

dec1: What are localization and detection:

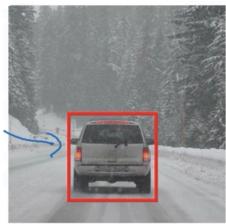
Image classification



car

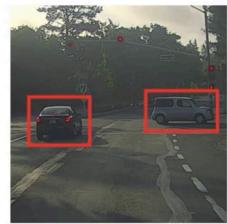
object

Classification with localization



car with
bounding box

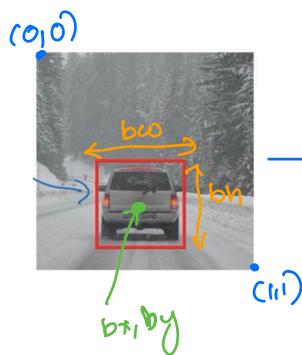
Detection



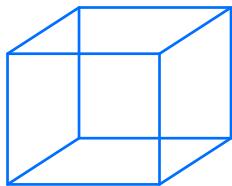
multi object

multi object, bounding box prediction.

Classification with localization:



$$\begin{aligned}bx &= 0.5 \\by &= 0.7 \\bh &= 0.3 \\bw &= 0.4\end{aligned}$$



- - -



→ softmax (4)

↳ b_x, b_y, b_h, b_w
bounding box (4)

1. pedestrian
2. car
3. motorcycle
4. background.



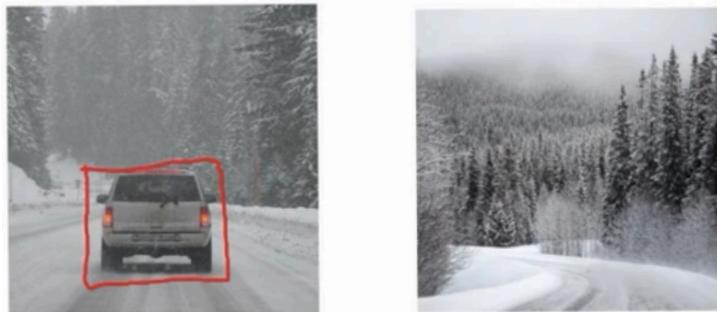
localization.

Target label:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \text{is there any object } (1/0)$$

assume that image has one object

if $x =$



$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

logistic regression loss.

squared error loss.

log likelihood loss

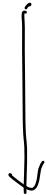
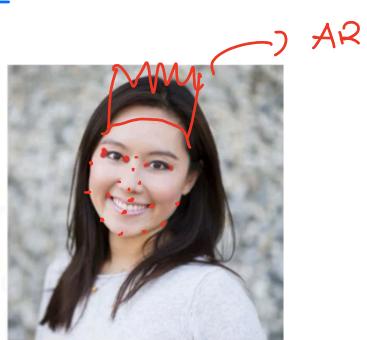
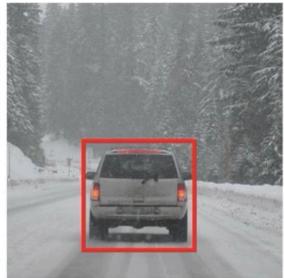
loss: if $y_i = p_c = 1$

$$\alpha(\hat{y}, y) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

loss if $y_i = p_c = 0$

$$\alpha(\hat{y}, y) = (\hat{y}_i - y_i)^2$$

dec2: landmark detection:



b_x, b_y, b_h, b_w

output = 4

l_{1x}, l_{1y}

l_{2x}, l_{2y}

l_{3x}, l_{3y}

l_{4x}, l_{4y}

:

l_{n_x}, l_{n_y} .

pose of person

l_{1x}, l_{1y}

:

l_{32x}, l_{32y} .

