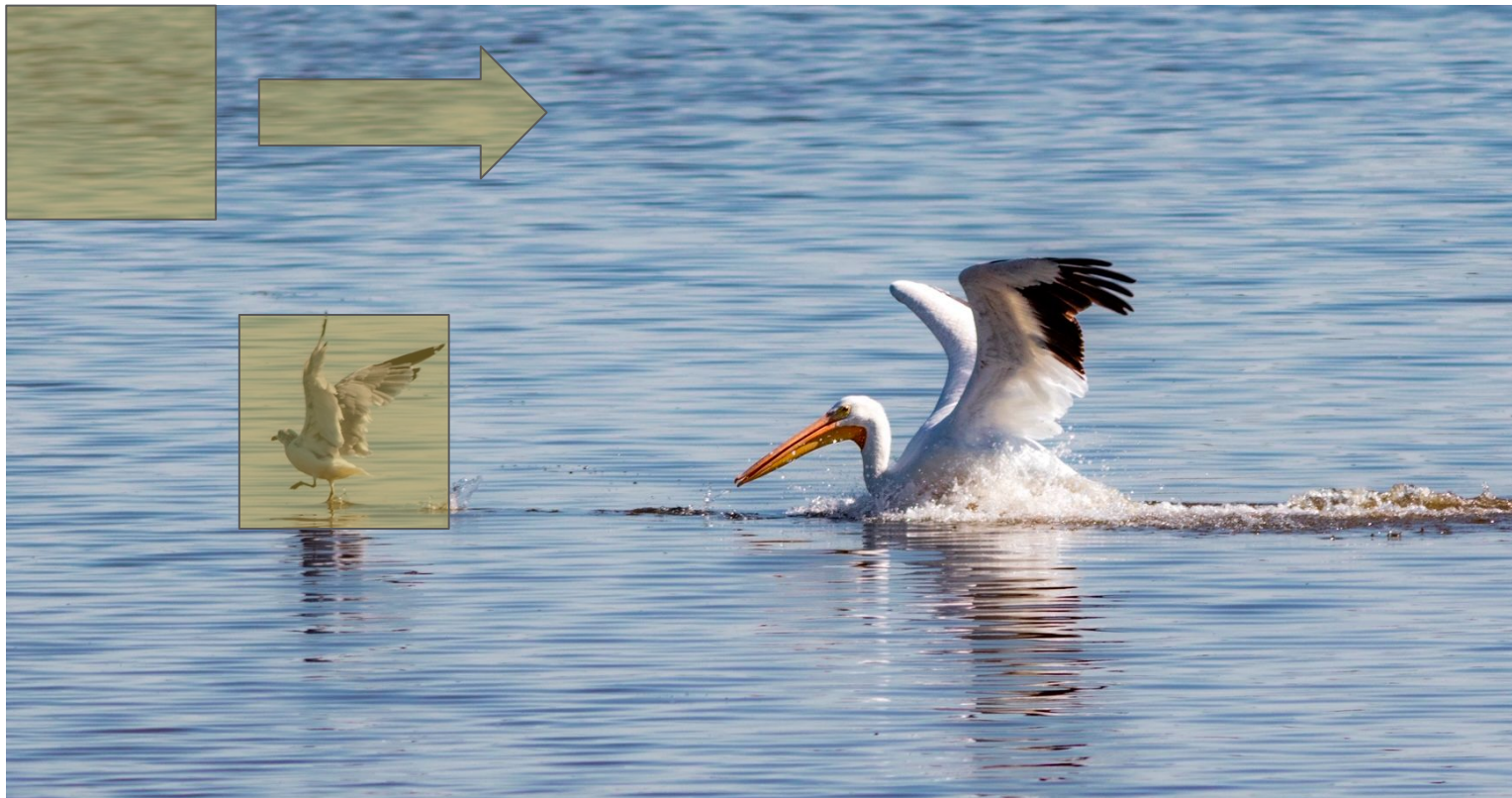
# Detection

# Detection

# Detection

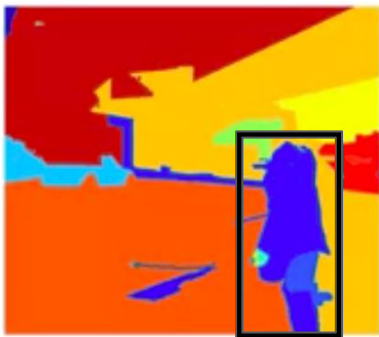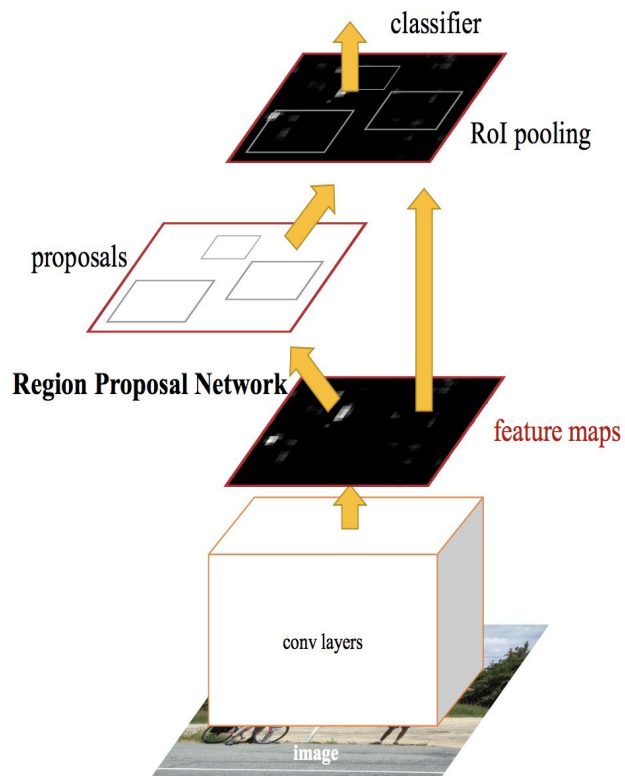# 2 Stage Detection: Regions with CNN (R-CNN)



- Narrow down candidate regions via semantic segmentation (Region Proposal Network)

- Train and inference times are slow
  - R-CNN (2014)
  - Fast R-CNN (2015)
  - Faster R-CNN (2015)

# **One Stage Detection**

- YOLO (You only look once 2016)
- SSD (Single Shot Detection 2016)
- Only one pass through the network
  - Simple
  - Fast
- Accuracy is ok

# **One Stage Detection**

- [YOLO](#) (You only look once)
- [SSD](#) (Single Shot Detection)
- Only one pass through the network
  - Simple
  - Fast
- Accuracy is ok

