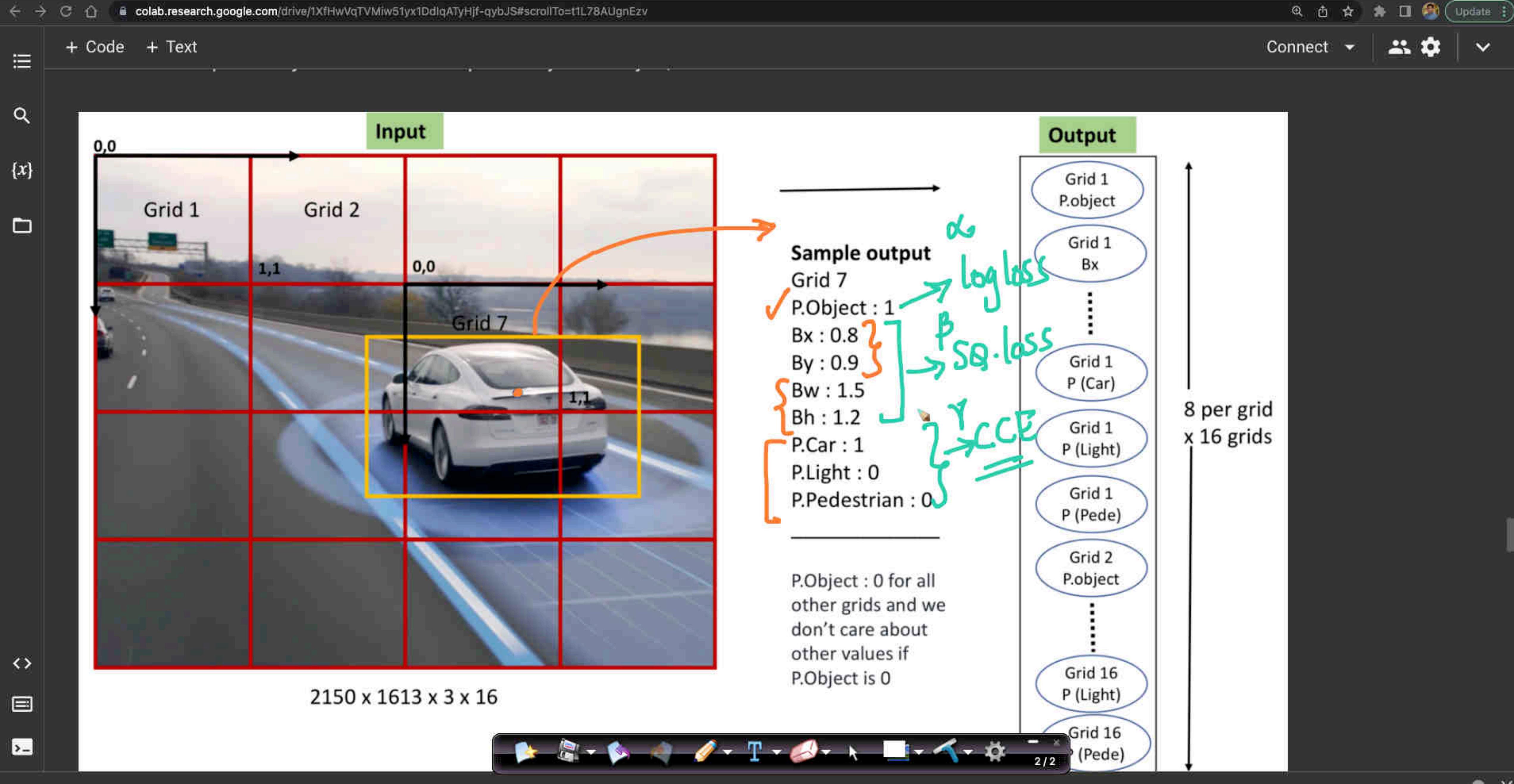


Agenda:

- ① YOLO , YOLO V3
- ② Overview of RetinaNet (Post read)
- ③ Segmentation & metrics
- ④ FCN + Transposed conv
- ⑤ U-net + Upsampling 2D
- ⑥ Dilated Conv ...



Network Architecture

The base model is similar to [GoogLeNet](#) with inception module replaced by 1x1 and 3x3 conv layers. The final prediction of shape $S \times S \times (5B + K)$ is produced by two fully connected layers over the whole conv feature map.

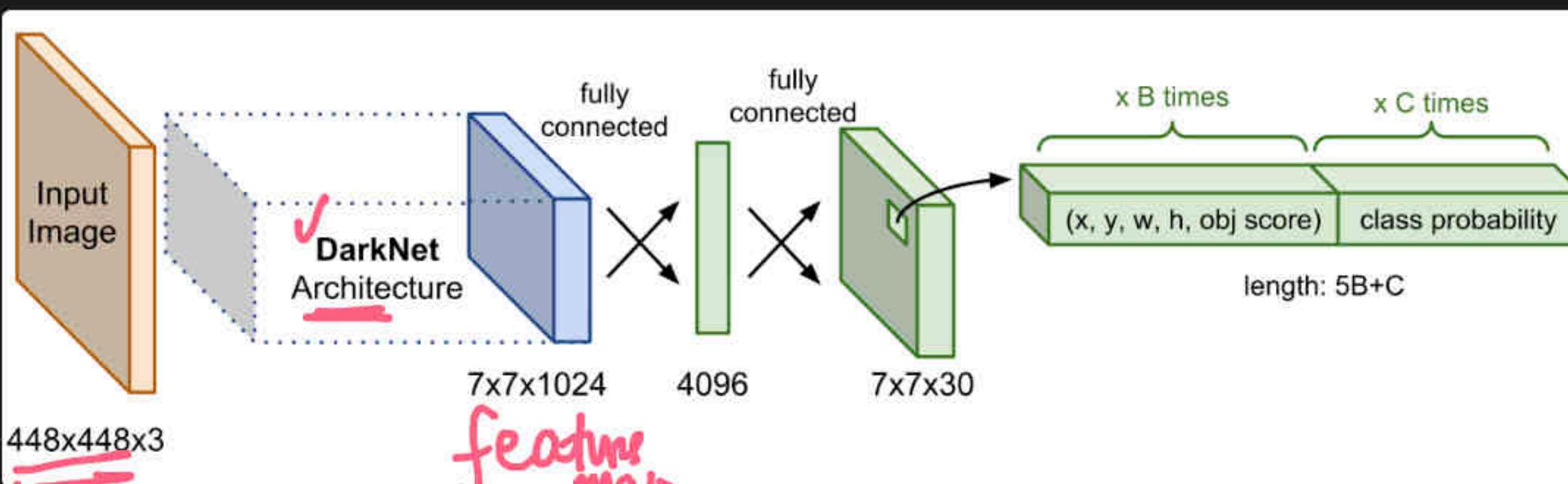


Fig. 2. The network architecture of YOLO.

Loss Function

The loss consists of two parts, the *localization loss* for bounding box offset prediction and the *classification loss* for conditional class probabilities. Both parts are computed as the sum of squared errors. Two scale parameters are used to control how much we want to increase the loss from bounding box coordinate predictions (λ_{coord}) and how much we want to decrease the loss of confidence score predictions for boxes without objects (λ_{noobj}). Down-weighting the loss contributed by background instances.

In the paper, the model sets $\lambda_{\text{coord}} = 5$ and $\lambda_{\text{noobj}} = 0.5$.

Network Architecture

The base model is similar to [GoogLeNet](#) with inception module replaced by 1x1 and 3x3 conv layers. The final prediction of shape $S \times S \times (5B + K)$ is produced by two fully connected layers over the whole conv feature map.

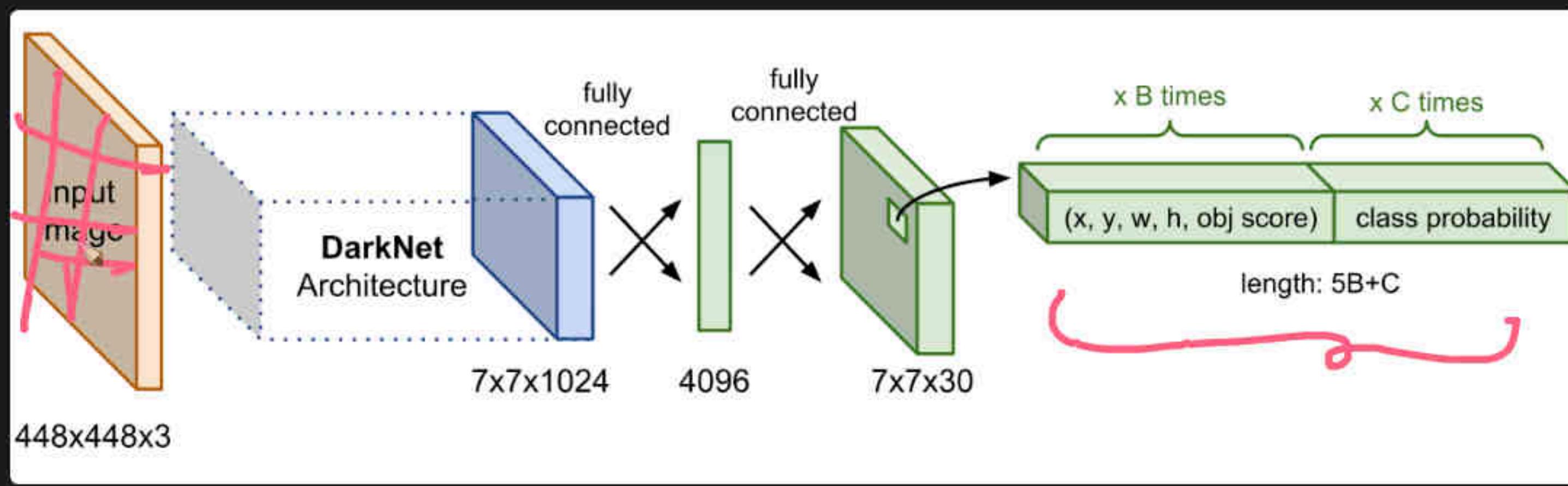


Fig. 2. The network architecture of YOLO.

Loss Function

The loss consists of two parts, the *localization loss* for bounding box offset prediction and the *classification loss* for conditional class probabilities. Both parts are computed as the sum of squared errors. Two scale parameters are used to control how much we want to increase the loss from bounding box coordinate predictions (λ_{coord}) and how much we want to decrease the loss of confidence score predictions for boxes without objects (λ_{noobj}). Down-weighting the loss contributed by background instances is also possible.



Network Architecture

RPN X
Faster RCNN
Yolo

The base model is similar to [GoogLeNet](#) with inception module replaced by 1x1 and 3x3 conv layers. The final prediction of shape $S \times S \times (5B + K)$ is produced by two fully connected layers over the whole conv feature map.

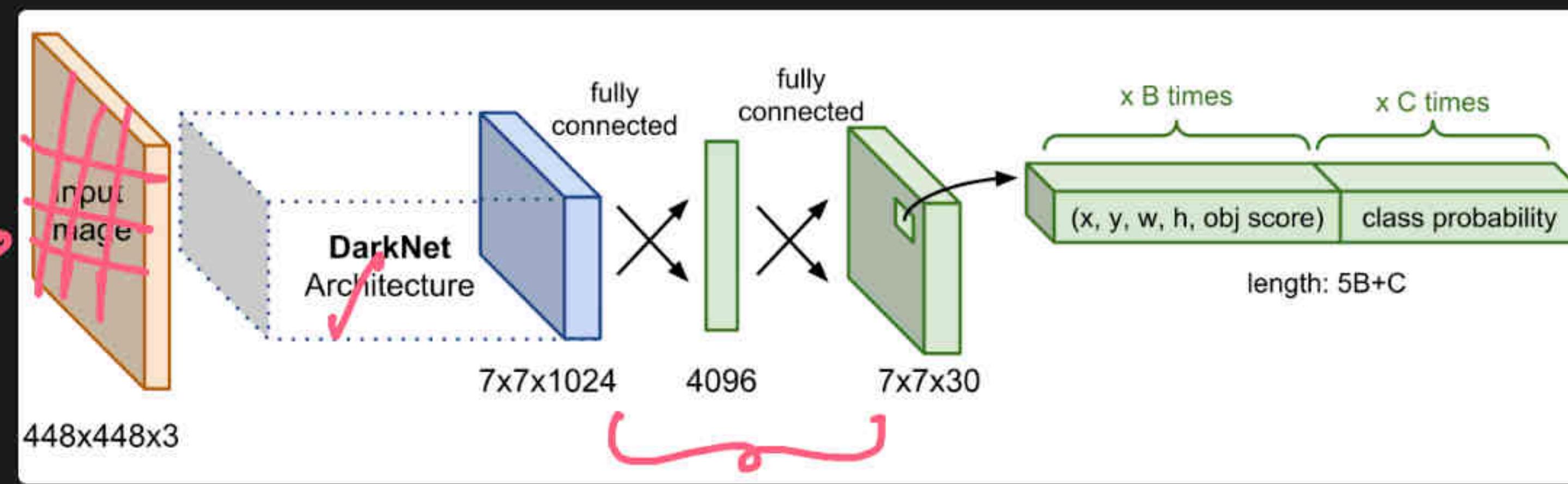
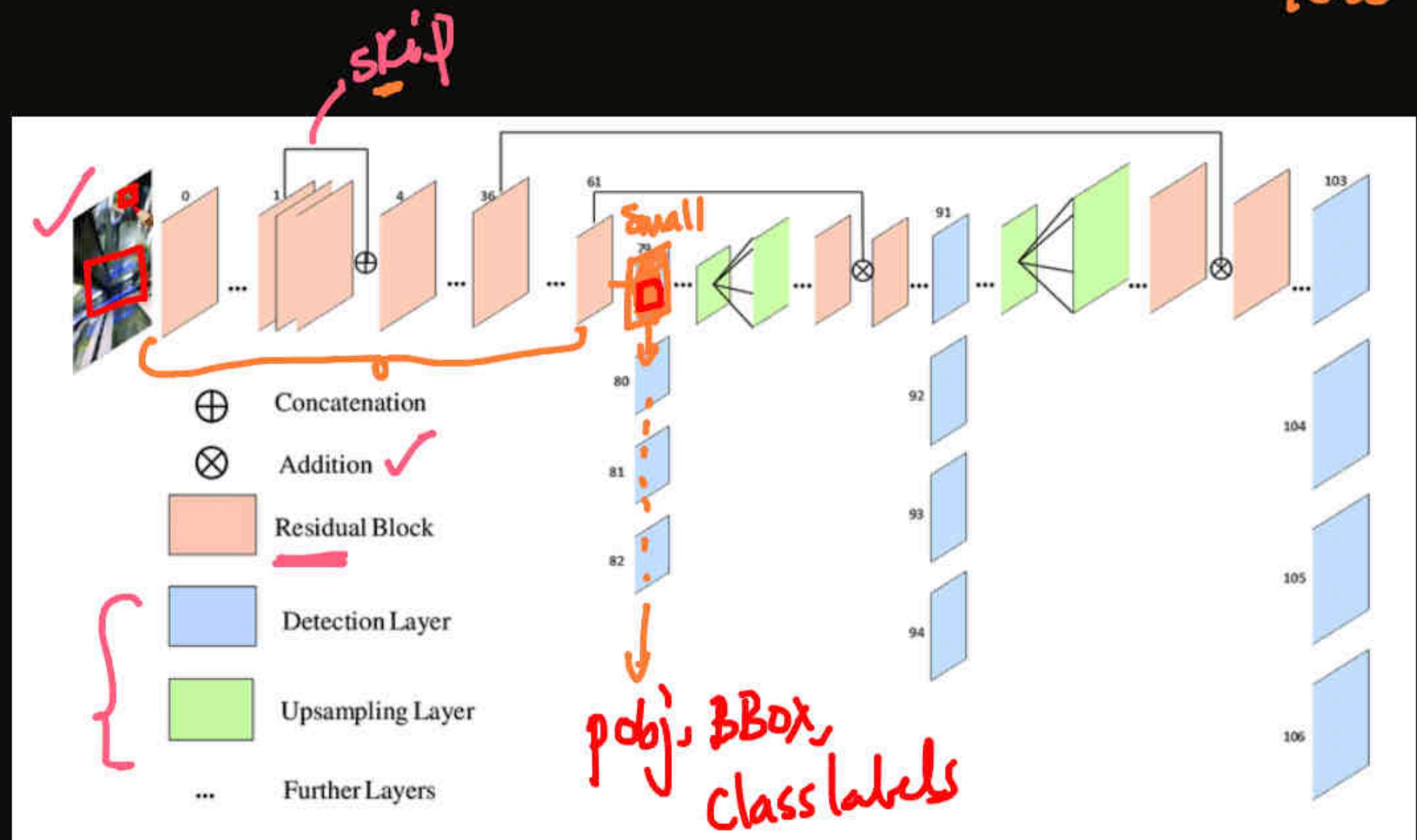


Fig. 2. The network architecture of YOLO.

Loss Function

The loss consists of two parts, the *localization loss* for bounding box offset prediction and the *classification loss* for conditional class probabilities. Both parts are computed as the sum of squared errors. Two scale parameters are used to control how much we want to increase the loss from bounding box coordinate predictions (λ_{coord}) and how much we want to decrease the loss of confidence score predictions for boxes without objects (λ_{noobj}). Down-weighting the loss contributed by background instances is also possible.





Network Architecture

The base model is similar to [GoogLeNet](#) with inception module replaced by 1×1 and 3×3 conv layers. The final prediction of shape $S \times S \times (5B + K)$ is produced by two fully connected layers over the whole conv feature map.

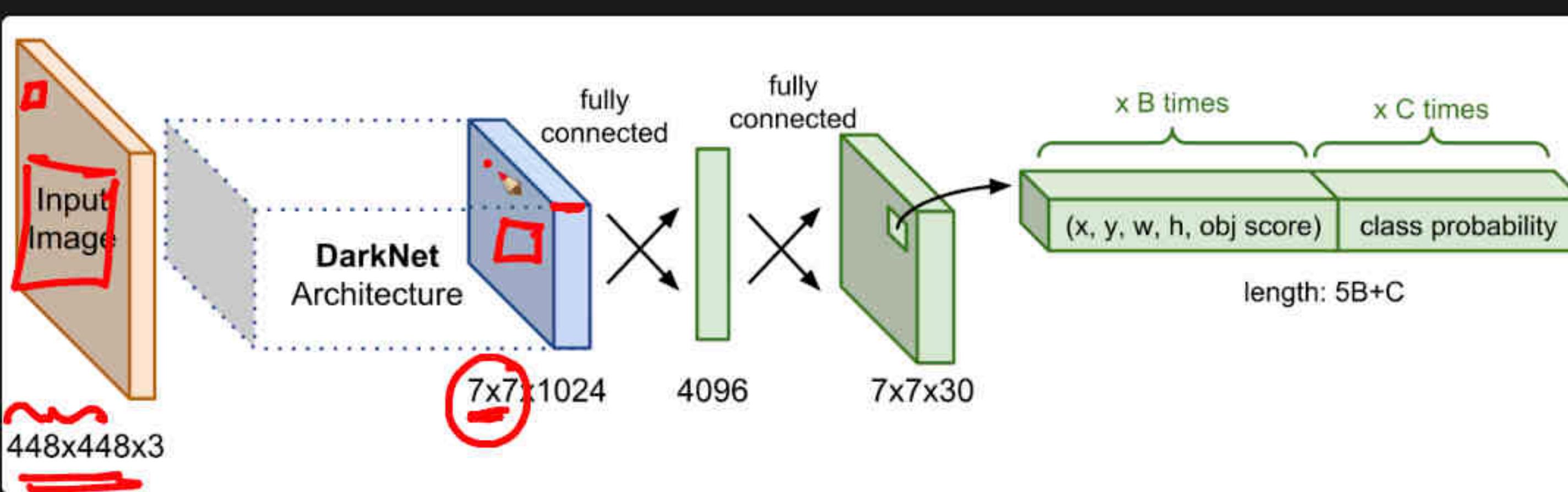
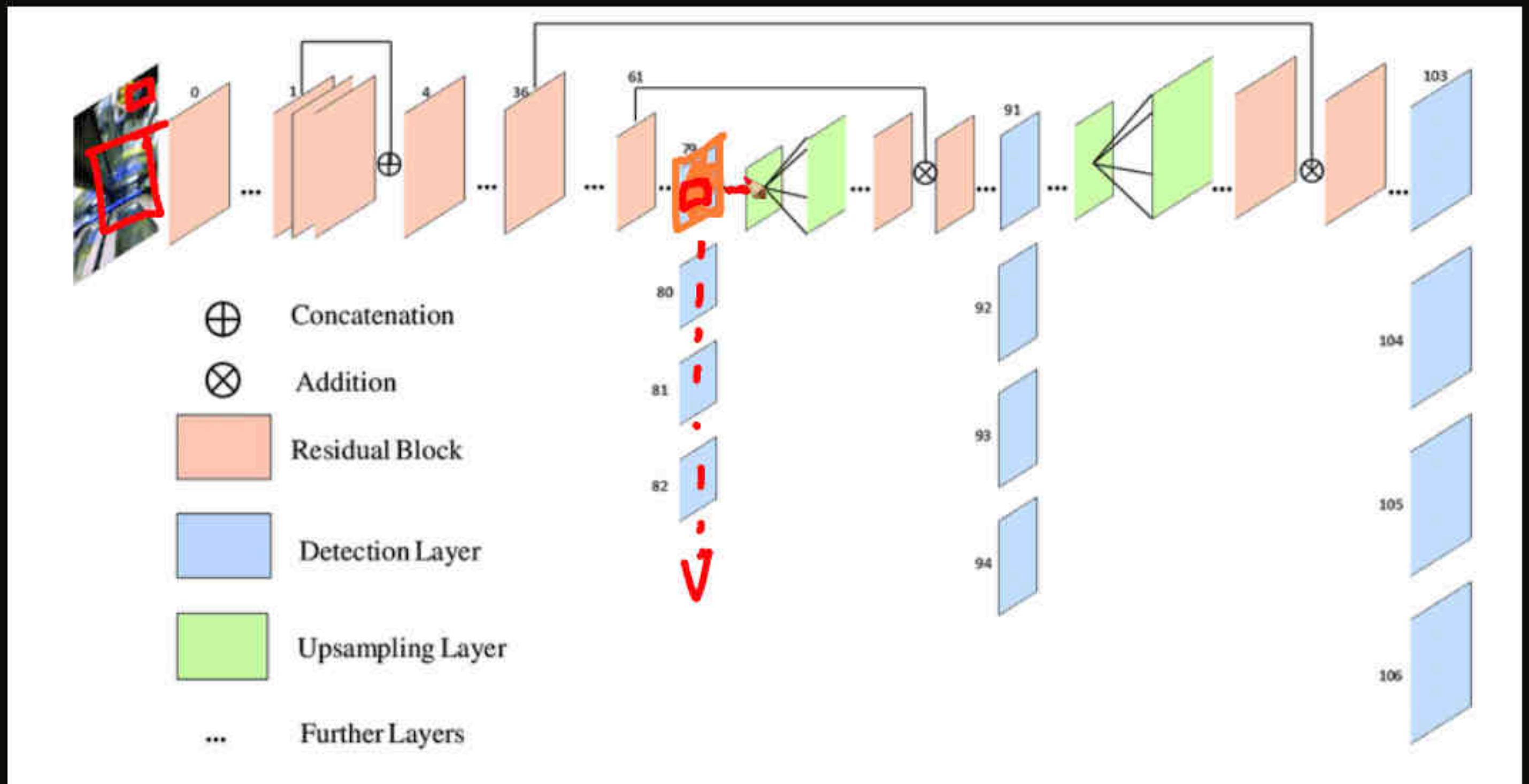


Fig. 2. The network architecture of YOLO.

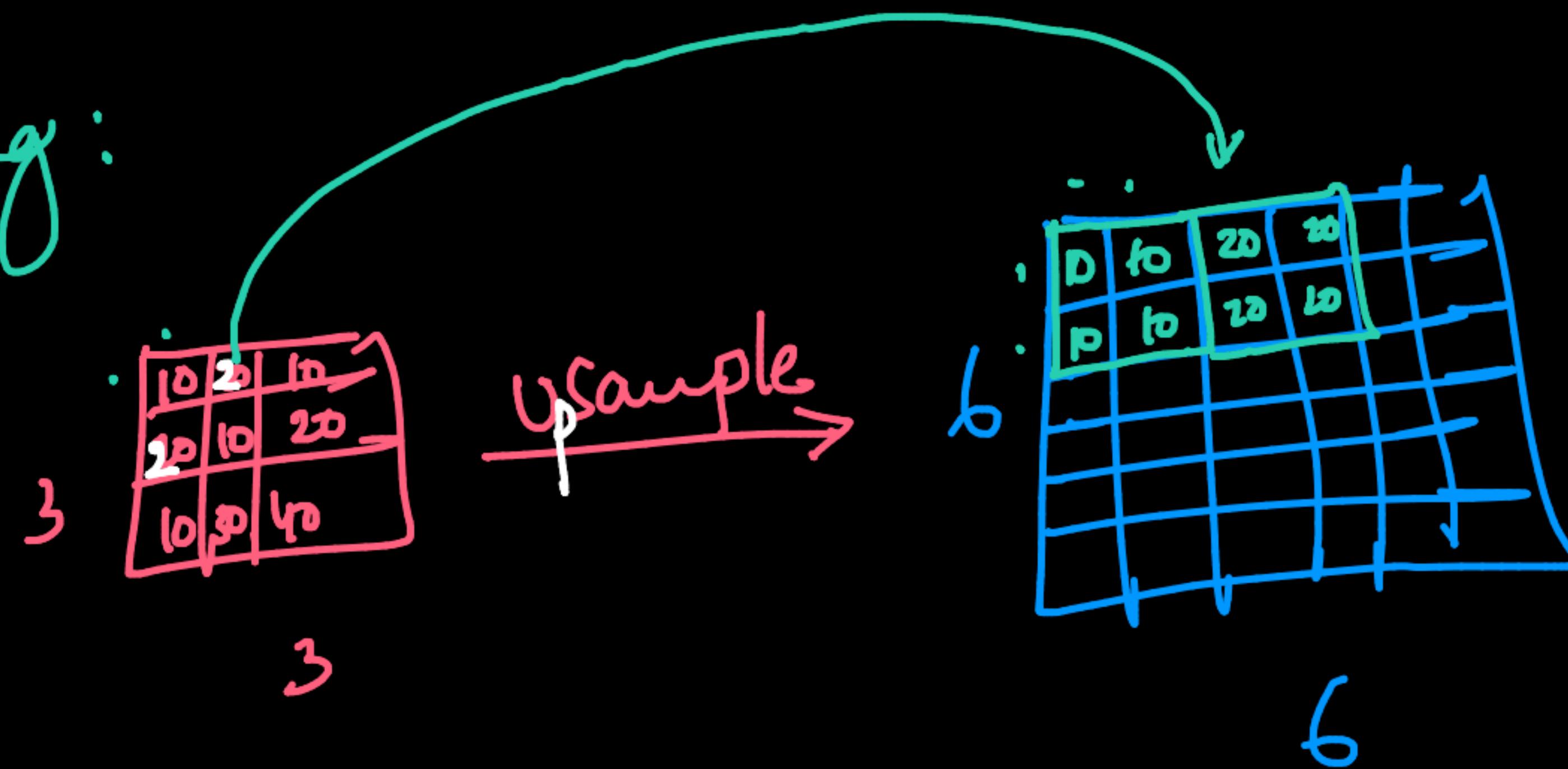
Loss Function

The loss consists of two parts, the *localization loss* for bounding box offset prediction and the *classification loss* for conditional class probabilities. Both parts are computed as the sum of squared errors. Two scale parameters are used to control how much we want to increase the loss from bounding box coordinate predictions (λ_{coord}) and how much we want to decrease the loss of confidence score predictions for boxes without objects ($\lambda_{\text{no obj}}$). Down-weighting the loss contributed by background instances.

multi-scale

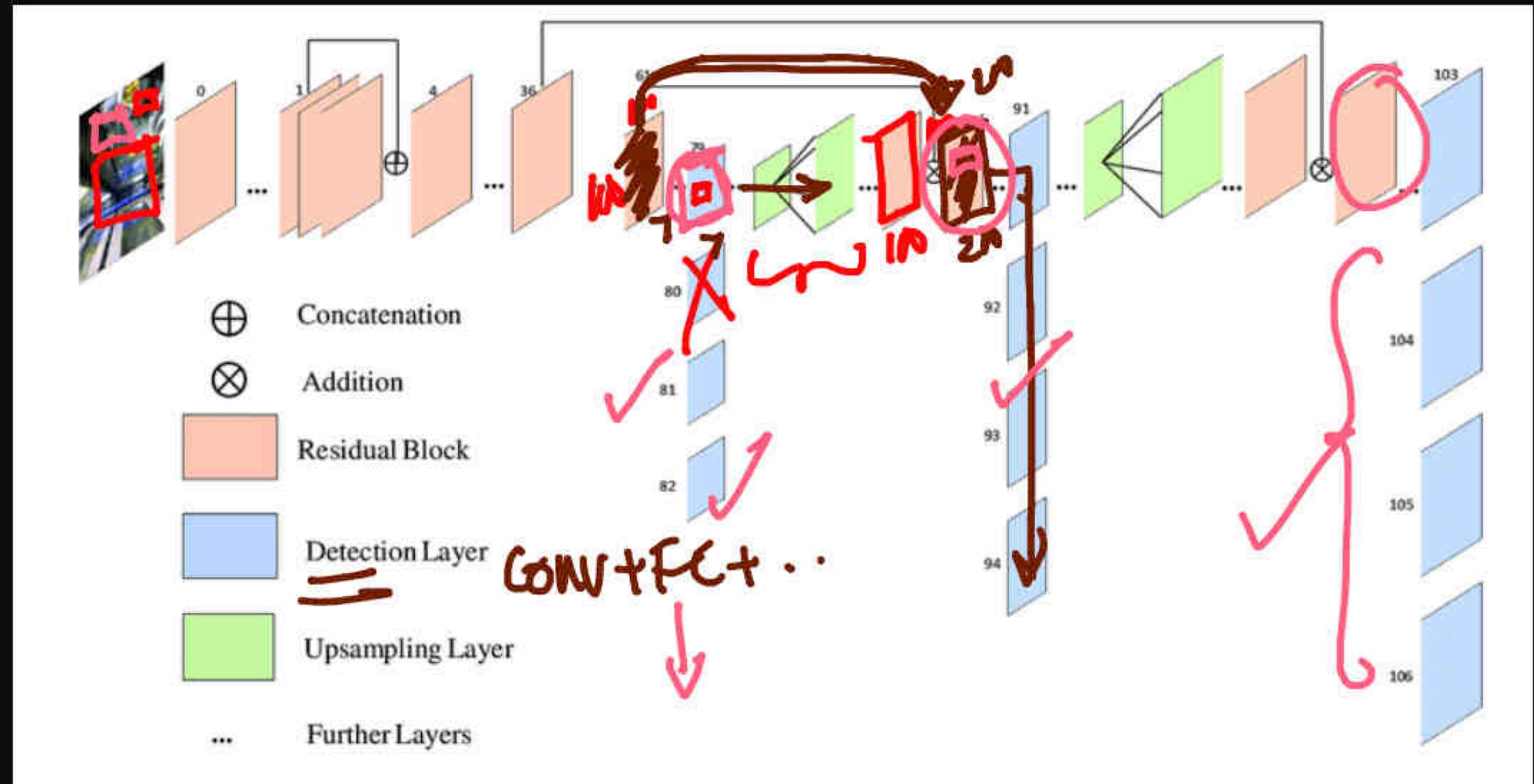


UpSampling:



$$128 \times 128 \xrightarrow{\text{US}} 256 \times 256$$

V3



L8: ObjectDetection X Object Detection Par X Network-architectur X Post_Read: RetinaNe X L9: Object Segments X GitHub - vdumoulin/c... X ROI pooling vs. ROI a... X An Introduction to di... X Li's Log X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=t1L78AUgnEzv

+ Code + Text Connect |  

Let's go back to our problem statement and see how we can use Pre-Trained YoloV5 to solve it:

Why YOLO V5:

- We are going to use pre-trained Yolo V5 models for our use case since they are trained using COCO-Dataset (<https://cocodataset.org/#home>) and already contain the Target-Classes which we want to Detect:
 - 'Car' , 'motorbike', 'aeroplane', 'bus', 'truck'
 - Traffic light
 - Person
 - so on..
- Due to its choice of framework being pytorch it's highly popular in Applied CV Community
- Pytorch models can be easily migrated to other frameworks such as onnx, tensorflow.
- It's faster than YoloV4 and has 5 checkpoints for different memory and speed requirements.
- For more details Check this blog: <https://blog.roboflow.com/yolov5-is-here/>
- We will be using ONNX Format of YoloV5

What is ONNX?