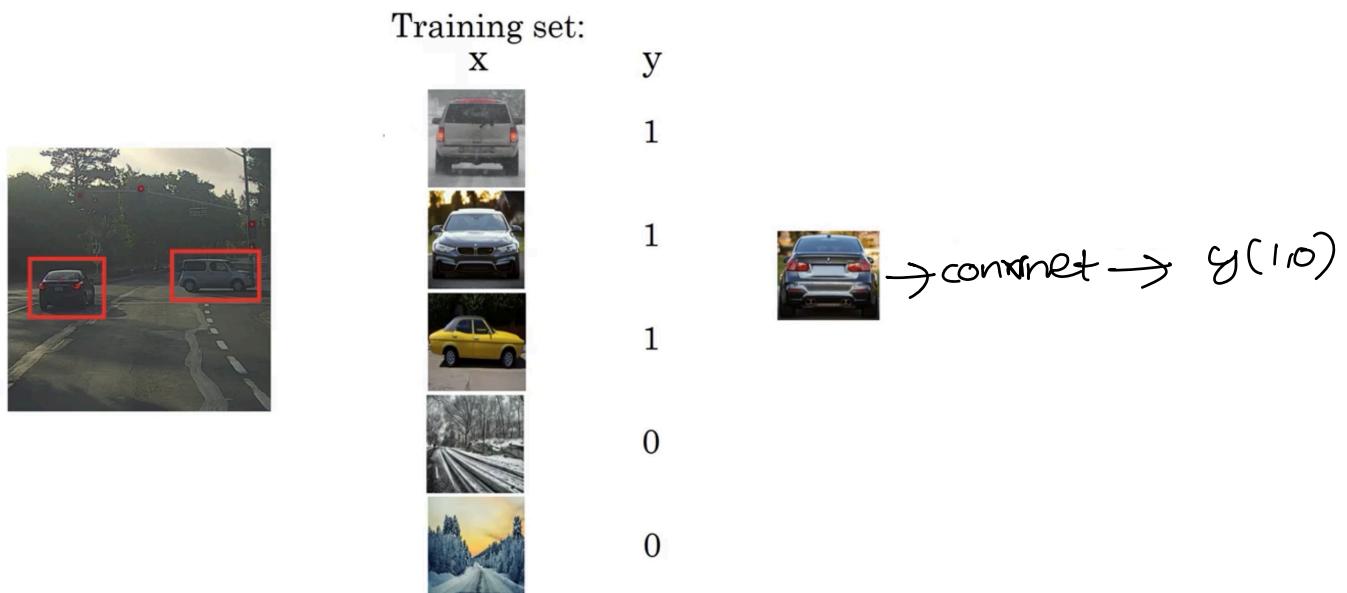
(1-landmark)

output = 129

= sum 2 + 1

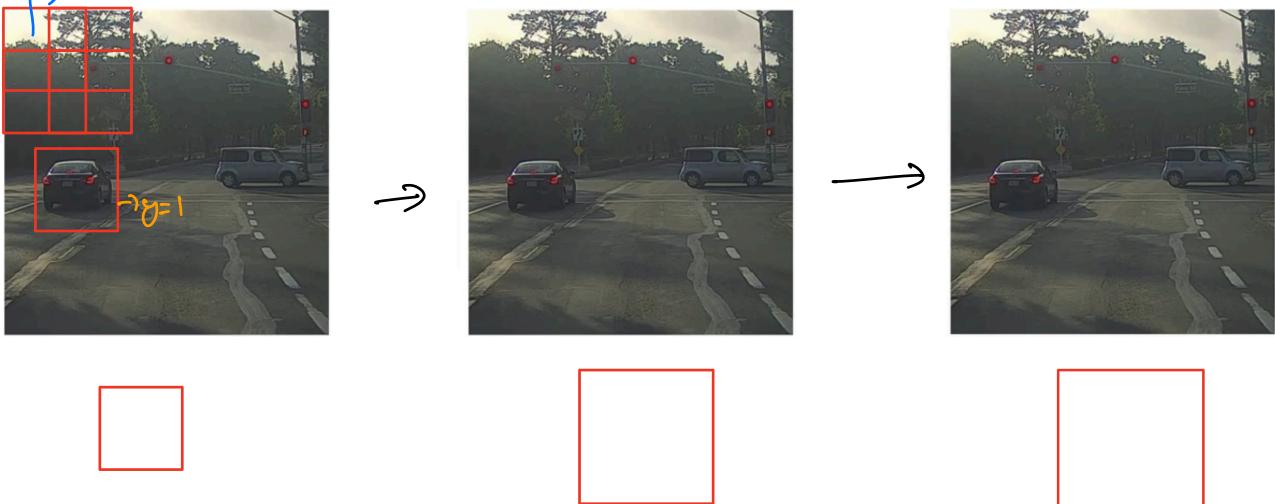
dec3: Object Detection:

Example: Car Detection:



Sliding window detection

\rightarrow convnet

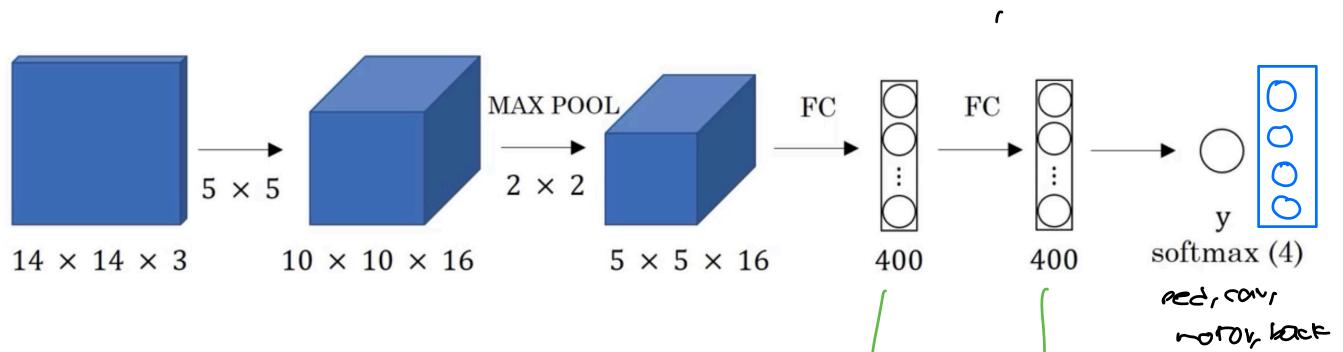


go through every region of input-image and pass through convnet

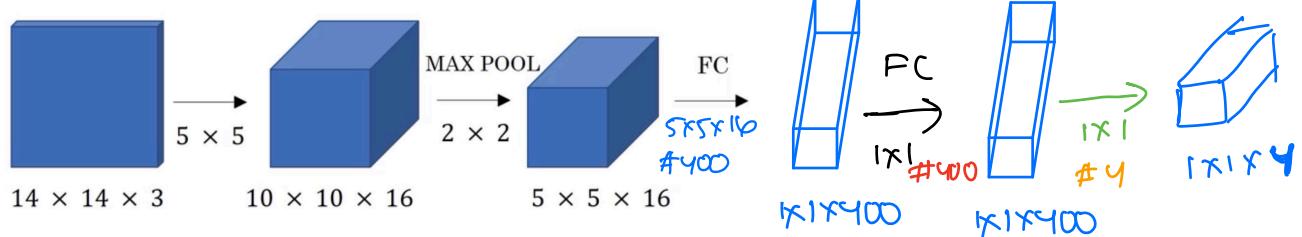
and increase the window size.

disadv: computational cost

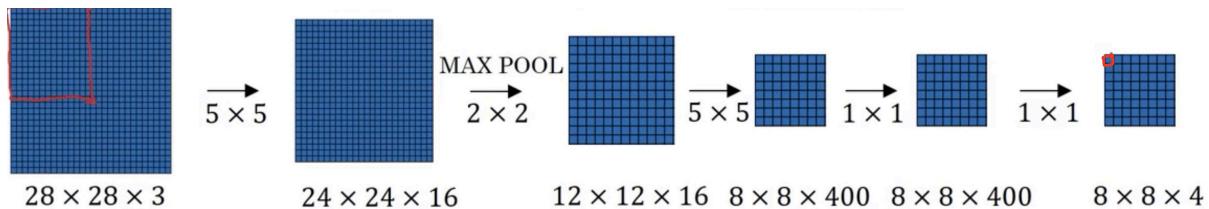
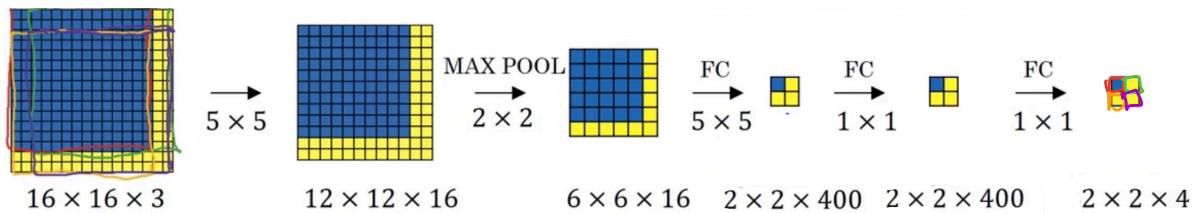
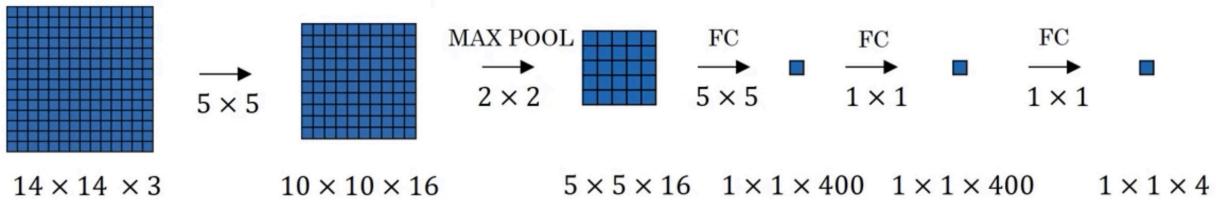
dec4: Turning FC layer into convolution layers:



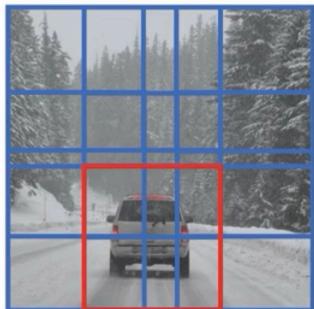
FC \rightarrow CONV



Convolution Implementation of Sliding Window:



pro:



sliding window problem can be solved
by a single forward pass convolution
network.