# One Stage Detection
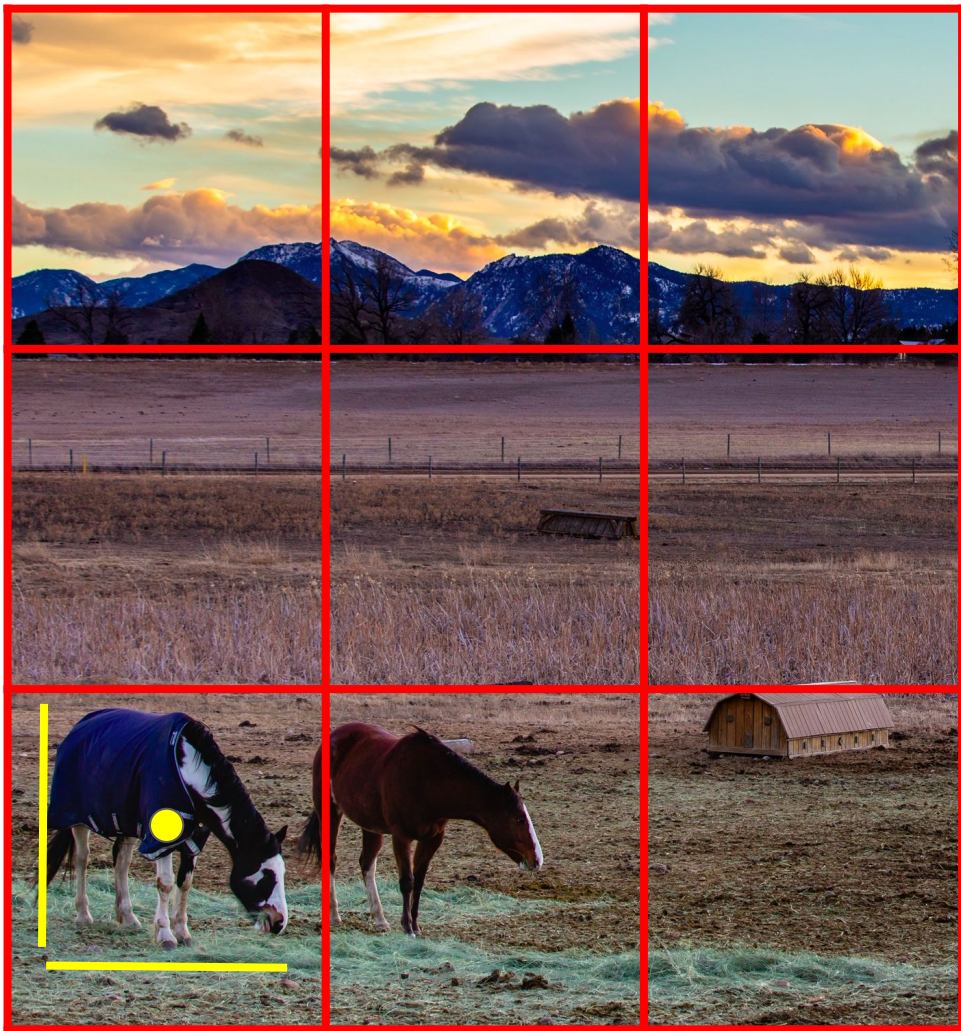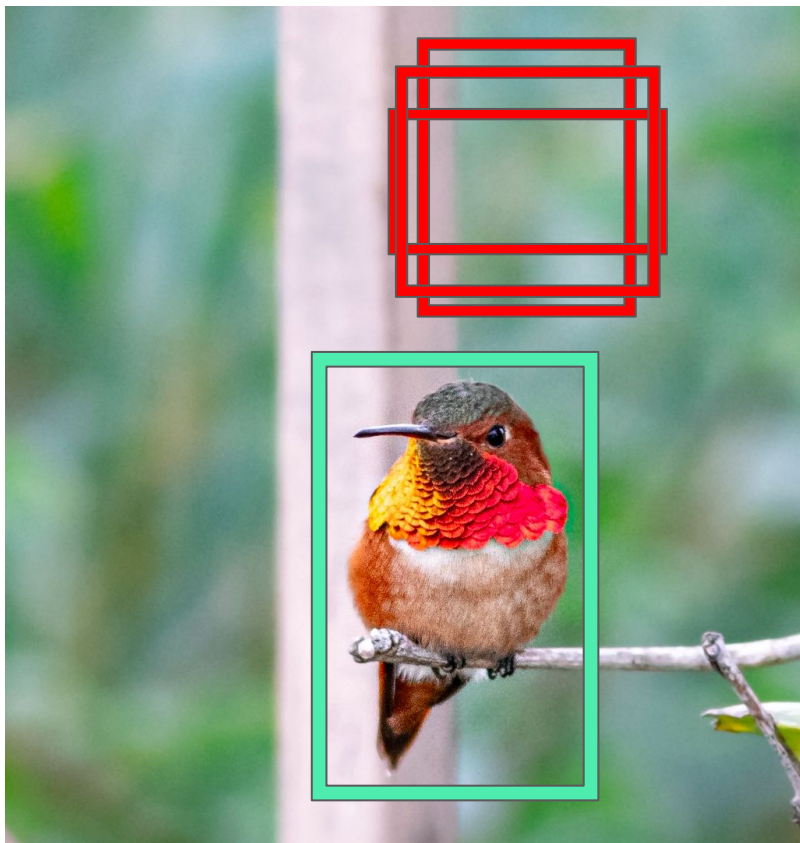
- YOLO (You only look once)
- SSD (Single Shot Detection)
- Only one pass through the network
  - Simple
  - Fast
- Accuracy is ok

# [RetinaNet](RetinaNet) Return of 1 stage



- Traditional problem with one-stage detectors: An extreme class imbalance causes "easy" background samples to overwhelm the loss function