ONNX is an open format built to represent machine learning models. ONNX defines a common set of operators - the building blocks of machine learning and deep learning models - and a common file format for saving them. ONNX is designed to be used by frameworks, tools, runtimes, and compilers.



11 / 11

The screenshot shows the official ONNX website (onnx.ai/index.html) with several handwritten annotations in red and green marker:

- A large green arrow points from the word "ONNX" to a green circle containing the text "Intel Xeon".
- A red arrow points from the "ONNX" logo to the central title area.
- A red bracket on the right side groups "TF2 [Keras]" and "PyTorch".
- A red bracket on the left side groups "Intel" and "Xeon".
- A red arrow points from the "GET STARTED" button to the explanatory text below.

NEWS: Onnx Python Package [Read more >](#)

ONNX

Open Neural Network Exchange

The open standard for machine learning interoperability

GET STARTED

ONNX is an open format built to represent machine learning models. ONNX defines a common set of operators - the building blocks of machine learning and deep learning models - and a common file format to enable AI developers to use models with a variety of frameworks, tools, runtimes, and compilers. [LEARN MORE >](#)

KEY BENEFITS

Interoperability
Develop in your preferred framework without worrying about downstream inferencing implications. ONNX enables you to use your preferred framework with your chosen inference engine.

[SUPPORTED FRAMEWORKS >](#)

Hardware Access
ONNX makes it easier to access hardware optimizations. Use ONNX-compatible runtimes and libraries designed to maximize performance across hardware.

[SUPPORTED ACCELERATORS >](#)

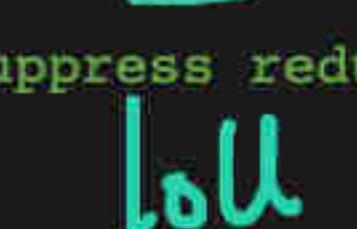
12 / 12

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |  

```
INPUT_HEIGHT = 640

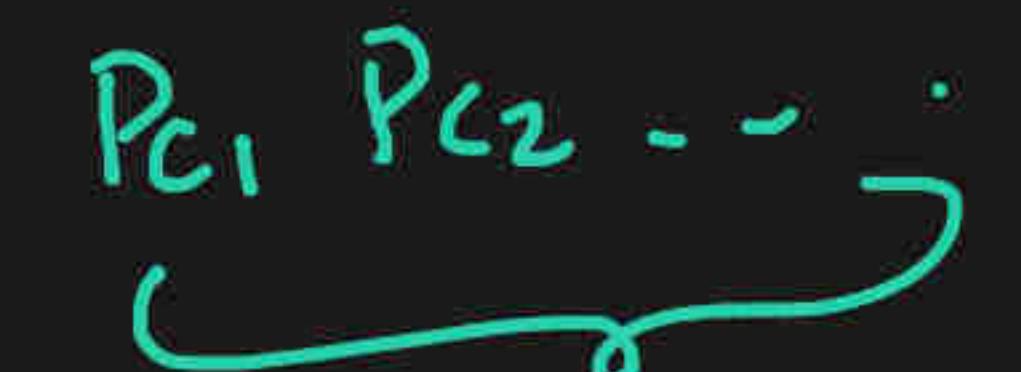
# probability threshold to filter boxes with object or no object
OBJECT_SCORE_THRESHOLD = 0.5

# # probability threshold to detect and Assign Class
CLASS_CONFIDENCE_THRESHOLD = 0.45
      
# IOU AREA Threshold to suppress redundant boxes using NMS
NMS_THRESHOLD = 0.45
      
# Text parameters used for annotating label on Image
FONT_FACE = cv2.FONT_HERSHEY_SIMPLEX
FONT_SCALE = 0.4
THICKNESS = 1
BOX_COLOR = (0,255,255)
FONT_COLOR= (0,0,0)

}

[ ] def draw_label(im, label, x, y):
    """Function used for Drawing text/label onto image at location."""

    # Get text size.
    text_size = cv2.getTextSize(label, FONT_FACE, FONT_SCALE, THICKNESS)
    dim, baseline = text_size[0], text_size[1]

    # Use text size to create a BLACK rectangle.
    cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);
          
# Display text inside the rectangle
```

```
BOX_COLOR = (0,255,255)
FONT_COLOR= (0,0,0)

[ ] def draw_label(im, label, x, y):
    """Function used for Drawing text/label onto image at location."""

    # Get text size.
    text_size = cv2.getTextSize(label, FONT_FACE, FONT_SCALE, THICKNESS)
    dim, baseline = text_size[0], text_size[1]

    # Use text size to create a BLACK rectangle.
    cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);

    # Display text inside the rectangle.
    cv2.putText(im, label, (x, y - dim[1]+baseline ), FONT_FACE, FONT_SCALE, FONT_COLOR, THICKNESS, cv2.LINE_AA)

[ ] def yolo_forward_pass(input_image, net):
    """Performs forward pass to generate prediction using image and Loaded Yolo Model as parameters """
    # Create a 4D blob from a frame
    blob = cv2.dnn.blobFromImage(input_image, 1/255, (INPUT_WIDTH, INPUT_HEIGHT), [0,0,0], 1, crop=False)

    # Sets the input to the network.
    net.setInput(blob)

    # Run the forward pass to get output of the output layers.
    outputs = net.forward(net.getUnconnectedOutLayersNames())

    return outputs
```

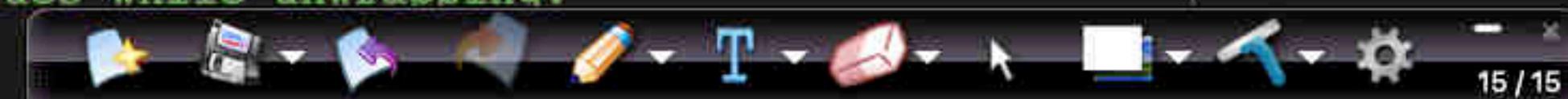
+ Code + Text

Connect  

```
[ ] # Use text size to create a BLACK rectangle.  
cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);  
  
{x} # Display text inside the rectangle.  
cv2.putText(im, label, (x, y - dim[1]+baseline ), FONT_FACE, FONT_SCALE, FONT_COLOR, THICKNESS, cv2.LINE_AA)
```

```
[ ] def yolo_forward_pass(input_image, net):  
    """Performs forward pass to generate prediction using image and Loaded Yolo Model as parameters """  
    # Create a 4D blob from a frame  
    blob = cv2.dnn.blobFromImage(input_image, 1/255, (INPUT_WIDTH, INPUT_HEIGHT), [0,0,0], 1, crop=False)  
    # Sets the input to the network.  
    net.setInput(blob)  
  
    # Run the forward pass to get output of the output layers.  
    outputs = net.forward(net.getUnconnectedOutLayersNames())  
  
    return outputs
```

```
[ ] def post_process_outputs(input_image, outputs):  
    """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:  
    1. Removing predictions with Low Confidence score  
    2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc  
    3. Perform NMS to suppress redundant or Duplicate boxes  
    4. Annotate/draw remaining boxes on the Image  
    """  
  
    # Lists to hold respective values while unwrapping.  
    class_ids = []
```



L8: Object Detect... X ONNX | Home X Object Detection X Network-architect... X Post_ Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X Li's Log X + Update

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |

```
[ ] # Get text size.  
text_size = cv2.getTextSize(label, FONT_FACE, FONT_SCALE, THICKNESS)  
dim, baseline = text_size[0], text_size[1]  
  
{x} # Use text size to create a BLACK rectangle.  
cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);  
  
# Display text inside the rectangle.  
cv2.putText(im, label, (x, y - dim[1]+baseline ), FONT_FACE, FONT_SCALE, FONT_COLOR, THICKNESS, cv2.LINE_AA)
```

```
[ ] def yolo_forward_pass(input_image, net):  
    """Performs forward pass to generate prediction using image and Loaded Yolo Model as parameters """  
    # Create a 4D blob from a frame  
    blob = cv2.dnn.blobFromImage(input_image, 1/255, (INPUT_WIDTH, INPUT_HEIGHT), [0,0,0], 1, crop=False)  
    # Sets the input to the network.  
    net.setInput(blob)  
  
    # Run the forward pass to get output of the output layers.  
    outputs = net.forward(net.getUnconnectedOutLayersNames())  
  
    return outputs
```

```
[ ] def post_process_outputs(input_image, outputs):  
    """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:  
    1. Removing predictions with Low Confidence score  
    2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc  
    3. Perform NMS to suppress redundant boxes  
    4. Annotate/draw remaining boxes
```

16 / 16

+ Code + Text

Connect | User Settings

```
[ ] # Get text size.  
text_size = cv2.getTextSize(label, FONT_FACE, FONT_SCALE, THICKNESS)  
dim, baseline = text_size[0], text_size[1]  
  
{x} # Use text size to create a BLACK rectangle.  
cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);  
  
# Display text inside the rectangle.  
cv2.putText(im, label, (x, y - dim[1]+baseline ), FONT_FACE, FONT_SCALE, FONT_COLOR, THICKNESS, cv2.LINE_AA)
```

```
[ ] def yolo_forward_pass(input_image, net):  
    """Performs forward pass to generate prediction using image and Loaded Yolo Model as parameters """  
    # Create a 4D blob from a frame  
    blob = cv2.dnn.blobFromImage(input_image, 1/255, (INPUT_WIDTH, INPUT_HEIGHT), [0,0,0], 1, crop=False)  
  
    # Sets the input to the network.  
    net.setInput(blob)  
  
    # Run the forward pass to get output of the output layers.  
    outputs = net.forward(net.getUnconnectedOutLayersNames())  
  
    return outputs
```

```
[ ] def post_process_outputs(input_image, outputs):  
    """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:  
    1. Removing predictions with Low Confidence score  
    2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc  
    3. Perform NMS to suppress redundant boxes  
    4. Annotate/draw remaining boxes  
    """
```



L8: ObjectDetect... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X Li's Log X + Update

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |

```
[ ] outputs = net.forward(net.getUnconnectedOutLayersNames())

return outputs
```

{x} [] def post_process_outputs(input_image, outputs):
 """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:
 1. Removing predictions with Low Confidence score
 2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc
 3. Perform NMS to suppress redundant or Duplicate boxes
 4. Annotate/draw remaining boxes on the Image
 """
 # Lists to hold respective values while unwrapping.
 class_ids = []
 confidences = []
 boxes = []

 # Rows
 rows = outputs[0].shape[1]
 image_height, image_width = input_image.shape[:2]

 # Resizing factor.
 x_factor = image_width / INPUT_WIDTH
 y_factor = image_height / INPUT_HEIGHT

 # Iterate through detections.
 for r in range(rows):
 row = outputs[0][0][r]
 confidence = row[4]

Discard bad detections

L8: ObjectDetect... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |  

```
[ ] # Run the forward pass to get output of the output layers.  
outputs = net.forward(net.getUnconnectedOutLayersNames())  
  
return outputs  
  
[ ] def post_process_outputs(input_image, outputs):  
    """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:  
    1. Removing predictions with Low Confidence score  
    2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc  
    3. Perform NMS to suppress redundant or Duplicate boxes  
    4. Annotate/draw remaining boxes on the Image  
    """  
  
    # Lists to hold respective values while unwrapping.  
    class_ids = []  
    confidences = []  
    boxes = []  
  
    # Rows  
    rows = outputs[0].shape[1]  
    image_height, image_width = input_image.shape[:2]  
  
    # Resizing factor.  
    x_factor = image_width / INPUT_WIDTH  
    y_factor = image_height / INPUT_HEIGHT  
  
    # Iterate through detections.  
    for r in range(rows):  
        row = outputs[0][0][r]  
        confidence = row[4]
```

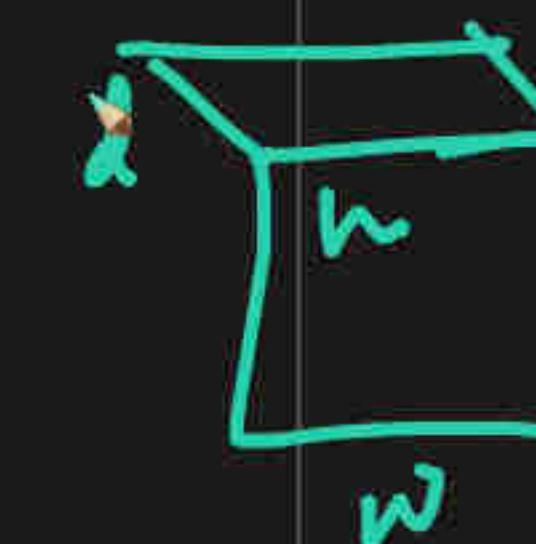
19 / 19

+ Code + Text

Connect  

```
[ ] # Use text size to create a BLACK rectangle.  
cv2.rectangle(im, (x,y), (x + dim[0], y - (dim[1]+baseline) ), (255,255,255), cv2.FILLED);  
  
{x} # Display text inside the rectangle.  
cv2.putText(im, label, (x, y - dim[1]+baseline ), FONT_FACE, FONT_SCALE, FONT_COLOR, THICKNESS, cv2.LINE_AA)
```

```
[ ] def yolo_forward_pass(input_image, net):  
    """Performs forward pass to generate prediction using image and Loaded Yolo Model as parameters """  
    # Create a 4D blob from a frame  
    blob = cv2.dnn.blobFromImage(input_image, 1/255, (INPUT_WIDTH, INPUT_HEIGHT), [0,0,0], 1, crop=False)  
  
    # Sets the input to the network.  
    net.setInput(blob)  
  
    # Run the forward pass to get output of the output layers.  
    outputs = net.forward(net.getUnconnectedOutLayersNames())  
  
    return outputs
```



```
[ ] def post_process_outputs(input_image, outputs):  
    """ Takes Model Prediction outputs from forward pass function and performs post-processing tasks:  
    1. Removing predictions with Low Confidence score  
    2. Filtering only classes required for our task such as Vehicles, Human, traffic light stop sgn etc  
    3. Perform NMS to suppress redundant or Duplicate boxes  
    4. Annotate/draw remaining boxes on the Image  
    """  
  
    # Lists to hold respective values while unwrapping.  
    class_ids = []  
    confidences = []
```

L8: Object Detecti... X ONNX | Home X Object Detection X Network-architect... X Post_ Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X Li's Log X + Update

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |

```
[ ] from google.colab.patches import cv2_imshow
import matplotlib.pyplot as plt

{x}
# Load class names.
classesFile = "/content/coco.names.txt"
classes = None

with open(classesFile, 'rt') as f:
    classes = f.read().rstrip('\n').split('\n')

# Load the model using cv2.dnn.readNet: Yolo V5 Nano Model
modelWeights = "/content/yolov5n.onnx"
net = cv2.dnn.readNet(modelWeights)

for image_path in sample_images:
    # Load image
    frame = cv2.imread(image_path)

    # Process image:
    detections = yolo_forward_pass(frame, net)
    pred_img = post_process_outputs(frame.copy(), detections)

    """
    Annotate efficiency information. The function getPerfProfile returns the overall time for inference(t)
    and the timings for each of the layers(in layersTimes).
    """

    t, _ = net.getPerfProfile()
    label = 'Inference time: %.2f ms' % t
    cv2.putText(pred_img, label, (10, 21), 21/21, LINE_AA)
```

L8: Object Detecti... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X LiLog X + Update

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |

```
ax2 = fig.add_subplot(1,2,2)
pred_img = cv2.cvtColor(pred_img, cv2.COLOR_BGR2RGB )
ax2.imshow(pred_img)
ax2.axis('off')

plt.show()
```

Inference time: 365.80 ms

car:0.74
car:0.69
car:0.60
car:0.69
car:0.61
truck:0.75
car:0.63

Inference time: 429.26 ms

L8: ObjectDetect... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW Update

+ Code + Text Connect |

```
ax2 = fig.add_subplot(1,2,2)
pred_img = cv2.cvtColor(pred_img, cv2.COLOR_BGR2RGB )
ax2.imshow(pred_img)
ax2.axis('off')

plt.show()
```

Inference time: 365.80 ms

car:0.74 car:0.60 car:0.69 car:0.61 car:0.63 truck:0.75

Inference time: 429.26 ms

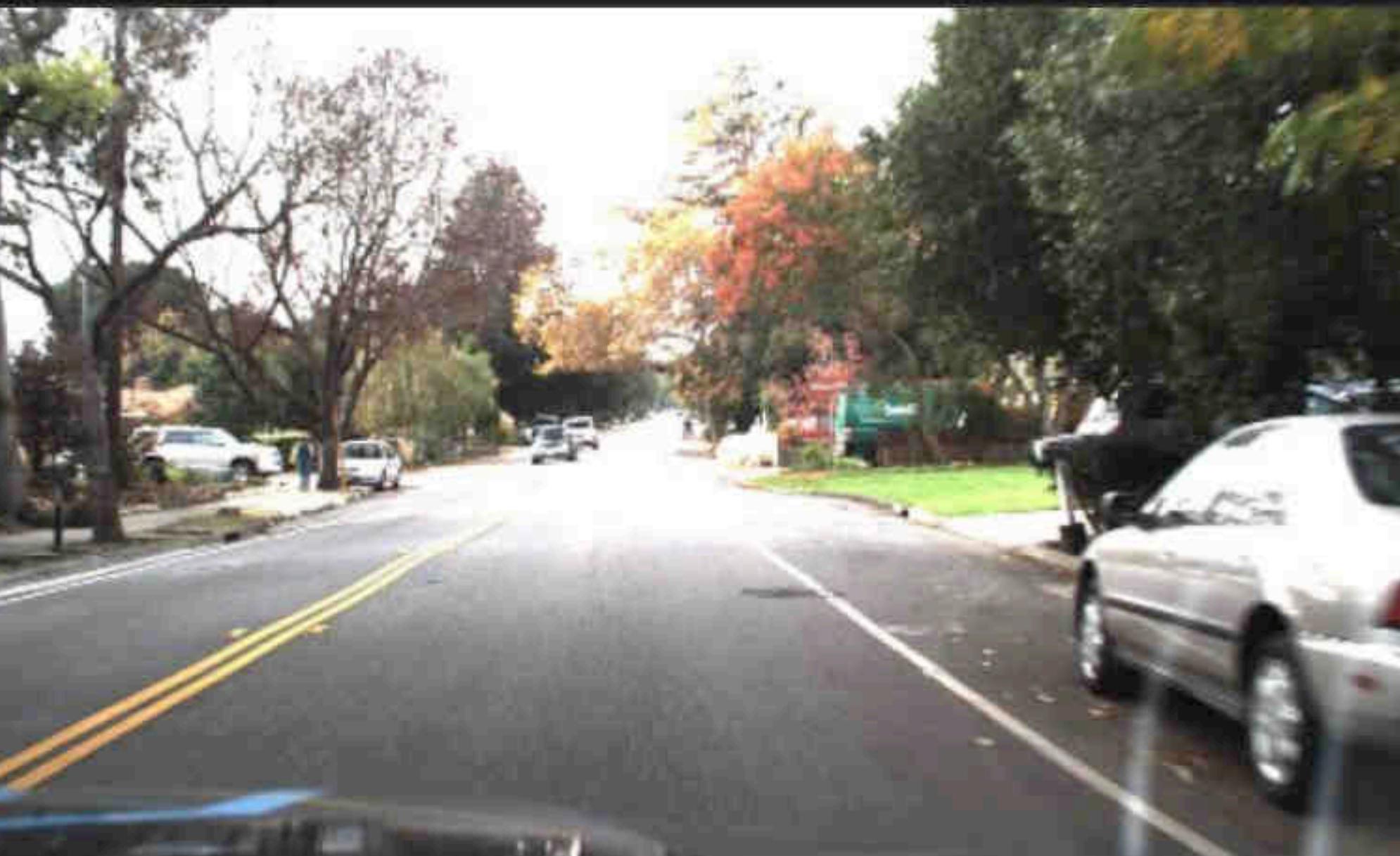
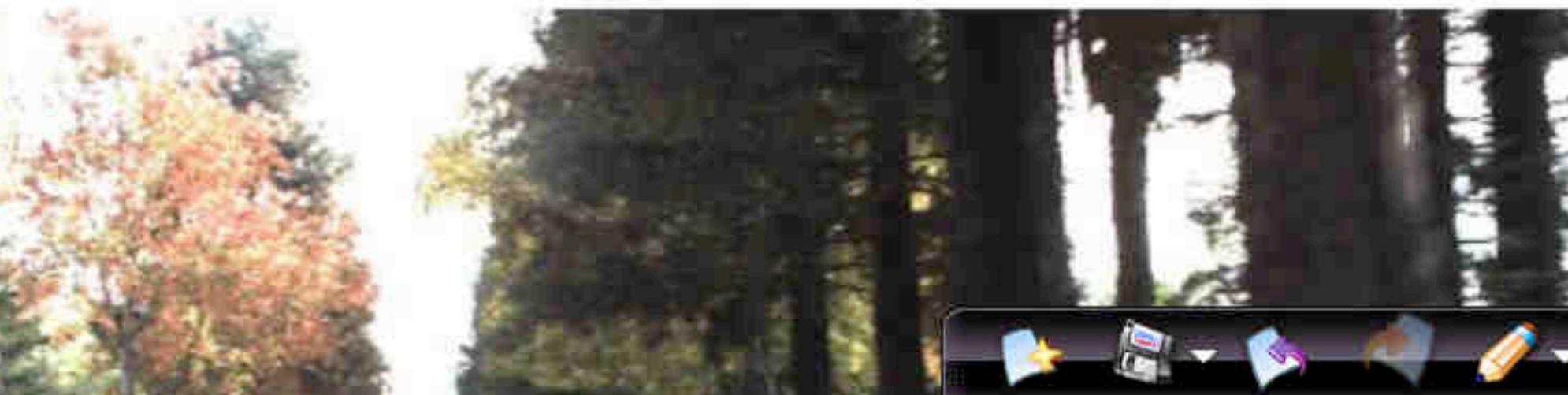
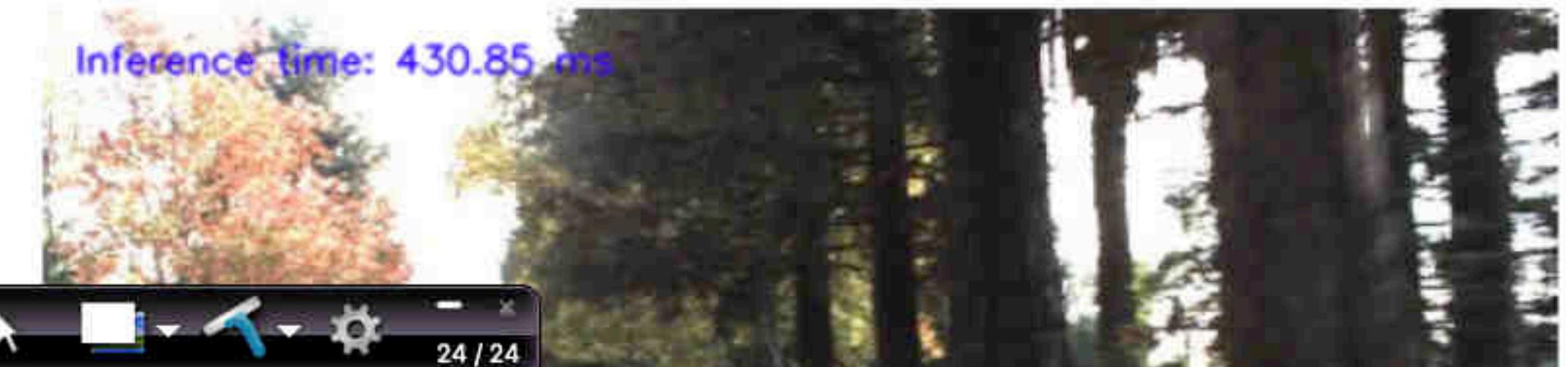
L8: ObjectDetect... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X Li's Log X + Update

colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |

```
[ ] ax2 = fig.add_subplot(1,2,2)
pred_img = cv2.cvtColor(pred_img, cv2.COLOR_BGR2RGB )
ax2.imshow(pred_img)
ax2.axis('off')

plt.show()
```

Inference time: 430.85 ms

24 / 24

L8: Object Detecti X ONNX | Home X Object Detection X Network-architect X Post_Read: Retin X L9: Object Segme X GitHub - vdumou X ROI pooling vs. RC X An Introduction to X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect |  

```
[ ] ax2 = fig.add_subplot(1,2,2)
pred_img = cv2.cvtColor(pred_img, cv2.COLOR_BGR2RGB )
ax2.imshow(pred_img)
ax2.axis('off')

plt.show()
```

