



## Deep Learning for Computer Vision (IV)

---

### Learning Objectives

- Learn the basics of advanced computer vision tasks beyond image classification: image segmentation and object detection.

# Beyond image classification: object detection and image segmentation

Single-label multi-class classification



- ☒ Biking -
- ☐ Running
- ☐ Swimming

Multi-label classification



- |  |  |
|--|--|
| <input checked="" type="checkbox"/> Bike -   | <input checked="" type="checkbox"/> Tree - |
| <input checked="" type="checkbox"/> Person - | <input type="checkbox"/> Car               |
| <input type="checkbox"/> Boat                | <input type="checkbox"/> House             |

Label  
非常简单

Image segmentation



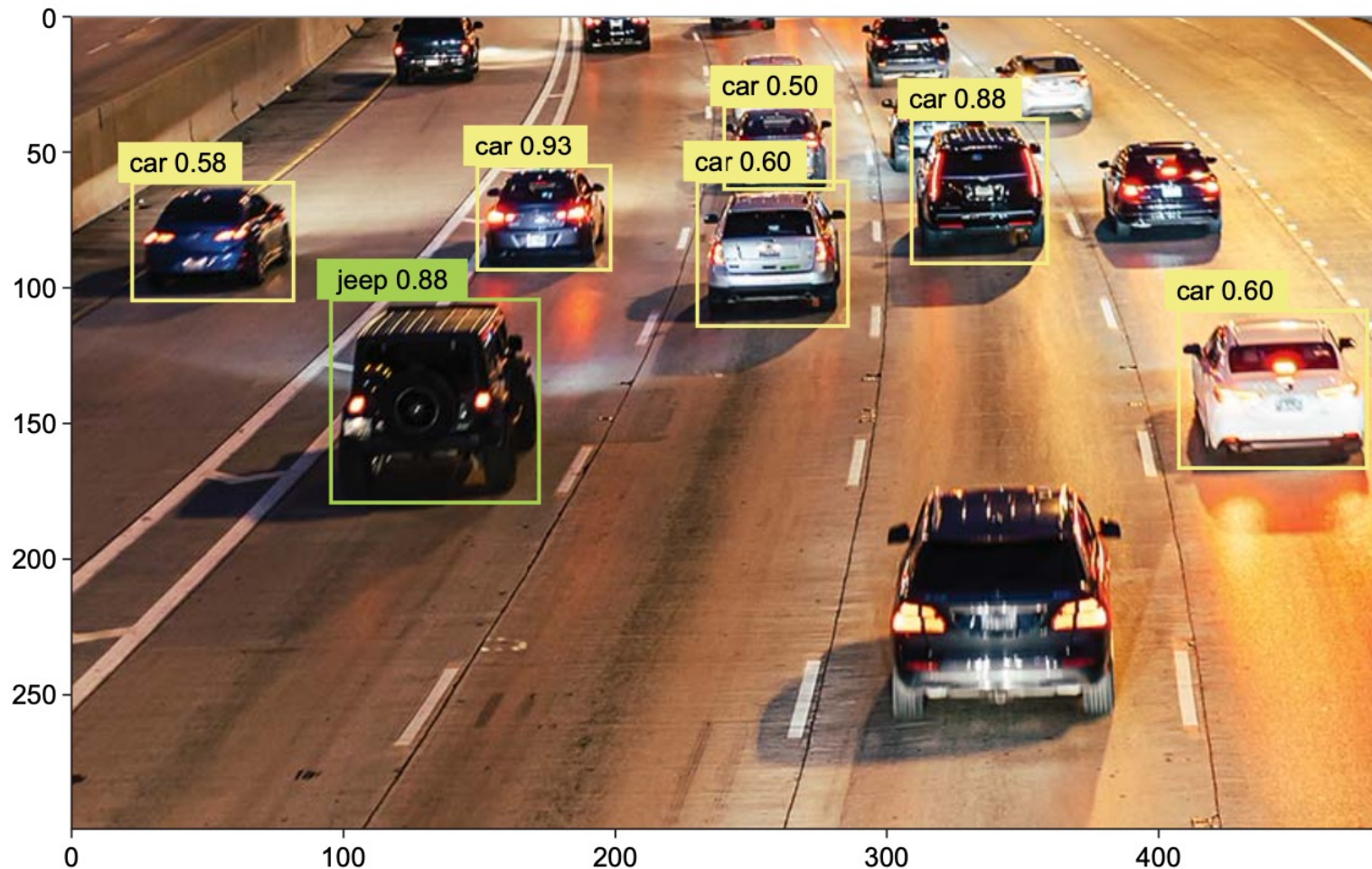
Object detection

x  
s



# Object Detection

- The goal is to draw rectangles (called bounding boxes) around objects of interest in an image, and associate each rectangle with a class.
- A self-driving car could use an object-detection model to monitor cars, pedestrians, and signs in view of its cameras, for instance.

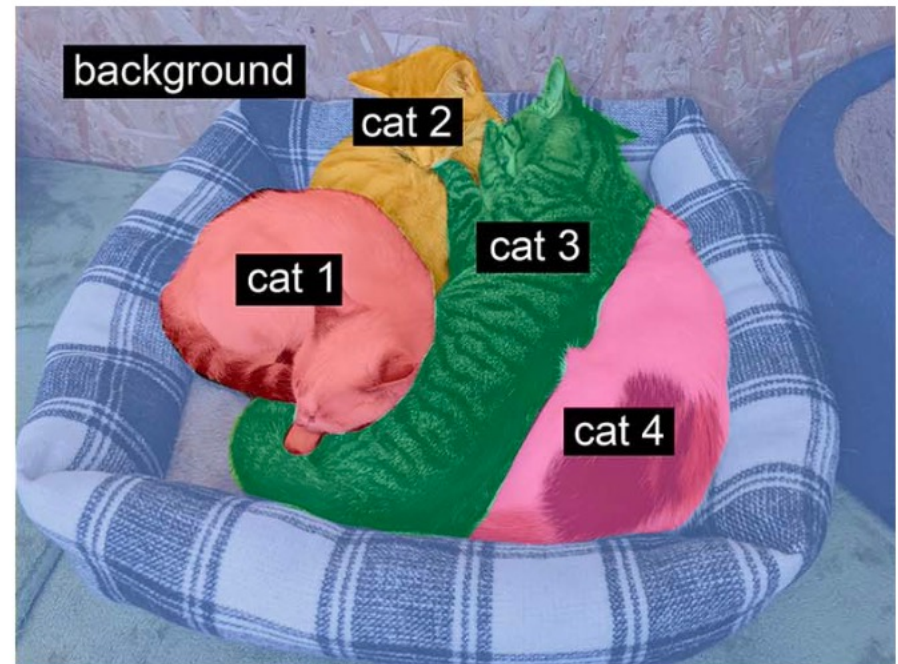


## Image Segmentation

- The goal is to “segment” or “partition” an image into different areas, with each area usually representing a category.