- **Focal Loss**: Keep the loss from easy background samples from overwhelming the loss from sparse hard samples
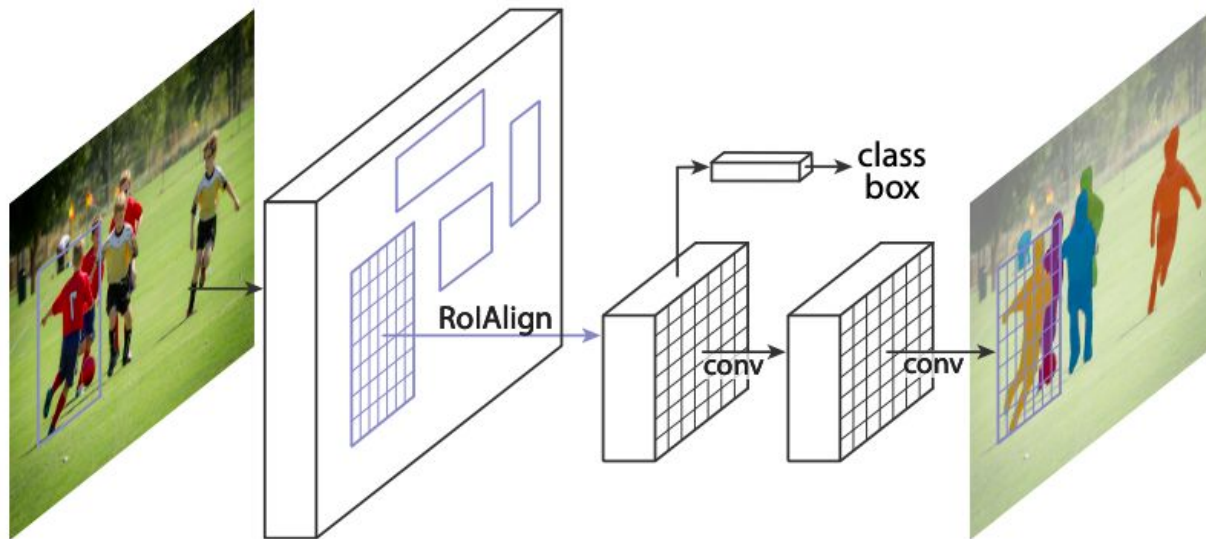
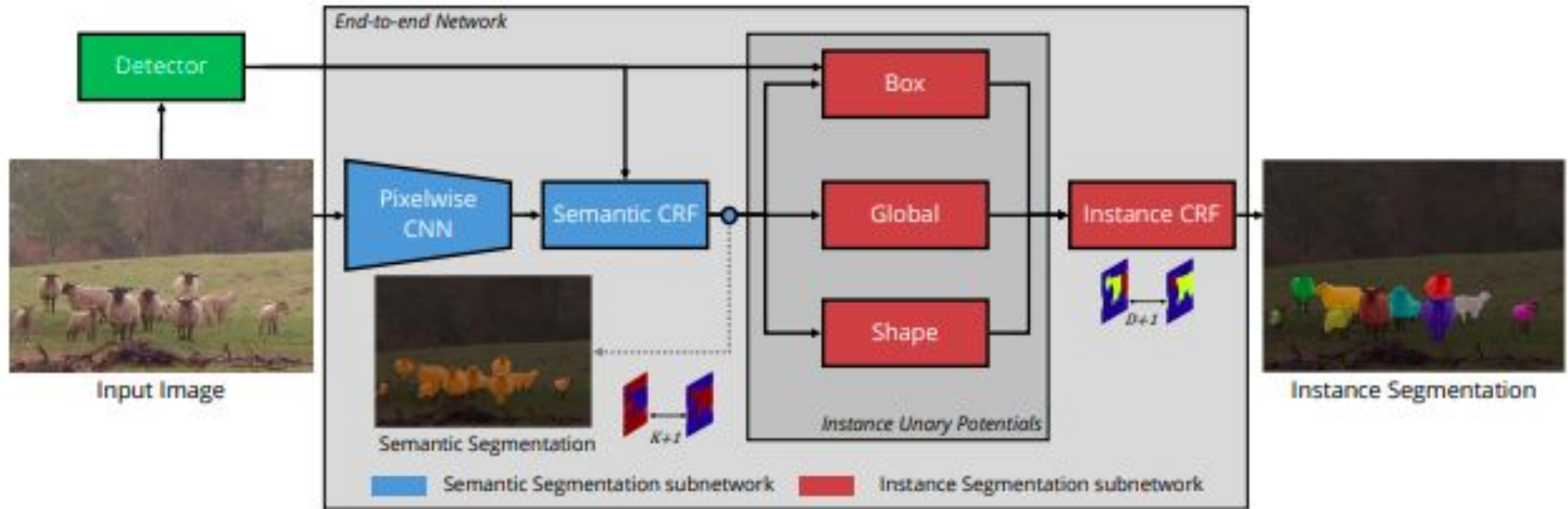# Instance segmentation

# Instance segmentation

# Detection-based approaches ("Top Down")



- [Mask R-CNN](#)
- [ShapeMask: Learning to Segment Novel Objects by Refining Shape Priors](#)

# Grouping-based approaches ("Bottom up")



- [Pixelwise Instance Segmentation with a Dynamically Instantiated Network](#) (above)
- [Semantic Instance Segmentation with a Discriminative Loss Function](#)
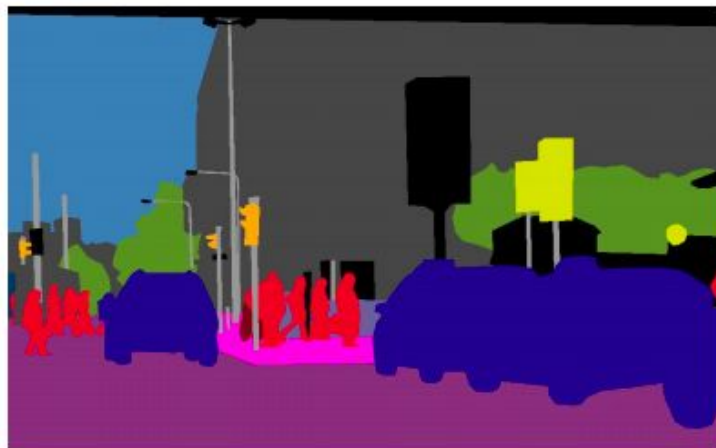- [InstanceCut: from Edges to Instances with MultiCut](#)

# Panoptic segmentation



Water?

Rocks?

- **Things**: Objects that come in discrete, countable instances
  - Well suited to instance segmentation
  - Examples: Cars, people, animals

- **Stuff**: Amorphous background regions
  - Better suited to semantic segmentation
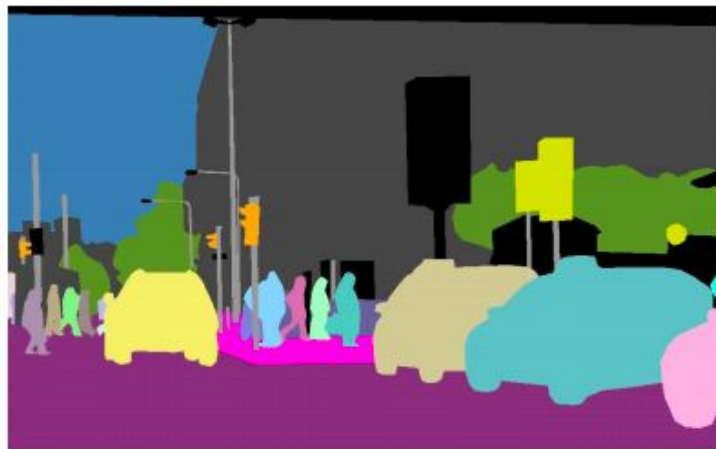  - Examples: Sky, grass, water