Inference time: 479.46 ms



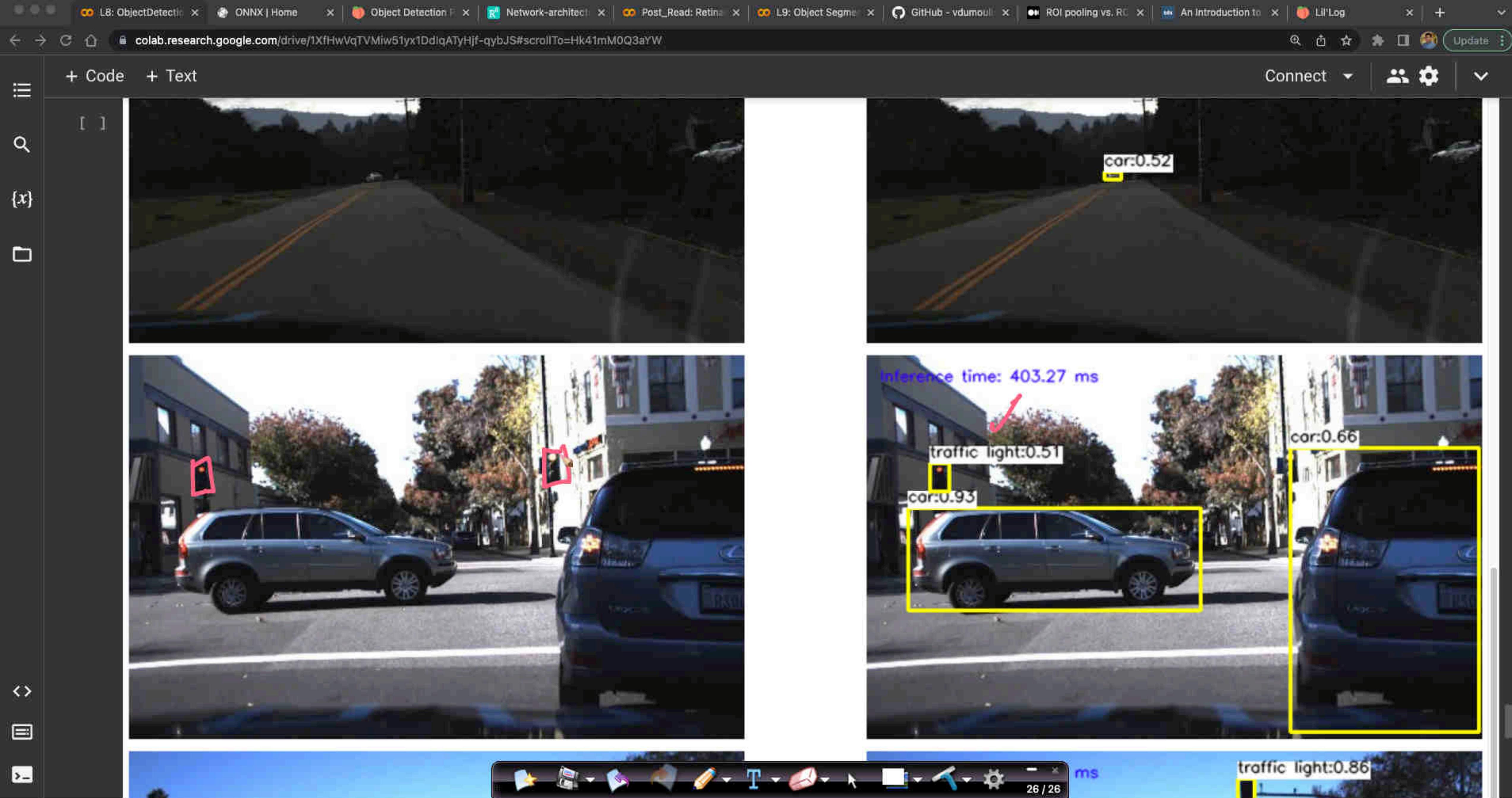
car:0.52

Inference time: 403.27 ms



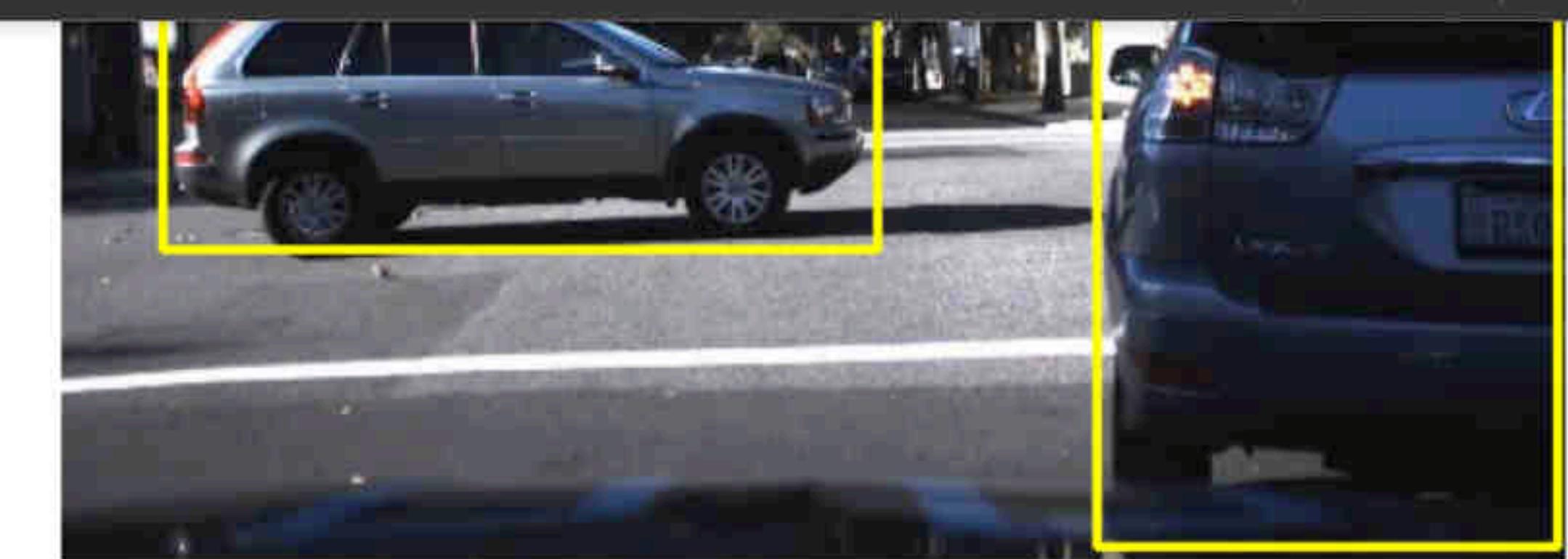
car:0.66

25 / 25



+ Code + Text

Connect ▾



- It did a great job of predicting Vehicles and

L8: Object Detecti X ONNX | Home X Object Detection X Network-architect X Post_Read: Retina X L9: Object Segme X GitHub - vdumou X ROI pooling vs. RC X An Introduction to X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect

iii {x} []

Inference time: 476.87 ms

traffic light:0.86

traffic light:0.57 t:0.61

car:0.50

car:0.73

car:0.76

car:0.55

car:0.51

• It did a great job of predicting Vehicles and

28 / 28

L8: Object Detecti X ONNX | Home X Object Detection X Network-architect X Post_Read: Retina X L9: Object Segme X GitHub - vdumou X ROI pooling vs. RC X An Introduction to X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjt-qybJS#scrollTo=Hk41mM0Q3aYW

+ Code + Text Connect

iii {x} []

Inference time: 476.87 ms

traffic light:0.86

traffic light:0.57

traffic light:0.61

car:0.50

car:0.73

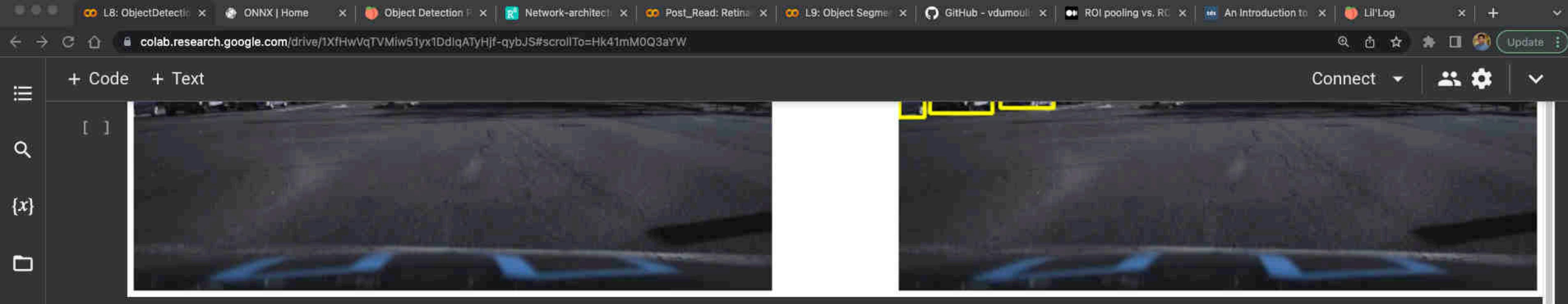
car:0.76

car:0.55

car:0.51

• It did a great job of predicting Vehicles and

29 / 29



- It did a great job of predicting Vehicles and Traffic light with just the pre-trained model and
- Every image was processed in less than a Second on CPU without any specialized hardware(1000 ms = 1 second).
- To Enhance performance we can finetune the model on our custom dataset as well.

Question:

- Try tuning Score threshold and Class threshold probabilities and check impact on predictions.
- Evaluate MAP of the pre-trained model on full dataset
- Run inference using Other Variants of Yolo Model and compare performance in terms of Latency and Accuracy/MAP
- Fine tune the model on dataset and compare performance with Pre-Trained Model in terms of MAP.

Conclusion:

- Using Single Stage Detection method specifically pre-trained YOLOV5 Model, we developed a very fast and accurate detector for our use case which can run inferences in Real Time

L8: ObjectDetect... X ONNX | Home X Object Detection X Network-architect... X Post_Read: Retin... X L9: Object Segme... X GitHub - vdumou... X ROI pooling vs. RC X An Introduction to... X LiLog X + colab.research.google.com/drive/1XfHwVqTVMiw51yx1DdlqATyHjf-qybJS#scrollTo=Hk41mM0Q3aYW

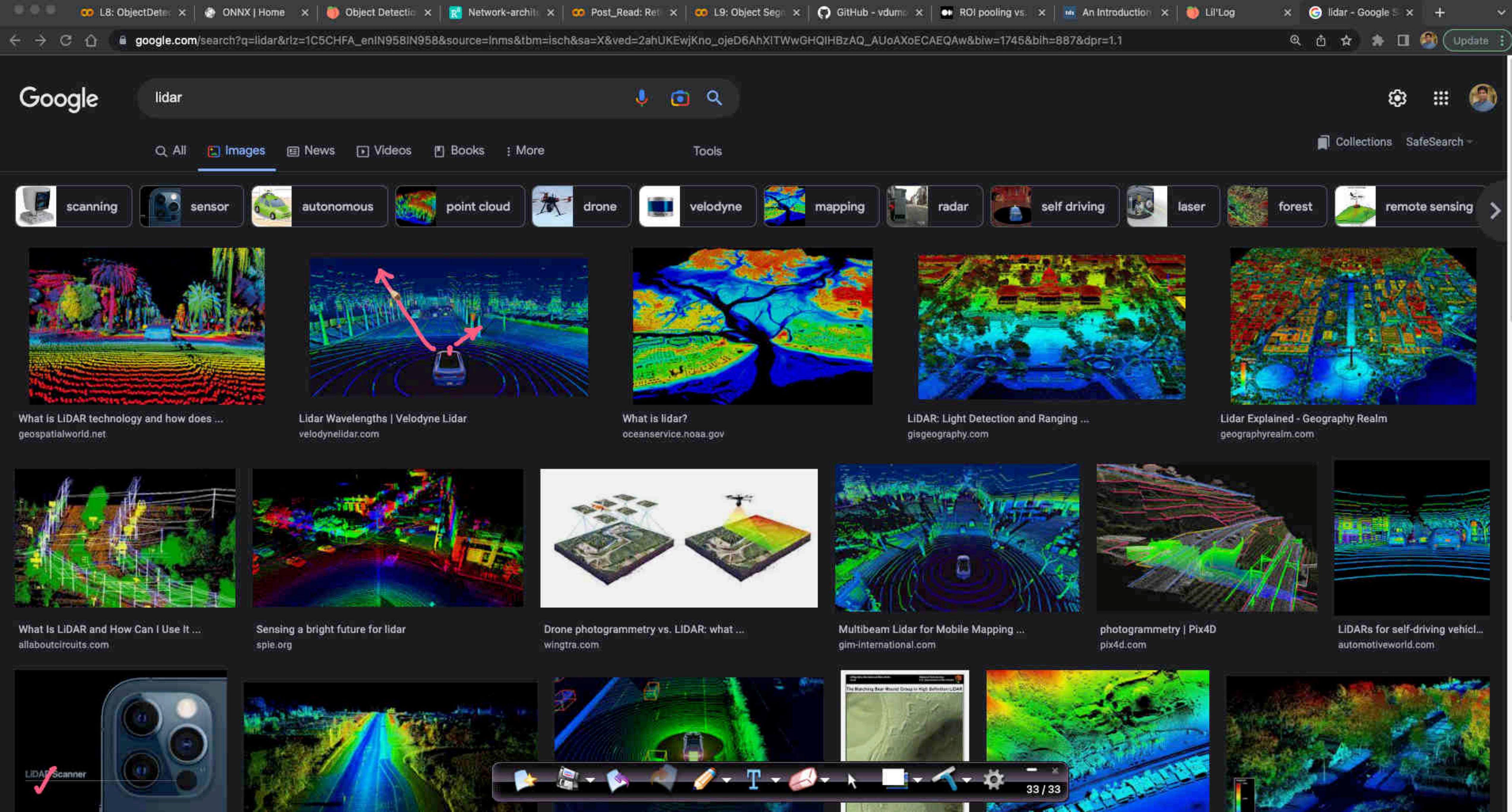
+ Code + Text Connect |

```
[ ] # Load the model using cv2.dnn.readNet: Yolo V5 Nano Model  
modelWeights = "/content/yolov5n.onnx"  
net = cv2.dnn.readNet(modelWeights)  
  
{x} for image_path in sample_images:  
    # Load image  
    frame = cv2.imread(image_path)  
  
    # Process image:  
    detections = yolo_forward_pass(frame, net)  
    pred_img = post_process_outputs(frame.copy(), detections)  
  
    """  
    Annotate efficiency information. The function getPerfProfile returns the overall time for inference(t)  
    and the timings for each of the layers(layersTimes).  
    """  
    t, _ = net.getPerfProfile()  
    label = 'Inference time: %.2f ms' % (t * 1000.0 / cv2.getTickFrequency())  
    cv2.putText(pred_img, label, (10, 20), FONT_FACE, FONT_SCALE, (255, 0, 50), THICKNESS, cv2.LINE_AA)  
  
    fig = plt.figure(figsize=(28,16))  
    ax1 = fig.add_subplot(1,2,1)  
    frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB )  
    ax1.imshow(frame)  
    ax1.axis('off')  
  
    ax2 = fig.add_subplot(1,2,2)  
    pred_img = cv2.cvtColor(pred_img, c
```

24 fps

Q

LIDAR



L8: ObjectDete X | ONNX | Home X | Object Detectio X | Network-archit X | Post_ Read: Re X | L9: Object Seg X | GitHub - vdum X | ROI pooling vs. X | An Introduction X | LiLog X | YOLOv5 is Here X +

blog.roboflow.com/yolov5-is-here/ Update

NEWS

YOLOv5 is Here: State-of-the-Art Object Detection at 140

FPS

Joseph Nelson, Jacob Solawetz

JUN 10, 2020 | 4 MIN READ



140 SEC

WE DELIVER

Get our latest content delivered directly to your inbox.

Enter email

Subscribe

Unsubscribe at any time. Review our [Privacy Policy](#).

TABLE OF CONTENTS

Less than 50 days after the release YOLOv4, YOLOv5 improves accessibility for realtime object detection.

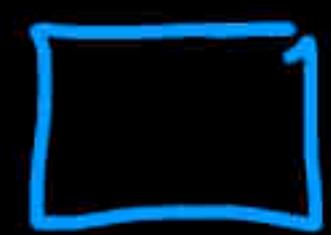
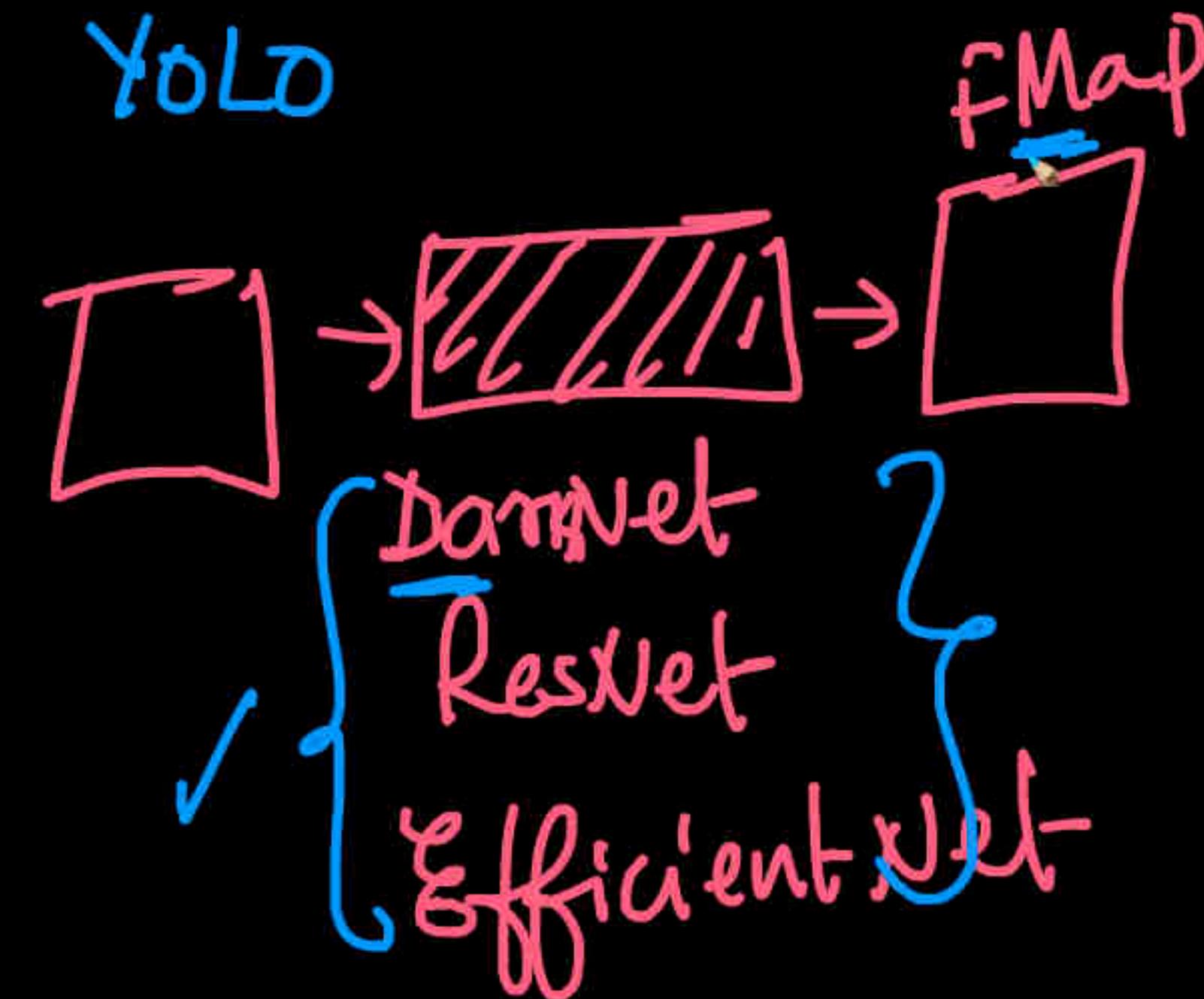
June 29, YOLOv5 has released the first official version of the repository. We wrote a new [deep dive](#) on YOLOv5.

Less than 50 days after the release YOLOv4, YOLOv5 improves accessibility for realtime object

35 / 35

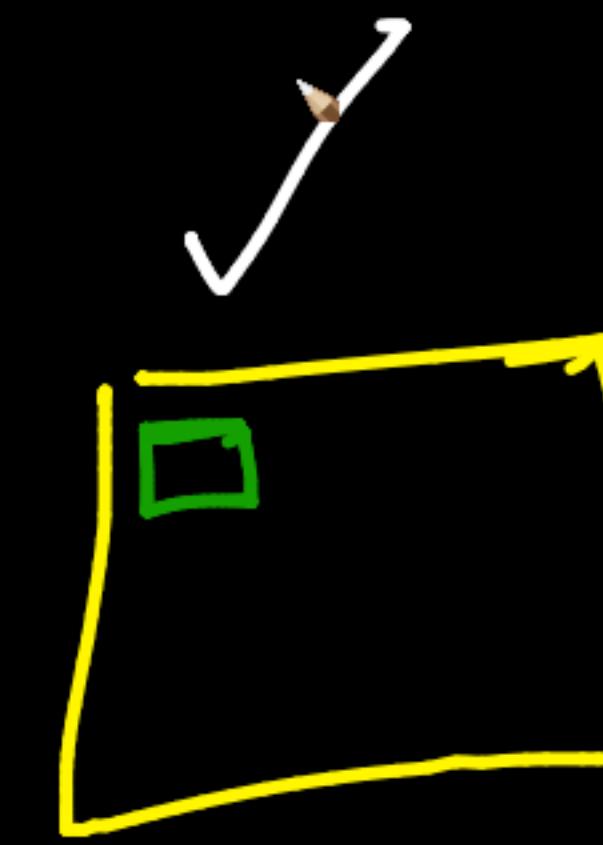
Q

YOLO



YOLO Limitations

①



small-objects

~
5
~~one~~ bounding boxes per cell [3 typically]

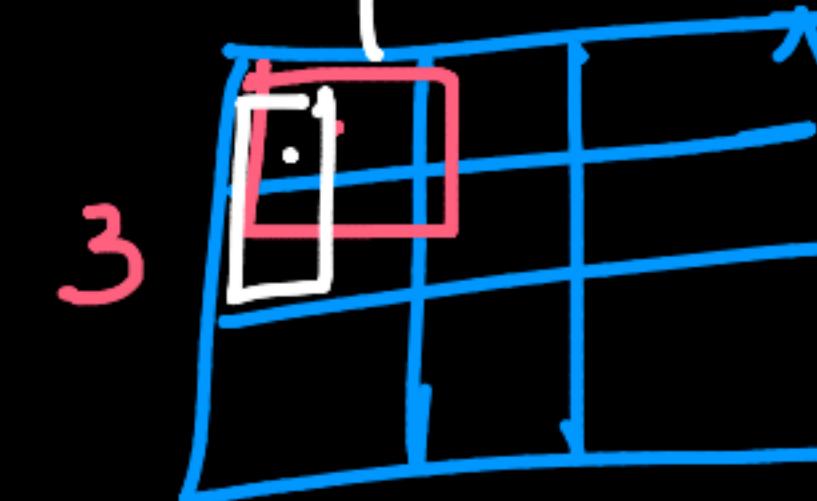
100x100

10x10

56x36

:

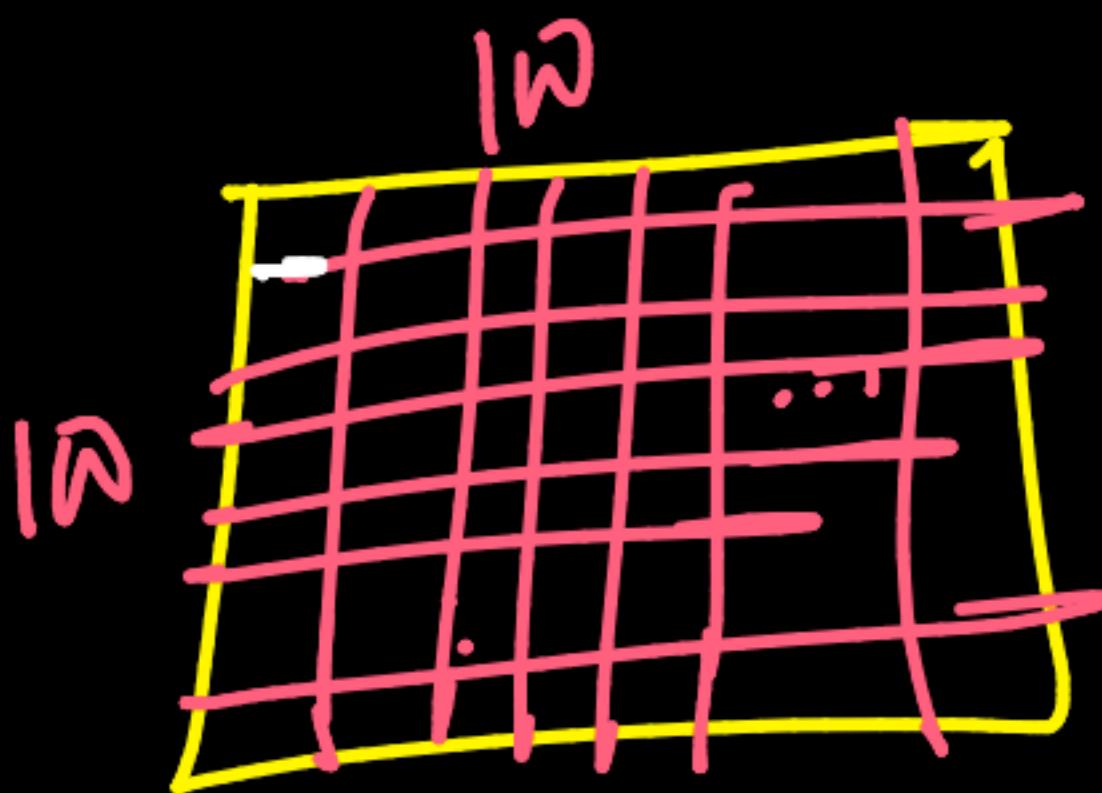
②



(let)

centroids
close to
one another

(3)



→ 10 Boxes
w centroid

lots of cells with ND
Centroid



Output Sparse



imbalanced classifier

L8: ObjectDet... X | ONNX | Home X | Object Detectio... X | Network-archit... X | Post_ Read: Re... X | L9: Object Seg... X | GitHub - vdum... X | ROI pooling vs... X | An Introduction X | LiLog X | YOLOv5 is Here... X | +

colab.research.google.com/drive/1cg42_iwvzeyC-3OPcvI4-1OF00_4snG#scrollTo=Ov0U0UXU8c_o

+ Code + Text Last edited on 12 September Connect | Update

F Map

(a) ResNet

(b) feature pyramid net

(c) class subnet (top)

(d) box subnet (bottom)

For better understanding, Let's understand each component of architecture separately

▼ The backbone Network

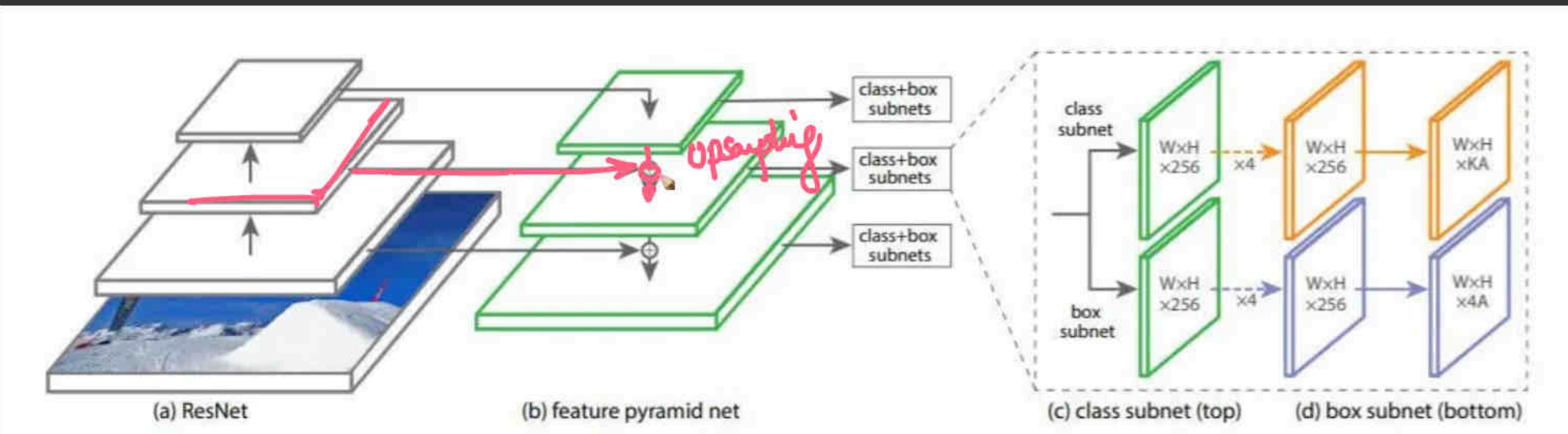
Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

Top down pathway with lateral connections:

+ Code + Text Last edited on 12 September

Connect |



For better understanding, Let's understand each component of architecture separately

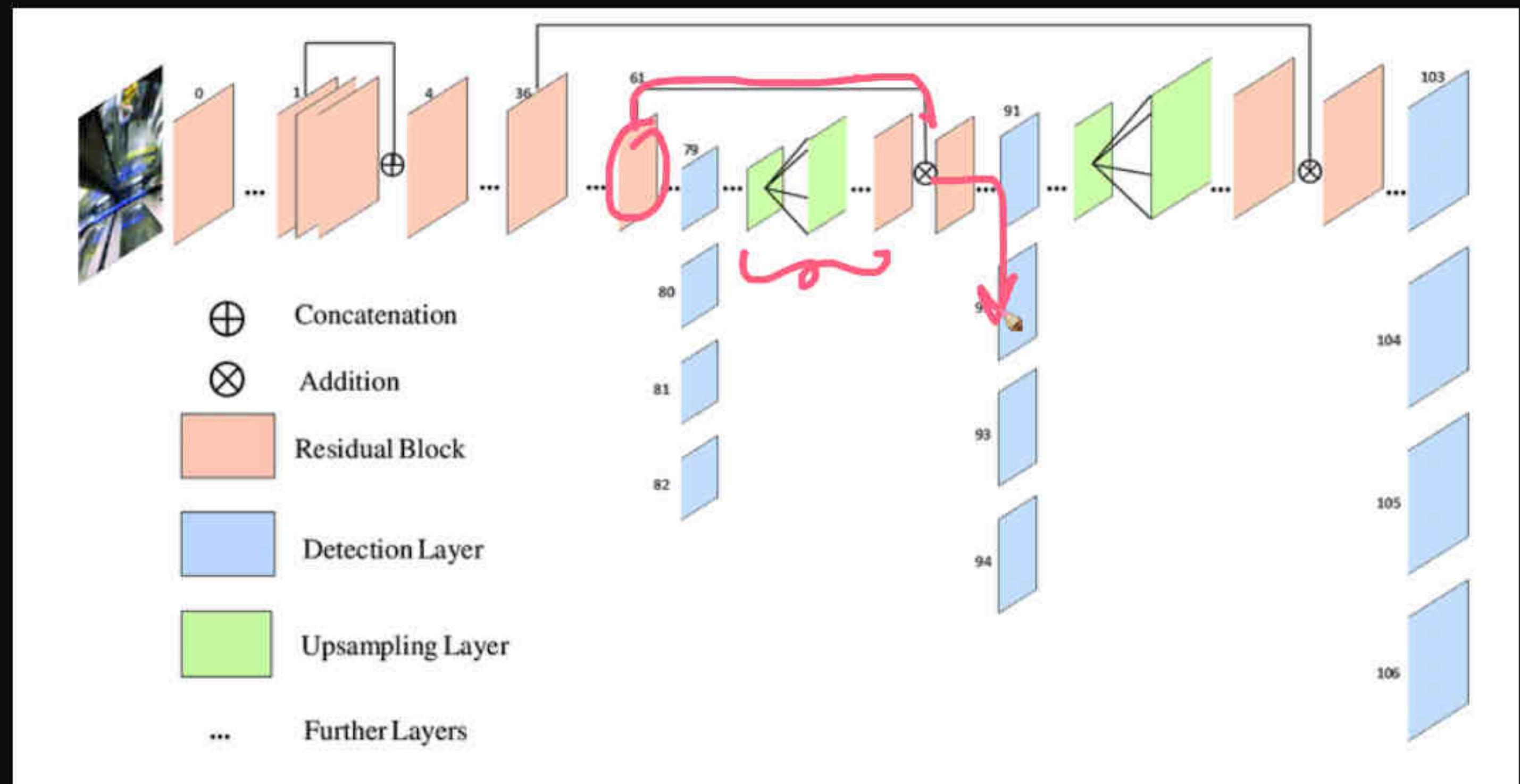
▼ The backbone Network

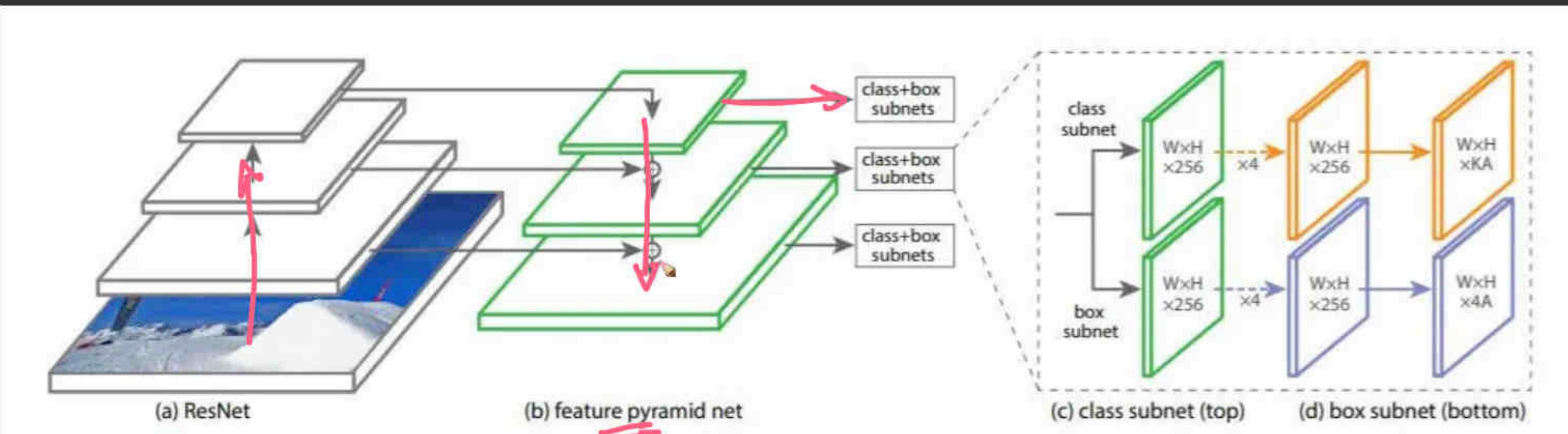
Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