Semantic segmentation

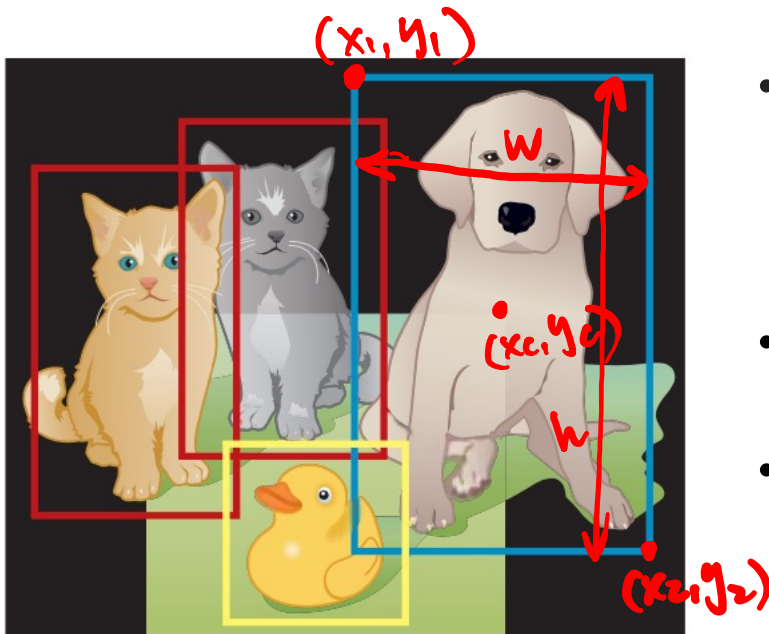


Instance segmentation

# **Learn the basics of object detections**

# Object Detection

- **Goal:** predict the location of objects in an image via bounding boxes and the classes of the located objects
- **Input:** an image with one or more objects
- **Output:** one or more bounding boxes (defined by coordinates  $[x_1, y_1, x_2, y_2]$  or  $[x_c, y_c, w, h, ]$ ) and a class label for each bounding box



Cat, Cat, Duck, Dog

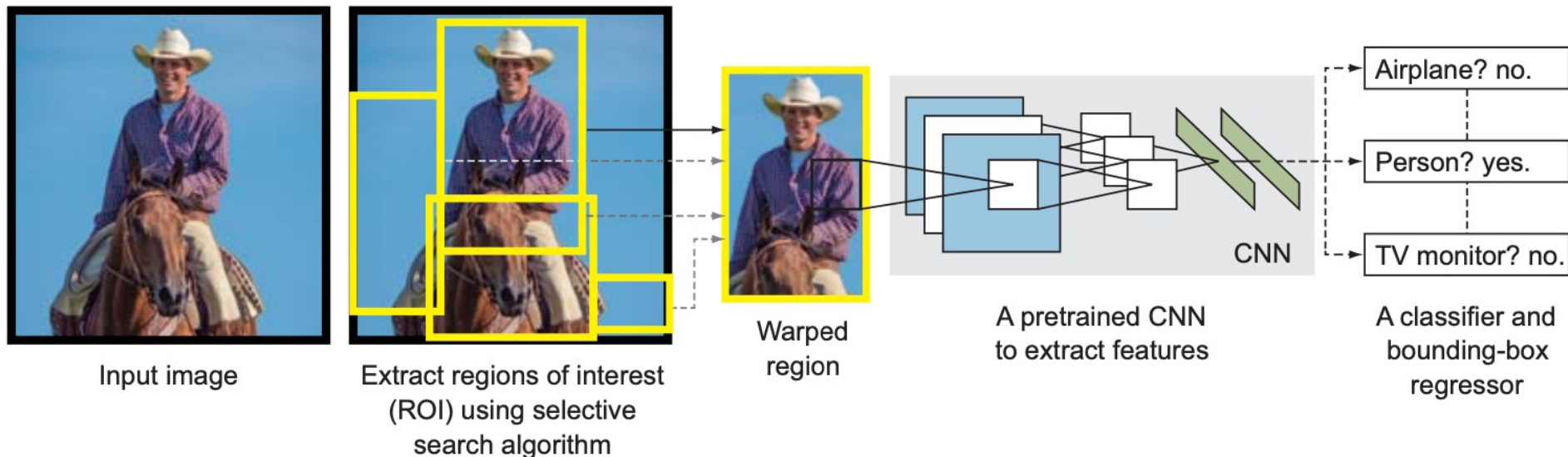
- The two most popular object detection systems are the **R-CNN** family of networks and the **YOLO** family of networks.
- R-CNN is a multi-stage detector and YOLO is a single-stage detector.
- In general, single-stage detectors tend to be less accurate than two-stage detectors but are significantly faster.

# Object Detection: R-CNN Family

- The R-CNN family of networks has three main variations: R-CNN, Fast R-CNN, and Faster R-CNN. R-CNN and Fast R-CNN use a selective search algorithm to propose Rols, whereas Faster R-CNN is an end-to-end DL system that uses a region proposal network to propose Rols.

