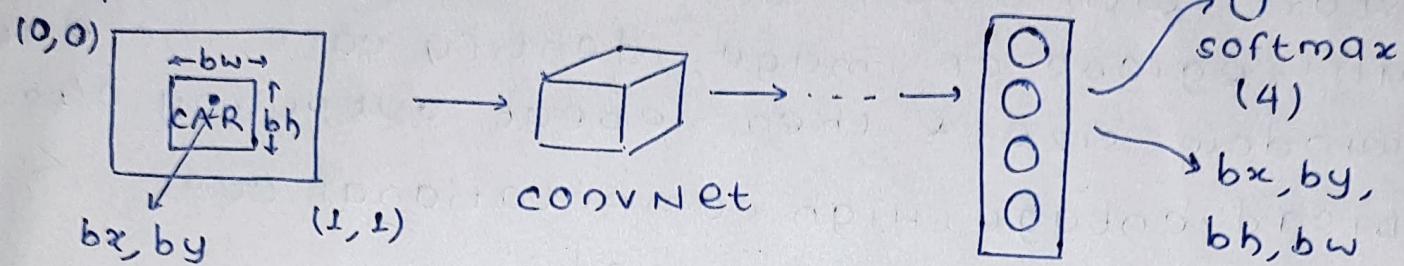


• Image classification

classification with localization → 1 object
Object detection → multiple objects

Input image



- 4 classes:
- 1 - pedestrian
 - 2 - car
 - 3 - motorcycle
 - 4 - background

→ is there an obj?
 $P_c = 1$ for classes (1, 2, 3)

$P_c = 0$ for class 4

* Target label (y) =

P_c
bx
by
bh
bw
c_1
c_2
c_3

IF $P_c = 1$, these will be assigned values.
 If $P_c = 0$, these will be "?" → don't care

* LOSS $\mathcal{L}(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 \Rightarrow y_i = 1$
 (squared error) $= (\hat{y}_1 - y_1)^2 \Rightarrow P_c = y_1 = 0$

• LANDMARK DETECTION (Landmark = imp. points)

e.g. FOR 64 landmarks, no. of output units = 128

1 = face/not face? & $128 = 64 \text{ landmarks}$

$(l_{1x}, l_{1y}), (l_{2x}, l_{2y}), \dots, (l_{64x}, l_{64y})$

* Landmark tables → consistent across images

* OBJECT DETECTION

x

y

Labeled training set \Rightarrow closely cropped images of cars

o/l

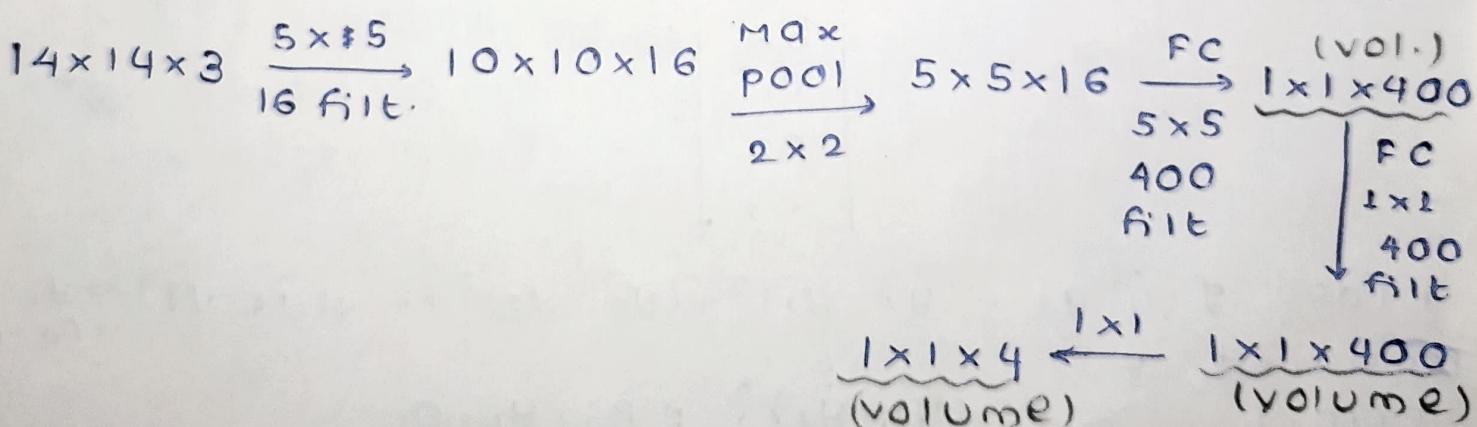
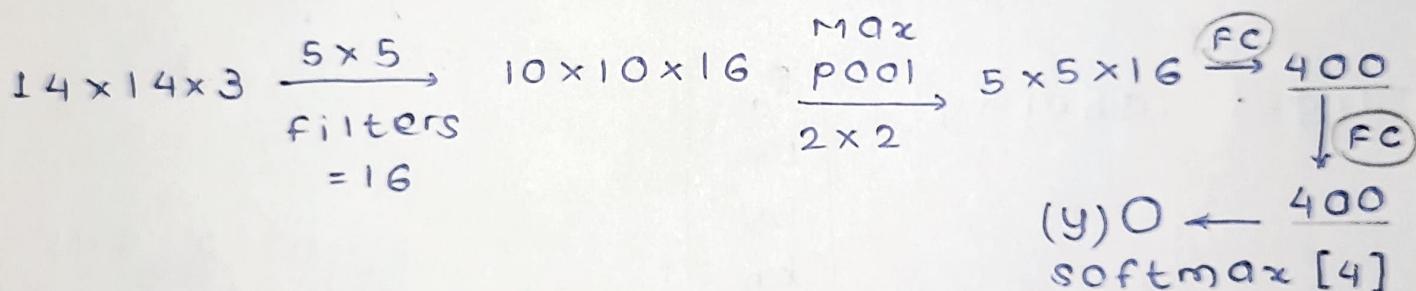
then input image \rightarrow convnet \rightarrow y