Top down pathway with lateral connections:







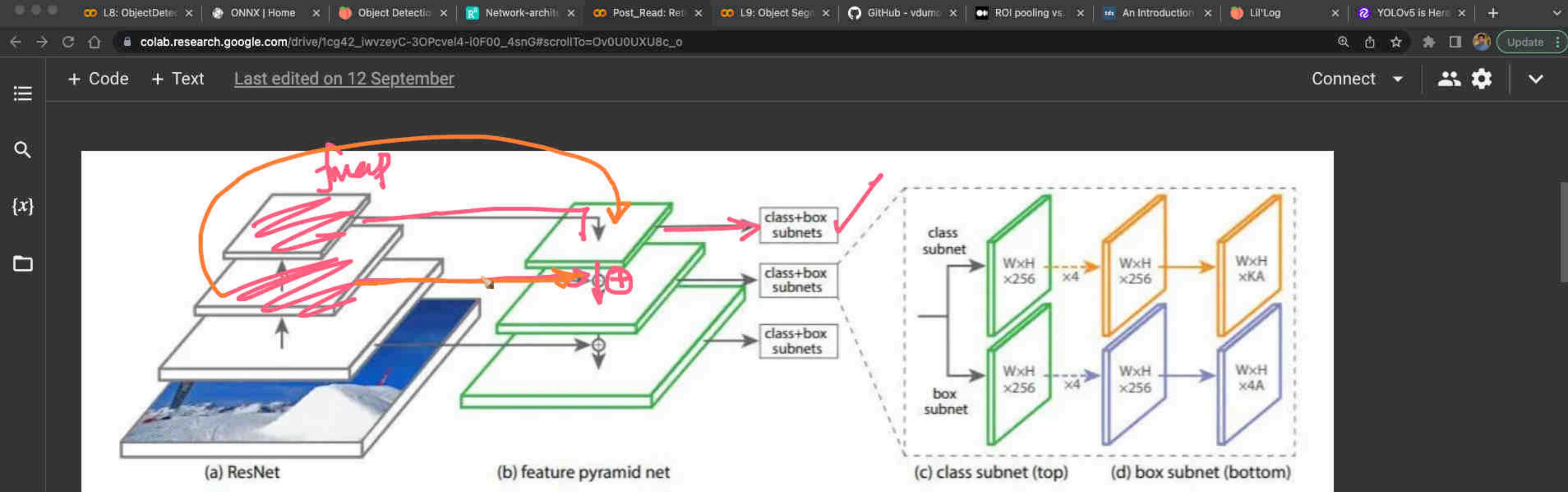
For better understanding, Let's understand each component of architecture separately

▼ The backbone Network

Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

Top down pathway with lateral connections:



For better understanding, Let's understand each component of architecture separately

▼ The backbone Network

Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

Top down pathway with lateral connections:

colab.research.google.com/drive/1cg42_iwvzeyC-3OPcvI4-1OF00_4snG#scrollTo=Ov0U0UXU8c_o

+ Code + Text Last edited on 12 September Connect | Update

{x}

(a) ResNet

(b) feature pyramid net

class+box subnets

class+box subnets

class+box subnets

class subnet

box subnet

WxH x256 x4

WxH x256 x4

WxH xKA

WxH x256

WxH x4A

(c) class subnet (top)

(d) box subnet (bottom)

For better understanding, Let's understand each component of architecture separately

▼ The backbone Network

Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

Top down pathway with lateral connections:

colab.research.google.com/drive/1cg42_iwvzeyC-3OPcvI4-1OF00_4snG#scrollTo=Ov0U0UXU8c_o

+ Code + Text Last edited on 12 September Connect | | ▾

{x}

(a) ResNet

(b) feature pyramid net

class+box subnets

class+box subnets

class+box subnets

class subnet

box subnet

(c) class subnet (top)

(d) box subnet (bottom)

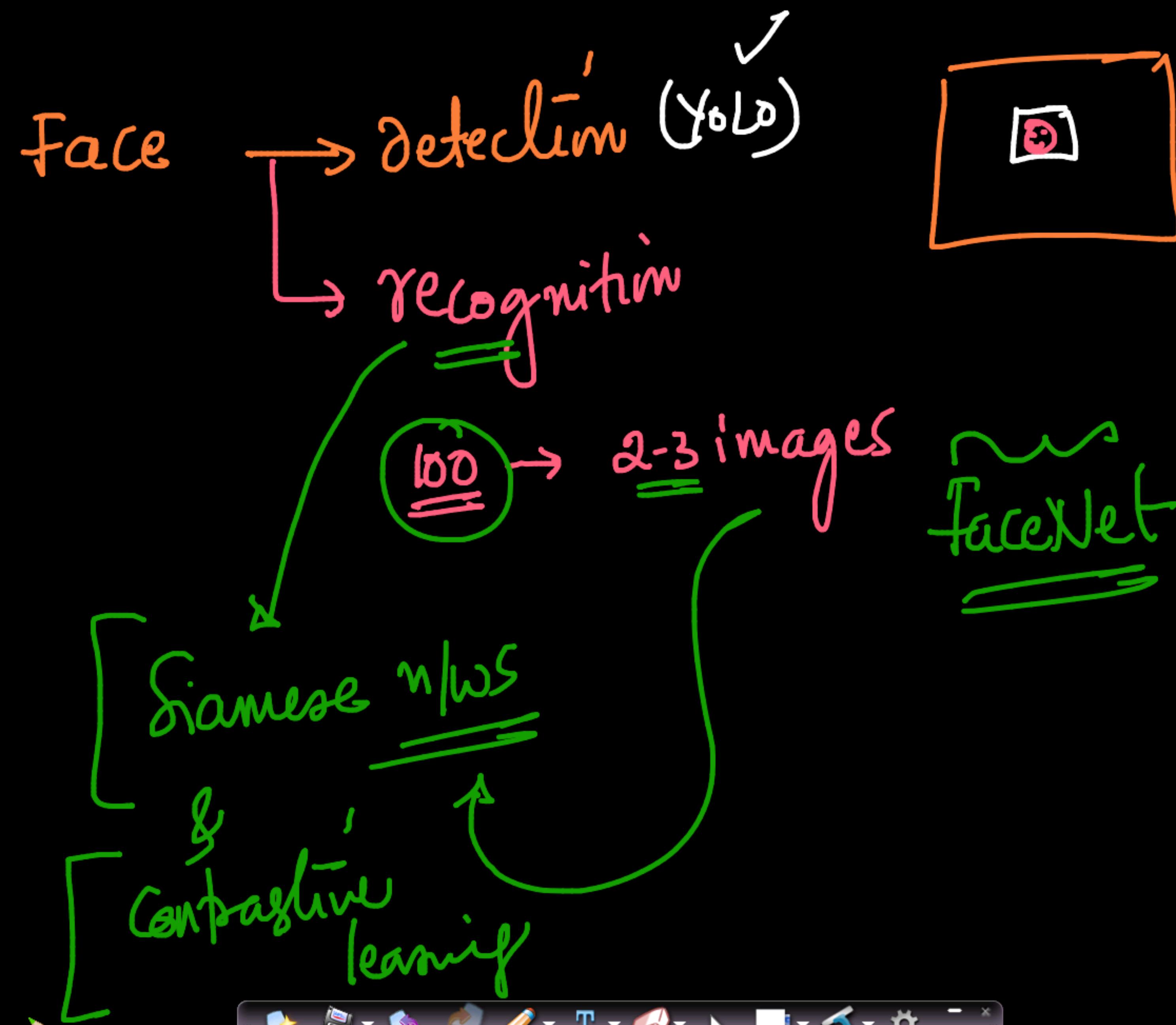
For better understanding, Let's understand each component of architecture separately

▼ The backbone Network

Bottom up pathway:

- Bottom up pathway (eg. ResNet) is used for **feature extraction**.
- So, It calculates the feature maps at different scales, irrespective of the input image size.

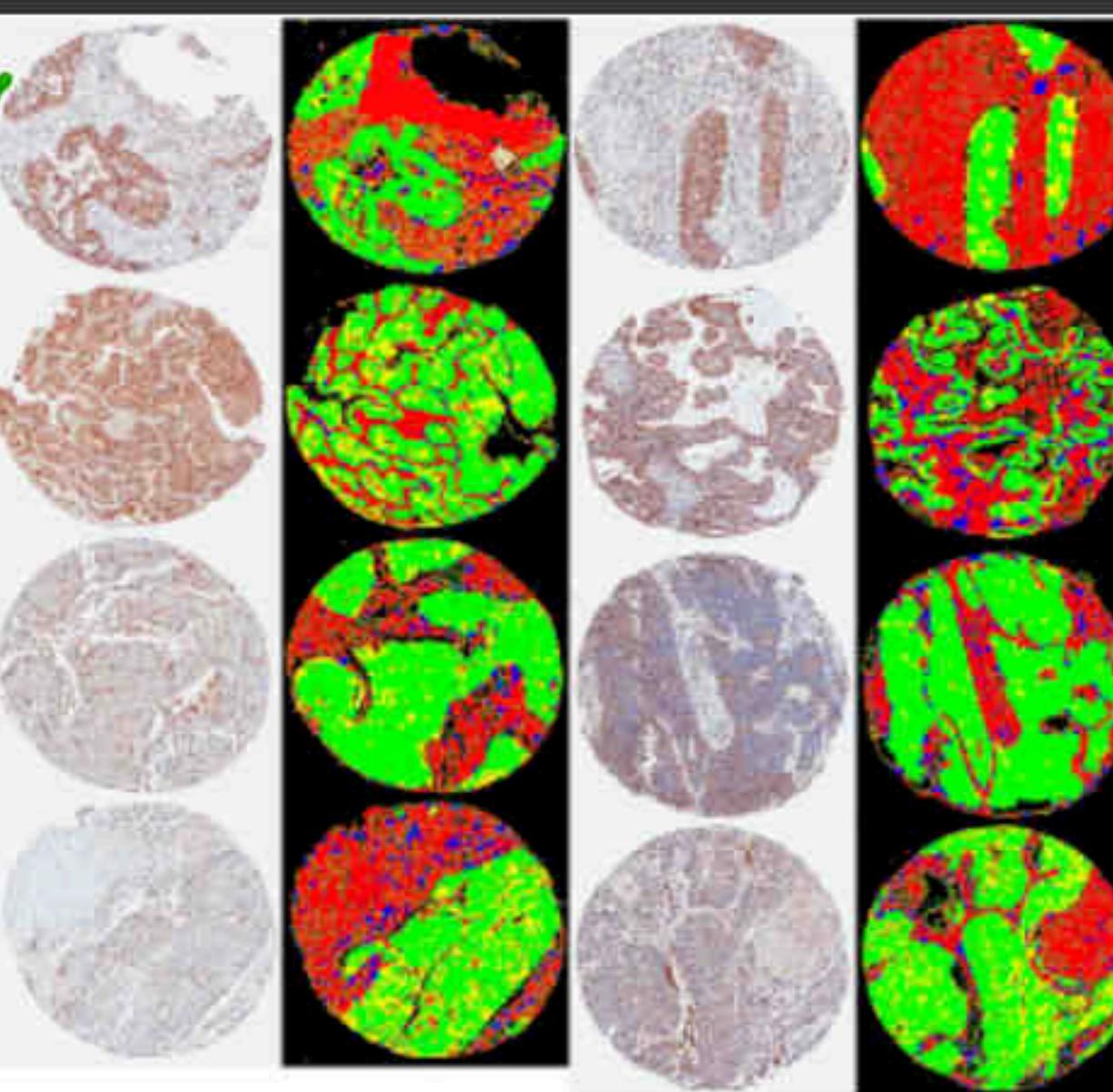
Top down pathway with lateral connections:



colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

- We started with image classification, which is very coarse and high level : whether an object is present in the image or not
- Then in Object Detection we learnt about object localization : where exactly the object is present in the image and we drew a rectangular bounding box around it
- But sometimes bounding boxes are not enough, we need more precision, extract boundary of the object
- For Example Tumor tissue detection, tissue in the red have quite complex structure and we cannot draw bounding boxes around them
- Today we are going to learn about image segmentation, where we do pixel level analysis of the object
- We will learn to form very complex boundary when required

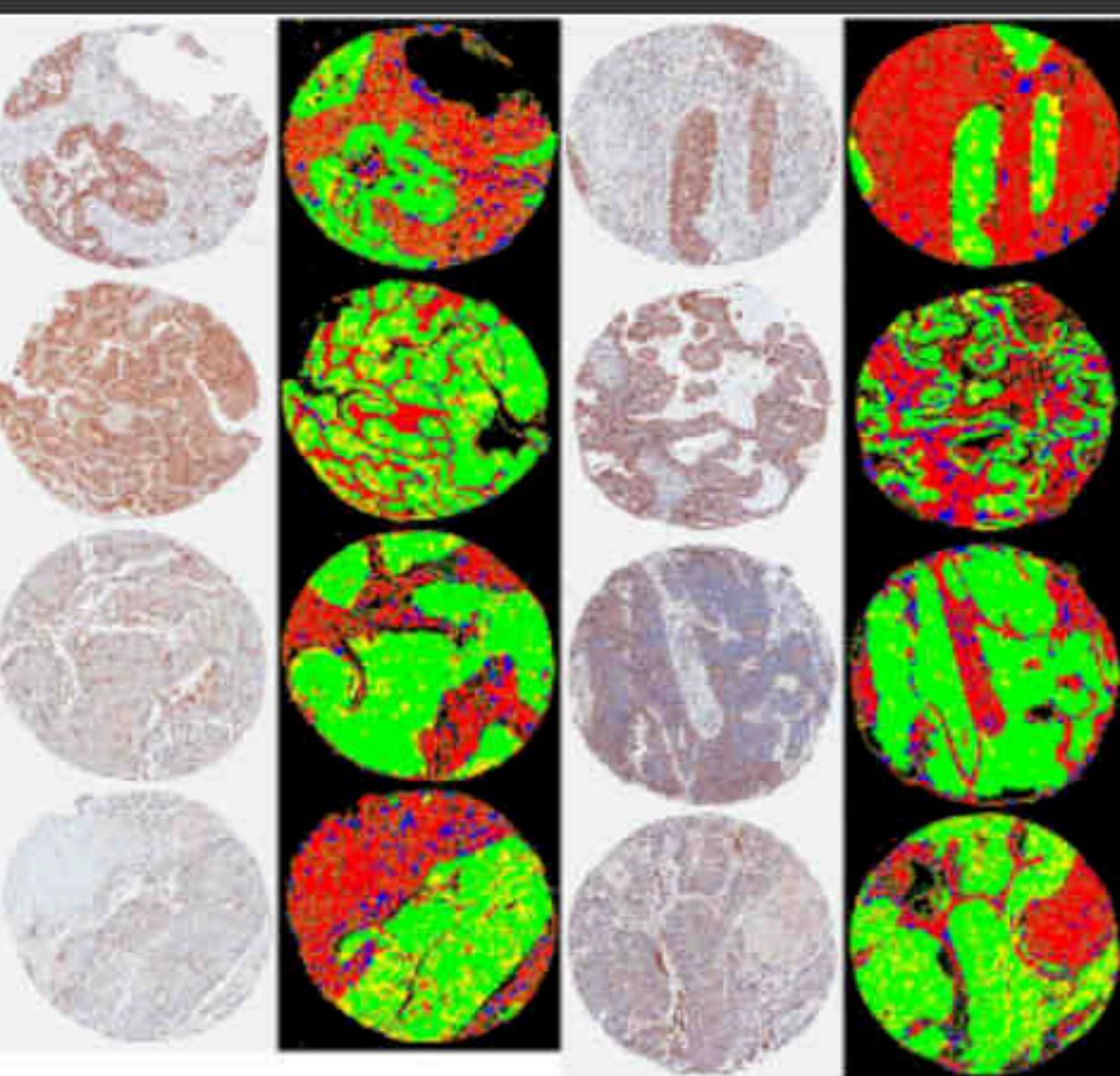


47 / 47

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

- We started with image classification, which is very coarse and high level : whether an object is present in the image or not
- Then in Object Detection we learnt about object localization : where exactly the object is present in the image and we drew a rectangular bounding box around it
- But sometimes bounding boxes are not enough, we need more precision, extract boundary of the object
- For Example Tumor tissue detection, tissue in the red have quite complex structure and we cannot draw bounding boxes around them
- Today we are going to learn about image segmentation, where we do pixel level analysis of the object
- We will learn to form very complex boundary when required



obj detection & locn ~

Segmentation

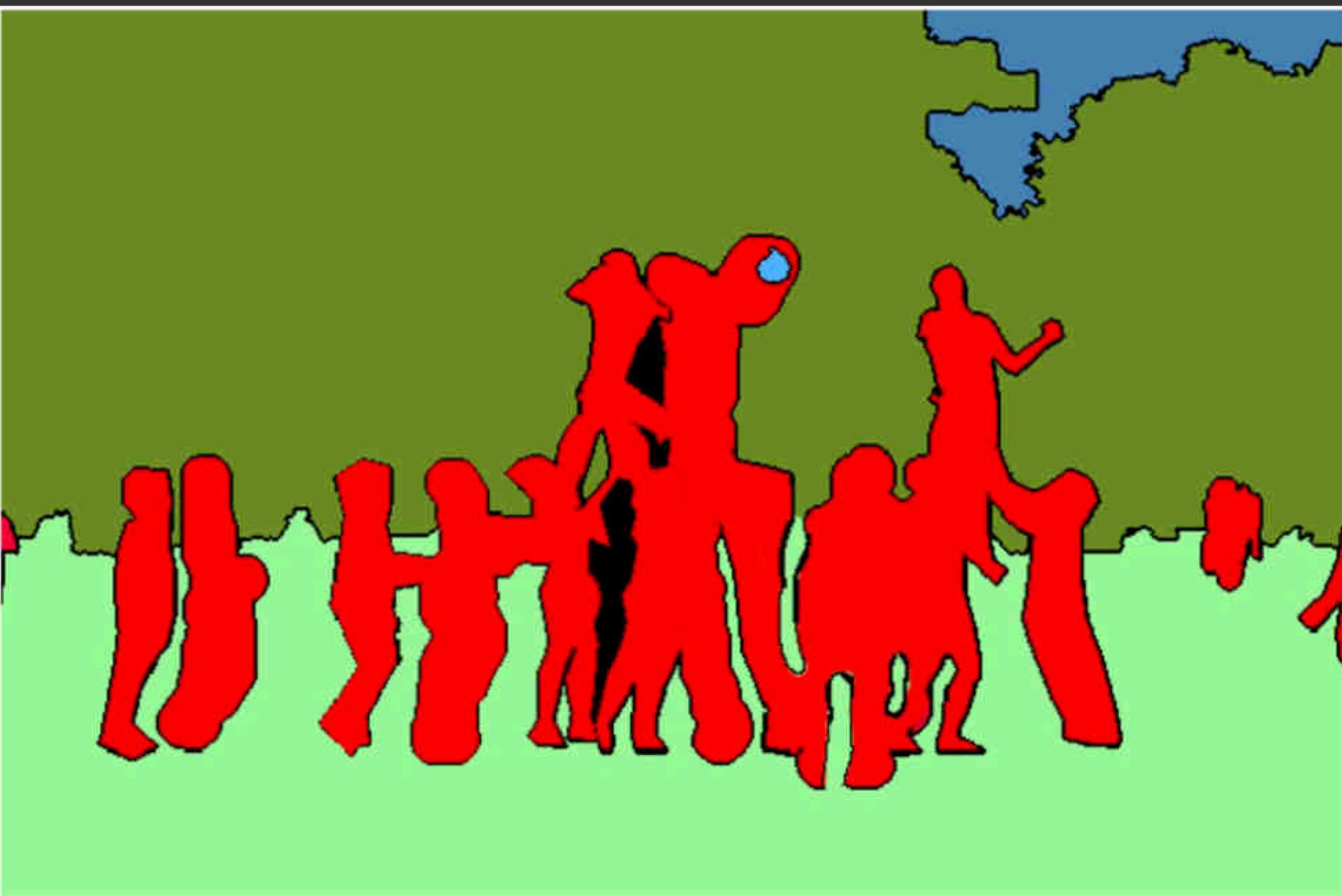


48 / 48

+ Code + Text

Connect |

- **Semantic Segmentation**: Each pixel is labeled in the image in categories. As shown in the following figure, the image is divided into human (red), tree (dark green), grass (light green), sky (blue).



+ Code + Text

Connect |  

How does the data looks?

```
[ ] import os  
import cv2  
import pandas as pd  
import matplotlib.pyplot as plt  
import numpy as np  
import tensorflow as tf
```

- Now lets look at the dataset of our bussniess case

```
[ ] from google.colab import drive  
drive.mount('/content/gdrive/', force_remount=True)  
path = '/content/gdrive/My Drive/Segmentation/'
```

Mounted at /content/gdrive/

```
[ ] images = np.load(path+'data/img_uint8.npy')  
  
mask = np.load(path+'data/msk_uint8.npy')  
  
print(f'Training X data: {images.shape}, Ground Truth: {mask.shape}')  
  
print(f'max pixel value: {mask.max()}, min pixel value: {mask.min()}')
```

Training X data: (18698, 128, 128, 3)

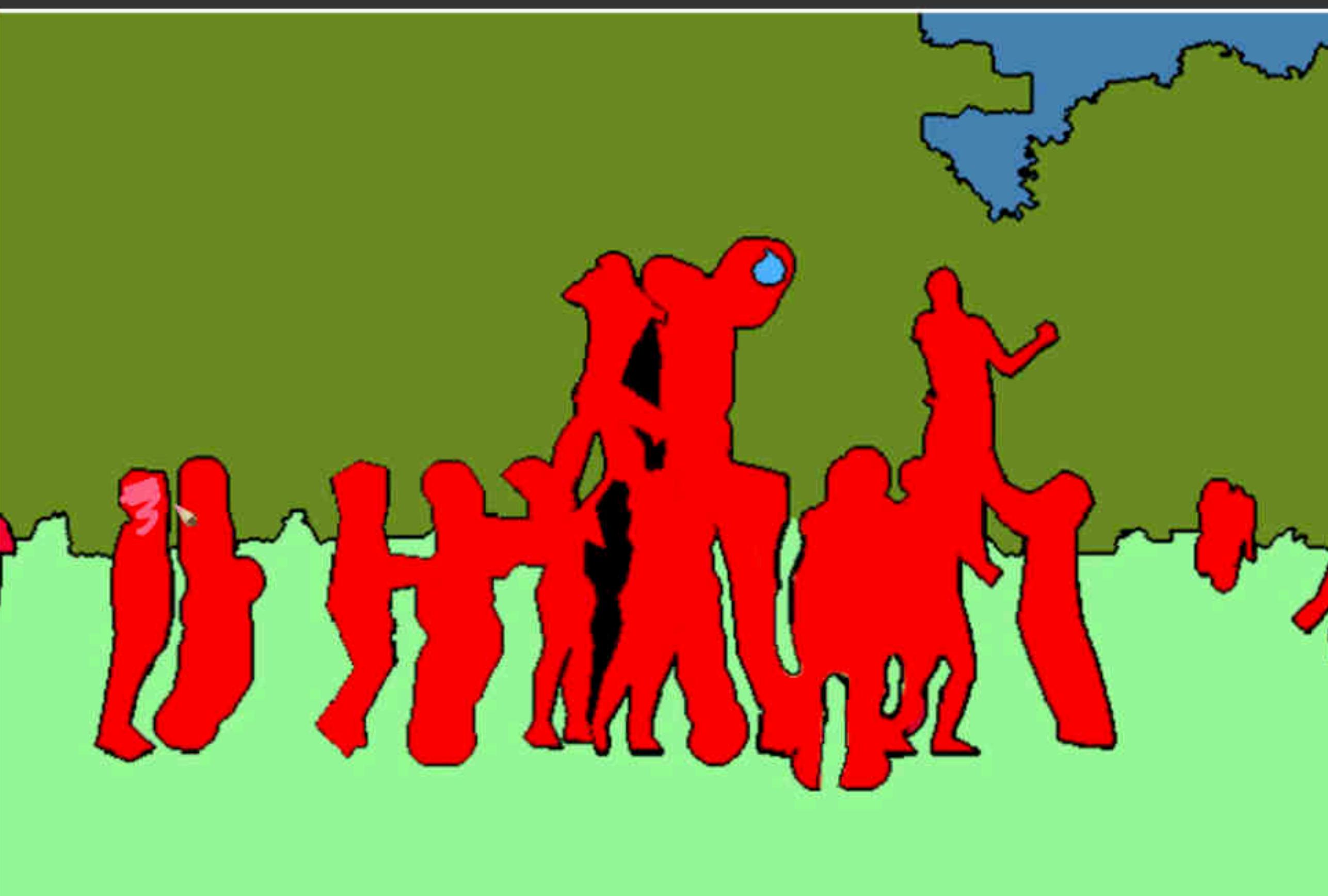


50/50

+ Code + Text

Connect |

- **Semantic Segmentation**: Each pixel is labeled in the image in categories. As shown in the following figure, the image is divided into human (red), tree (dark green), grass (light green), sky (blue).



+ Code + Text

Connect |  

How does the data looks?

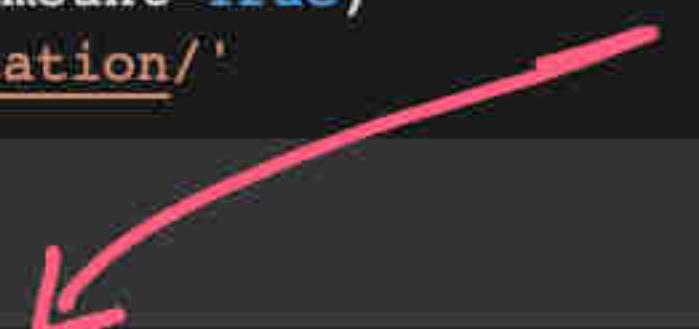
```
[ ] import os  
import cv2  
import pandas as pd  
import matplotlib.pyplot as plt  
import numpy as np  
import tensorflow as tf
```

- Now lets look at the dataset of our bussniess case

```
[ ] from google.colab import drive  
drive.mount('/content/gdrive/', force_remount=True)  
path = '/content/gdrive/My Drive/Segmentation/'
```

Mounted at /content/gdrive/

```
[ ] images = np.load(path+'data/img_uint8.npy')  
mask = np.load(path+'data/msk_uint8.npy')  
  
print(f'Training X data: {images.shape}, Ground Truth: {mask.shape}')  
  
print(f'max pixel value: {mask.max()}, min pixel value: {mask.min()}')
```



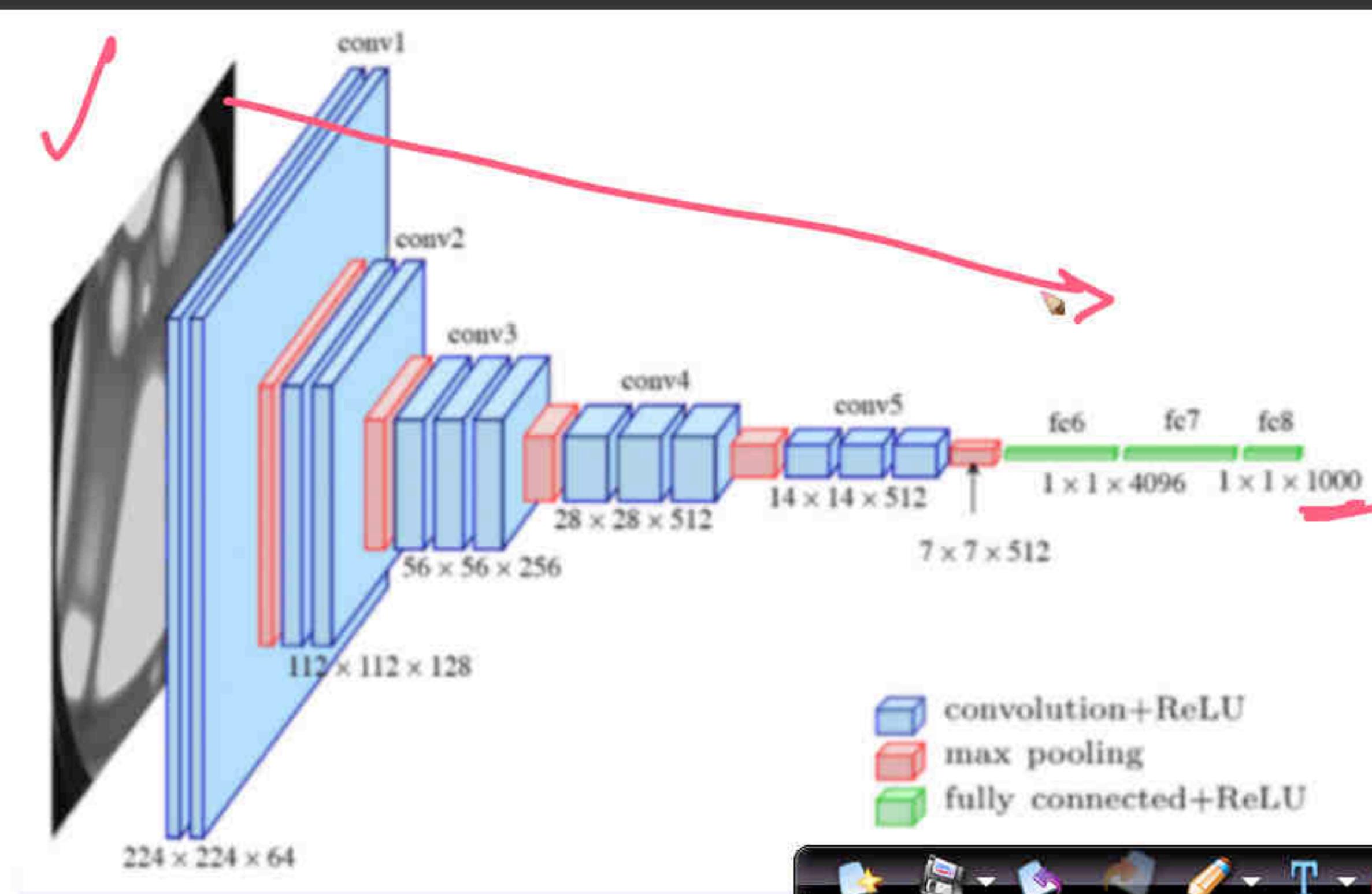
colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

- Instead we can go through every pixel of the image and classify if it is background or not

How should we design the architecture ?