## R-CNN

*Region of Interest*

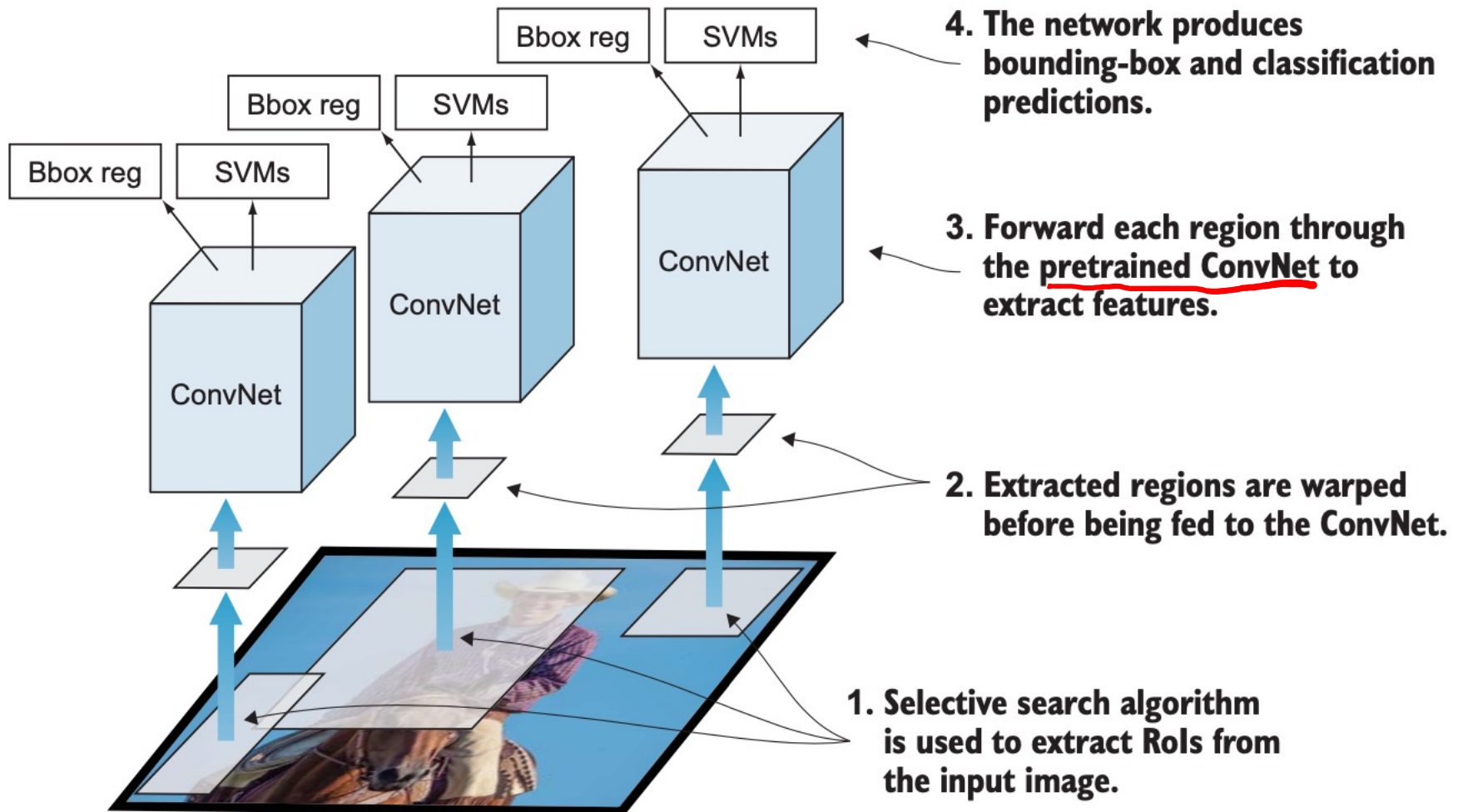


# Selective Search Algorithm

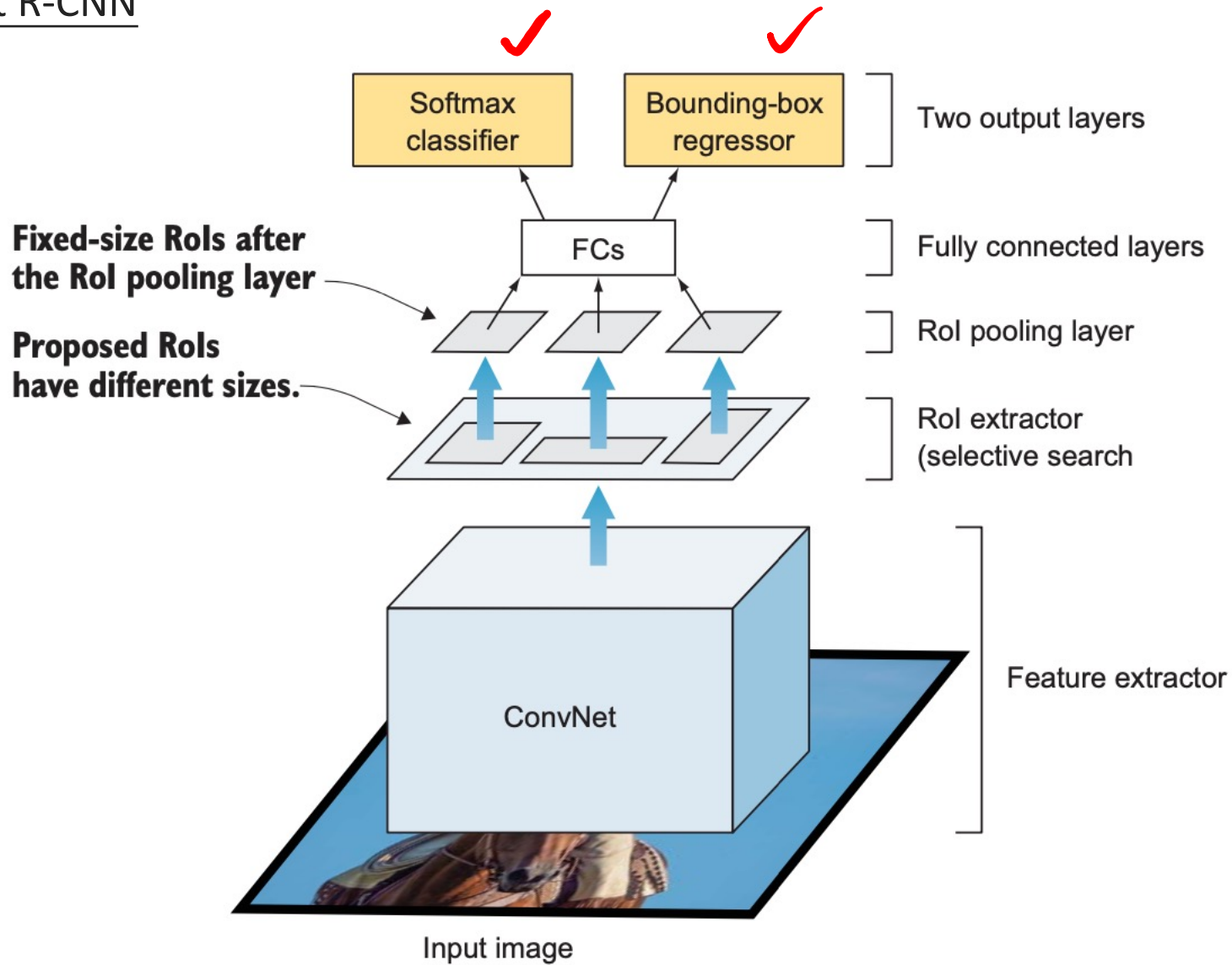
- A greedy search algorithm that is used to provide region proposals that potentially contain objects.
- It tries to find areas that might contain an object by combining similar pixels and textures into rectangular boxes.