• Sliding window detection

start with small window size, go through all regions of image, identify cars, increase window size & then repeat outputting o/l's

• Disadvantage: High computational cost
slow speed

* Fully connected layer representation as convolutional layer

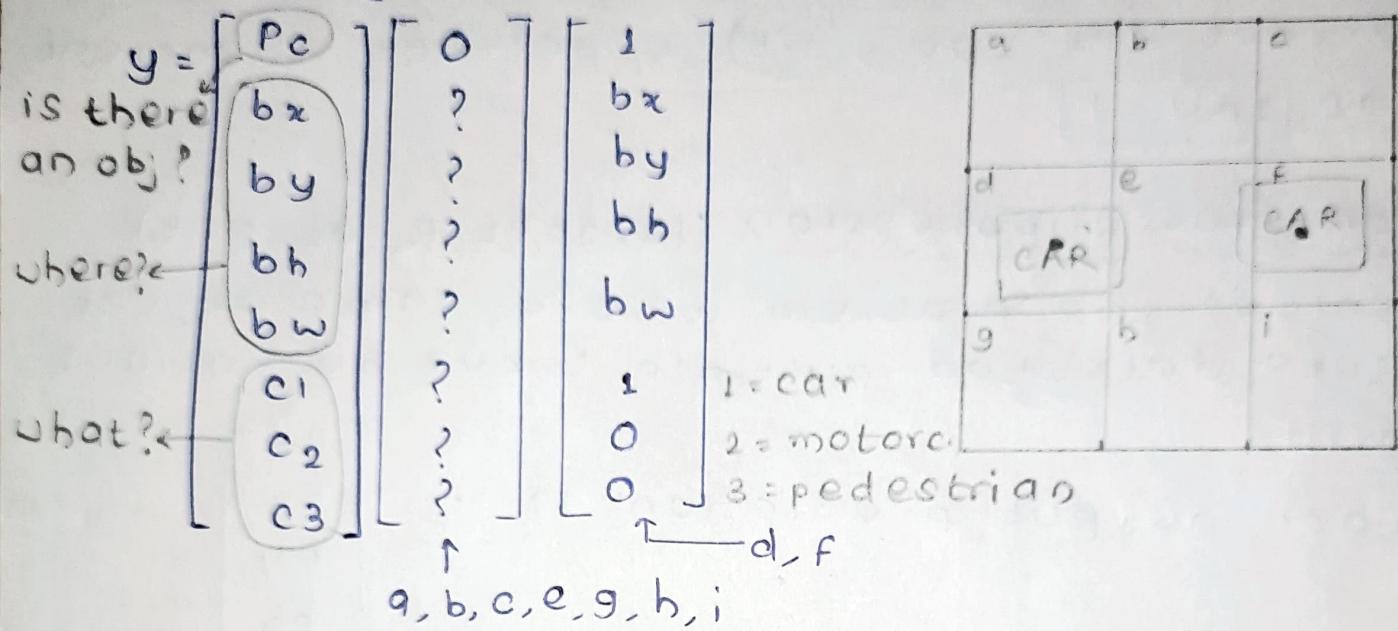


• Sliding windows using convolution
 \Rightarrow process entire image at once.