- Here we are predicting labels for each pixel, therefore we need to use a CNN layer with same HxW of the input img shape
- Each class will require one channel, therefore output layer will have dim (bs, H, W, #oflabels)



conv1
conv2
conv3
conv4
conv5
fc6
fc7
fc8

112 × 112 × 128
56 × 56 × 256
28 × 28 × 512
 $14 \times 14 \times 512$
 $7 \times 7 \times 512$

224 × 224 × 64

convolution+ReLU
max pooling
fully connected+ReLU

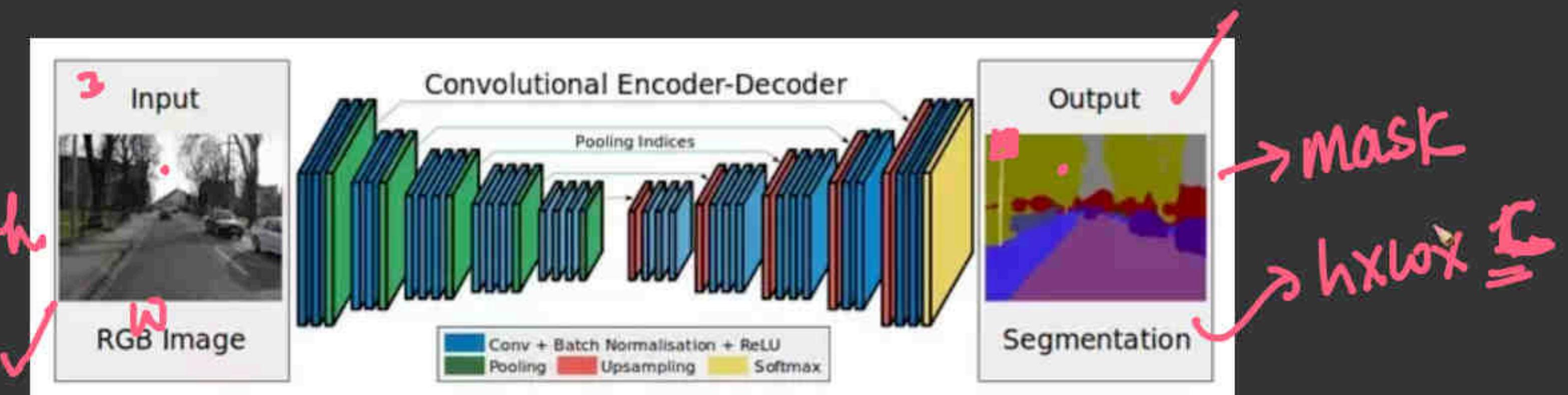
53 / 53

+ Code + Text

Connect



- Final CNN layers conveys Deep, coarse, semantic information: which Answers the "What"
- First few layers conveys Shallow, fine, appearance information: which Answers the "Where"
- In image segmentation we need both "What" and "Where", "What" helps in detecting the label(face/background)
- "Where" helps in detecting the boundary of the "What", in image classification models we lack the "Where" information
- Therefore we skip connect initial layers with final layer (to infuse "Where")

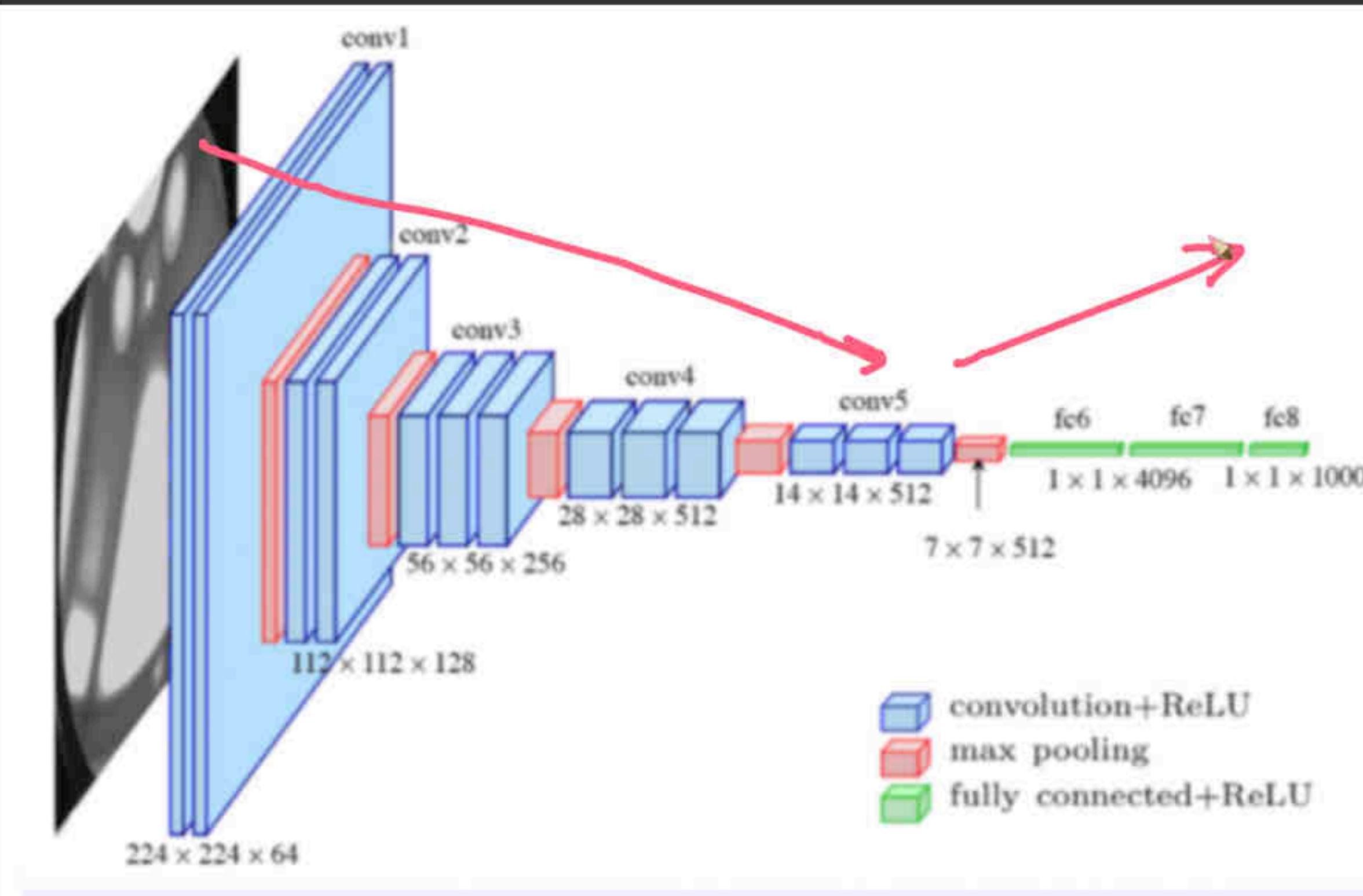


- But How would we increase HxW 7x7 to 224x224?
- [Q] What if we don't use it? Ans. We will lack the "What"
- We use something called Encoder Decoder network
- Encoder part is same as designed for object classification without Flatten/FC layers
- Each encoder layer has a corresponding decoder layers
- Decoder network up-samples its input feature map using transposed Convolutional Layer

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

- Each class will require one channel, therefore output layer will have dim (bs, H, W, #oflabels)



conv1

conv2

conv3

conv4

conv5

112 × 112 × 128

56 × 56 × 256

28 × 28 × 512

14 × 14 × 512

7 × 7 × 512

fc6

fc7

fc8

1 × 1 × 4096

1 × 1 × 1000

convolution+ReLU

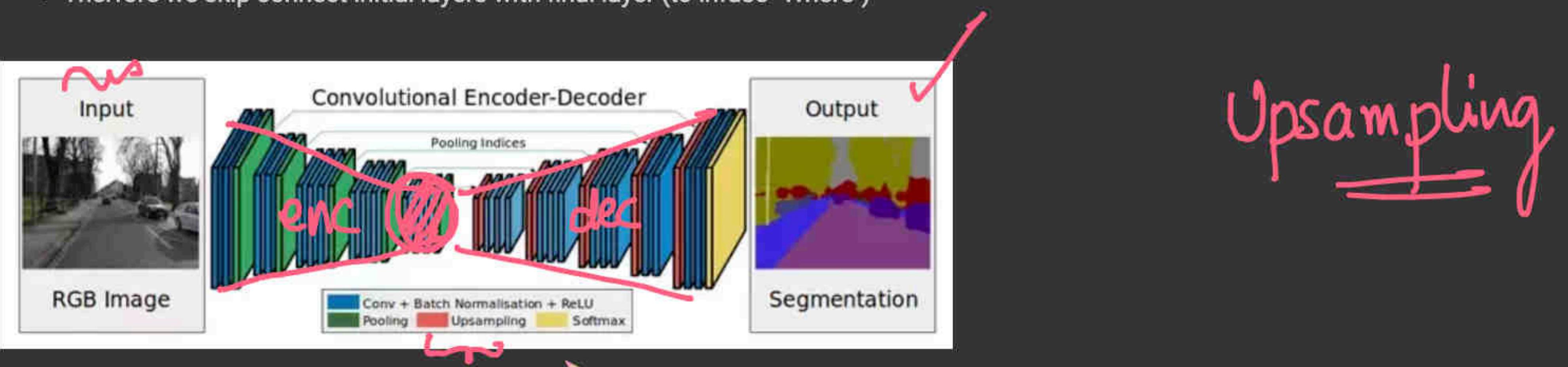
max pooling

fully connected+ReLU

224 × 224 × 64

- Till now we have seen models where HxW dim decreases as we go down the model and number of features increases
- For ex here, HxW → 224 → 112 → 56 → 28 → 14 → 7 and #offeatures 3 → 64 → 128 → 256 → 512
- But we would require 224x224 instead of 7x7 / 1 instead of 512 for our final layer

- Final CNN layers conveys Deep, coarse, semantic information: which Answers the "What"
- First few layers conveys Shallow, fine, appearance information: which Answers the "Where"
- In image segmentation we need both "What" and "Where", "What" helps in detecting the label(face/background)
- "Where" helps in detecting the boundary of the "What", in image classification models we lack the "Where" information
- Therefore we skip connect initial layers with final layer (to infuse "Where")



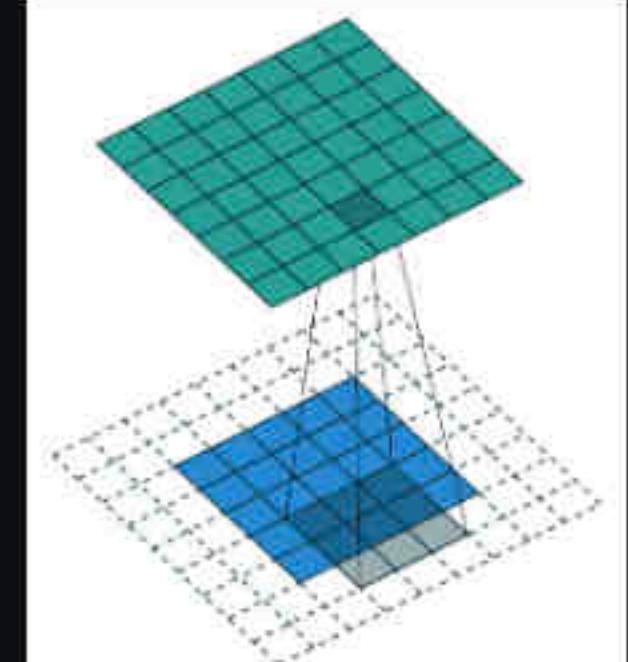
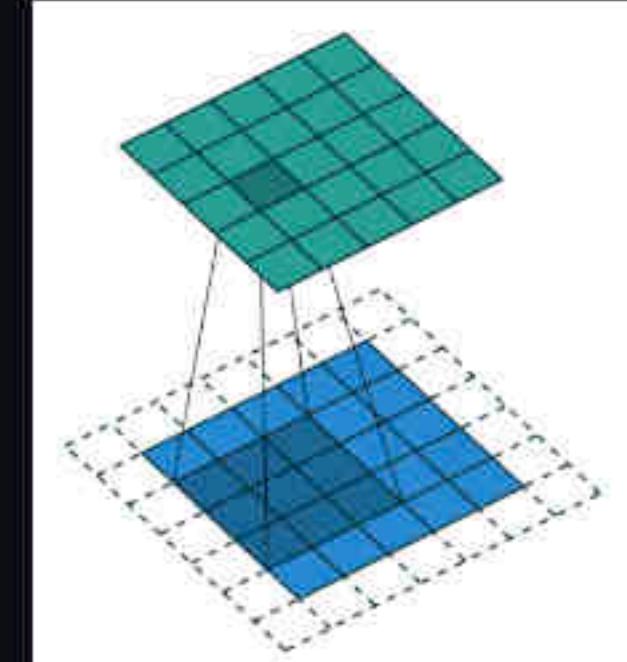
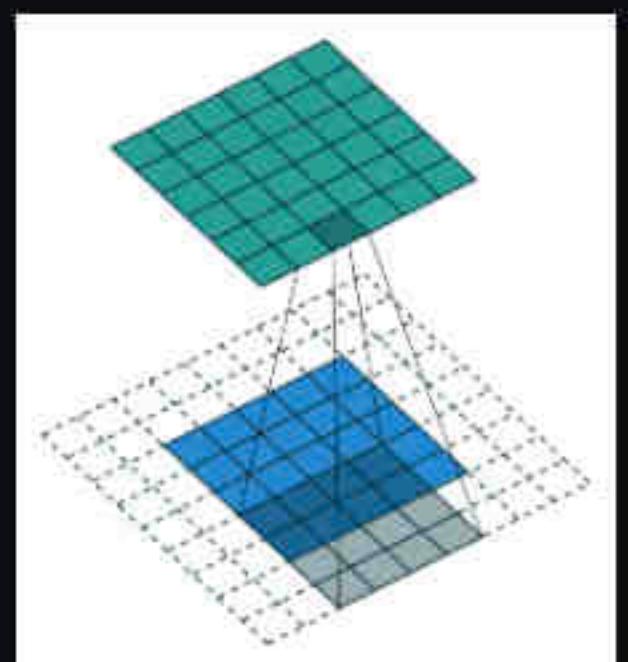
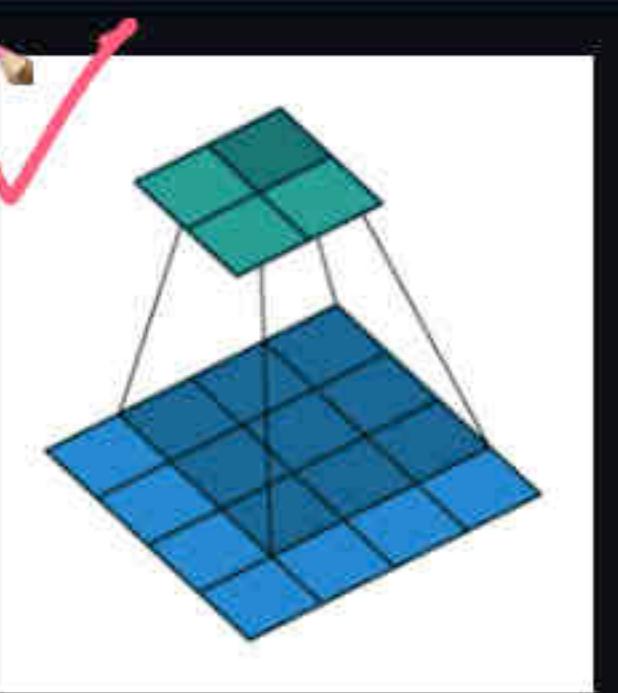
- But How would we increase HxW 7x7 to 224x224?
- [Q] What if we don't use it? Ans . We will lack the "What"
- We use something called Encoder Decoder network
- Encoder part is same as designed for object classification without Flatten/FC layers
- Each encoder layer has a corresponding decoder layers
- Decoder network up-samples its input feature map using transposed Convolutional Layer

README.md

TeX 79.7% Python 20.3%

Convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

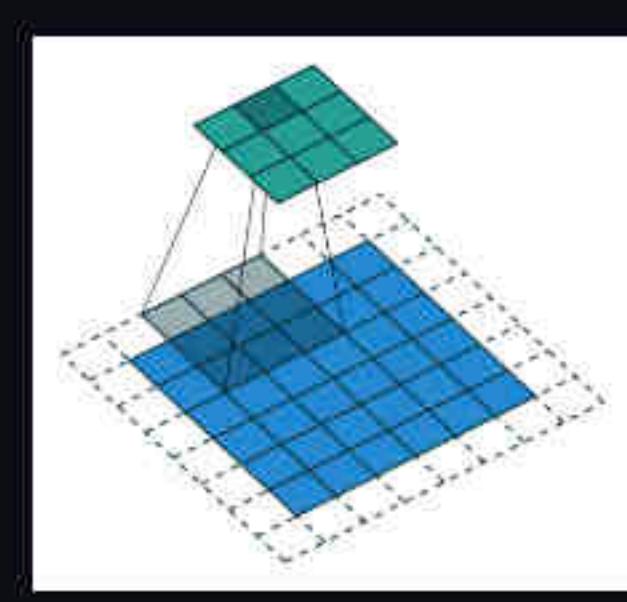
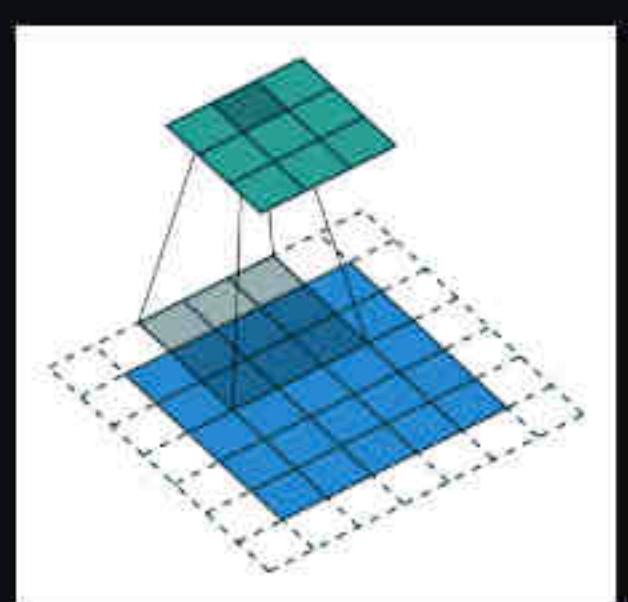
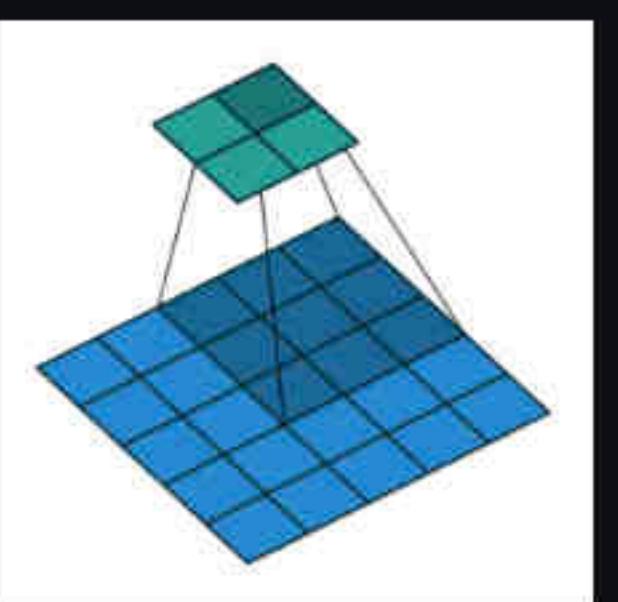


No padding, no
strides

Arbitrary padding, no
strides

Half padding, no
strides

Full padding, no
strides



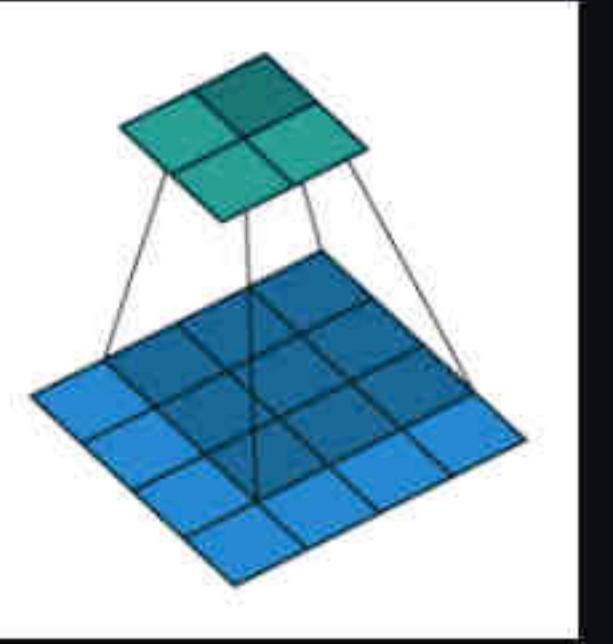
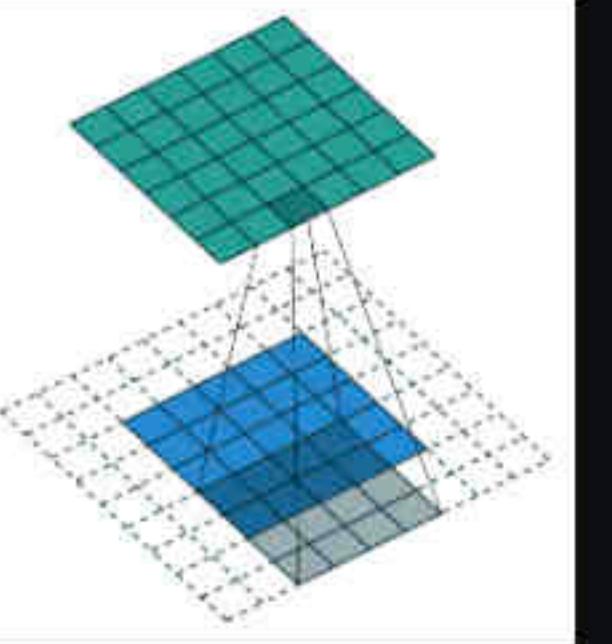
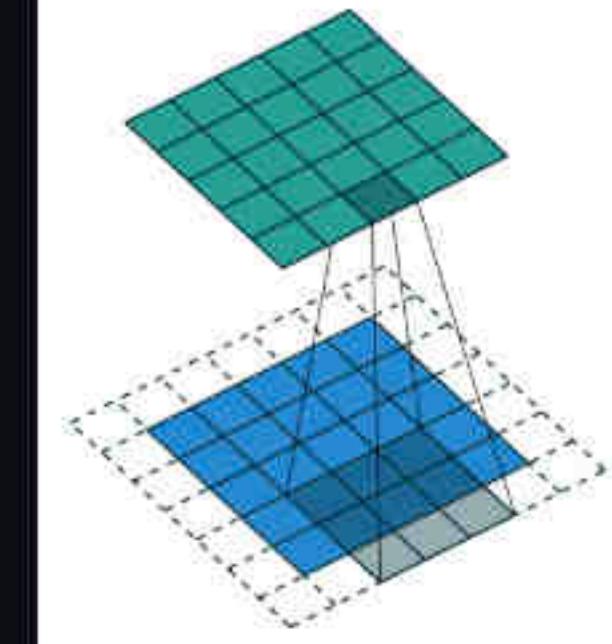
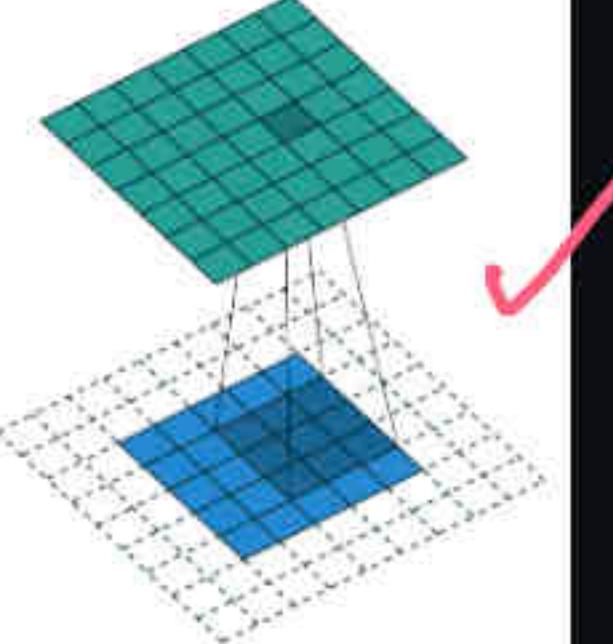
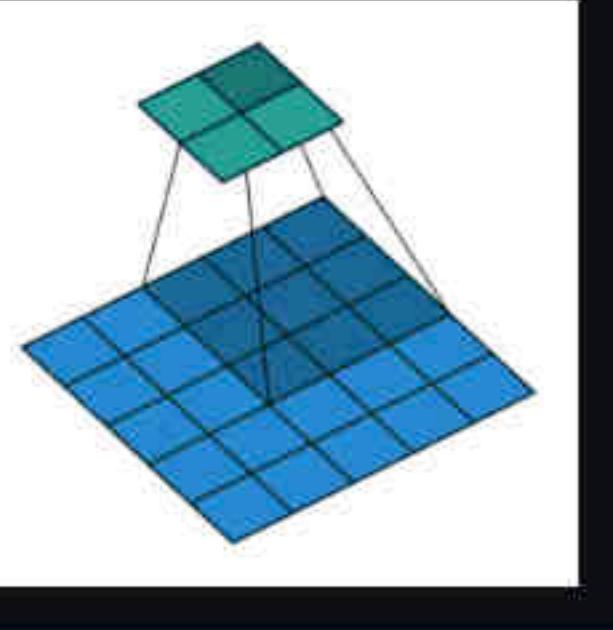
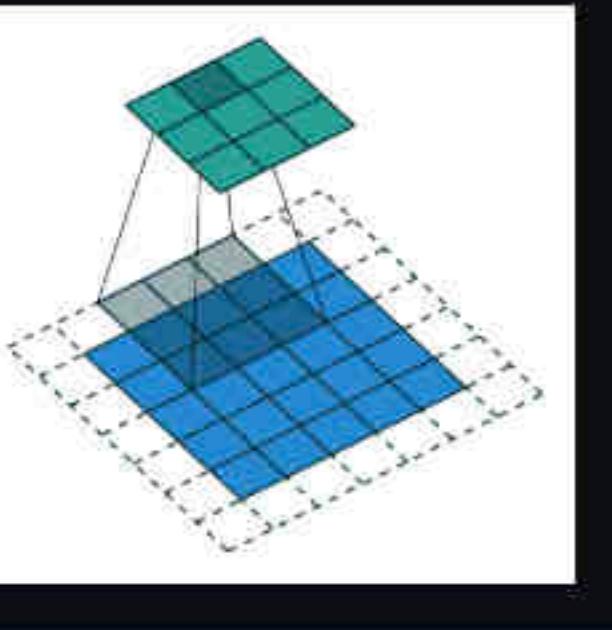
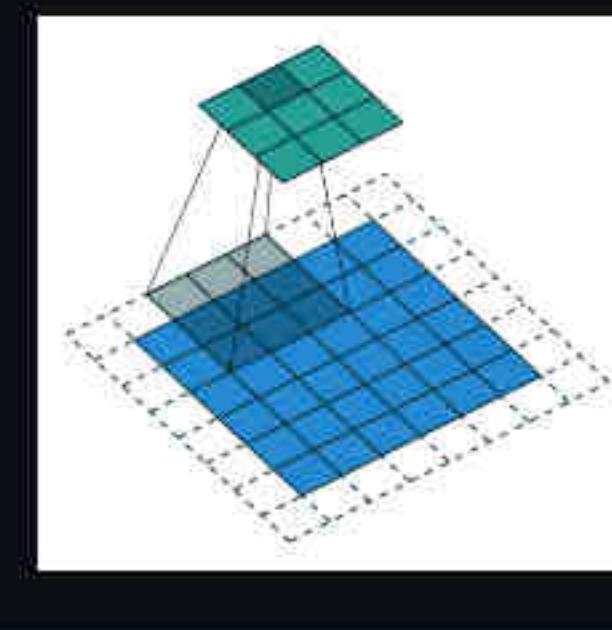
No padding, strides

Padding, strides

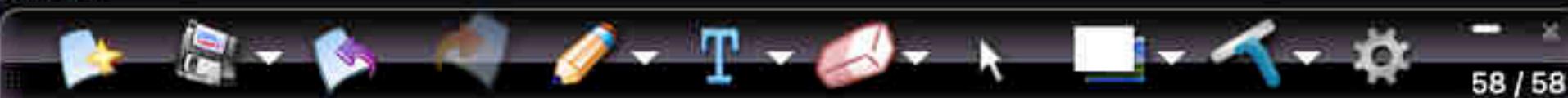
Padding, strides (odd)

README.md

N.B.: Blue maps are inputs, and cyan maps are outputs.