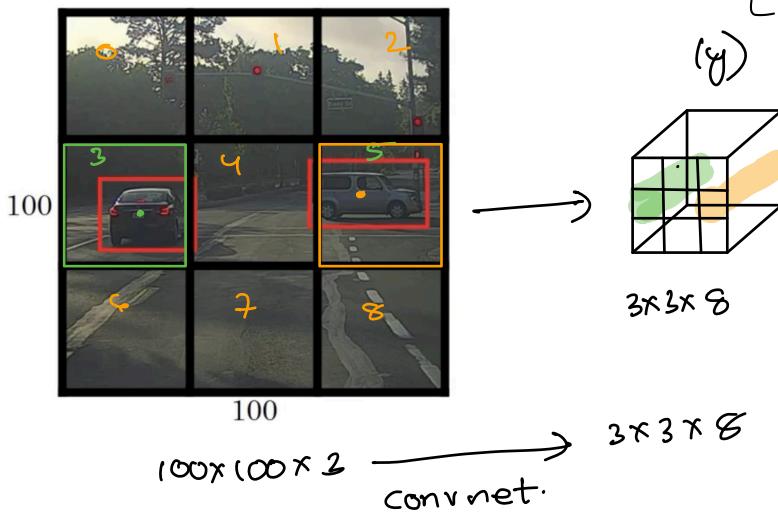
cons:

- bounding box position, not accurate.

dec 5: Output accurate bounding boxes:

YOLO Algorithm: (2015)

you only look once.
- real time object detection.



tables for training

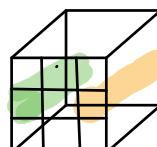
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_w \\ b_h \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_w \\ b_h \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

(y)



3x3x8

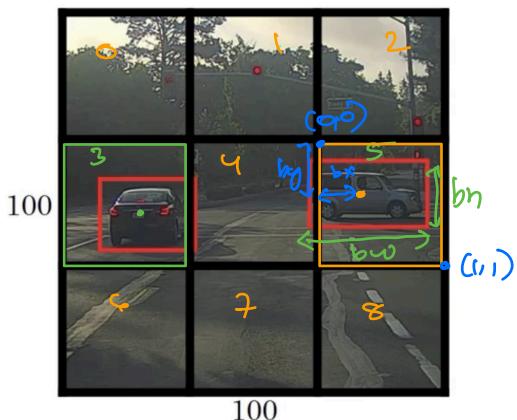
outputs:

→ conv implementation, much less compute and efficient
→ precise bounding box instead relying on the stride and filter size.

→ use bit tinner grid (19x19)
→ outputs (19x19x8)

problem:
multiple object in a single grid. (use smaller grid).

Specify the bounding boxes:

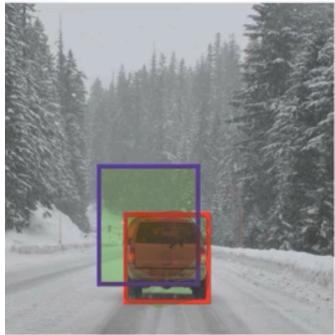


$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_w \\ b_h \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} r \\ 0.4 \\ 0.3 \\ 0.5 \\ 0.9 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

(fraction of cell grid height)
(fraction of cell grid weight)

$b_x, b_y < 1$
 $b_h, b_w \text{ could be } \geq 1$

dec 6: Intersection Over Union: (IoU)



$$IoU = \frac{\text{size of intersection}}{\text{size of union}}$$

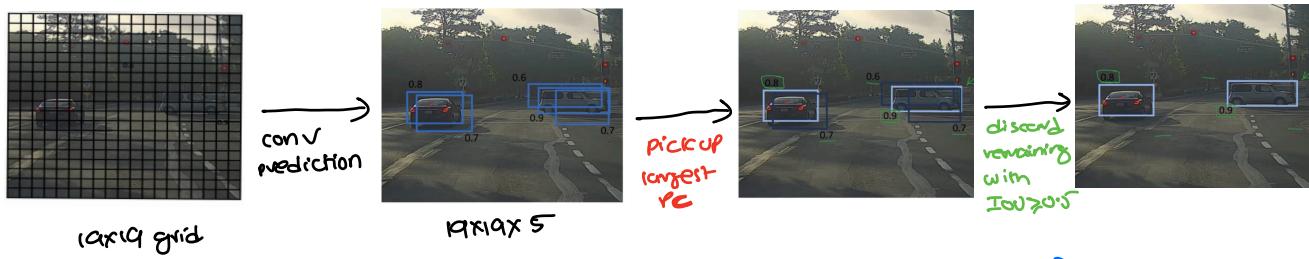
correct if $IoU \geq 0.5$ (0.6)