[problem] \Rightarrow bounding boxes are not accurate

• YOLO ALGORITHM (you only look) once 2

- To achieve accurate bounding boxes
- Input images are layered with a grid
- Training labels for each grid:



Target output = $3 \times 3 \times 8$ dimensional target vector
grid cells $y = \begin{bmatrix} . \\ . \\ . \end{bmatrix}$

* **TRAINING** FAST & EFFICIENT (useful for real time.)

Input [x] \rightarrow conv \rightarrow pool $\rightarrow \dots \rightarrow$ Output [y]
 $100 \times 100 \times 3$ $3 \times 3 \times 8$

* **Specify bounding box** $y = \begin{bmatrix} 1 \\ bx \\ by \\ bh \\ bw \\ 0 \\ 1 \end{bmatrix}$

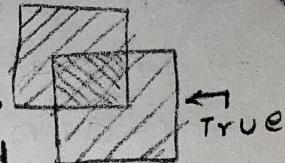
$\left. \begin{array}{l} bx \\ by \\ bh \\ bw \end{array} \right\} 0.4 \quad \left. \begin{array}{l} 0.4 \\ 0.3 \\ 0.5 \\ 0.9 \end{array} \right\} b/w \text{ or } 1$

$\left. \begin{array}{l} 0 \\ 1 \end{array} \right\} \text{can be } > 1$

YOLO

NOTE → only one forward propagation pass

- Intersection over union → accuracy algorithm
- to evaluate object detection algorithm
- compute $\text{IOU} = \frac{\text{size of intersection}}{\text{size of union}}$
- correct for $\text{IOU} > 0.5$ & perfect alignment has $\boxed{\text{IOU} = 1}$



- Non-max suppression (largest P_c output)
- solves the problem where same object gets detected multiple times by various grids
- Each output prediction is $y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$

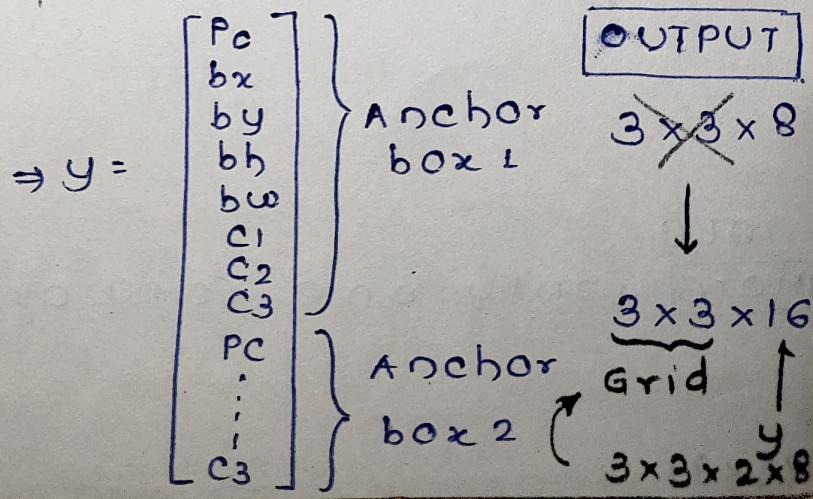
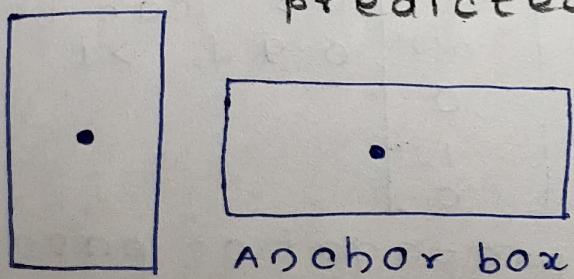
$\bullet C_1, C_2, C_3$
 not included as we are detecting car only

discard boxes with $\boxed{P_c \leq 0.6}$

For remaining boxes:

- pick box with largest P_c , output that as a prediction
- discard any remaining box with $\text{IOU} > 0.5$ with box output in previous step. (largest P_c)

- Anchor boxes (to detect 2 objects in one grid) → suitable shape w.r.t output to be predicted.