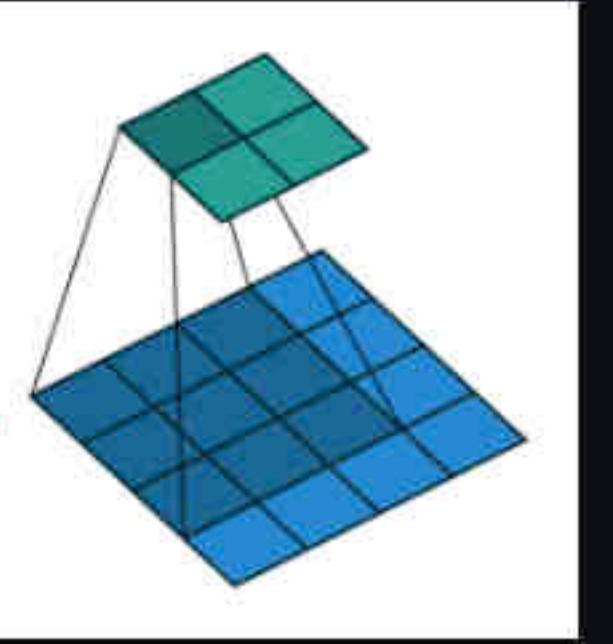
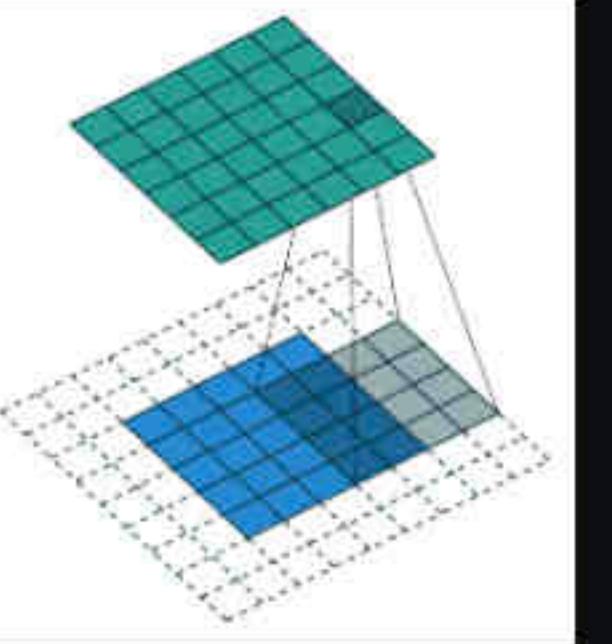
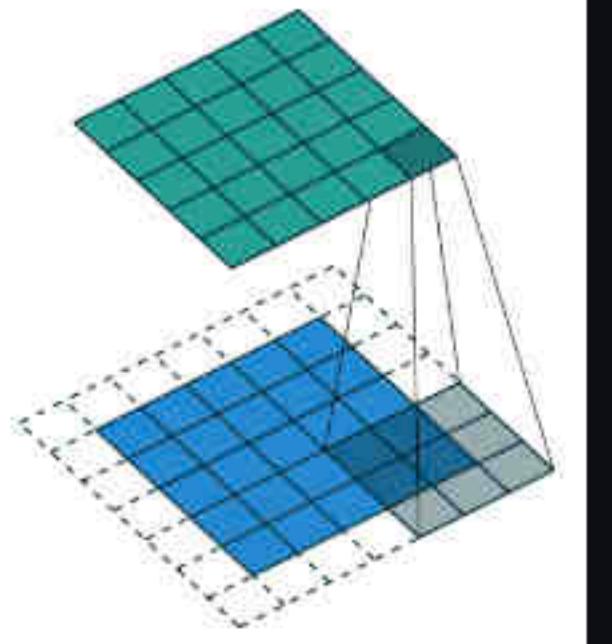
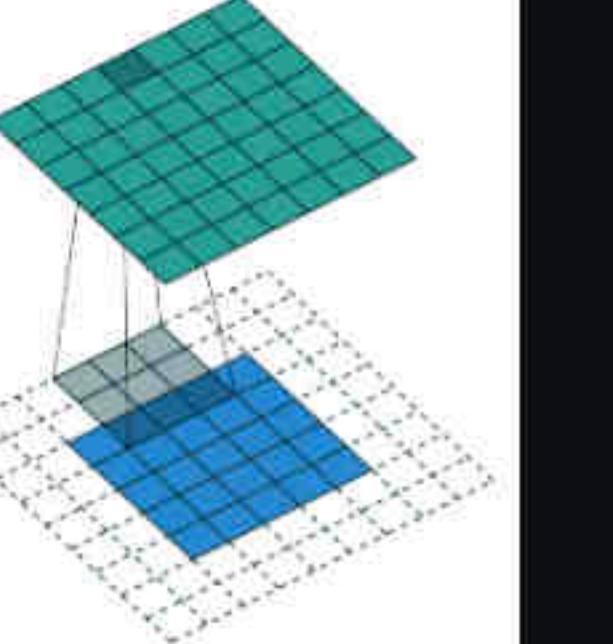
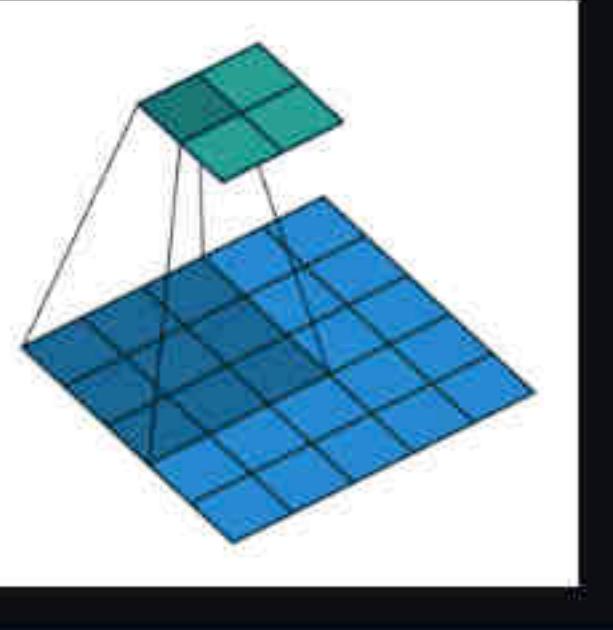
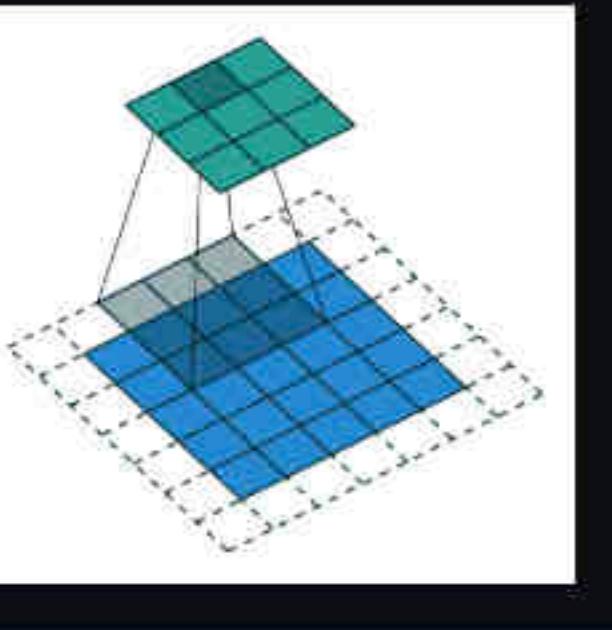
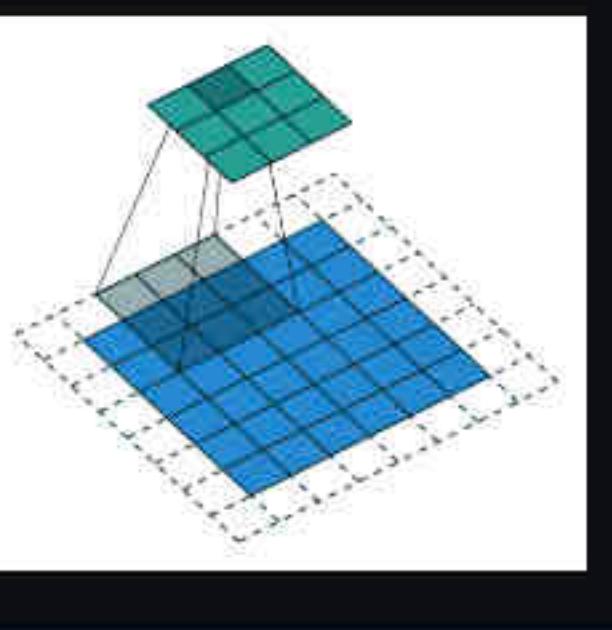
			
No padding, no strides	Arbitrary padding, no strides	Half padding, no strides	Full padding, no strides
			
No padding, strides	Padding, strides	Padding, strides (odd)	

Transposed convolution animations



README.md

N.B.: Blue maps are inputs, and cyan maps are outputs.

 3x3 Conv			
No padding, no strides	Arbitrary padding, no strides	Half padding, no strides	Full padding, no strides
			
No padding, strides	Padding, strides	Padding, strides (odd)	

Transposed convolution animations



README.md

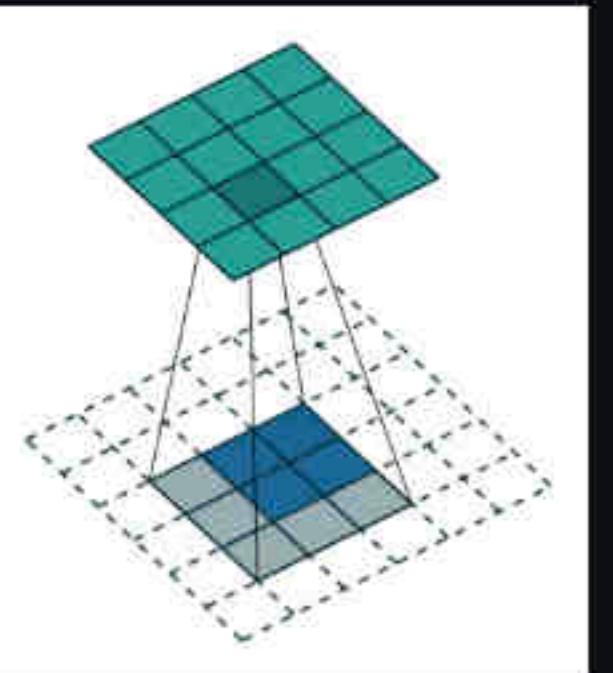
No padding, strides

Padding, strides

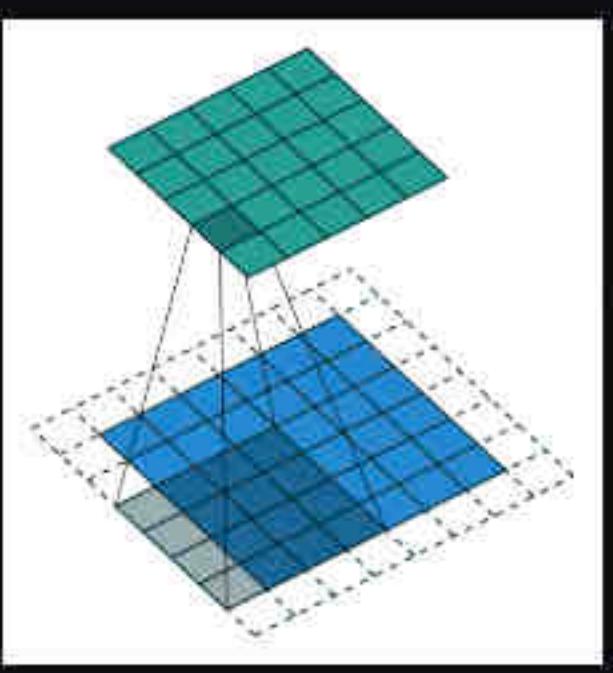
Padding, strides (odd)

Transposed convolution animations

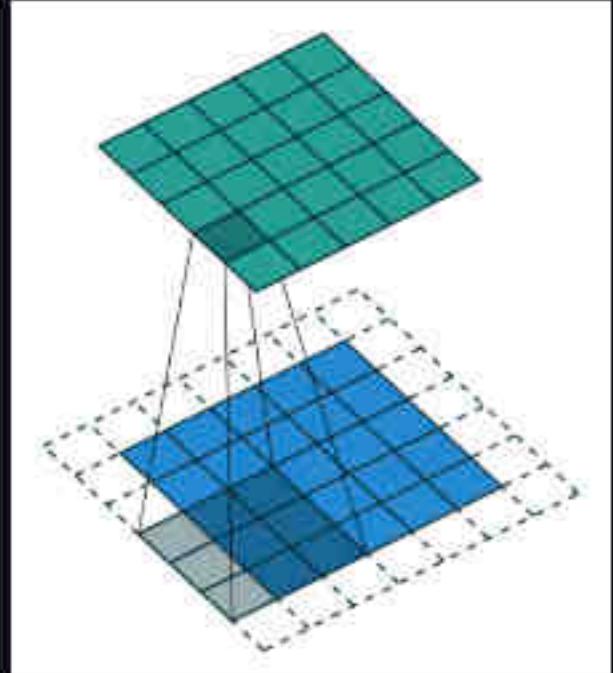
N.B.: Blue maps are inputs, and cyan maps are outputs.



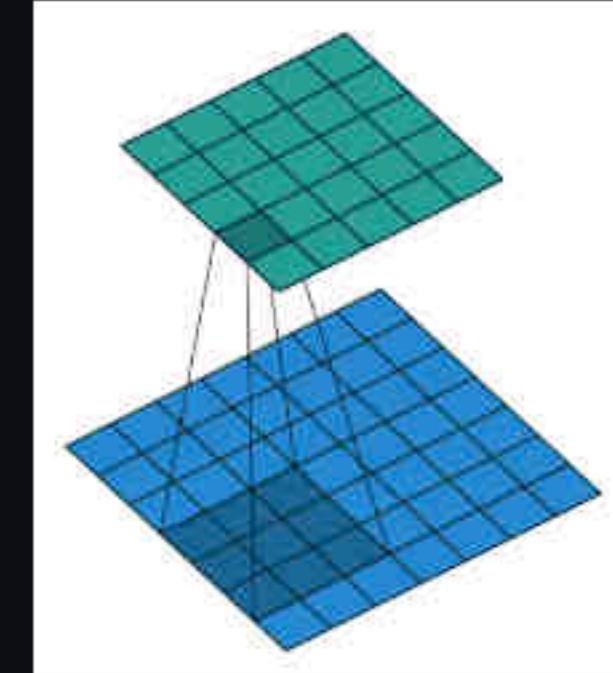
No padding, no strides,
transposed



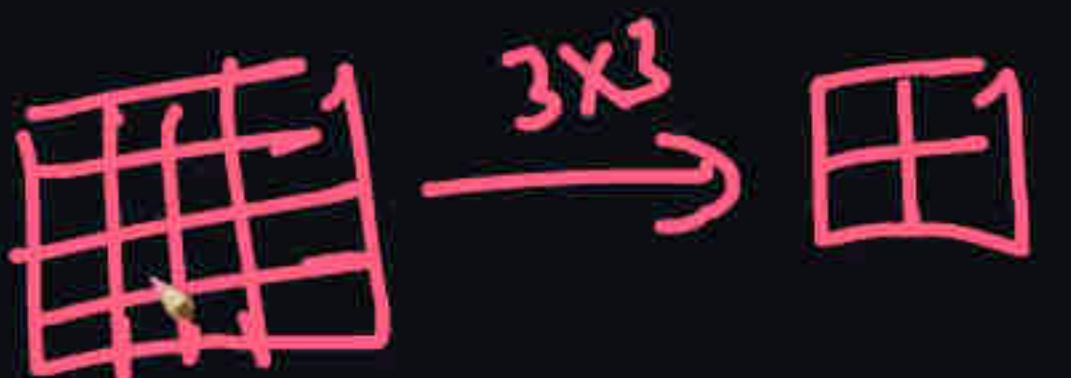
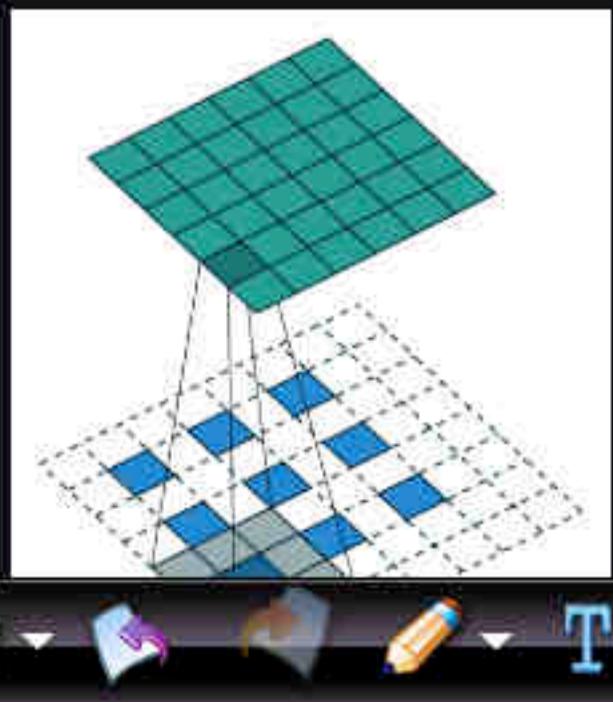
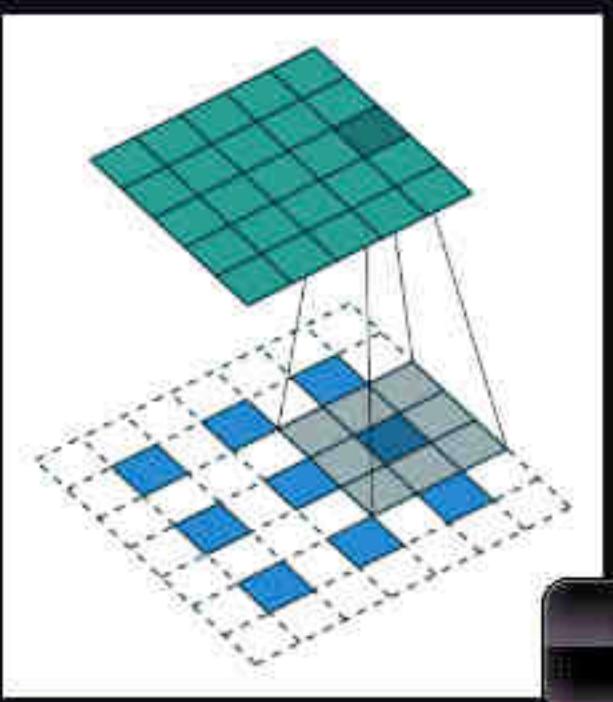
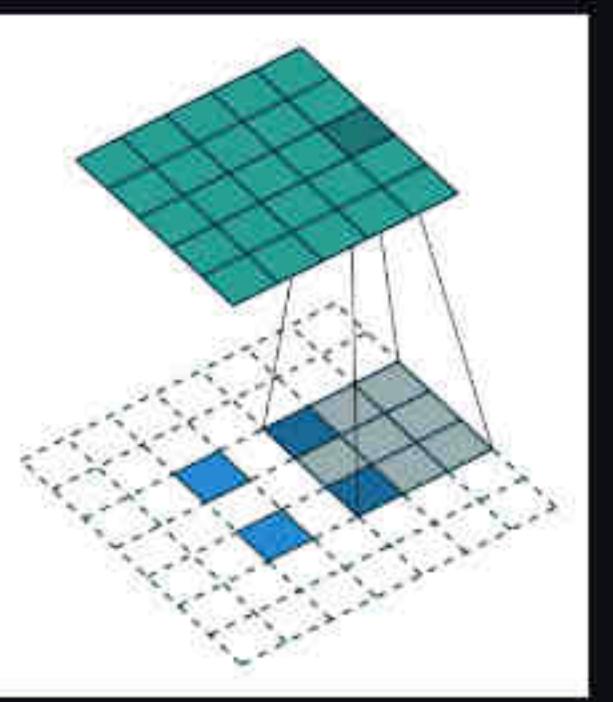
Arbitrary padding, no
strides, transposed



Half padding, no
strides, transposed



Full padding, no strides,
transposed



README.md

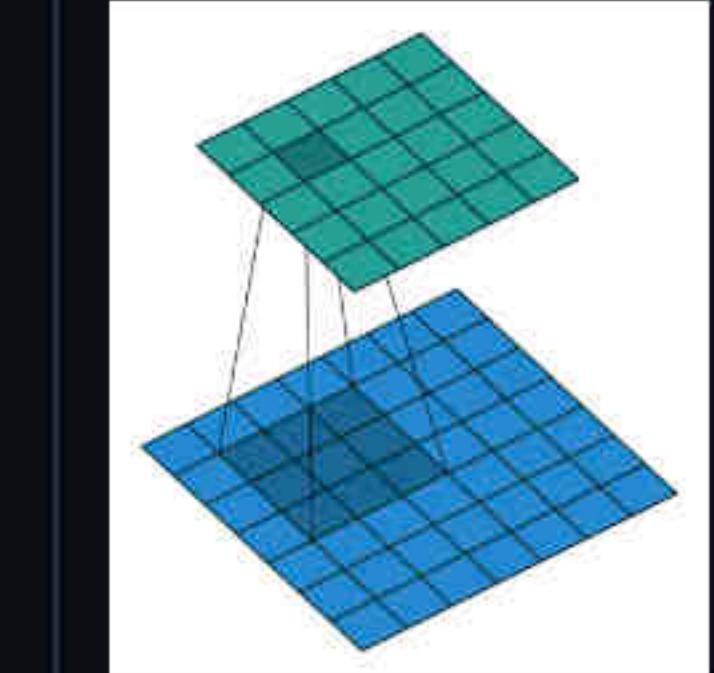
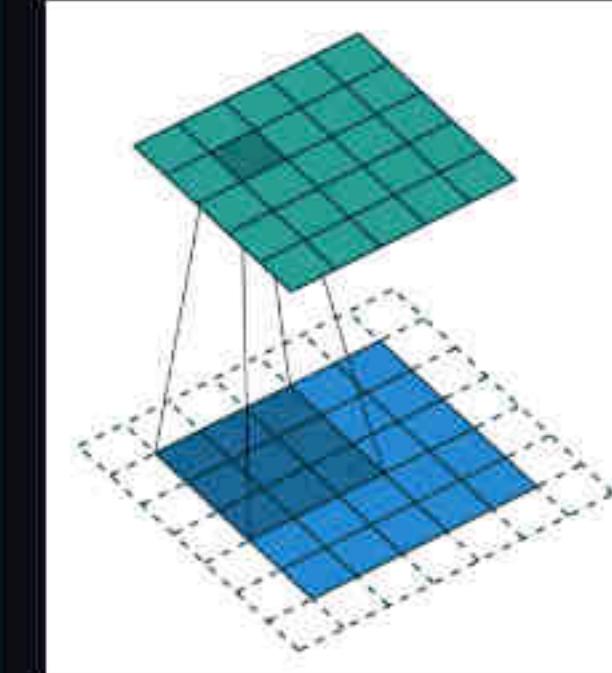
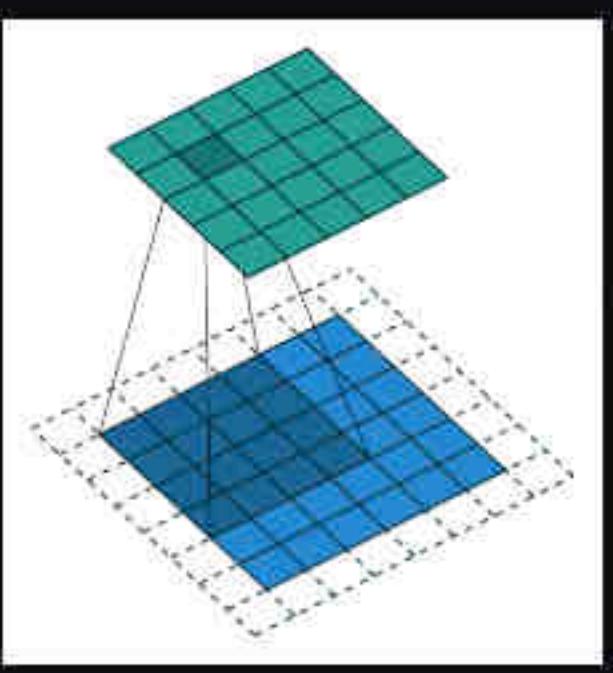
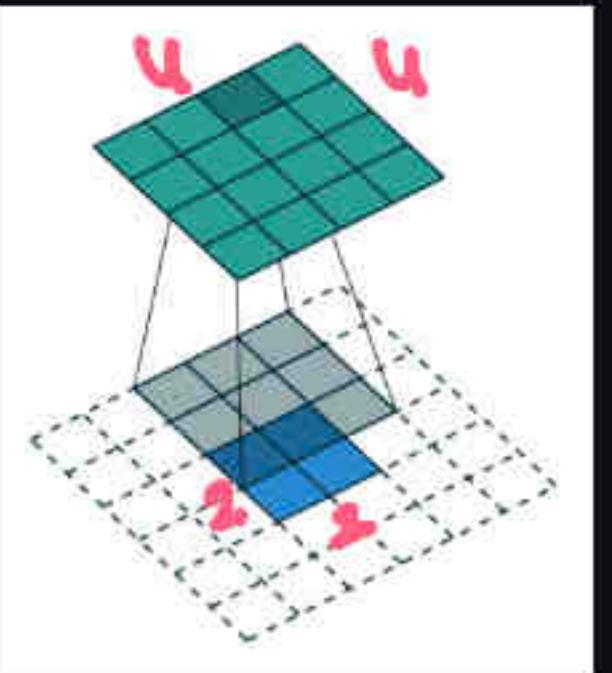
No padding, strides

Padding, strides

Padding, strides (odd)

Transposed convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

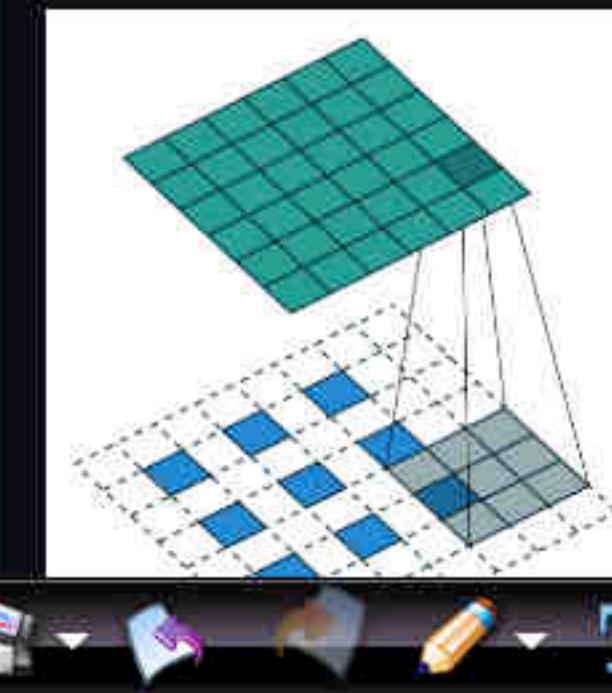
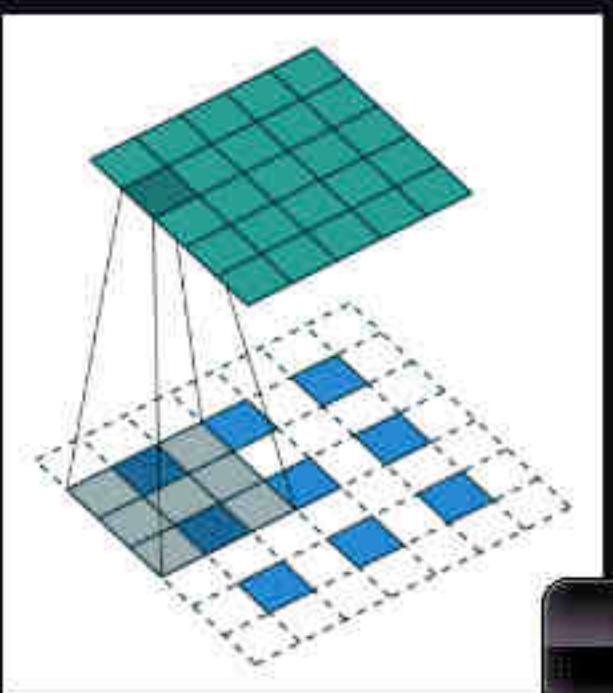
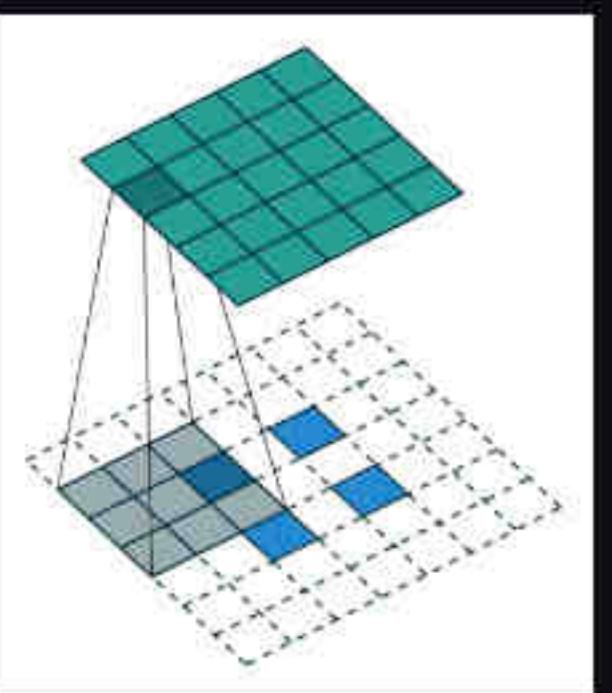