- gt
- predicted

→ measure or overlap between two bounding boxes

1

dec 7: Non Max Suppression:



Each prediction:

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

if multiple objects (C_1, C_2, C_3)
do non max suppression for all the classes.

① Discard all boxes with $p_c \leq 0.6$

② while there are any remaining boxes:

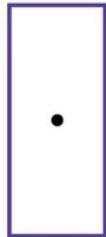
- pick the box with largest p_c output as prediction.

- discard any remaining box with $IoU \geq 0.6$ with the box output in the previous step.

dec 8: Anchor Box



Anchor box 1:



Anchor box 2:

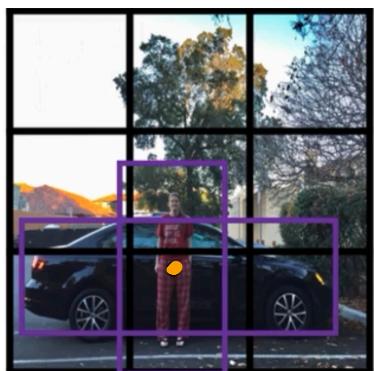


→ Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IOU.
(grid cell, anchor box, object).

output: $(3 \times 3 \times 16) \rightarrow (3 \times 3 \times 2 \times 8)$

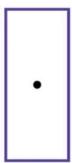
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}_{16 \times 1}$$

Anchor box 1
Anchor box 2



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}_{16 \times 1} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}_{16 \times 1} \xrightarrow{\text{ped & car}} \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}_{16 \times 1} \xrightarrow{\text{only car}}$$

Anchor box 1:

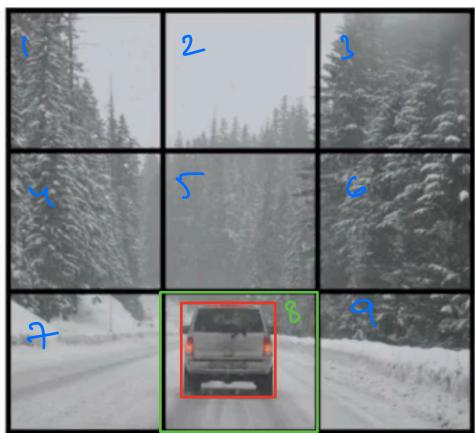


Anchor box 2:



deca: complete YOLO algorithm:

Training



$$x = 100 \times 100 \times 3$$

classes:

1. pedestrian
 2. car
 3. motorcycle

Anchor box 1:

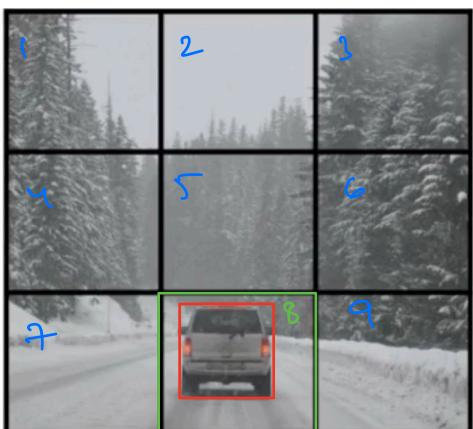
Anchor box:

$$\begin{array}{c}
 \left(1, 2, 3, 4, 5, 6, 7, 9 \right) \\
 \left[\begin{array}{c} PC \\ bY \\ bG \\ bH \\ bW \\ C_1 \\ C_2 \\ C_3 \end{array} \right] = \left[\begin{array}{c} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{array} \right] = \left(8 \right) \\
 \left[\begin{array}{c} PC \\ bY \\ bG \\ bH \\ bW \\ C_1 \\ C_2 \\ C_3 \end{array} \right] = \left[\begin{array}{c} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{array} \right] = \left[\begin{array}{c} bR \\ bG \\ bH \\ bW \\ O \\ I \\ O \end{array} \right]
 \end{array}$$

$$y = 3 \times 3 \times 16$$

output size

Predictions:



3x 3x16

Handwritten notes for the 3rd and 5th chromosomes:

3₄, 5_{6,7,9}) (8)

Black vertical line (3rd chromosome):

- Top circle: 0
- Middle circle: 0
- Bottom circle: 0

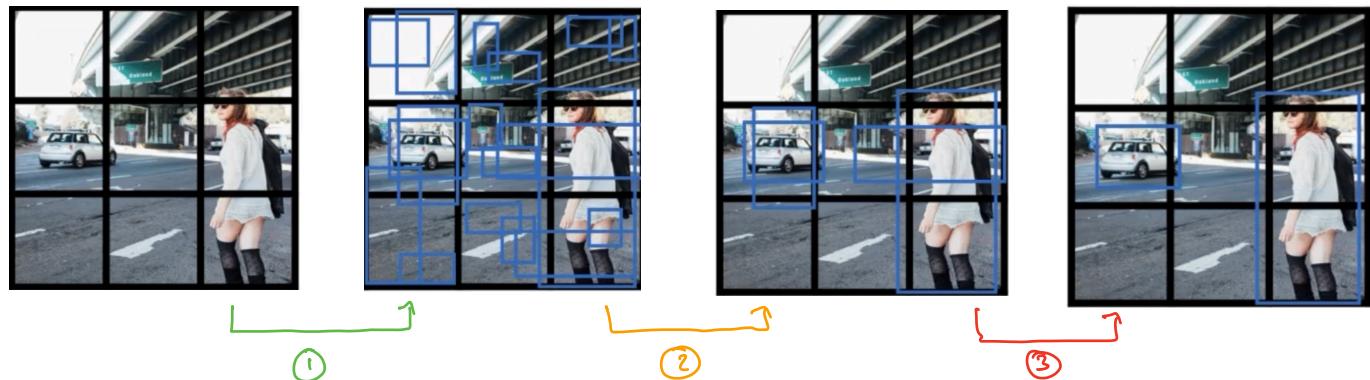
Green vertical line (5th chromosome):

- Top circle: 0
- Middle circle: 0
- Bottom circle: 0

Labels on the green line:

- bx
- by
- bn
- bw
- 0
- 0

Outputting the non-max suppressed outputs:



- ① For each grid cell, get 2 predictions bounding boxes
- ② Get rid of low probability predictions
- ③ for each class (ped, car, motorcycle) use non-max suppression to generate final prediction.

dec10: Region Proposals: (R-CNN)



cNN run's on unwanted regions.

unwanted regions



segmentation
→



run the classifier on segmentation
output blobs.

RCNN : propose regions. classify proposed regions one at a time.
output label + bounding box.

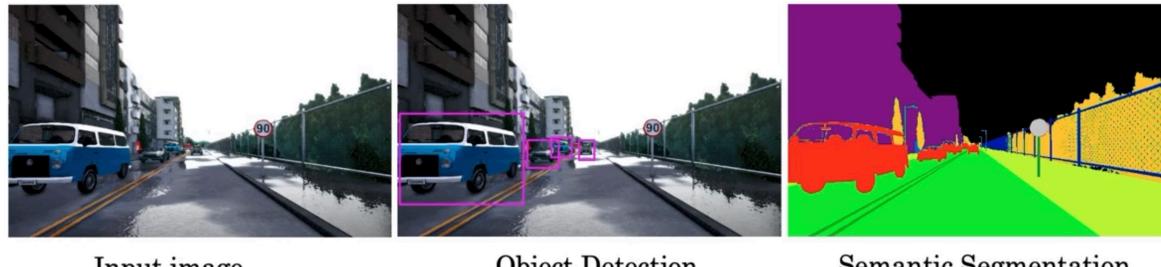
cons: quiet slow.

Fast R-CNN: propose regions. use convolution implementation of sliding windows to classify all the proposed regions.

cons: clustering step in propose regions is slow.

Faster R-CNN: Use convolutional network to propose regions.
(2016)

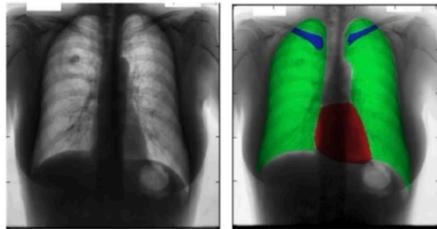
de11: Semantic Segmentation with UNet:



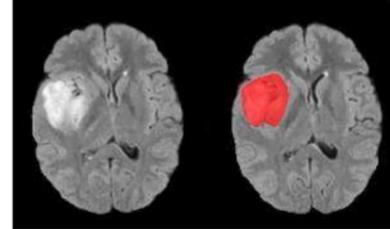
Input image

Object Detection

Semantic Segmentation



Chest X-Ray



Brain MRI

Per Pixel Class labels: (segmentation Map)

