

**Convolutional layer**Input  $C_{in} \times H \times W$ hyp: Kernel size  $K_H \times K_W$ 通常  $K_H = K_W$ #filters  $C_{out}$ 32, 64, 128, 256  $2^n$ Padding  $P$ 图像边界额外加  $P$  圈  $0 \rightarrow$  preserving more spatial info  $\rightarrow$  imp performance

$$P = \frac{K-1}{2}$$

$$\begin{cases} P=0, \text{ No (Zero) Padding} \\ P=1, \text{ One pixel (Zero) } \sim, \text{ Zero-padding} \\ P = \lfloor \frac{K-1}{2} \rfloor, \text{ Same } \sim (\text{Zero}) \\ P = K-1, \text{ FULL Zero } \sim \end{cases}$$

Stride  $S$ 

$$\begin{cases} S=1, \text{ Unit stride, highest resolution} \\ > \text{ Strided conv, } \downarrow \\ < \text{ Fractional stride, used in transposed conv. } \text{Size}_{out} > \text{Size}_{in} \end{cases}$$

Weight matrix  $C_{out} \times C_{in} \times K_H \times K_W$   
#filtersBias vector  $C_{out}$ Output size (#output ele)  $C_{out} \times H' \times W'$   
 $H' = \frac{(H - K + 2P)}{S} + 1$ 

$$W' = \frac{(W - K + 2P)}{S} + 1$$

无法整除则 Floor

Memory usage For 32-bit floating point, bytes per ele: 4 (unit)

$$\text{Weights } m = C_{in} \times C_{out} \times K_H \times K_W \times 4 \quad ①$$

$$\text{biased } m = C_{out} \times 4 \quad ②$$

feature maps  $m$ :

$$\text{input } m = H_{in} \times W_{in} \times C_{in} \times 4 \quad ③$$

$$\text{out } m = H' \times W' \times C_{out} \times 4 \quad ④$$

$$\text{Total mem} = \frac{①+②+③+④}{0}$$

Full Form	Units	Bytes
1 Bit	Binary Digit (0/1)	
1 Nibble	4 bits	
1 Byte	8 bits	
1 kilobyte(KB)	1024 byte	$2^{10}$ bytes
1 Megabyte(MB)	1024 KB	$2^{20}$ bytes
1 Gigabyte (GB)	1024 MB	$2^{30}$ bytes
1 Terabyte(TB)	1024 GB	$2^{40}$ bytes
1 Petabyte(PB)	1024 TB	$2^{50}$ bytes
1 Exabyte(EB)	1024 PB	$2^{60}$ bytes
1 Zettabyte(ZB)	1024 EB	$2^{70}$ bytes
1 Yottabyte(YB)	1024 ZB	$2^{80}$ bytes
1 Brontobyte	1024 YB	$2^{90}$ bytes
1 Geopbyte	1024 Brontobyte	$2^{100}$ bytes

$$\begin{aligned} \# \text{Para} (\# \text{weight}) &= \text{weight shape} + \text{bias shape} \\ &= C_{in} \times C_{out} \times K_W \times K_H + C_{out} \end{aligned}$$

$$\# \text{FLOP} = \# \text{Output size} \times \text{Oper per out ele}$$

$$\text{floating pt. oper.} = C_{out} \times H' \times W' \times C_{in} \times K_H \times K_W$$

**Pooling,**无学习参数,  $C_{in} = C_{out}$ , 有 S.P.  $H'W'$  计算同 conv.

$$\text{FLOP} = C_{out} \times H' \times W' \times \begin{cases} (K \times K - 1) \\ K \times K \end{cases} \quad \begin{matrix} \text{max-pooling} \\ \text{avg} - \sim \end{matrix}$$

**Flatten**

无学习参数

$$\text{Output size} = C_{in} \times H \times W$$

$$\text{FLOP} = 0$$

无可学习参数

Output size  $C_{in} \times H \times W$   
 FLOP 0  
 Memory  $= C_{in} \times H \times W \times 4$

For each class: **confidence score** the probability that the prediction is correct.

eg. person

**IOU** intersection-to-parallel ratio (Jaccard)

1° 给定坐标  $\rightarrow$  画图  $\rightarrow$  面积交集除以并集  
 2° 给定矩阵  $\rightarrow$  pred 标记  $\rightarrow$   $\frac{TP}{TP+FP+FN}$   
 TP, FP, FN 点 重叠 面积