No padding, no strides,
transposed

Arbitrary padding, no
strides, transposed

Half padding, no
strides, transposed

Full padding, no strides,
transposed



github.com/vdumoulin/conv_arithmetic

README.md

No padding, strides Padding, strides Padding, strides (odd)

Transposed convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

3x3 CONV ✓

No padding, no strides, transposed

Arbitrary padding, no strides, transposed

Half padding, no strides, transposed

Full padding, no strides, transposed

Tr. conv

62 / 62

README.md

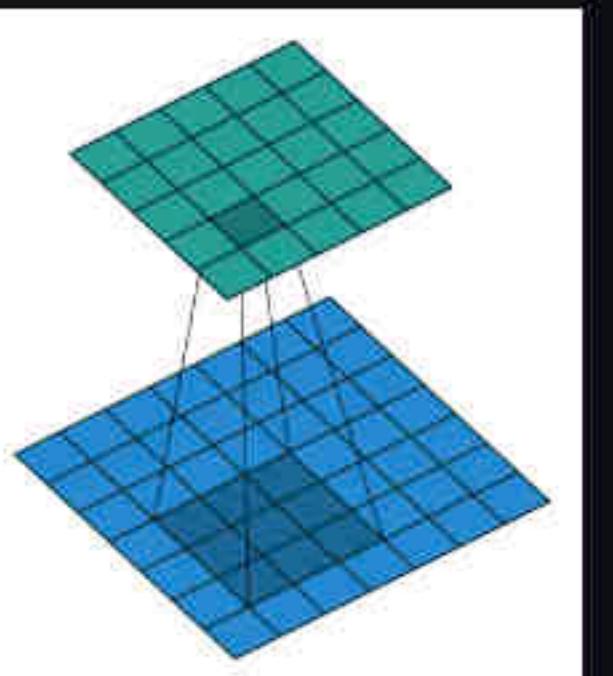
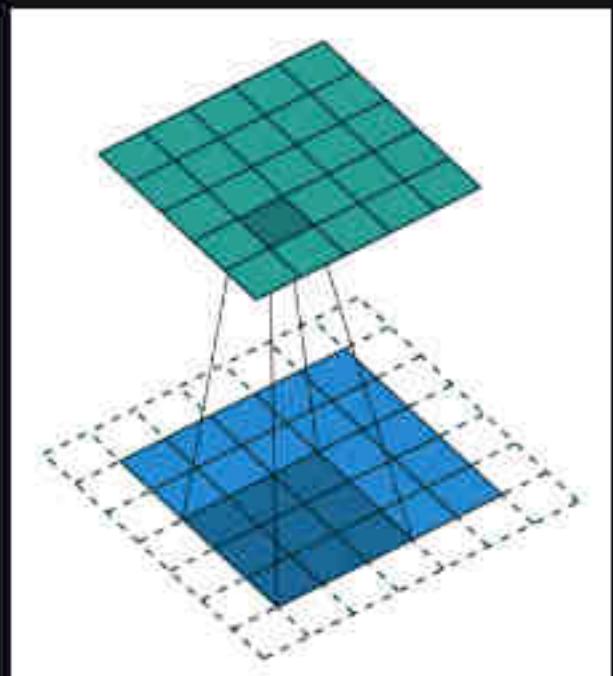
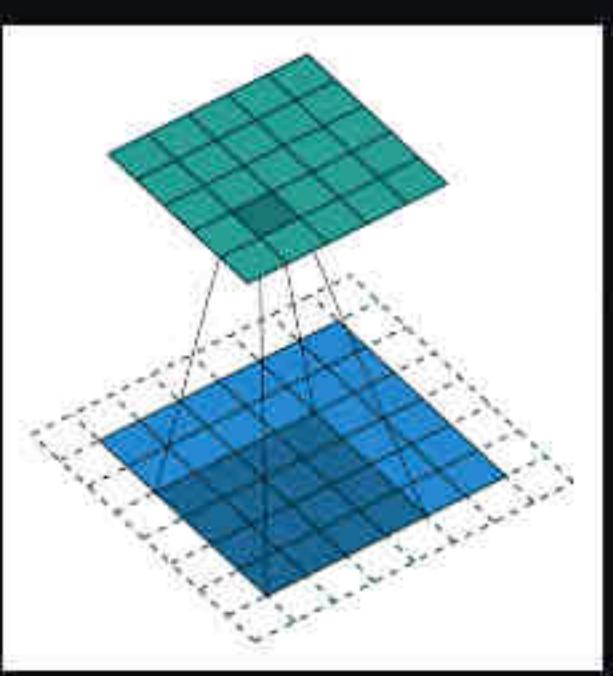
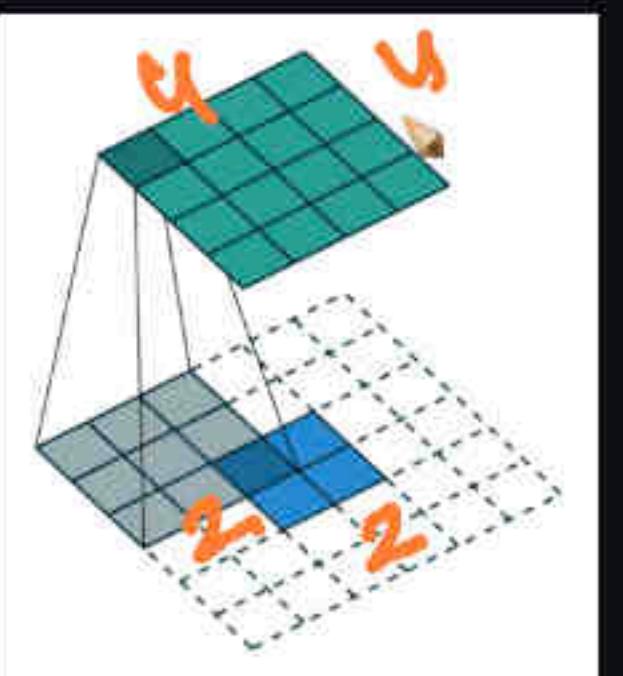
No padding, strides

Padding, strides

Padding, strides (odd)

Transposed convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

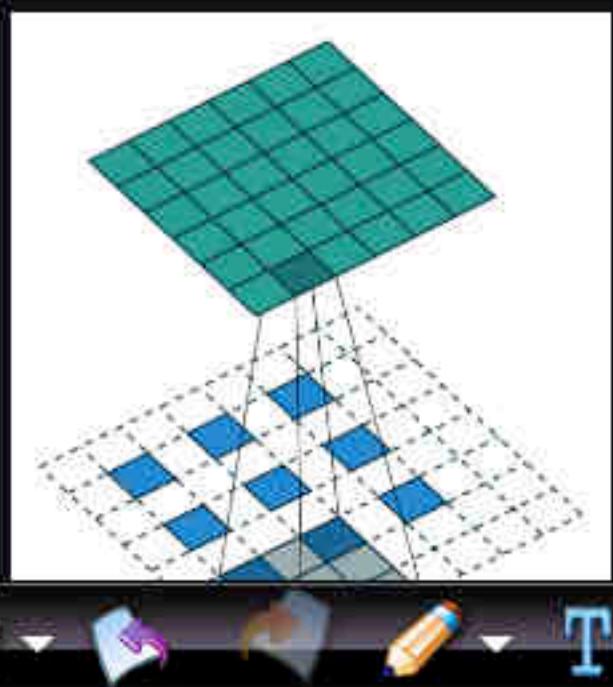
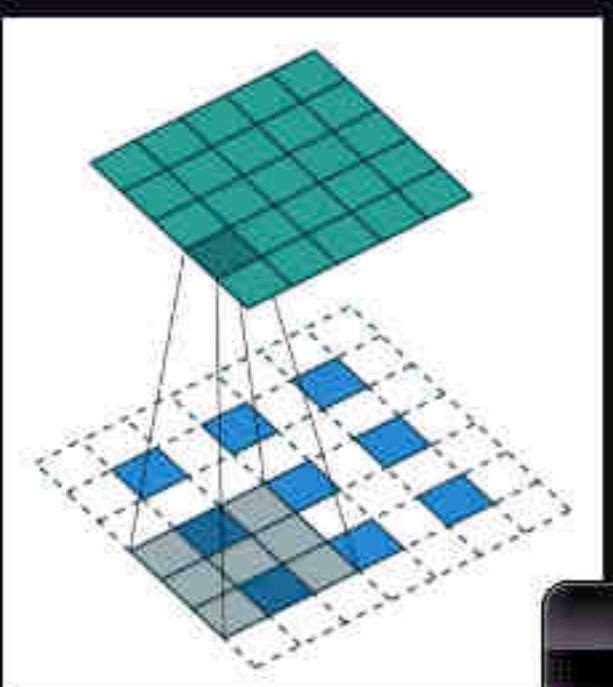
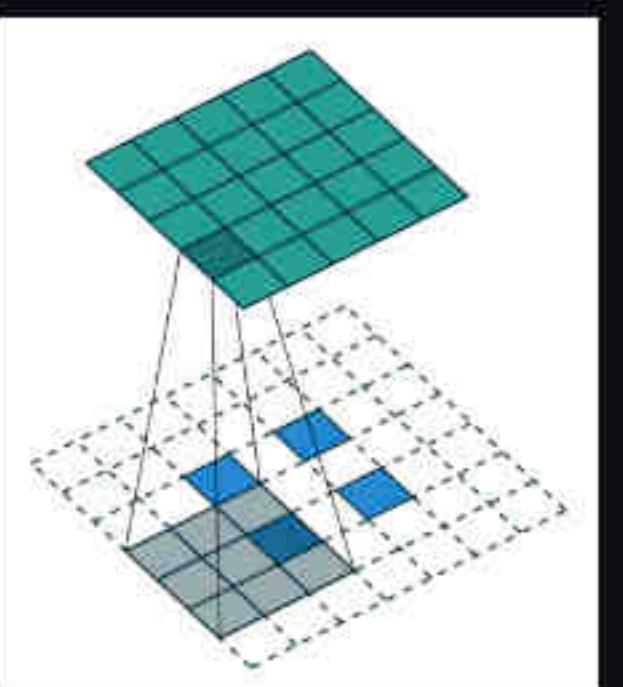


No padding, no strides,
transposed

Arbitrary padding, no strides, transposed

Half padding, no
strides, transposed

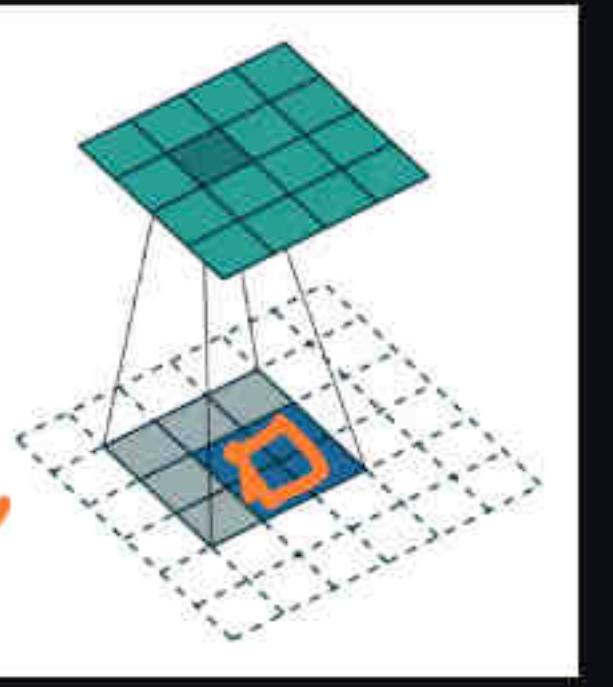
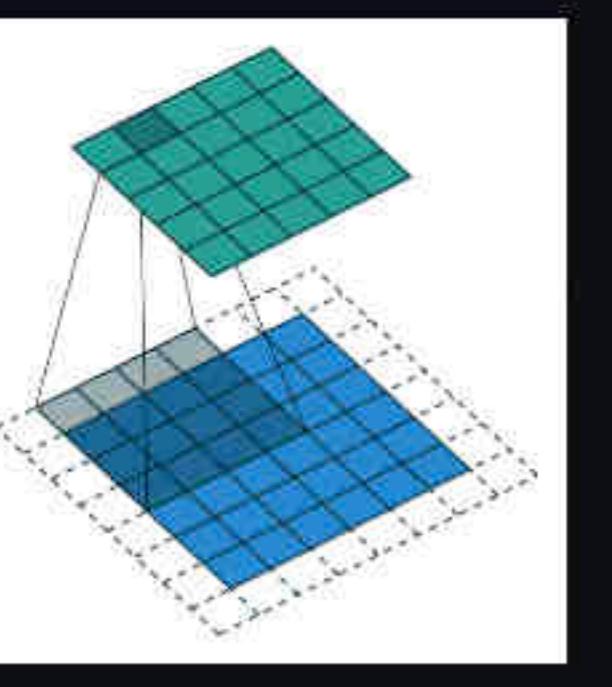
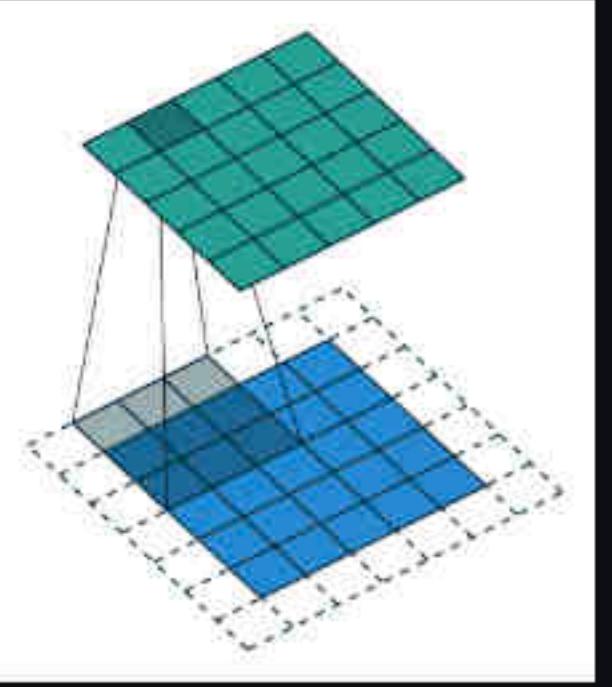
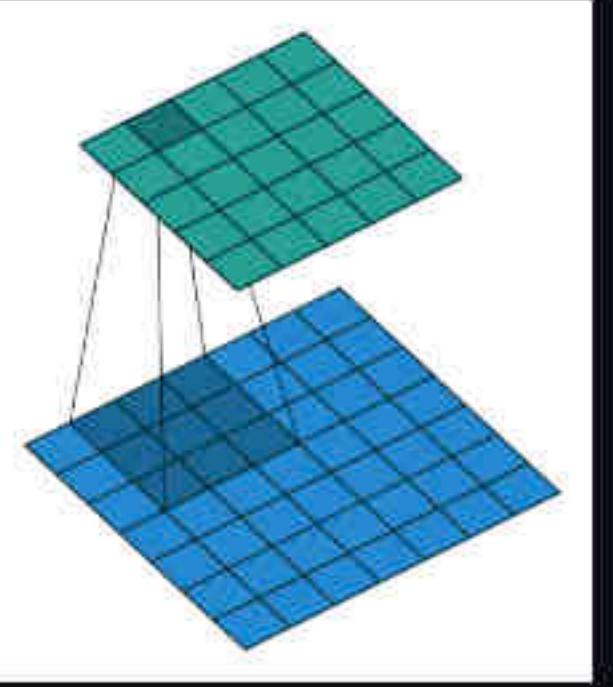
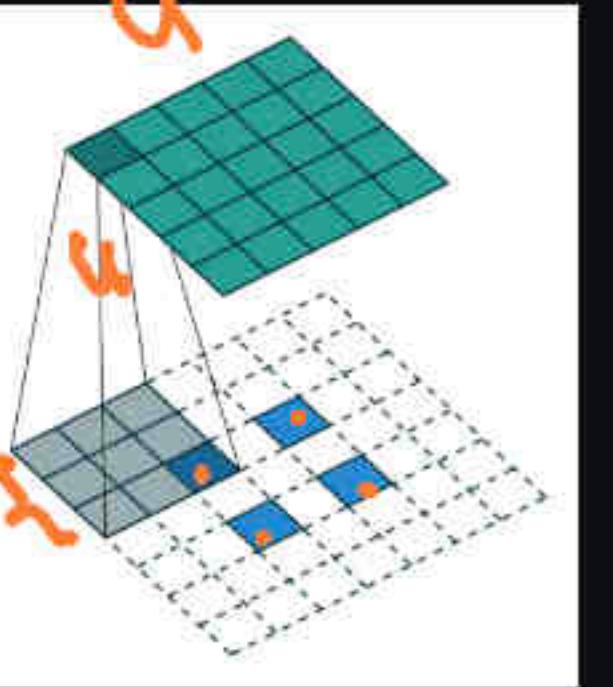
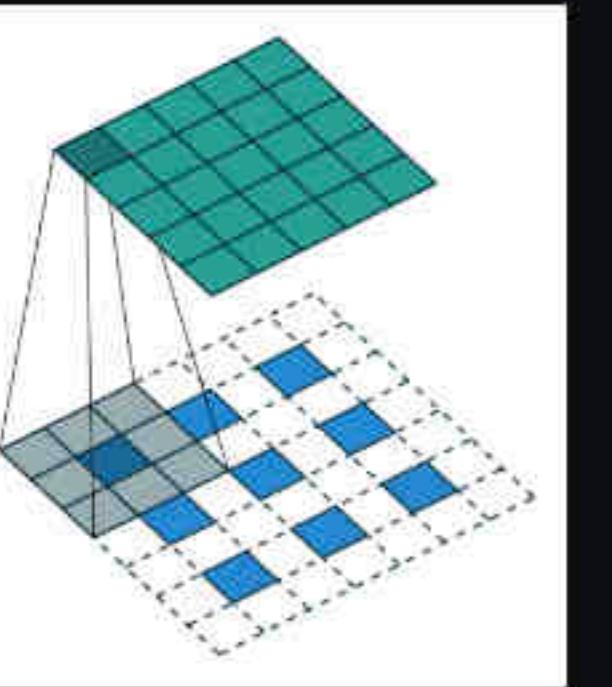
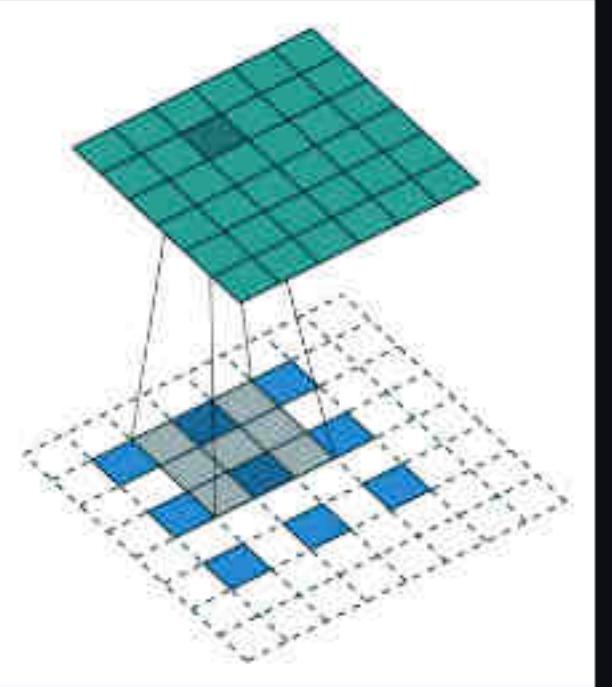
Full padding, no strides,
transposed



github.com/vdumoulin/conv_arithmetic

README.md

N.B.: Blue maps are inputs, and cyan maps are outputs.

			
No padding, no strides, transposed	Arbitrary padding, no strides, transposed	Half padding, no strides, transposed	Full padding, no strides, transposed
			
No padding, strides, transposed	Padding, strides, transposed	Padding, strides,	

3x3

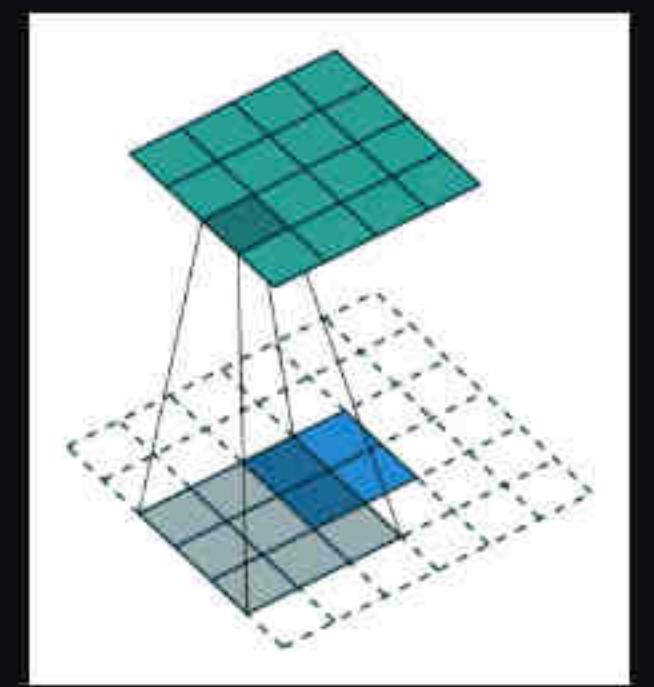
Conv $t \rightarrow s$

Tr·Conv $s \rightarrow l$

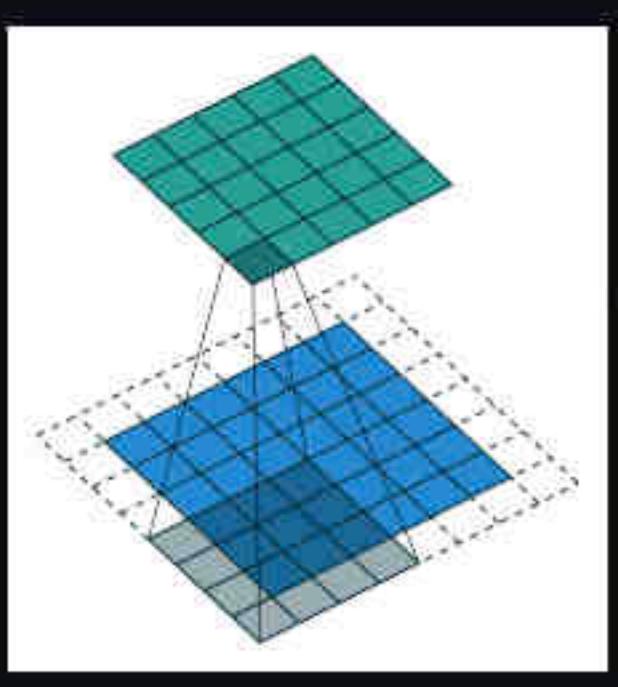
64 / 64

README.md

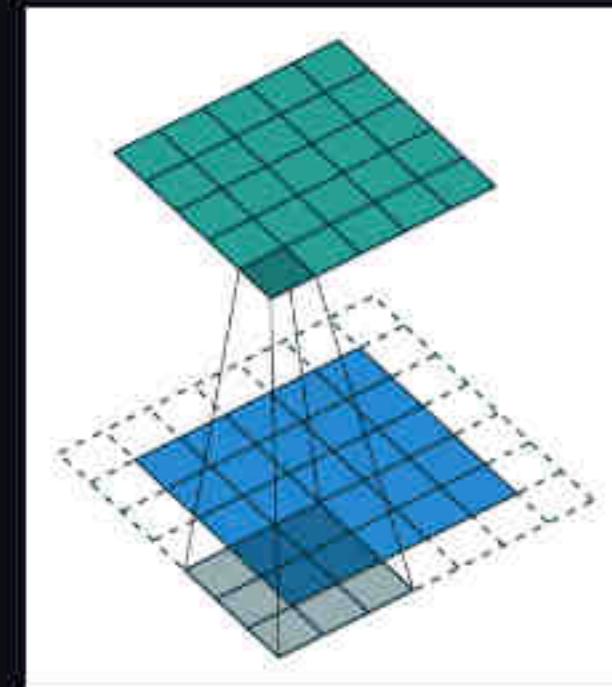
N.B.: Blue maps are inputs, and cyan maps are outputs.



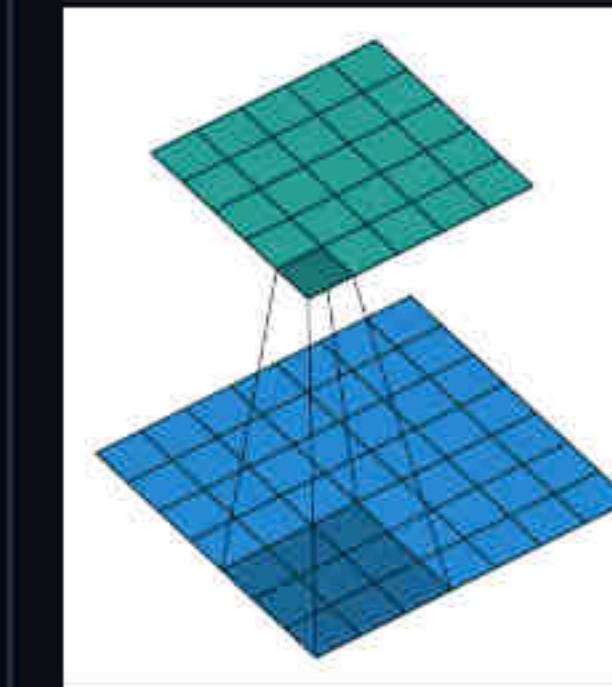
No padding, no strides,
transposed



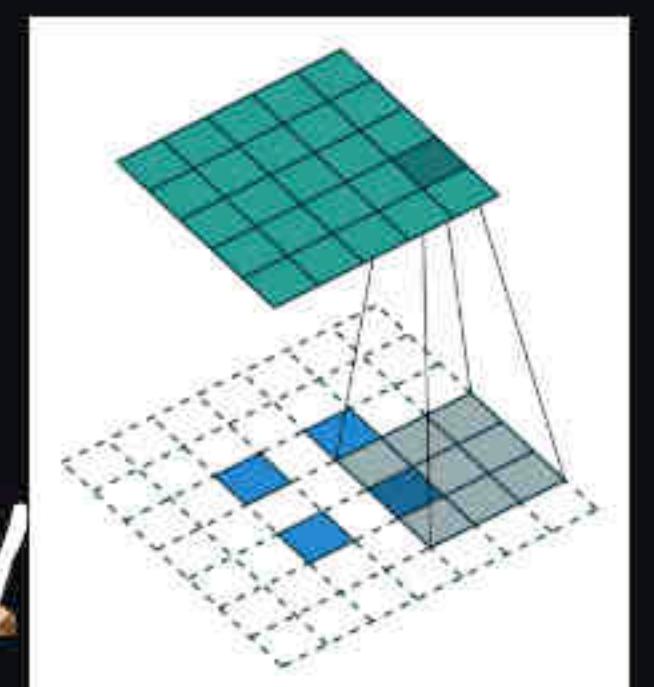
Arbitrary padding, no
strides, transposed



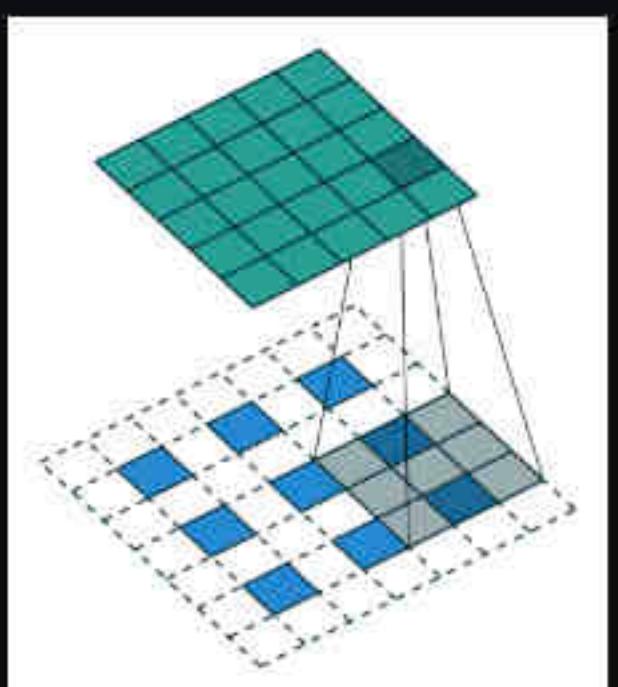
Half padding, no
strides, transposed



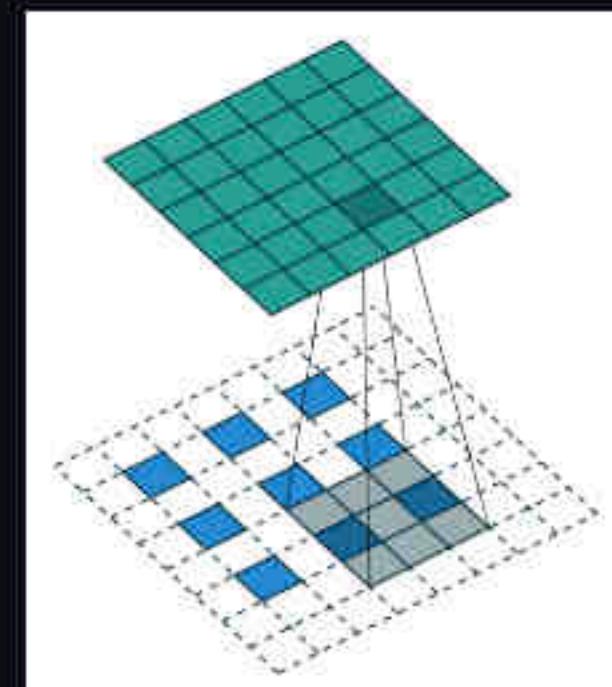
Full padding, no strides,
transposed



No padding, strides,
transposed



Padding, strides,
transposed



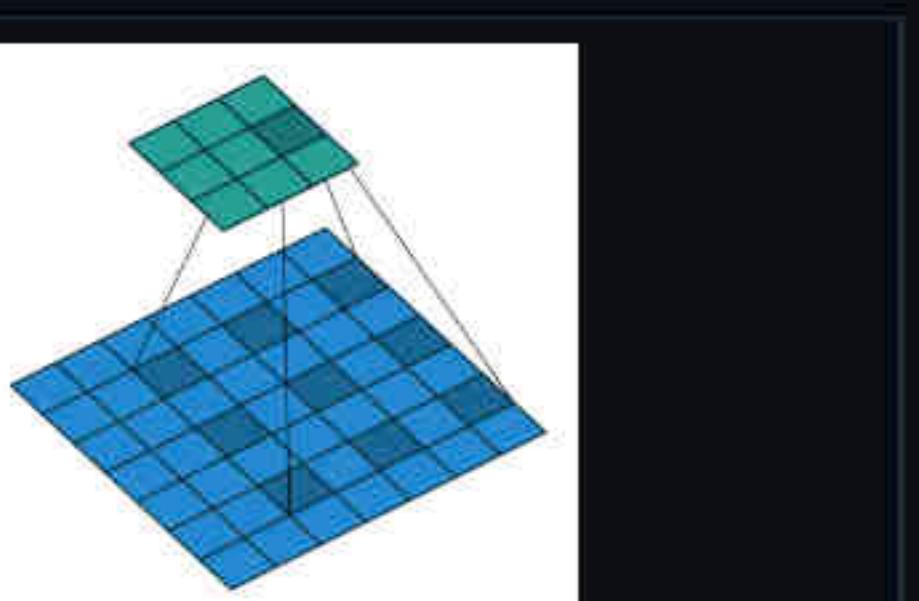
Padding, strides,
transposed

github.com/vdumoulin/conv_arithmetic

README.md

Dilated convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.