\uparrow
 \uparrow

• YOLO ALGORITHM

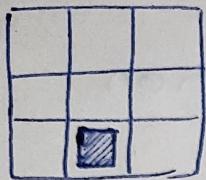
? \Rightarrow Random nos.
or noise

* Training

INPUT

Image & 3×3

Grid

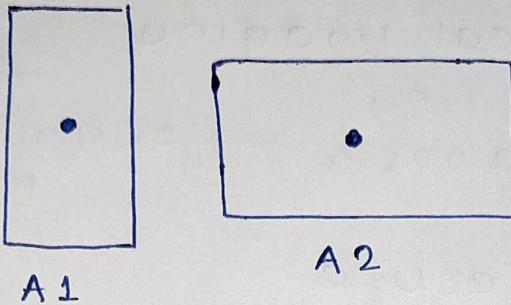


1 - pedestrian

2 - car

3 - motorcycle

$y =$



p_c	A ₁
b_x	
b_y	
b_b	
b_w	
c_1	
c_2	
c_3	
p_c	
c_3	

p_c	A ₂
b_x	
b_y	
b_b	
b_w	
c_1	
c_2	
c_3	
p_c	
c_3	

OUTPUT: $3 \times 3 \times 16$ or $3 \times 3 \times 2 \times 8$

* Outputting non-max suppressed outputs

- each grid cell = 2 predicted bounding boxes
- eliminate low probability predictions ($p_c \leq 0.5$)
- use non-max suppression to generate final output prediction for each class

* Region proposal (R-CNN)

- perform segmentation & run classifier (sliding window) on only those blobs

R-CNN : propose regions, classify proposed regions one at a time, output
→ label + bounding box

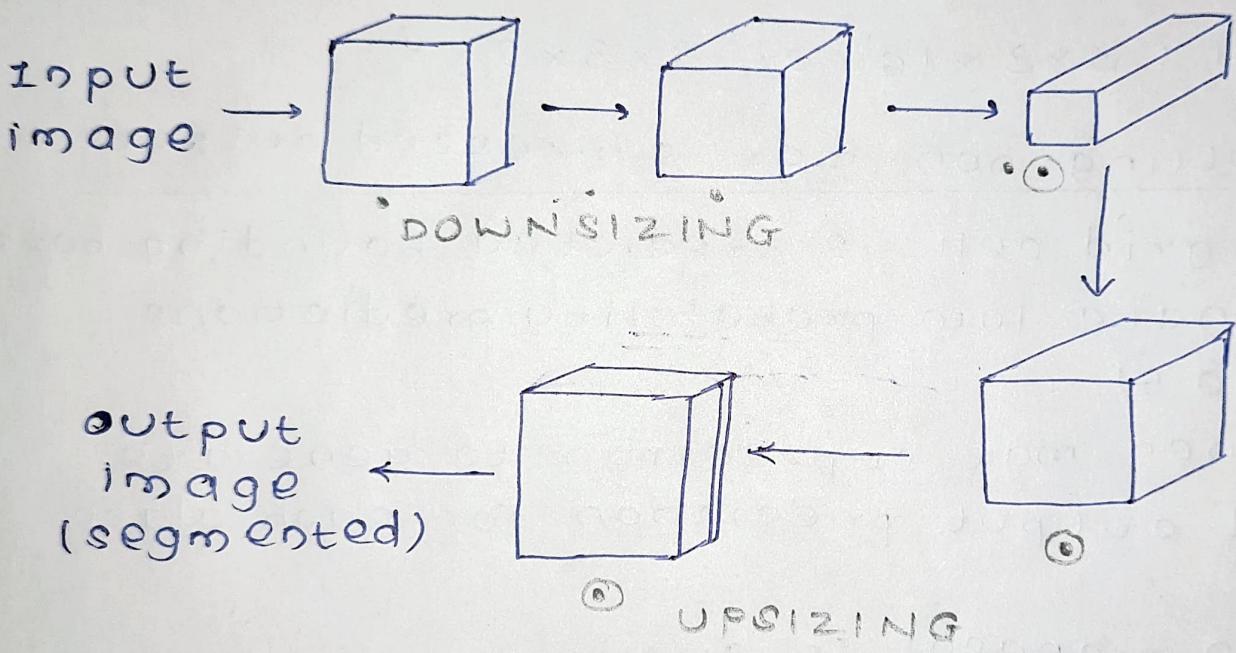
Fast R-CNN : propose regions, classify using convolutional implementation of sliding windows