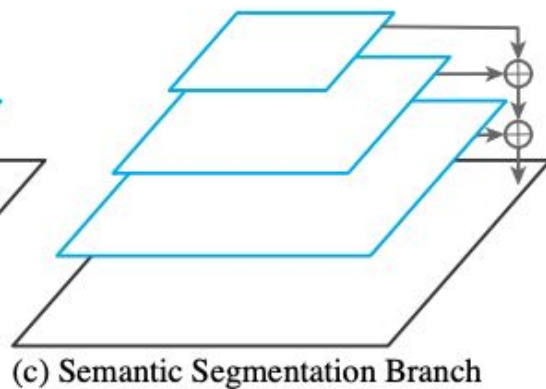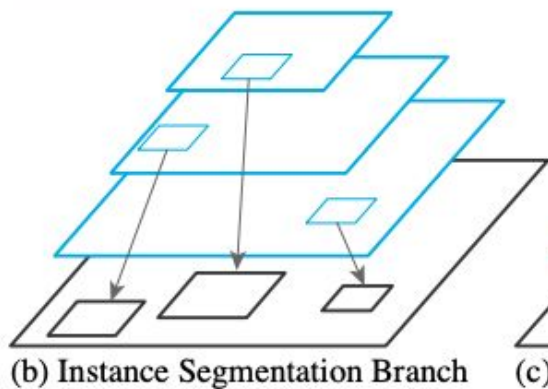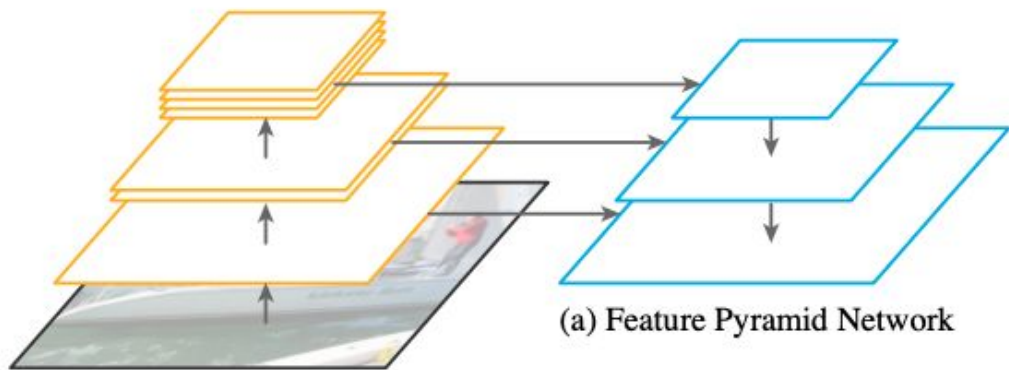(a) image

(b) semantic segmentation

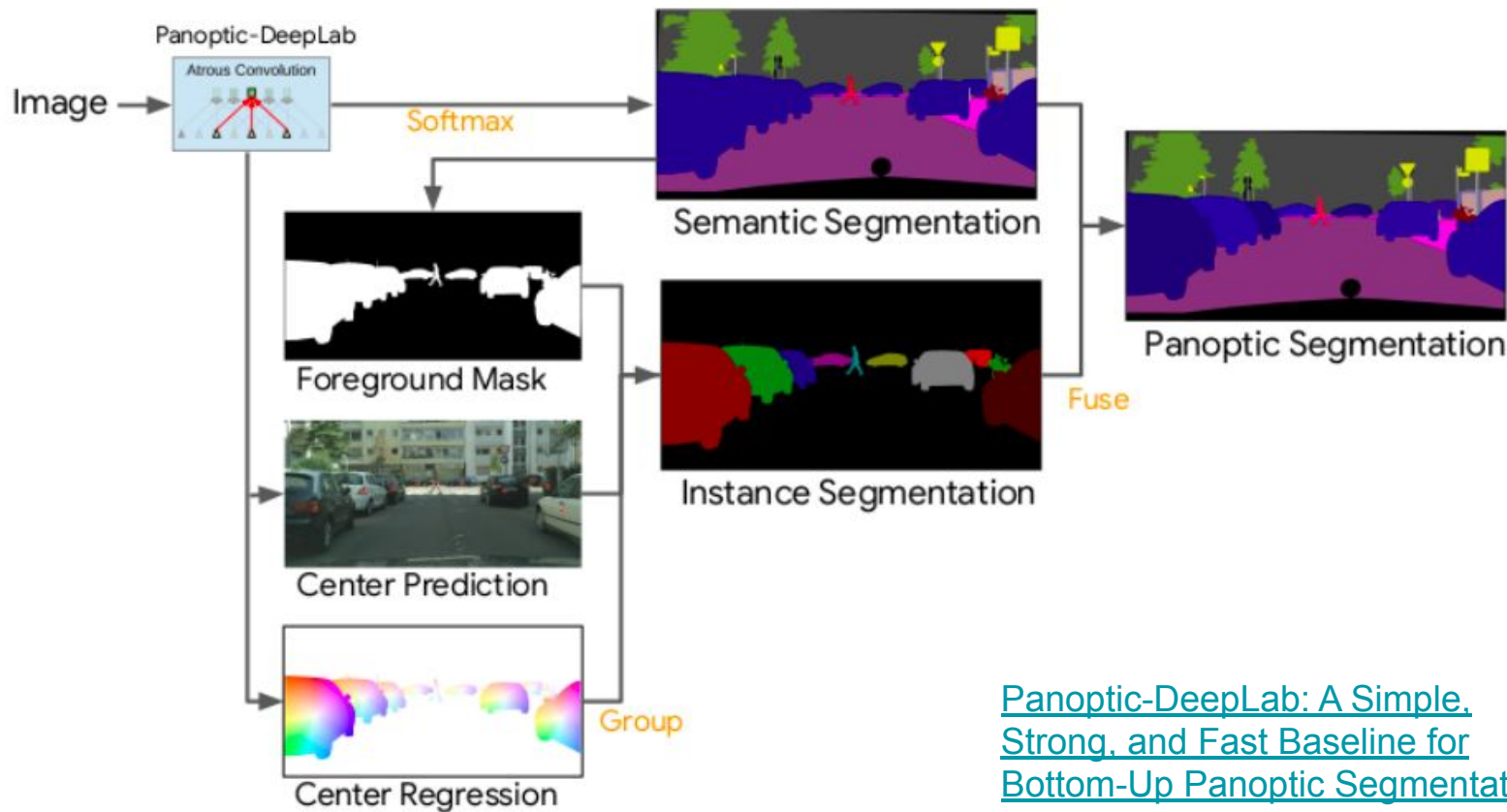(c) instance segmentation

(d) panoptic segmentation

# Top Down



(a) Feature Pyramid Network

(b) Instance Segmentation Branch

(c) Semantic Segmentation Branch

[Panoptic Feature Pyramid Networks](#)