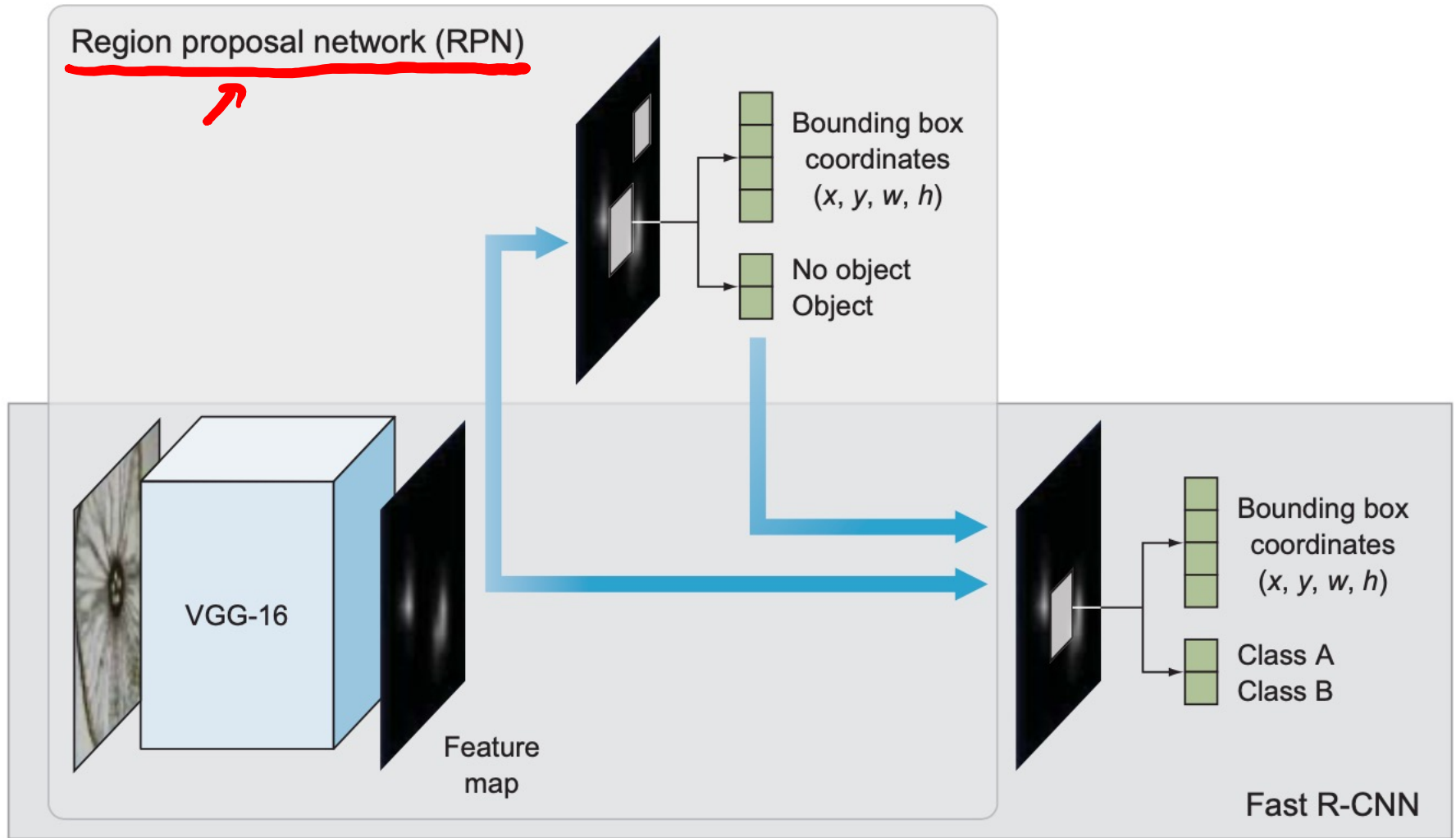
## R-CNN



## Fast R-CNN

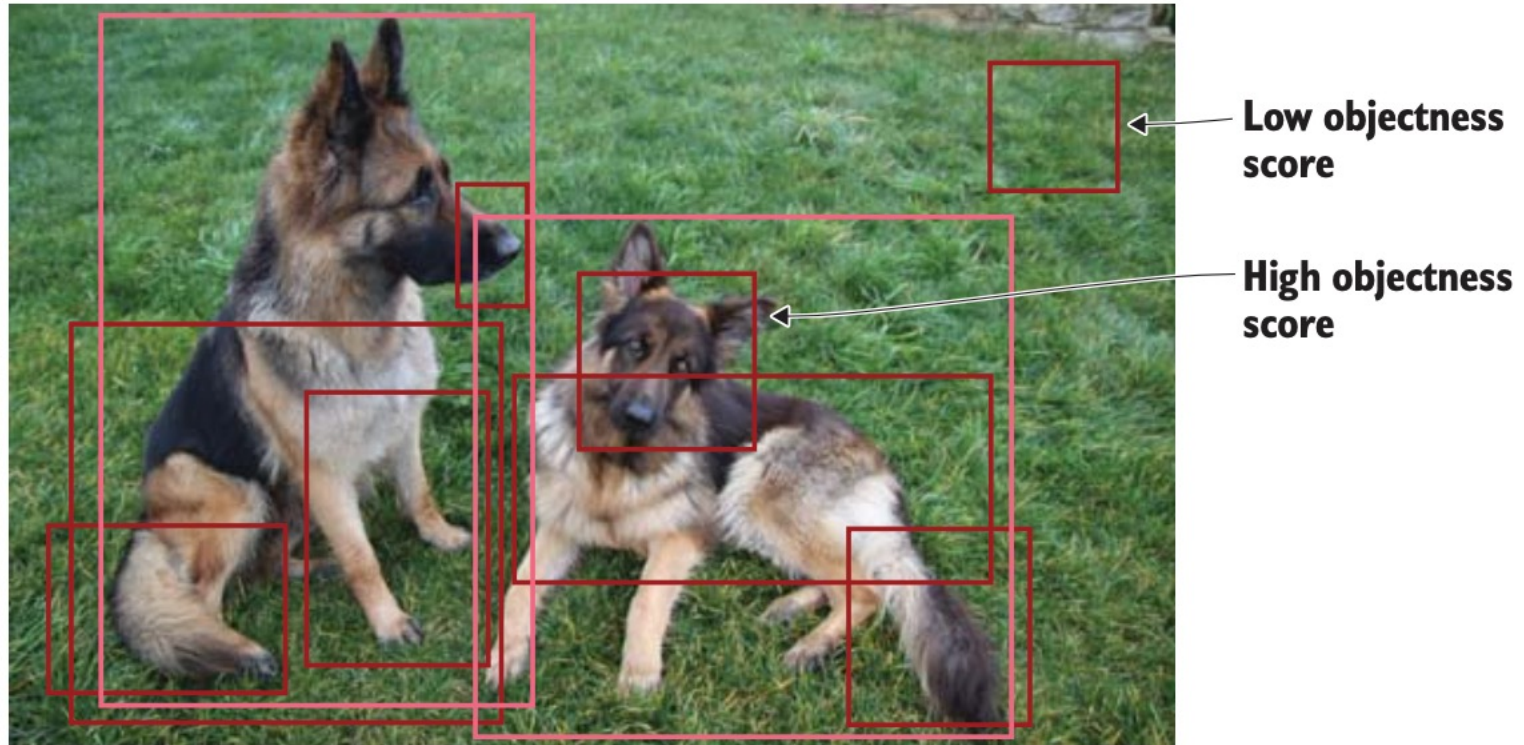


## Faster R-CNN (RPN + Fast R-CNN)



# Region Proposals

Regions of interest (Rols) proposed by the system. Regions with high objectness score represent areas of high likelihood to contain objects (foreground), and the ones with low objectness score are ignored because they have a low likelihood of containing objects (background).