Faster R-CNN : use CNN to propose regions.

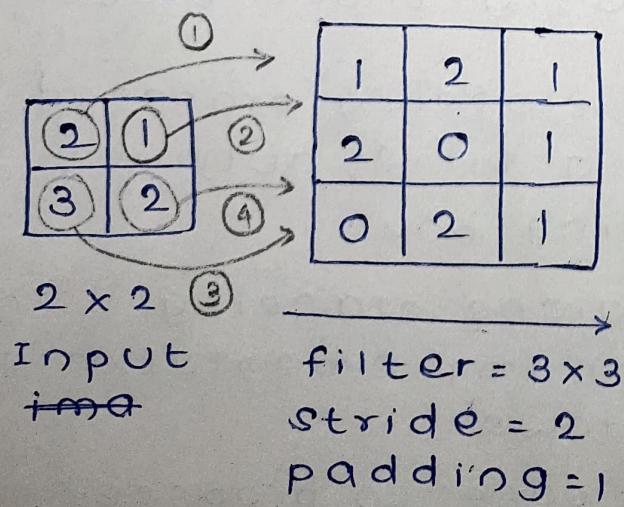
SEMANTIC SEGMENTATION WITH U-NET

- pixel categorization / labelling
- draw exact outline around objects over use of bounding boxes.
- uses : (i) self driving cars
(ii) medical imaging
- input image → specify classes → segmentation map.

→ U-Net Architecture



* TRANSPOSE CONVOLUTION



0	0	2+2	0	1
4+6	2+0	2+4	1+2	
10	+8+2	6	3	
0	3+4	0	2	
0	7	0	2	
6	3+0	4	2	
6	3	4	2	

4x4 output (pad=1)

• U-NET

Input image → convolution → ENCODER → transpose convolutions → DECODER → output image

