# CV / VLM

Unit 2: Introduction to Object
Detection (OD)

TIL-AI
TODAY I LEARNED AI

# 2.3.1

## Single Shot Detectors (SSD) and YOLO

Introduction to SSD and YOLO architectures

TIL-AI
TODAY I LEARNED AI

# OD Architectures: R-CNN Vs. SSD

Algorithms in the R-CNN family use a 2
step approach (Region proposal + CNN
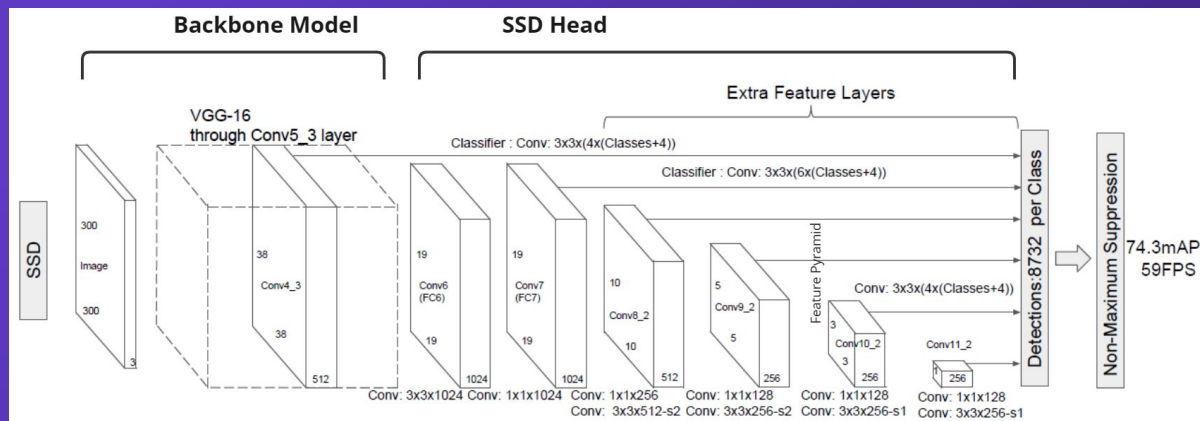detection)
  + Better accuracy
  - Sub-seconds prediction

Single-shot detector (SSD) approaches
utilize a deeper/custom neural network
but only require a single pass (1 step)
  + More efficient
  + Real-time (predictions in milliseconds)
  - Moderate accuracy

TIL-3
TODAY I LEARNED 1
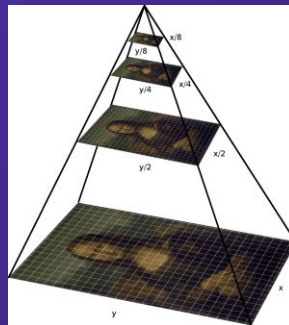
# Single-Shot Detector (SSD)

Two Components

1. **Backbone model**: Pre-trained image CNN (Resnet, VGG)
2. **SSD Head**: More convolutional layers added to the backbone, whereby the outputs are bounding boxes and classes of the objects in the spatial locations.
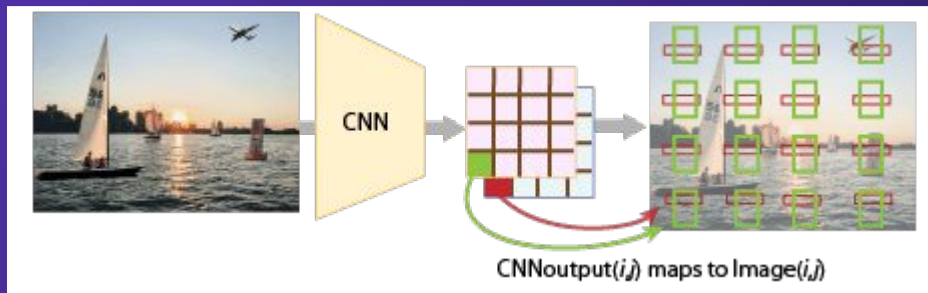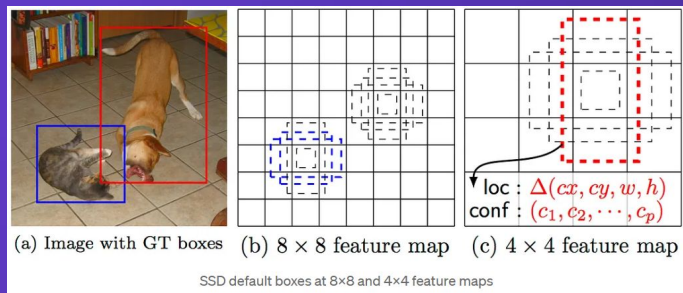


Wei Liu et al. in the paper SSD: Single Shot MultiBox Detector.

# Single-Shot Detector (SSD)
## Innovations

- Applies various feature map grid cell sizes (e.g. 8x8, 6x6, ..., 1x1) to detect objects of different sizes [image pyramid]



- Anchor Boxes

This is all done in the SSD Head network!