## Object Detection: R-CNN Family

$[x_c, y_c, w, h]$

- *R-CNN*—Bounding boxes are proposed by the selective search algorithm. Each is warped, and features are extracted via a deep convolutional neural network such as AlexNet, before a final set of object classifications and bounding box predictions is made with linear SVMs and linear regressors.
- *Fast R-CNN*—A simplified design with a single model. An RoI pooling layer is used *after* the CNN to consolidate regions. The model predicts both class labels and Rols directly.
- *Faster R-CNN*—A fully end-to-end DL object detector. It replaces the selective search algorithm to propose Rols with a region proposal network that interprets features extracted from the deep CNN and learns to propose Rols directly.

# Detectron2

## 1. Detectron2

- a. **Data preparation (dataset registration)**

DatasetCatalog  
MetadataCatalog

- b. **Configuration setup (hyper-parameters)**

cfg = get\_cfg()

- c. **Build a trainer**

DefaultTrainer(cfg)

- d. **COCO evaluation**

COCOEvaluator

- e. **Inference by predictor**

DefaultPredictor

## 2. Rebar hook angle detection

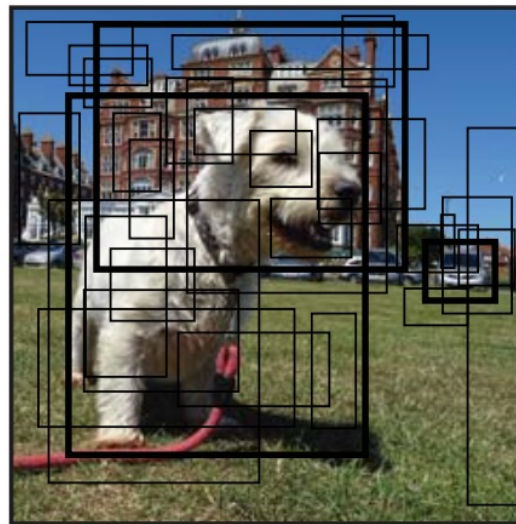
Train your first Detectron2 detector

# Object Detection: YOLO Family

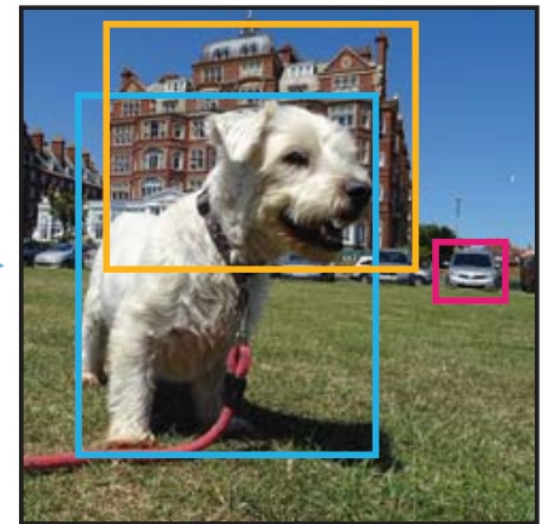
- The YOLO family is a series of end-to-end DL models designed for fast object detection, and it was among the first attempts to build a fast real-time object detector.
- Although the accuracy of the models is close but not as good as R-CNNs, they are popular for object detection because of their detection speed, often demonstrated in real-time video or camera feed input.
- For a good introduction, see <https://aiacademy.tw/yolo-v4-intro>



Splits the image into grids



Predicts bounding boxes  
and classifications



Final predictions after  
non-maximum suppression

# Darknet (YOLO v4)

## 1. Darknet

### a. Data preparation (dataset registration)

`.jpg/.txt (images/annotations)`  
`train.txt/valid.txt`

### b. Configuration setup (hyper-parameters)

`Obj.data`  
`Obj.names`  
`yolov4.cfg`

### c. Training

`./darknet detector train obj.data yolov4.cfg yo.weights -map`

### d. Evaluation

`./darknet detector map obj.data yolov4.cfg yo.weights -iou_thresh 0.5`

### e. Image prediction

`./darknet detector test yolov4.cfg yo.weights image`

## 2. Construction machine detection

[Train your first YOLO v4 detector](#)

**Learn the basics of image segmentation**

# Image Segmentation

- The goal is to “segment” or “partition” an image into different areas, with each area usually representing a category.



## Semantic segmentation

each pixel is independently classified into a semantic category, like “cat.” If there are two cats in the image, the corresponding pixels are all mapped to the same generic “cat” category