$$Dice = \frac{2TP}{2TP + FP + FN}$$

$$precision = \frac{TP}{TP+FP} \rightarrow \text{预测的 1 预测精度}$$

$$recall = \frac{TP}{TP+FN} \rightarrow \text{真实 1 中哪些被检测到, 召回}$$

## Non-Maximum Suppression

removes **redundant** bounding boxes from a set of detected objects

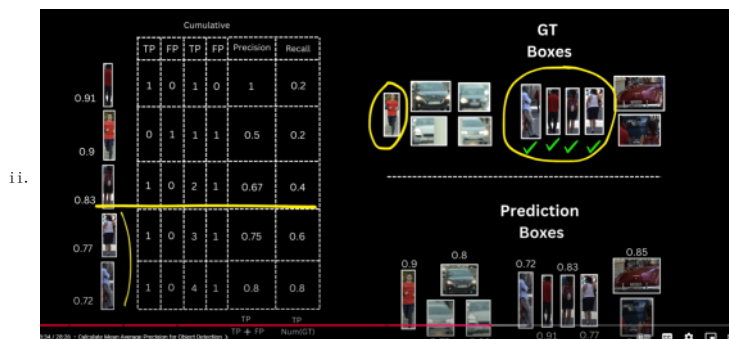
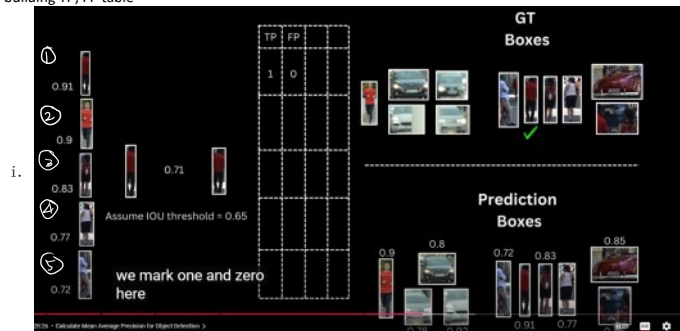
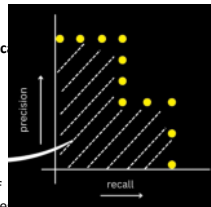
- ① Bounding box generation
- ② Sorting by cs (conf score)
- ③ Suppression:
  - 1) Select the box with the highest cs
  - 2) Cal IOU with all other boxes
  - 3) Suppress overlapping boxes(exceeds a predefined IOU threshold)
  - 4) Repeat This continues until all boxes have been either selected as a detection or suppressed.

**average precision (AP)** src: Mean Average Precision (mAP) | Explanation and Implementation for Object Detection

The average of precision **at different recall levels**. It is calculated as the **area under the precision-recall curve**.  
 Improving one often comes at the cost of reducing the other.

AP for person class:

- a. take all the predicted boxes that have person as the pred class
- b. sort them in the decreasing order of confidence scores
- c. For each, find the best matching ground truth box for matching (based on overlap(Area of this GT box has not been matched before, and its IOU is greater than some predefined threshold, mark the pred as a TP, GT as matched.
- d. building TP, FP table

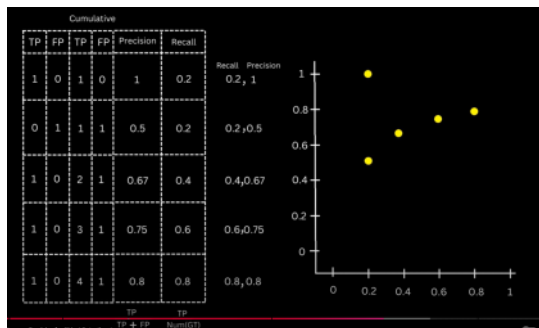


We got the precision and recall values using different confidence threshold. e.g. 0.83. 2 detection below it will be ignored.?

For all classes/  
queries.

$$mAP = \frac{1}{N} \sum_{i=1}^N AP_i$$

↑  
mean



link them -> precision recall curve -> Find the area under the curve (AUC-PR).  
For object detection tasks, AP (Average Precision) is equivalent to the area under the Precision-Recall curve.

In VOC 2007, AP is calculated using the 11-point interpolation method, but modern object detection frameworks often use a finer, continuous approximation to compute AP.

In VOC2012, mAP = meanPrecision \* meanRecall

Doing this among all classes and doing the mean of the APs, will get mAP.

For all classes  
mean average precision

Activation function

## Activation Functions

### Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

### tanh

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

### ReLU

$$\max(0, x)$$

### Leaky ReLU

$$\max(0.2x, x)$$

### Softplus

$$\log(1 + \exp(x))$$

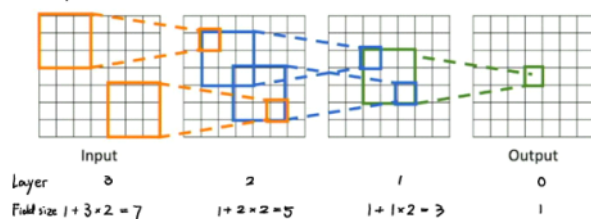
### ELU

$$f(x) = \begin{cases} x & x > 0 \\ \alpha(\exp(x) - 1) & x \leq 0 \end{cases}$$

ReLU is a good default choice for most problems

Receptive fields

## 3\* Receptive Fields



$$L \text{ Layer receptive field size} = 1 + L \times (K - 1)$$

↑  
f1 size

e.g. input  $1000 \times 1000$ ,  $K=3$ ,  $L=?$

$$1000 = 1 + L \times (3 - 1)$$

$$L = \frac{1000 - 1}{2} = 499.5$$

Large imgs need many layers for each outputs to "see" the whole img, sol: (global context)