# Bottom Up



Panoptic-DeepLab

Image → Atrous Convolution

Softmax

Semantic Segmentation

Foreground Mask

Center Prediction

Center Regression

Group

Instance Segmentation

Fuse

Panoptic Segmentation
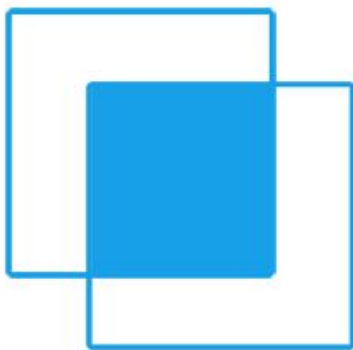
Panoptic-DeepLab: A Simple, Strong, and Fast Baseline for Bottom-Up Panoptic Segmentation

# Intersection over Union



$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

- For **stuff** IoU is a standard evaluation metric

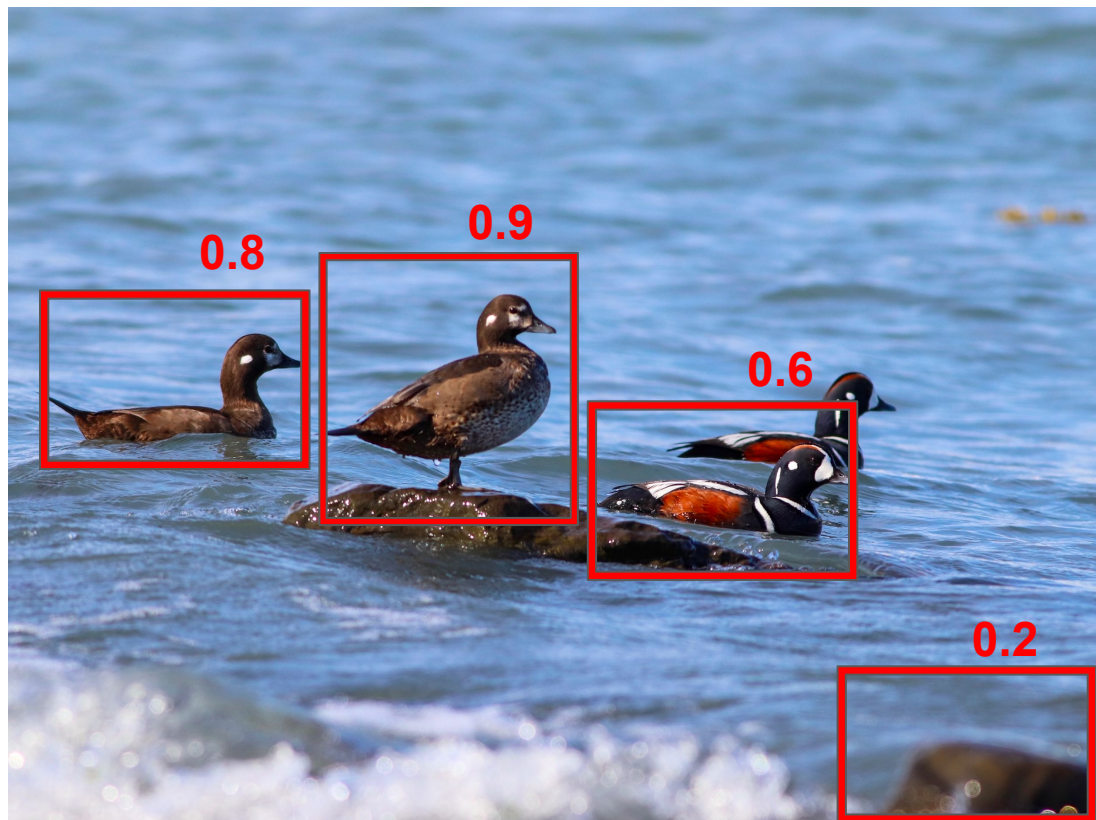- For **things,** you decide what IoU is a true positive (IoU = 0.5 is common)

# Actual Values

|  | Positive | Negative |
|---|---|---|
| **Predicted Values** Positive | True Positive (TP) | False Positive (FP, Type I Error) |
| Negative | False Negative (FN, Type II Error) | True Negative (TN) |

$$\text{Precision} = \frac{\# \text{ TP}}{\# \text{ TP} + \# \text{ FP}}$$
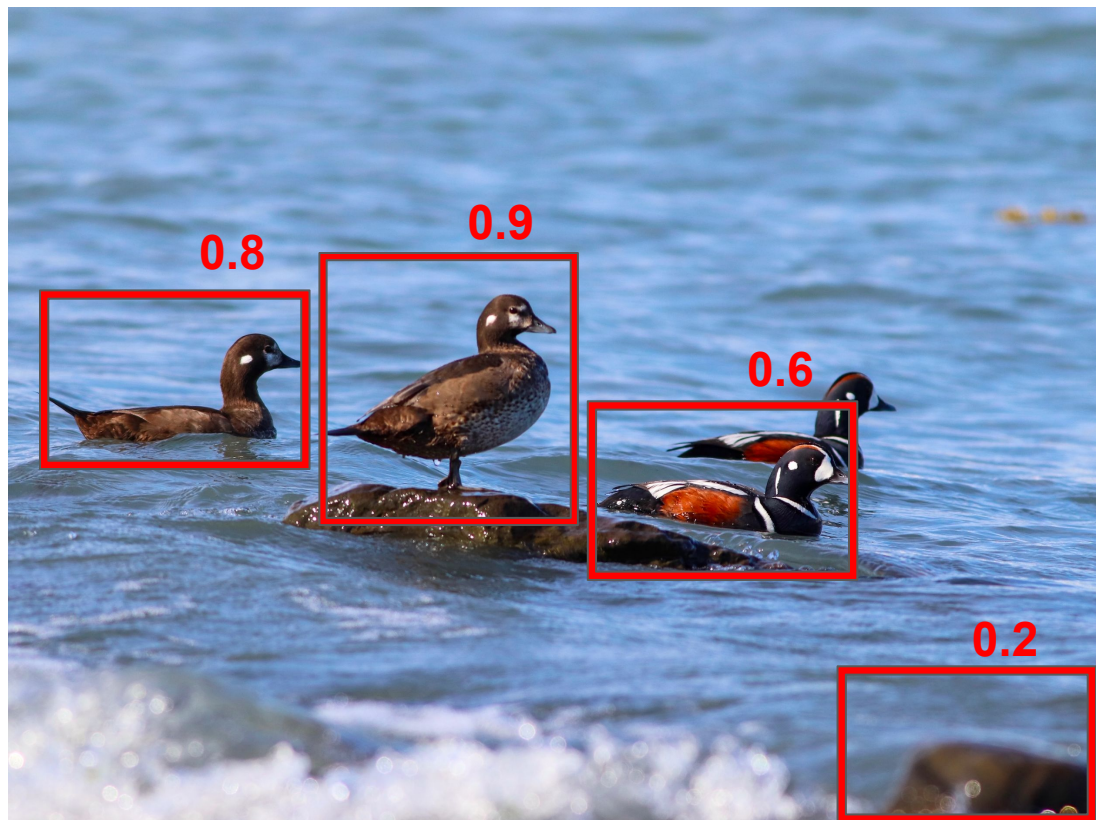
With confidence threshold of 0:

PR = 3 / (3 + 1) = 0.75



Model predictions with confidences

Precision = $\dfrac{\text{\# TP}}{\text{\# TP + \# FP}}$

With confidence threshold of 0.5:
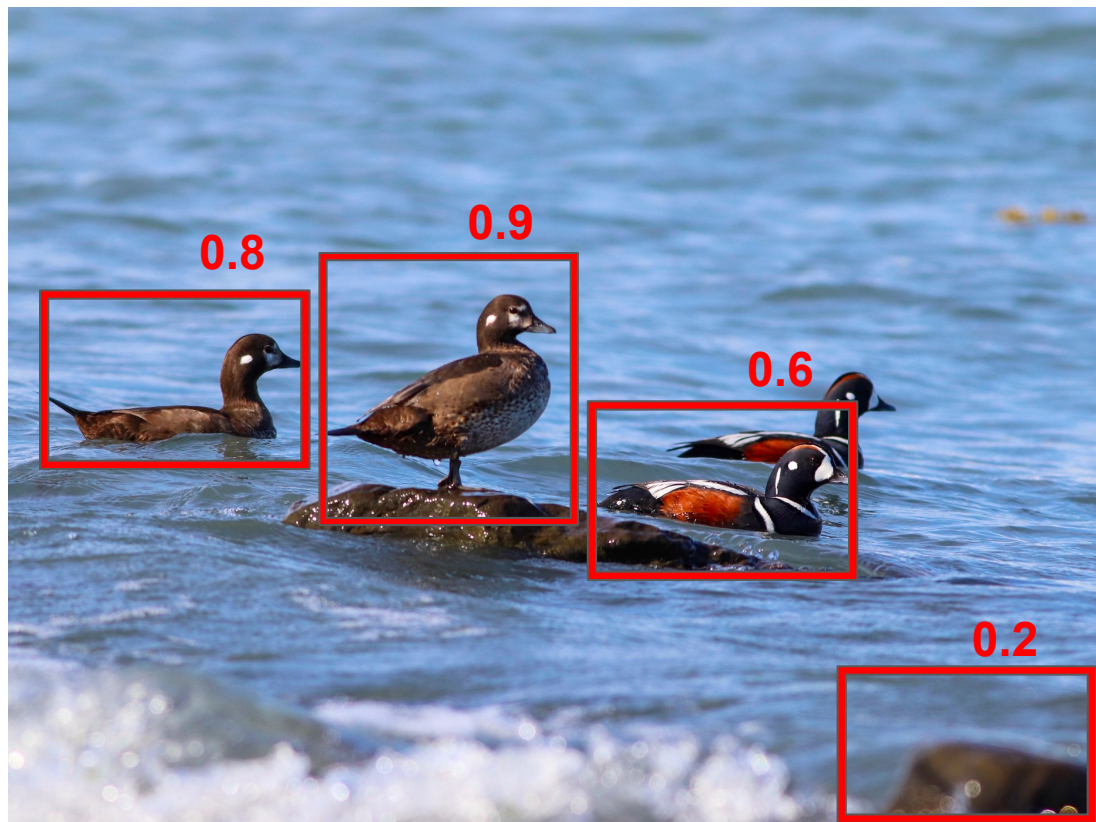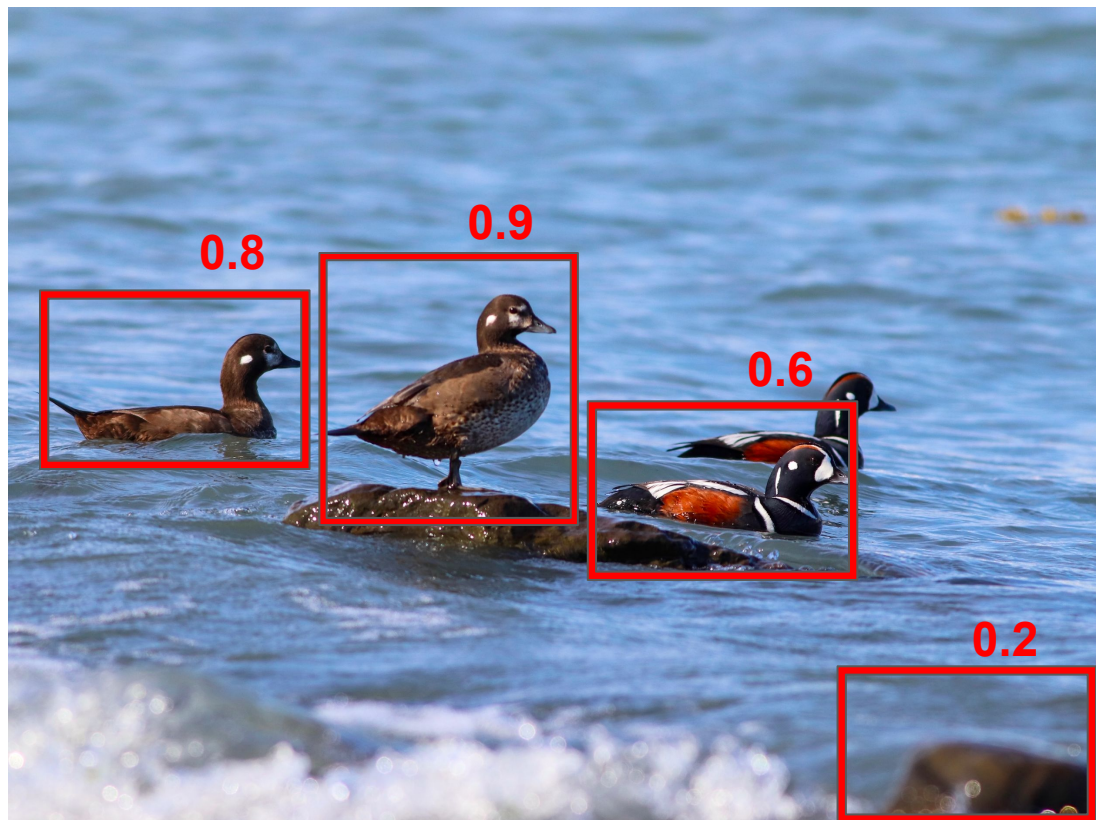
PR = 3 / (3 + 0) = 1