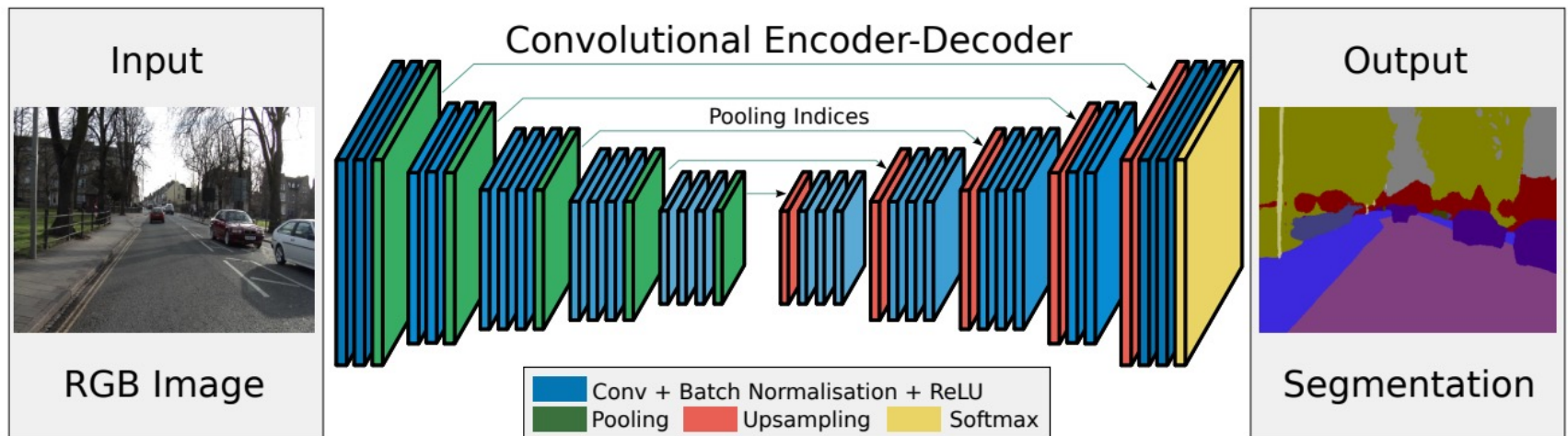


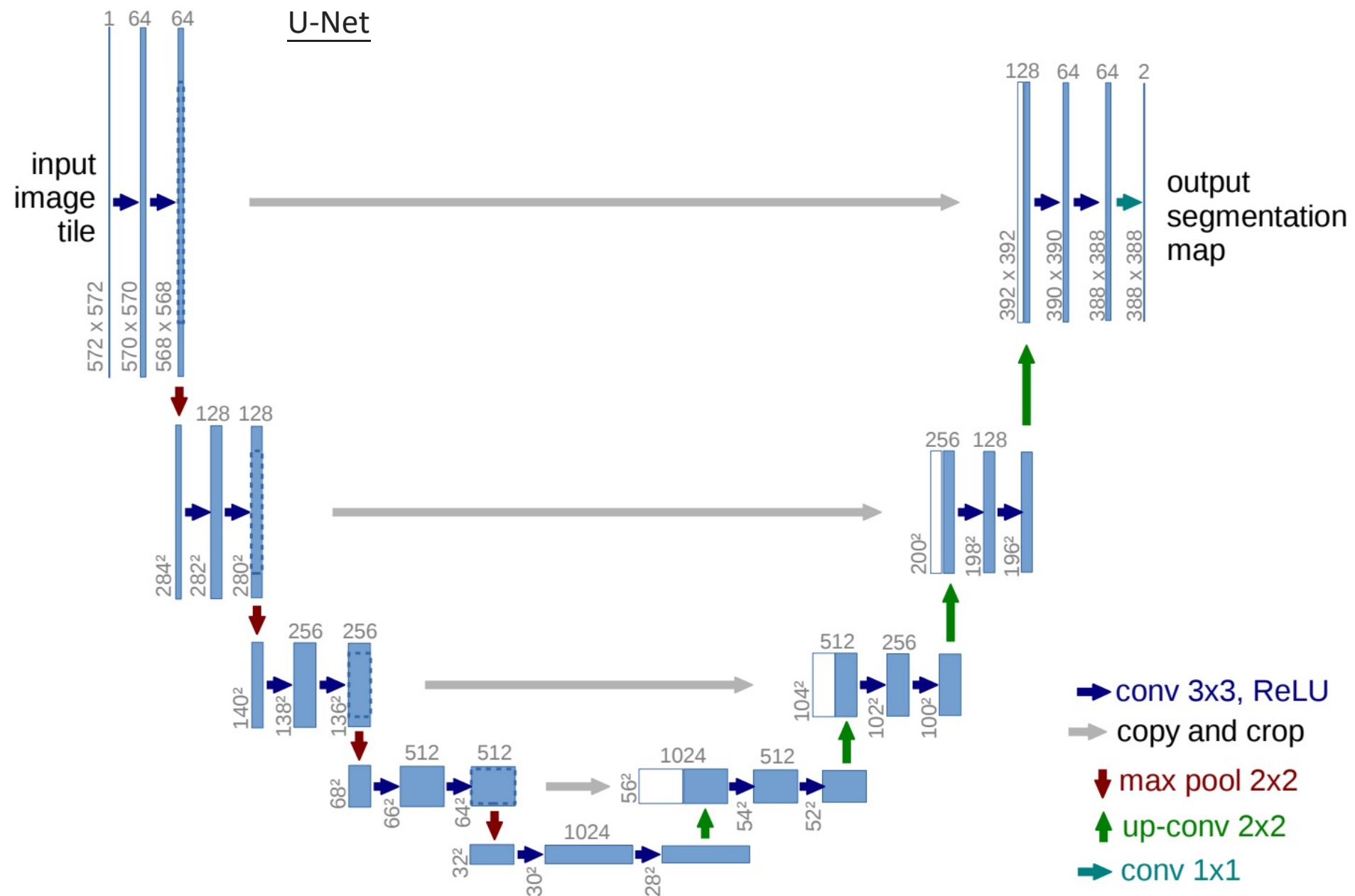
## Instance segmentation

seeks not only to classify image pixels by category, but also to parse out individual object instances. In an image with two cats in it, instance segmentation would treat “cat 1” and “cat 2” as two separate classes of pixels

# Semantic Segmentation

- Convolution **Encoder**: closely resembles the kind of convnet you'd use for image classification: a stack of Conv2D layers, with gradually increasing filter sizes. The purpose is to encode the images into smaller feature maps, where each spatial location (or pixel) contains information about a large spatial chunk of the original image. You can understand it as a kind of compression.
- Convolution **Decoder**: apply a kind of inverse of the transformations we've applied so far—something that will upsample the feature maps instead of downsampling them.
- [A Simple Encoder-Decoder Example for Semantic Segmentation](#)





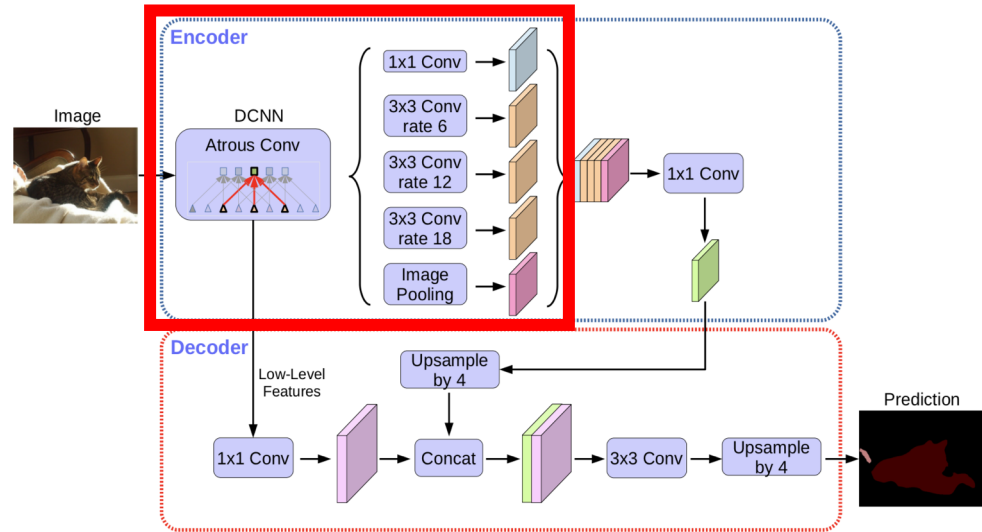
# Semantic Segmentation (Pixel Level Recognition)