(a) Image with GT boxes   (b) 8 × 8 feature map   (c) 4 × 4 feature map

$loc : \Delta(cx, cy, w, h)$
$conf : (c_1, c_2, \cdots, c_p)$

SSD default boxes at 8×8 and 4×4 feature maps
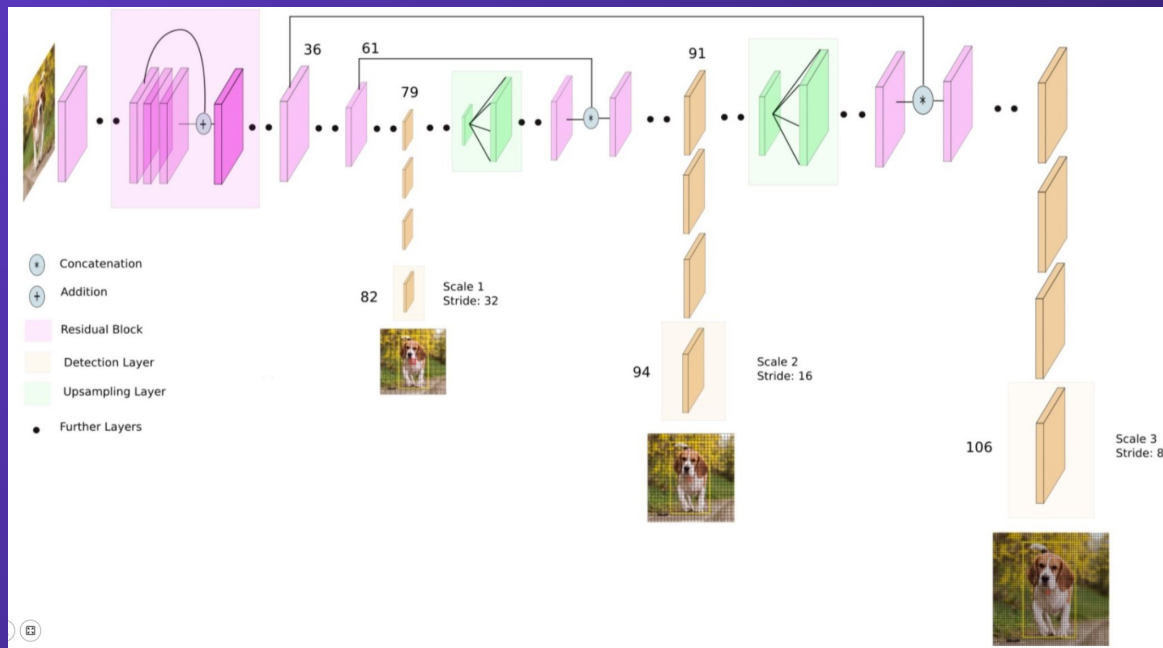


CNNoutput(i,j) maps to Image(i,j)

# You Only Look Once (YOLO) v3
## Architecture

Backbone model:
Pre-trained CNN (Darknet)

YOLO v3 unique customization:

- Adapted ResNet-style residual blocks.
- Upsampling and concatenation of feature layers with earlier feature layers which preserve fine-grained features
- Three scales for detection