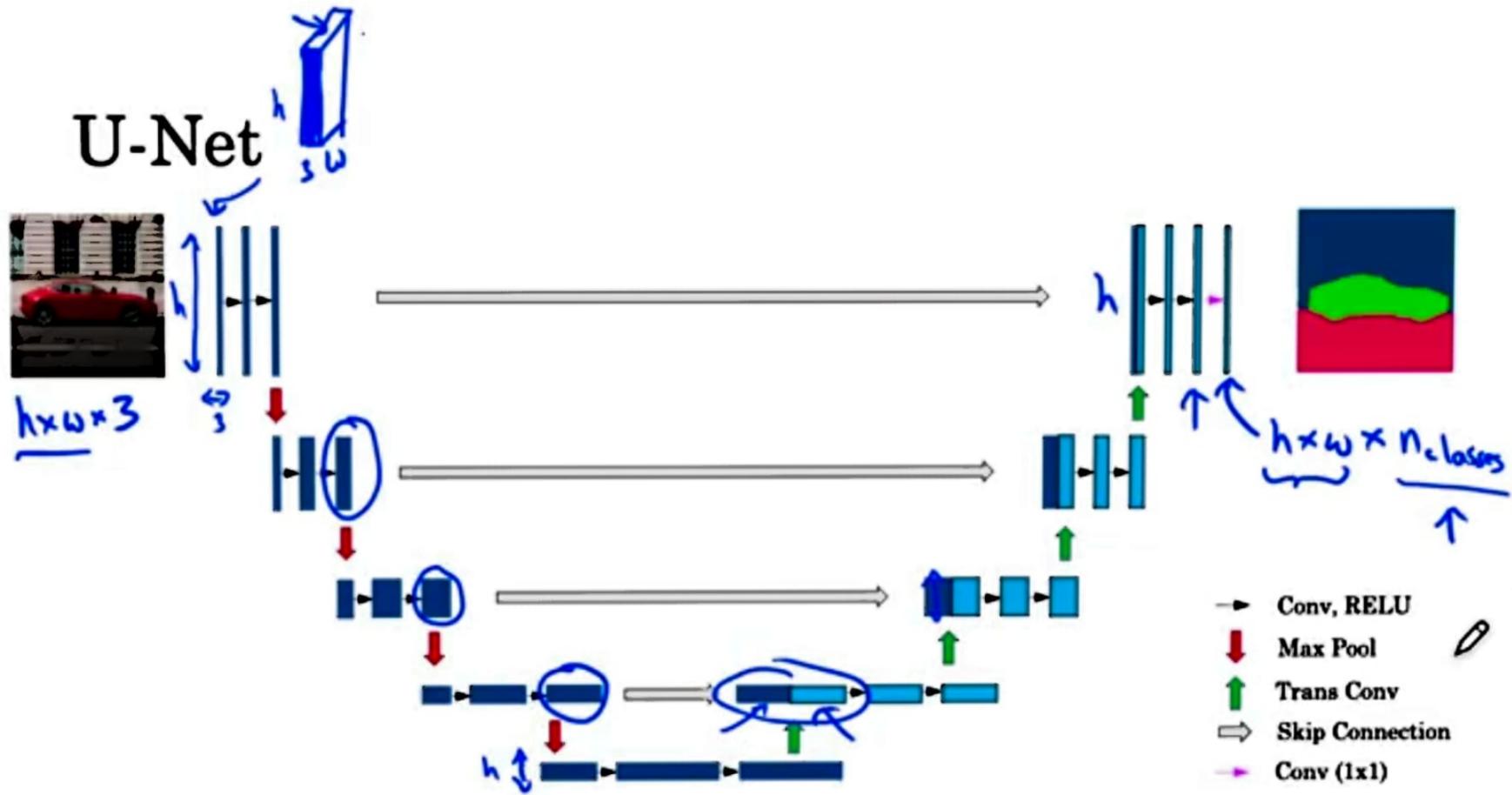
* Skip connections (gives info from encoder to decoder)

Input image $b \times w \times 3$ → Encoder → Decoder → segmentation map $b \times w \times n$ classes

- for 3 classes
 $O/P = b \times w \times 3$

- Steps:
- (i) conv, RELU →
 - (ii) Max pool ↓
 - (iii) transconv ↑
 - (iv) skip connection ⇒
 - (v) conv ($L \times L$) → ... (just before O/P)