## Google DeepLab V3+

### Encoder-Decoder Structure

**Encoder:** gradually reduces the feature maps and captures higher semantic information

**Decoder:** gradually recovers the spatial information



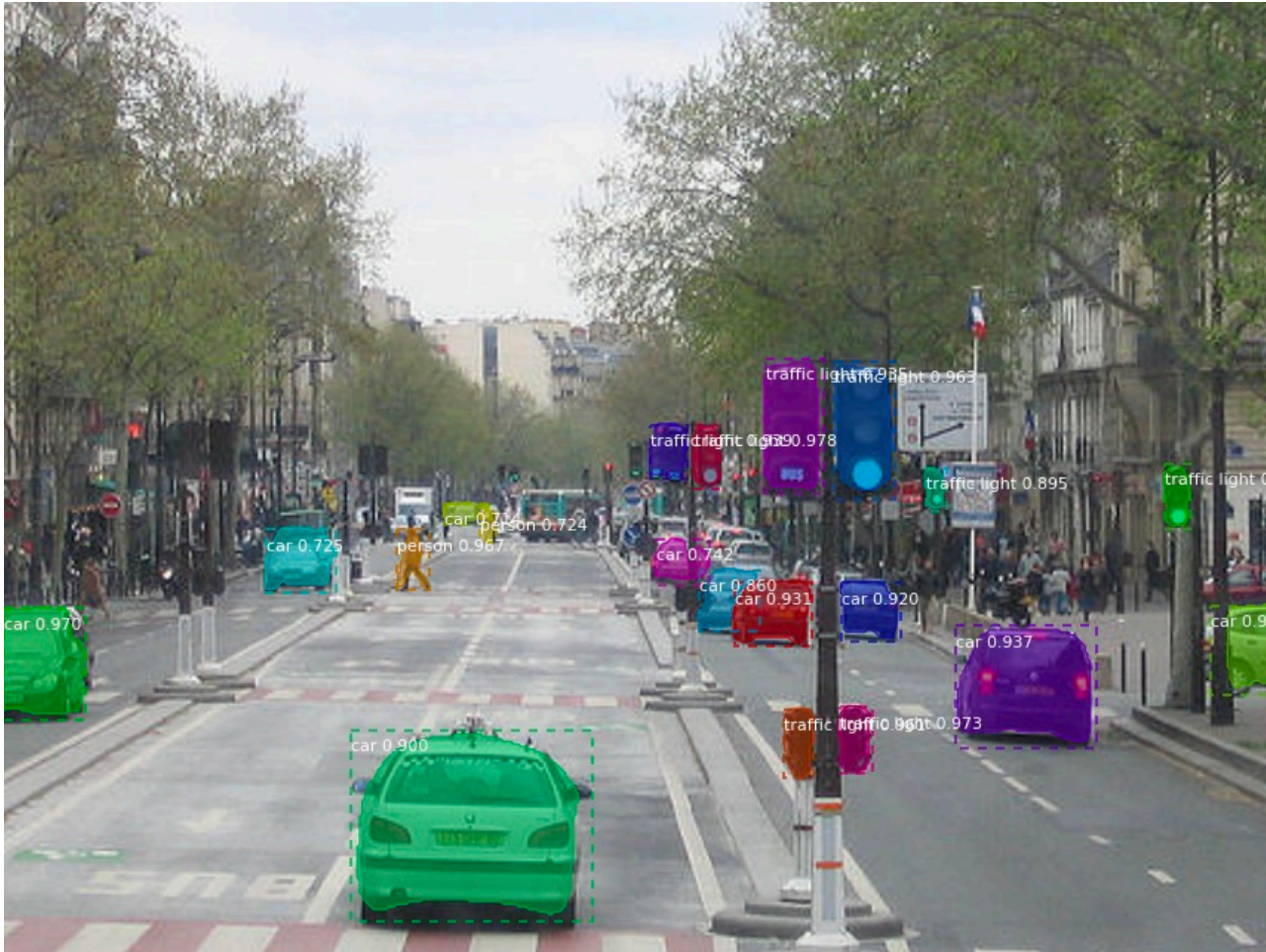
### Atrous Spatial Pyramid Pooling

Apply several parallel atrous convolution with different rates

### Depthwise separable convolution

Reduce the computation cost and number of parameters while maintaining similar performance.

# Instance Segmentation (Pixel Level Recognition) Mask R-CNN



Kaiming He et al. (2018). Mask R-CNN. [arXiv:1703.06870](https://arxiv.org/abs/1703.06870)

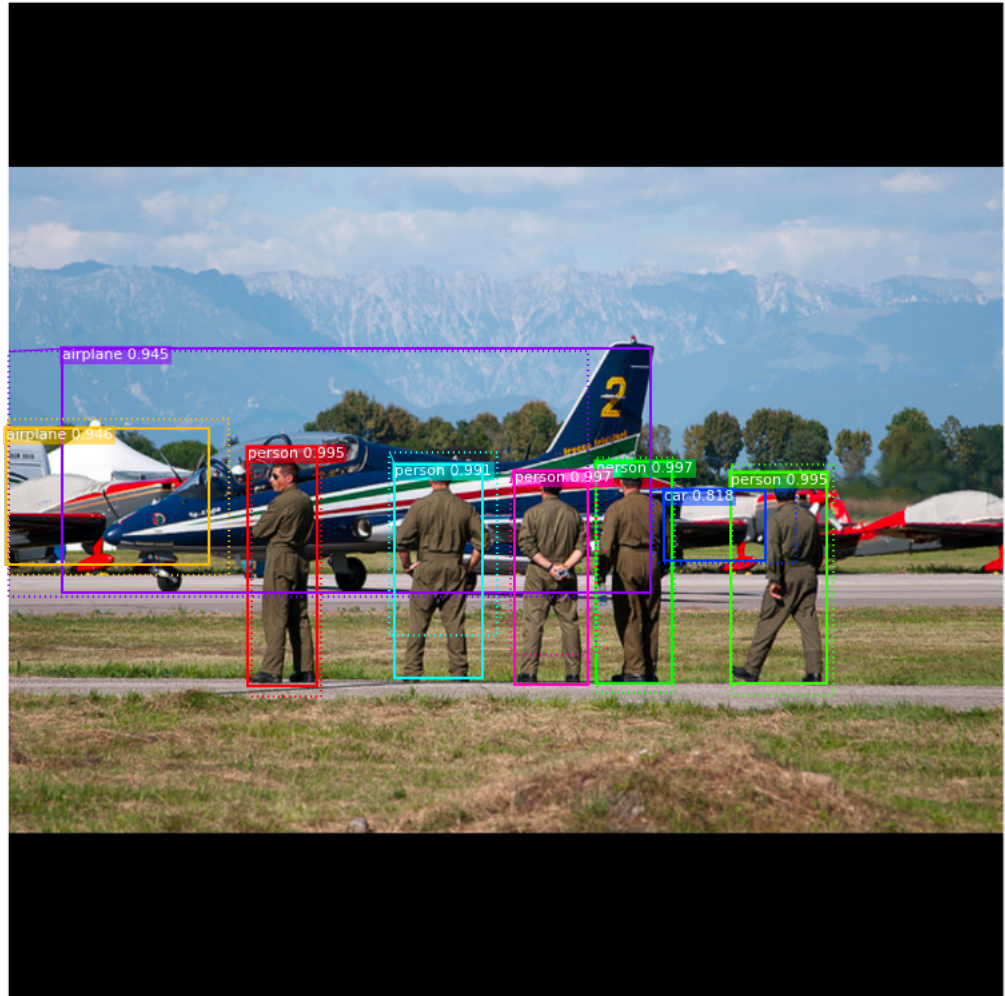
# Instance Segmentation (Pixel Level Recognition) Mask R-CNN