Model predictions with confidences

$$\text{Recall} \quad = \quad \frac{\# \text{ TP}}{\# \text{ TP} + \# \text{ FN}}$$

With confidence threshold of 0:

PR = 3 / (3 + 1) = 0.75



Model predictions with confidences

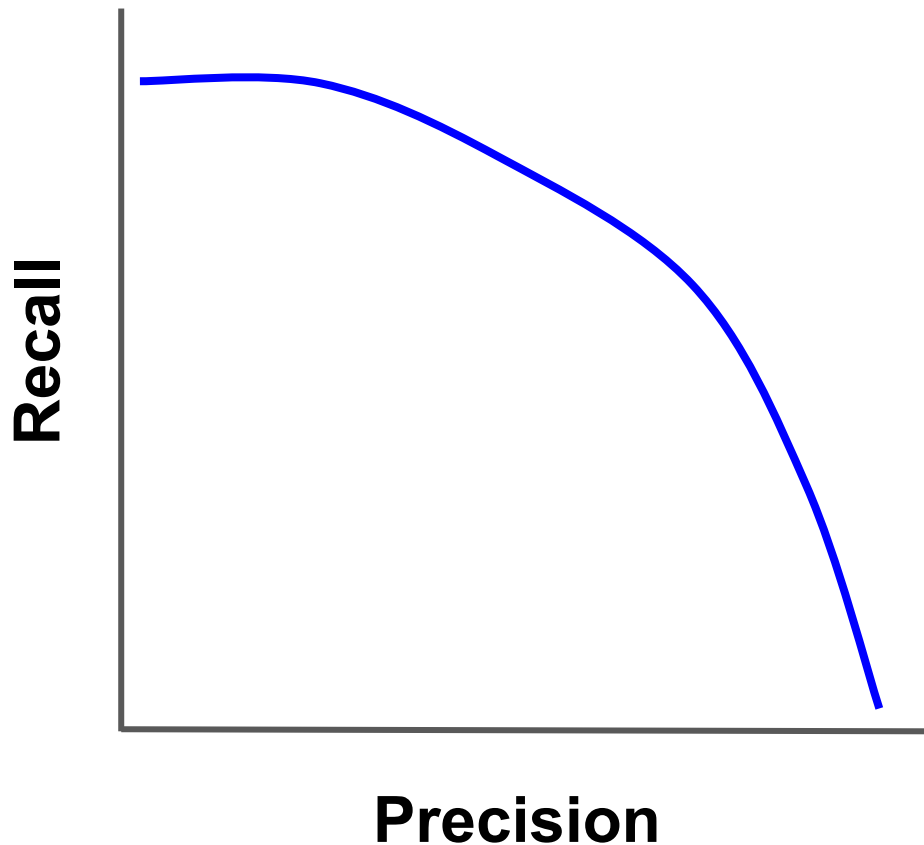Recall $\quad = \quad \dfrac{\text{\# TP}}{\text{\# TP + \# FN}}$

With confidence threshold of 0.7:

PR = 2 / (2 + 2) = 0.5



Model predictions with confidences

# Precision Recall Curves



- Plot precision vs recall at all the confidence thresholds

- Average Precision (AP) is the area under the curve

- Mean Average Precision (mAP) is the mean of all the APs for each class

- Similarly, you will see mean IOU (mIOU) for semantic segmentation metrics

# Public Benchmarks

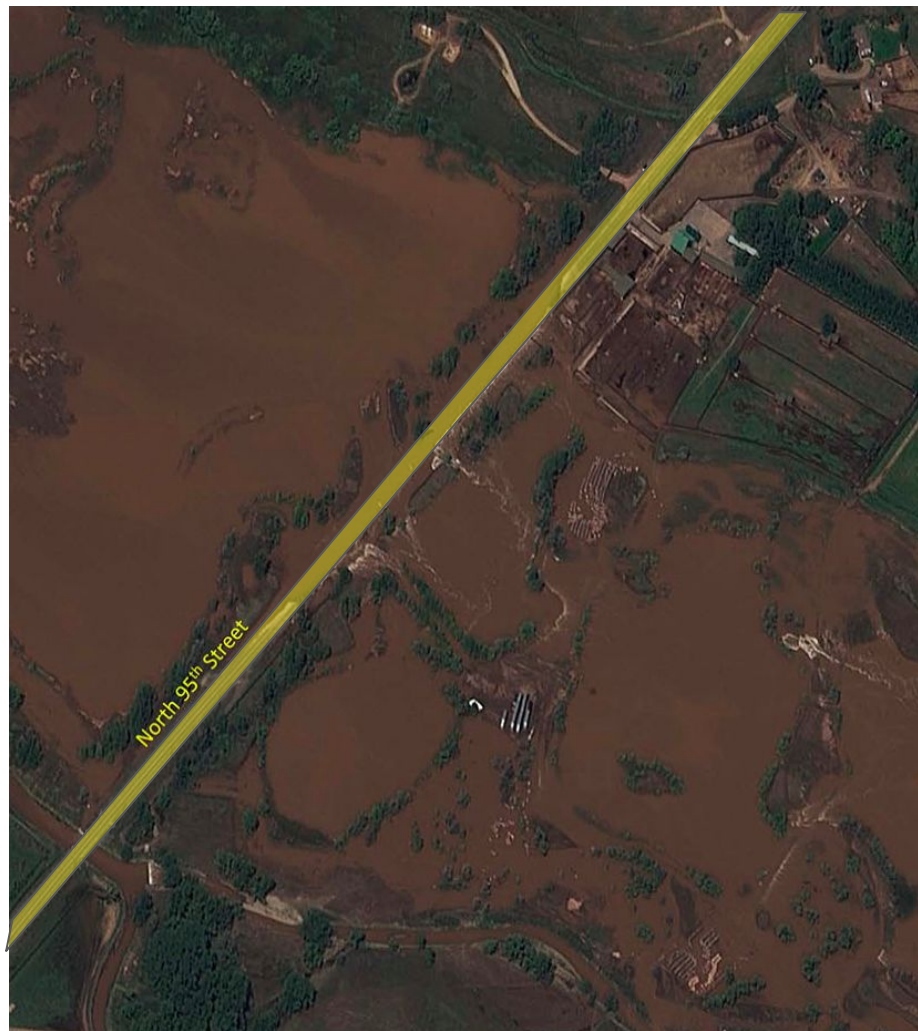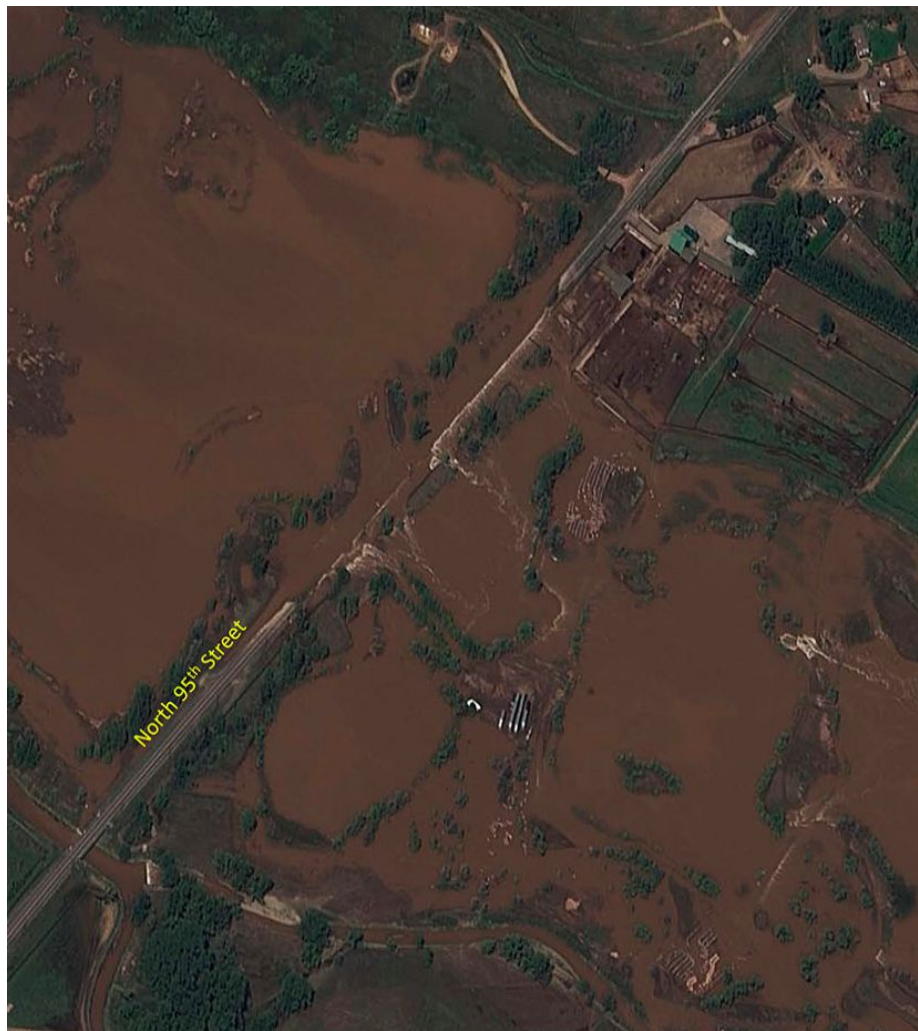- COCO
- Cityscapes
- Mapillary Vistas

# Evaluating models for real world applications

- Mapping:
  - False positive: Add a non-existent building to the map
  - False negative: Missing a building on the map
- Autonomous Vehicles:
  - False positive: Vehicle detects a non-existent stop sign, stops, and gets rear ended
  - False negative: Vehicle drives through a stop sign and causes an accident
- Medicine:
  - False positive: Unnecessary procedures -> higher healthcare costs, strain health care system
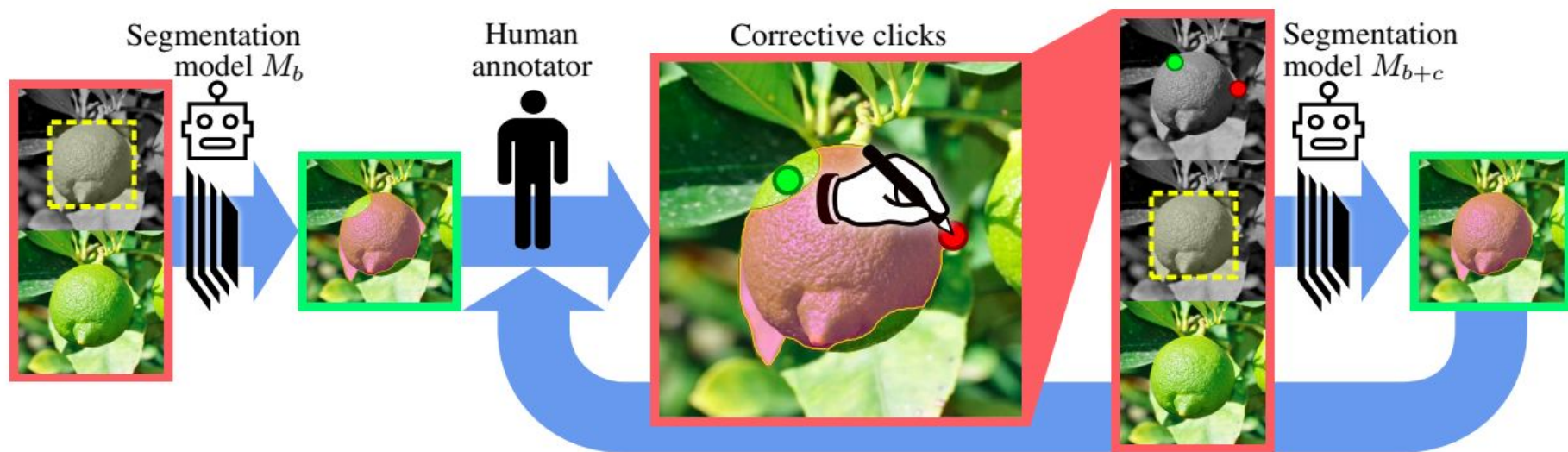  - False negative: A serious disease goes untreated

North 95th Street

# Resource Constraints: Data and Labeling



- [Large-scale interactive object segmentation with human annotators](#)
- [Interactive Full Image Segmentation by Considering All Regions Jointly](#)

# Resource Constraints:  Computation

- Autonomous vehicles need to do inference on device
    - [MobileNet](#)
- Very deep backbones (e.g. ResNet) are expensive at train and inference time
    - Start with a pre-trained backbone
- High training cost to search for your architecture
- Bigger data sets yield better results but require disk space