Latest version: YOLO v8

# YOLO v8 Architecture
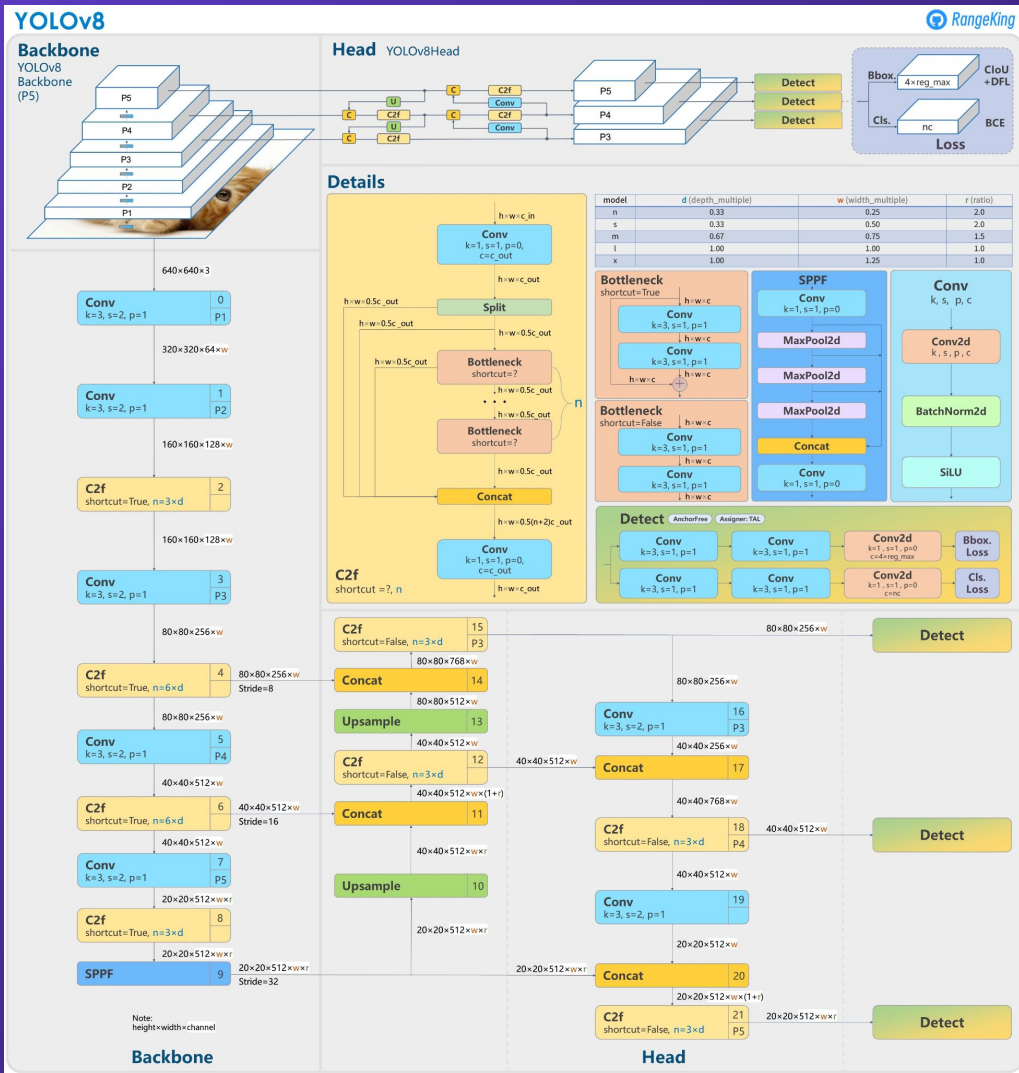
High-level view of the architecture and its components.

## Backbone

- with Feature Pyramid (downscale through Conv from 640x640 to 20x20)

## Head

- With C2F (bottleneck, SPPF) (Pooling/Reduce Dimension) and CoV[SiLU] (Activation)
- Bounding Box (without Anchor) and Classifications

- Latest version: YOLO v8
- Adapted from Brief summary of YOLOv8 model structure · GitHub

# YOLO v8 Architecture

YOLOv8 has also integrated other submodules

- <u>Classify</u> models pretrained on the <u>ImageNet</u> dataset.
- <u>Detect</u>, <u>Segment</u> and <u>Pose</u> models pretrained on the <u>COCO</u> dataset (with track mode).
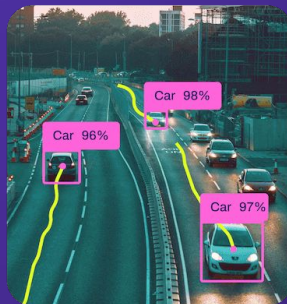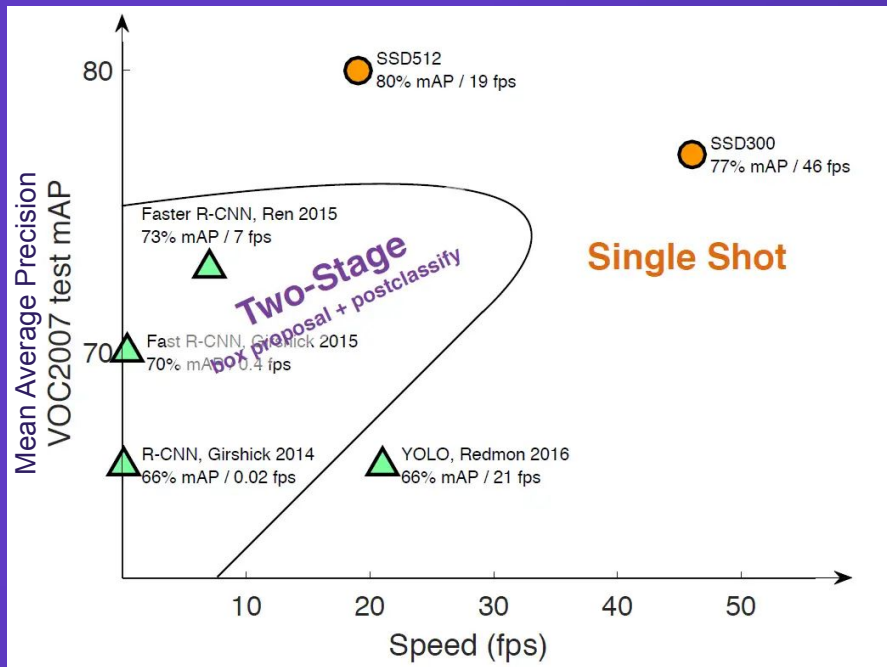


- <u>Latest version: YOLO v8</u>
- Adapted from <u>Brief summary of YOLOv8 model structure · GitHub</u>

# Comparison Between One Stage and two-Stage Object detection



## Single Stage Detectors (SSD)

- Focus on speed
- Harder on training

## Two Stage Detectors

- Better precision
- More flexible

However, It is worthwhile to also note that SSD such as YOLO has undergone quite a bit of developments since.