## 1. Anchor sorting and filtering



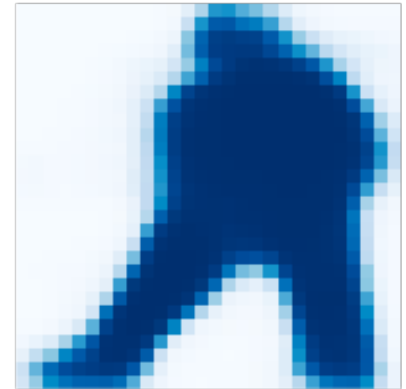
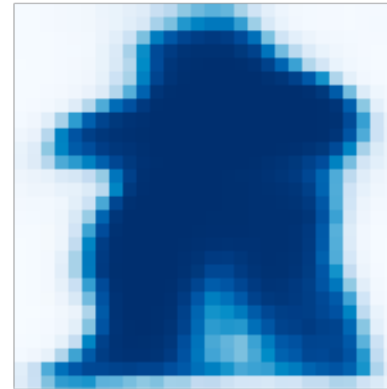
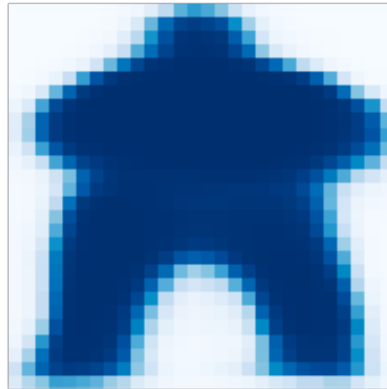
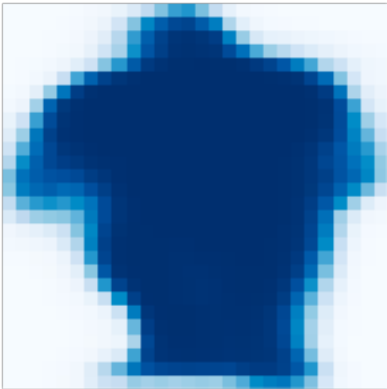
# Instance Segmentation (Pixel Level Recognition) Mask R-CNN

## 2. Bounding Box Refinement



# Instance Segmentation (Pixel Level Recognition) Mask R-CNN

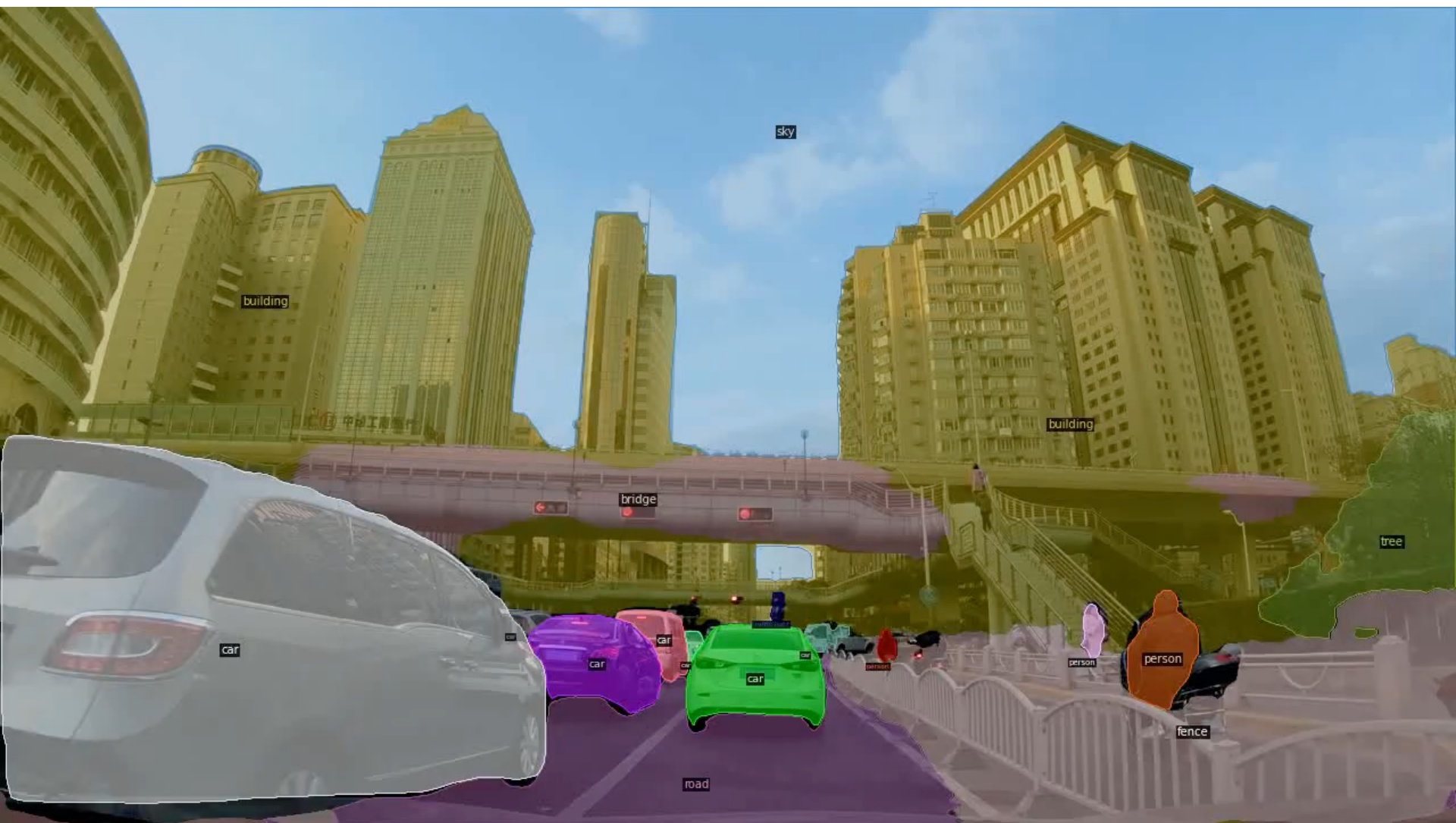
## 3. Mask Generation



# Instance Segmentation (Pixel Level Recognition) Mask R-CNN

4. Composing the different pieces into a final result





# HW7