1. cav

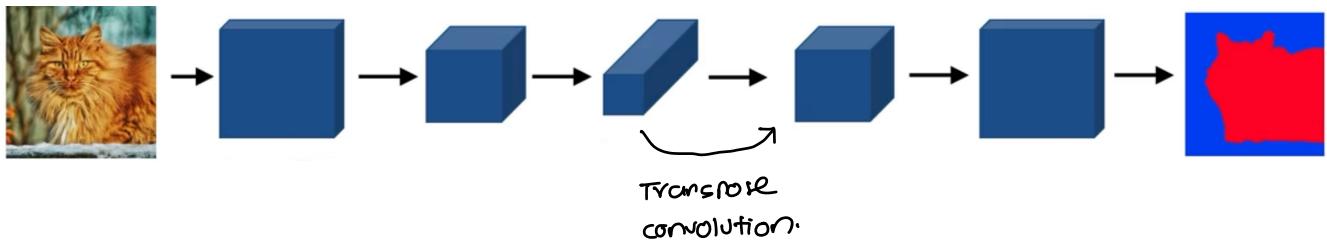
0. background

A red car is parked in front of a building with a grid overlay consisting of 1s and 0s.

1. cav

- 2. building
- 3. road

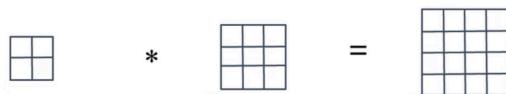
Architecture of Semantic Segmentation:



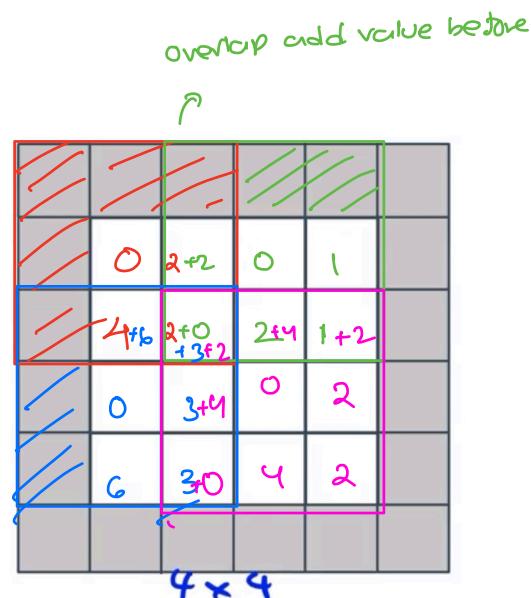
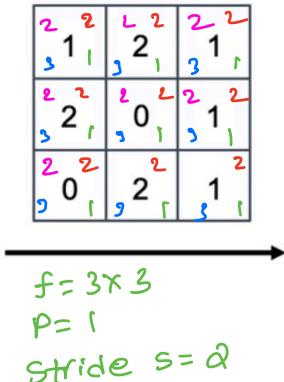
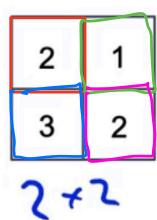
Normal Convolution



Transpose Convolution



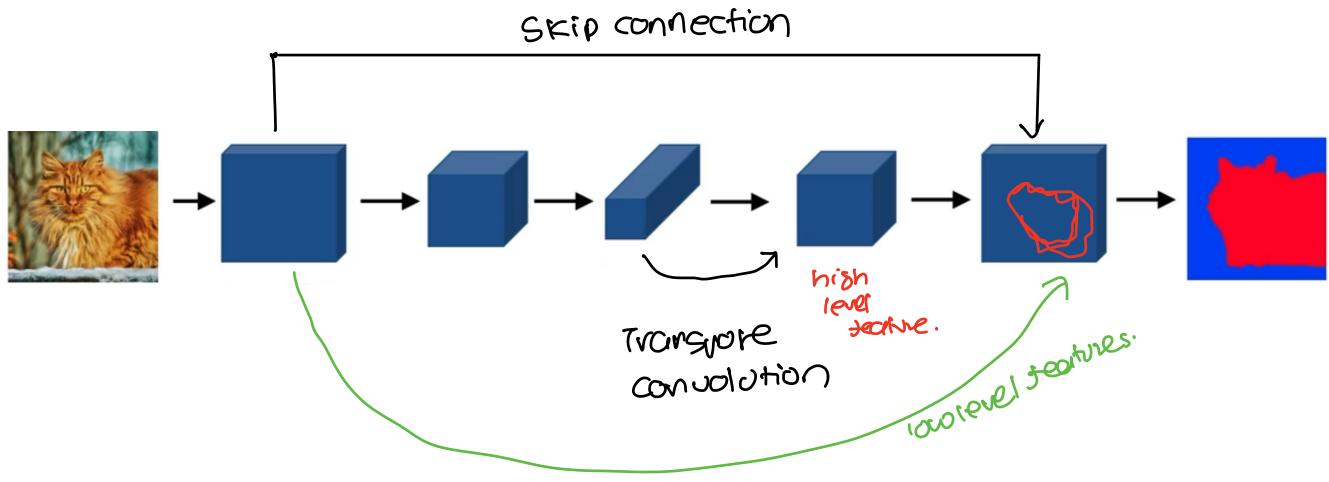
Transpose convolution:



Output:

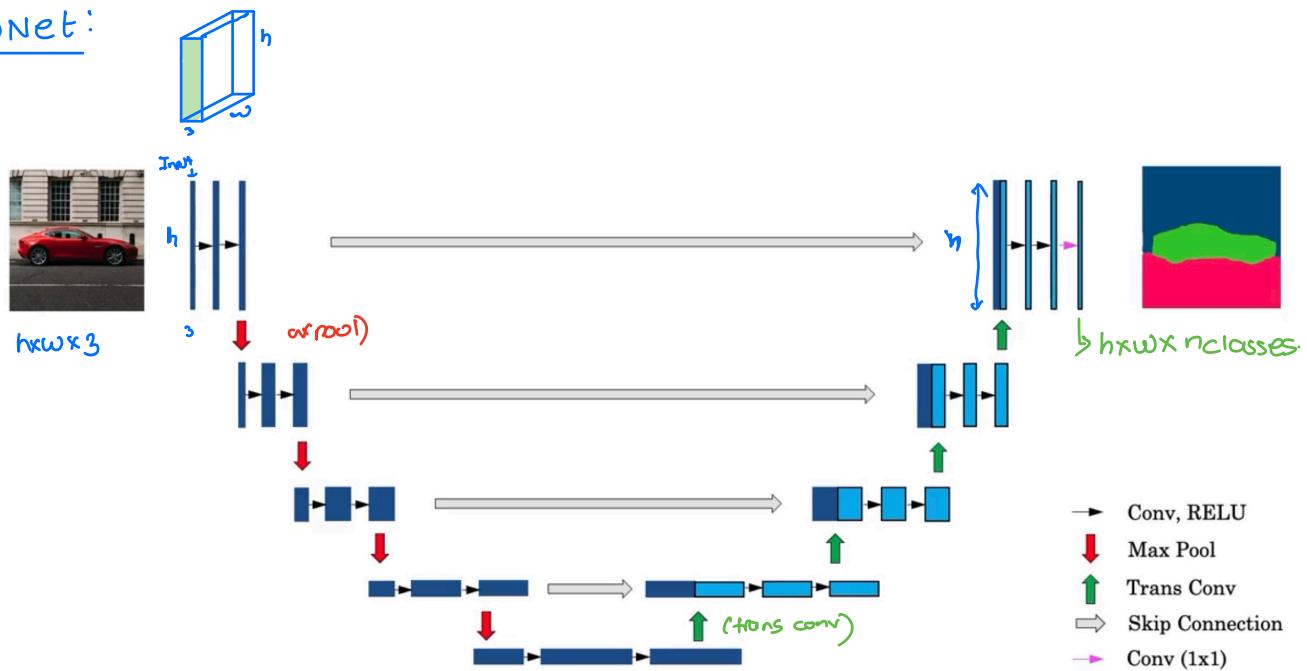
$$\begin{bmatrix} 0 & 4 & 0 & 1 \\ 10 & 7 & 6 & 3 \\ 0 & 7 & 0 & 2 \\ 6 & 3 & 4 & 2 \end{bmatrix} \quad (4 \times 4)$$

learns parameters for filters.



gives the NN, low features into directly, low level features - texture
 high level features - where the cat is.

UNet:



$$\begin{array}{ccc|ccc} & & & 1 & 0 & 1 \\ 1 & 3 & & 0 & 0 & 0 \\ 2 & 4 & & 1 & 0 & 1 \end{array}$$

404

10

۲۰۱

$$x=10 \quad y=0 \quad z=6$$

$$\begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} \quad \begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{matrix}$$

$$\bar{z} = -z$$

$$\cancel{x} = 0$$

$$d = 2$$

	0	0	0	0
	+0	-2+4	-4	+4
	-1	-1+3	-2	-2
	+3	+3		
	0	0	0	0
	-3	-3	-4	-4
	4	-4		

$$\frac{12}{20+8} = \frac{12}{\cancel{2}8} \overset{3}{\cancel{2}} \quad \frac{4}{16+12}$$

4'
287