U-Net



U-Net

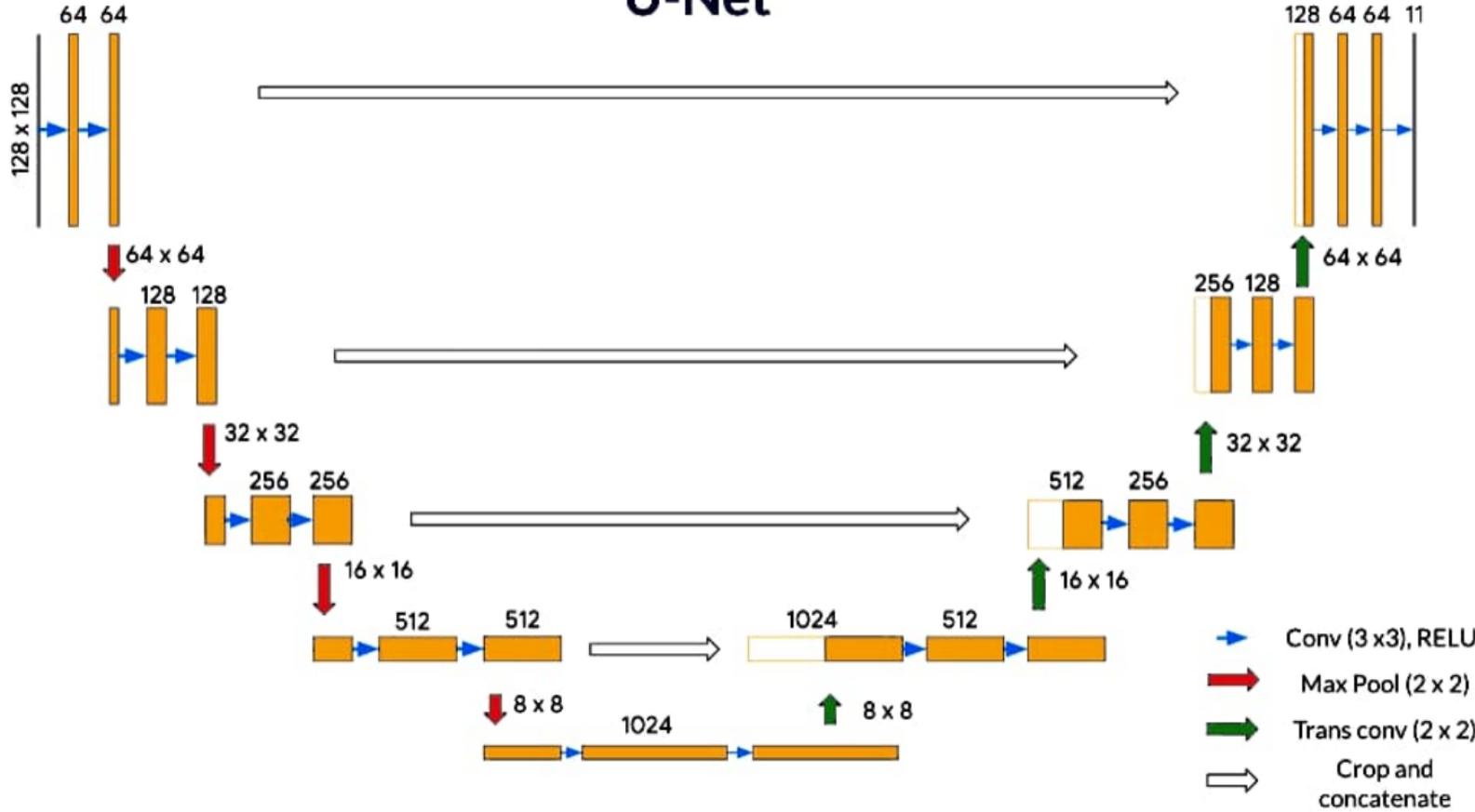


Figure 2 : U-Net Architecture

Encoder

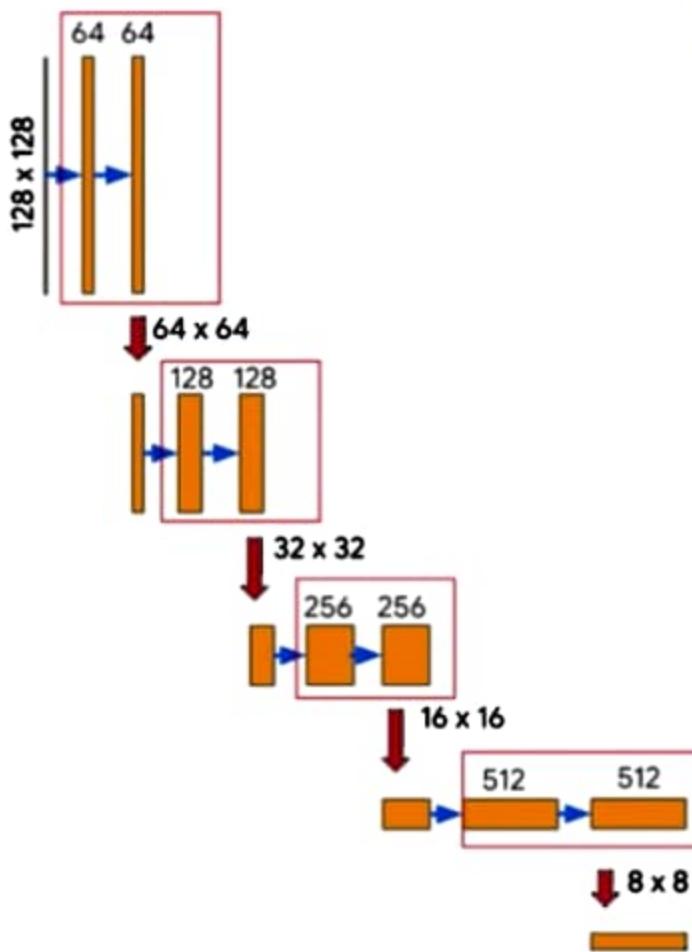


Figure 3: The U-Net Encoder up close