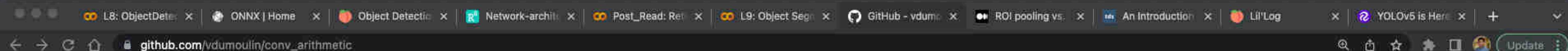
No padding, no stride, dilation

Generating the Makefile

From the repository's root directory:

```
$ ./bin/generate_makefile
```

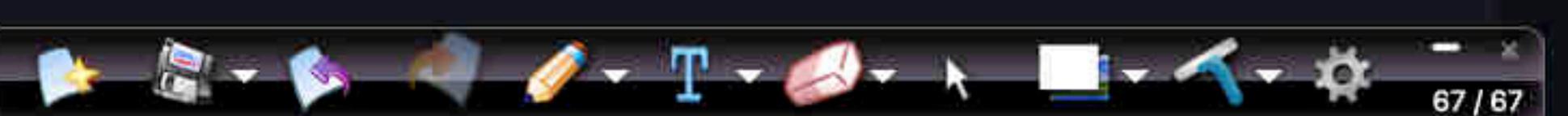
Generating the animations



Generating the Makefile

From the repository's root directory:

```
$ ./bin/generate_makefile
```



L8: ObjectDet... X | ONNX | Home X | Object Detectio... X | Network-archi... X | Post_ Read: Re... X | L9: Object Seg... X | GitHub - vdum... X | ROI pooling vs... X | An Introduction X | LiLog X | YOLOv5 is Here... X | +

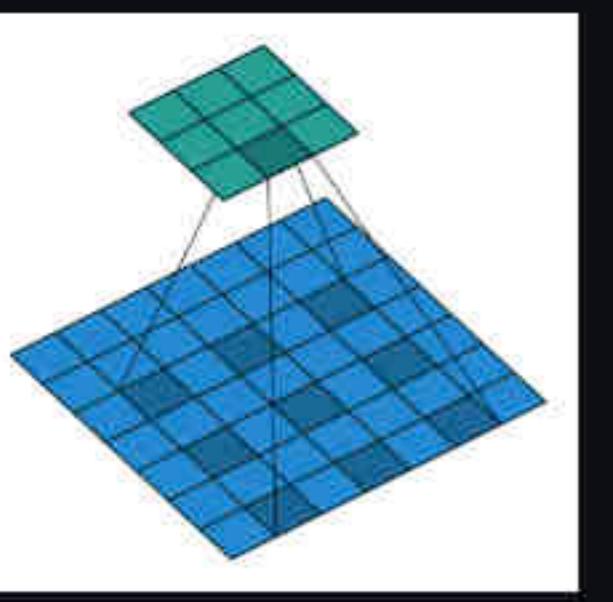
github.com/vdumoulin/conv_arithmetic

README.md

No padding, strides, transposed	Padding, strides, transposed	Padding, strides, transposed (odd)
---------------------------------	------------------------------	------------------------------------

Dilated convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

 A diagram showing a 3x3 blue convolution kernel applying to a 5x5 blue input map. The output is a 3x3 cyan map. Handwritten annotations in pink read "3x3 conv" and "dilated".

No padding, no stride, dilation

Generating the Makefile

From the repository's root directory:

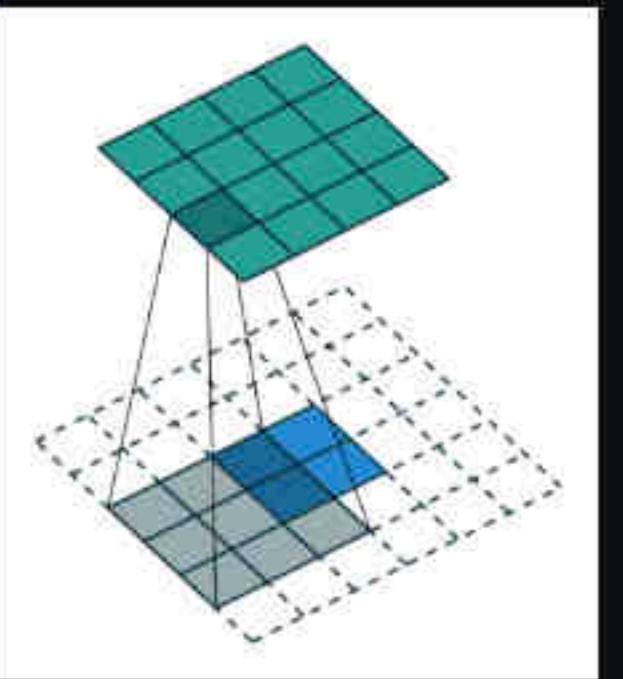
```
$ ./bin/generate_makefile
```

68 / 68

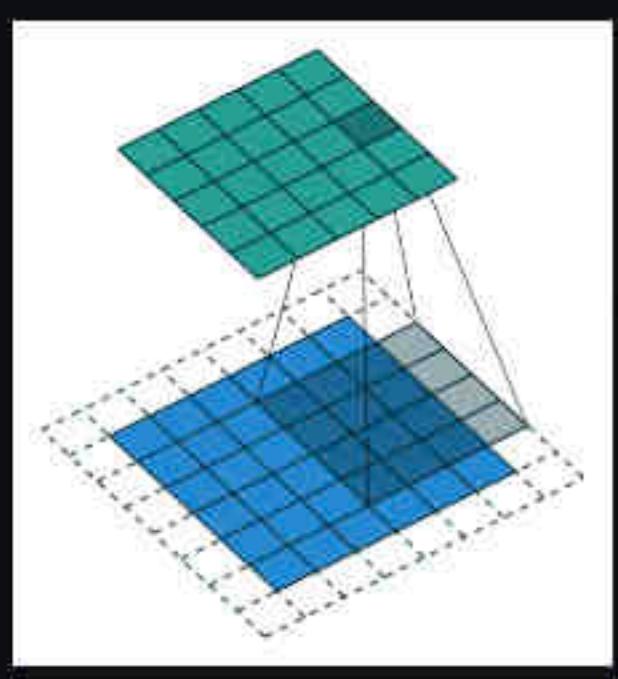
README.md

Transposed convolution animations

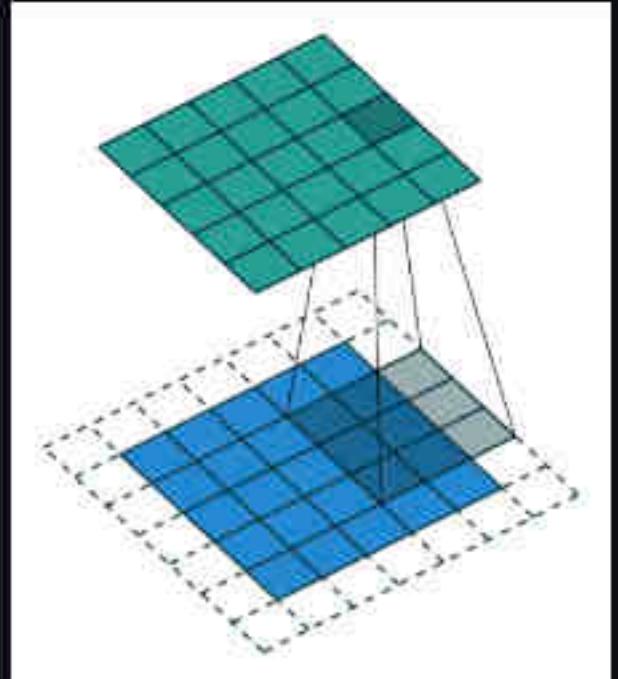
N.B.: Blue maps are inputs, and cyan maps are outputs.



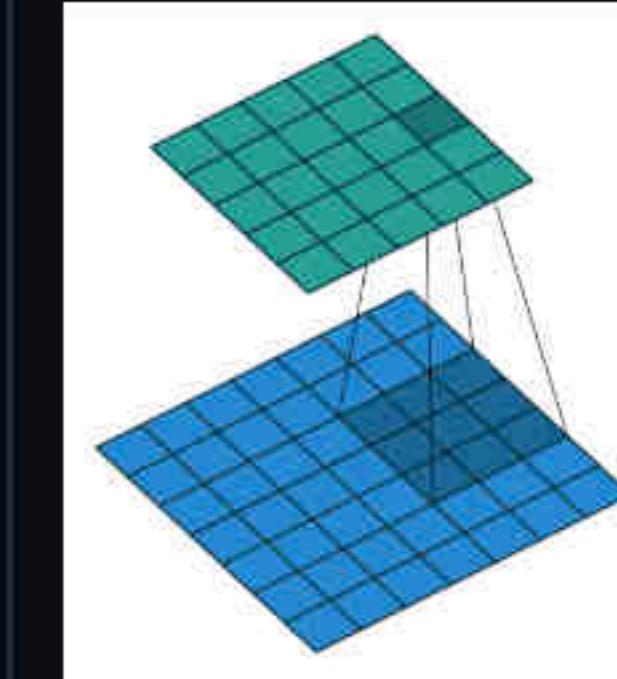
No padding, no strides,
transposed



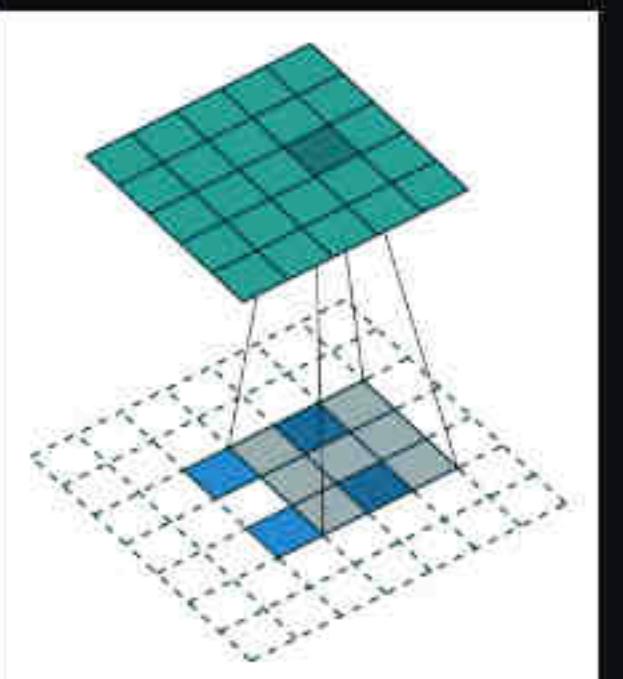
Arbitrary padding, no
strides, transposed



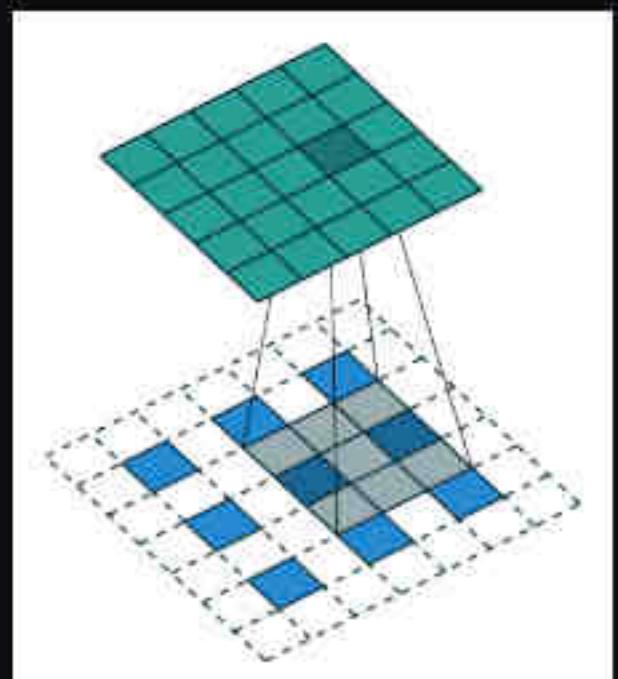
Half padding, no
strides, transposed



Full padding, no strides,
transposed



No padding, strides,
transposed



Padding, strides,
transposed



How we will evaluate Image Segmentation Models?

- We classify pixels of an image as 1 or 0. If there is a mask in a pixel we state 1, if there is not a mask (background) we state 0.
 - Making pixelwise binary classification of image
 - we can formulate it a two-class semantic segmentation problem, with binary cross entropy loss as the loss function.
 - . However, it works badly because of two reasons:
 - highly unbalanced label distribution and
 - per-pixel intrinsic issue of cross entropy loss.
 - When using cross entropy loss, the statistical distributions of labels play a big role in training accuracy.
 - The more unbalanced the label distributions are, the more difficult the training will be.
 - the loss is calculated as the average of per-pixel loss, without knowing whether its adjacent pixels are boundaries or not.
 - As a result, cross entropy loss only considers loss in a micro sense rather than considering it globally, which is not enough for image level prediction.
 - So we need a new loss that can handle class imbalance and consider a macro sense of prediction distribution
 - IOU

L8: ObjectDet... X | ONNX | Home X | Object Detectio... X | Network-archit... X | Post_ Read: Re... X | L9: Object Seg... X | GitHub - vdum... X | ROI pooling vs... X | An Introduction X | LiLog X | YOLOv5 is Here... X | +

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect | Update

scribble

How we will evaluate Image Segmentation Models?

per pixel

metric

~ loss

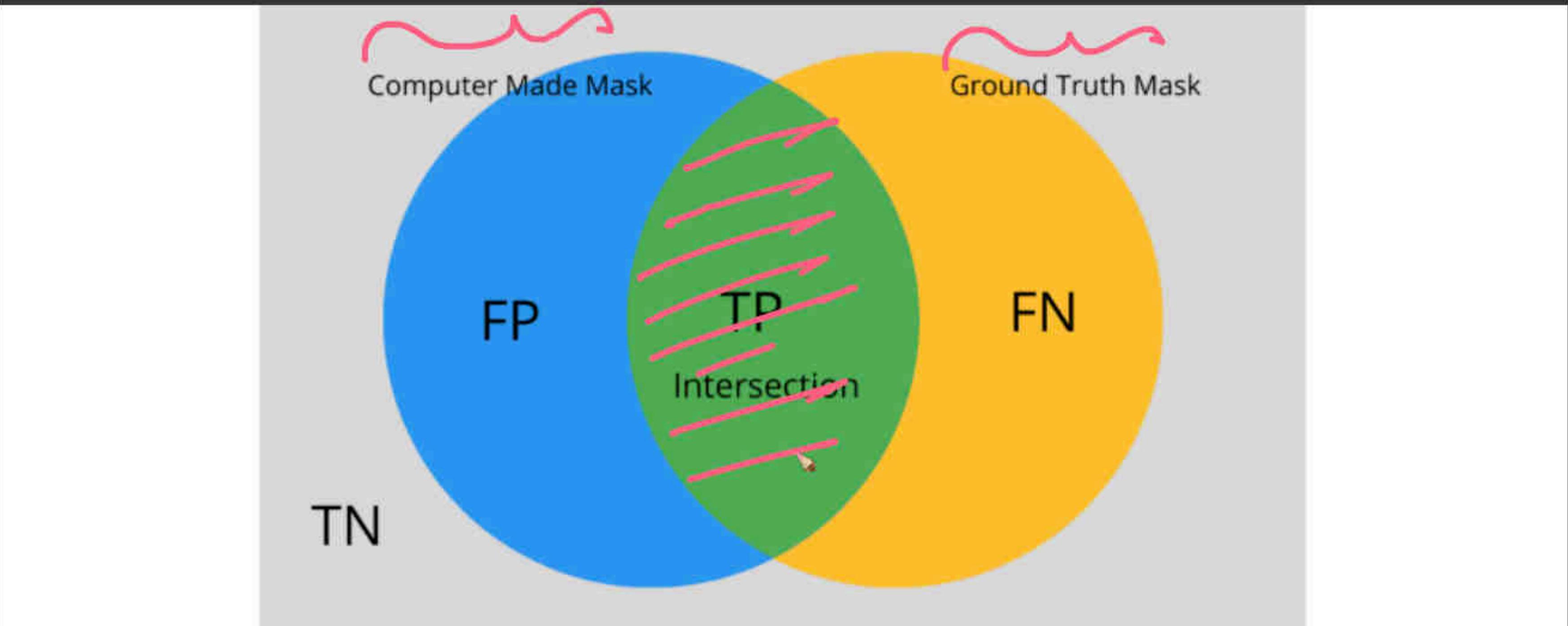
- We classify pixels of an image as 1 or 0. If there is a mask in a pixel we state 1, if there is not a mask (background) we state 0.
- Making pixelwise binary classification of image
- we can formulate it a two-class semantic segmentation problem, with binary cross entropy loss as the loss function.
- . However, it works badly because of two reasons:
 - highly unbalanced label distribution and
 - per-pixel intrinsic issue of cross entropy loss.
- When using cross entropy loss, the statistical distributions of labels play a big role in training accuracy.
- The more unbalanced the label distributions are, the more difficult the training will be.
- the loss is calculated as the average of per-pixel loss, without knowing whether its adjacent pixels are boundaries or not.
- As a result, cross entropy loss only considers loss in a micro sense rather than considering it globally, which is not enough for image level prediction.
- So we need a new loss that can handle class imbalance and consider a macro sense of prediction distribution

71 / 71

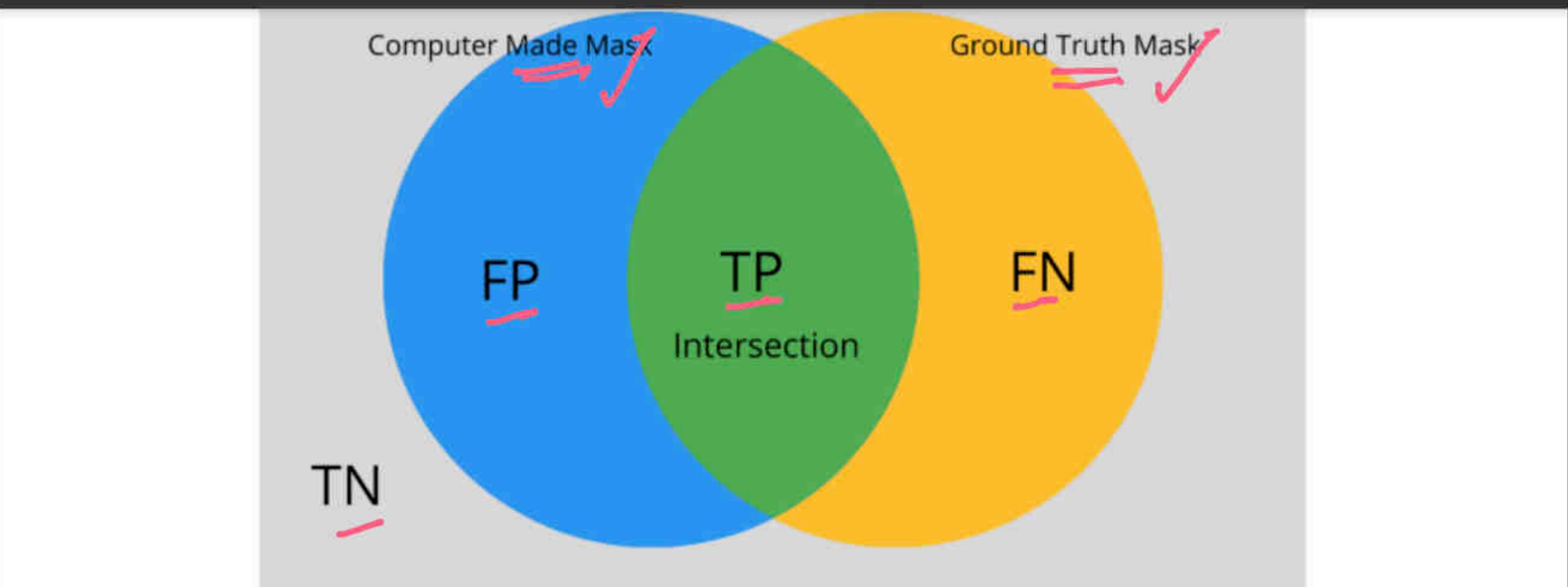
+ Code + Text

Connect | + | Update

- **IOU**
- We have discussed IOU in object detection



$$\frac{\text{Intersection}}{\text{Union}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$$



$$\frac{\text{Intersection}}{\text{Union}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$$

= loll ✓

- IOU can be used in image segmentation, but we use a slightly different metric **Dice**
- Dice coefficient is very similar to IOU. Dice co

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

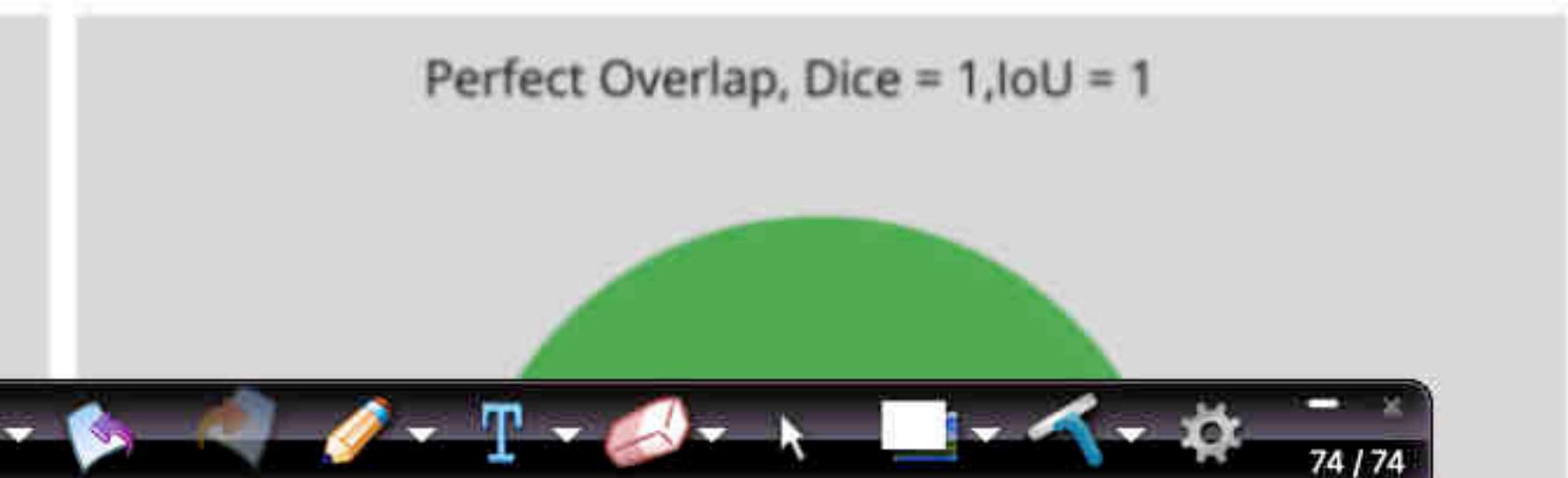
$$\frac{\text{Intersection}}{\text{Union}} = \frac{2\text{TP}}{2\text{TP} + \text{FN} + \text{FP}}$$

• IOU can be used in image segmentation, but we use a slightly different metric **Dice**
• Dice coefficient is very similar to IOU. Dice coefficient double counts the intersection(TP).

$$\text{Dice Coefficient} = \frac{2 \times \text{Intersection}}{\text{Union} + \text{Intersection}} = \frac{2\text{TP}}{2\text{TP} + \text{FN} + \text{FP}}$$



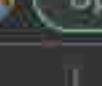
No Overlap, Dice = 0, IoU = 0



Perfect Overlap, Dice = 1, IoU = 1

74 / 74

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

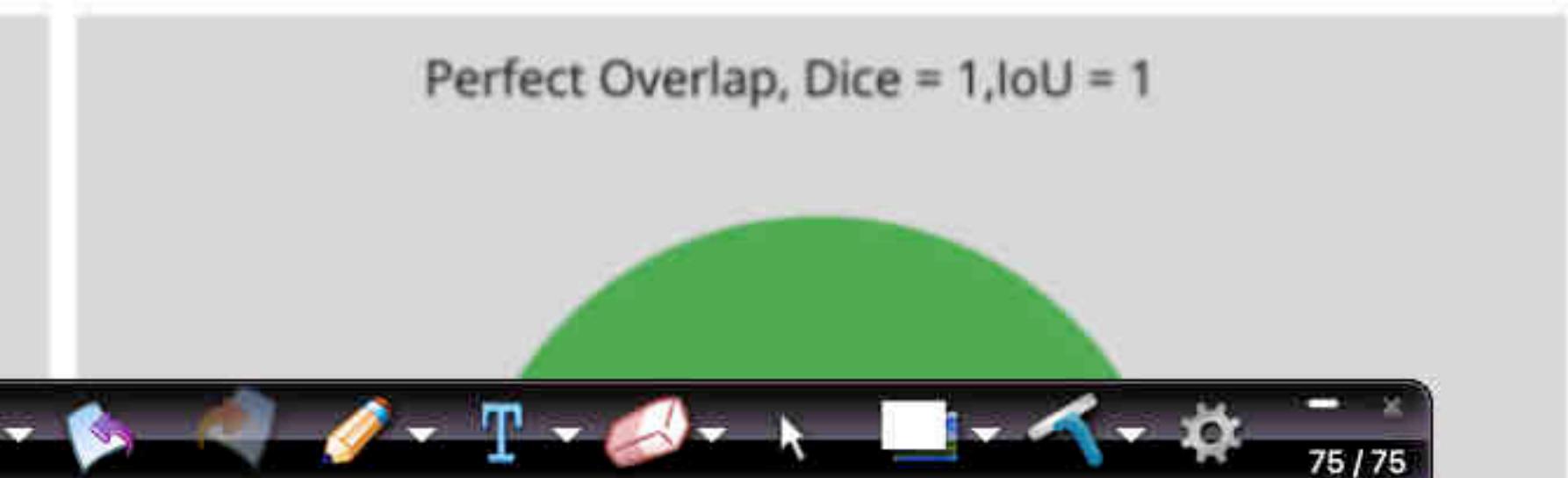
$$\frac{\text{Intersection}}{\text{Union}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$$

• IOU can be used in image segmentation, but we use a slightly different metric **Dice**
• Dice coefficient is very similar to IOU. Dice coefficient double counts the intersection(TP).

$$\text{Dice Coefficient} = \frac{2 \times \text{Intersection}}{\text{Union} + \text{Intersection}} = \frac{2\text{TP}}{2\text{TP} + |\text{FN}| + |\text{FP}|}$$



No Overlap, Dice = 0, IoU = 0



Perfect Overlap, Dice = 1, IoU = 1

75 / 75

L8: ObjectDet... X | ONNX | Home X | Object Detectio... X | Network-archit... X | Post_ Read: Re... X | L9: Object Seg... X | GitHub - vdum... X | ROI pooling vs... X | An Introduction X | LiLog X | YOLOv5 is Here... X | +

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect | | ▾

FCN Coding Tutorial

- With This, we will dive deeper into the Convolutional Encoder-Decoder models by implementing a Fully Convolutional Neural Networks from scratch

The diagram illustrates the Fully Connected Network (FCN) architecture. It shows an input image on the left being processed by a series of layers (96, 256, 384, 384, 256, 4096, 4096, 21) to produce a segmentation map on the right. A double-headed arrow labeled "forward/inference" points from the input to the output. Another double-headed arrow labeled "backward/learning" points from the output back to the input, indicating the flow of gradients during training.

How does the model understand the key features of the image ?

- The FCN model uses convolutional layers as feature extractor
- When used in the Encoder part of the FCN model by downsampling the image

Now, the last layers contains all the key features we pass it to a decoder part of the model instead of Fully Connected Layers

- The decoder consists of Deconvolutional layers that upsample these key features to the original size of the image

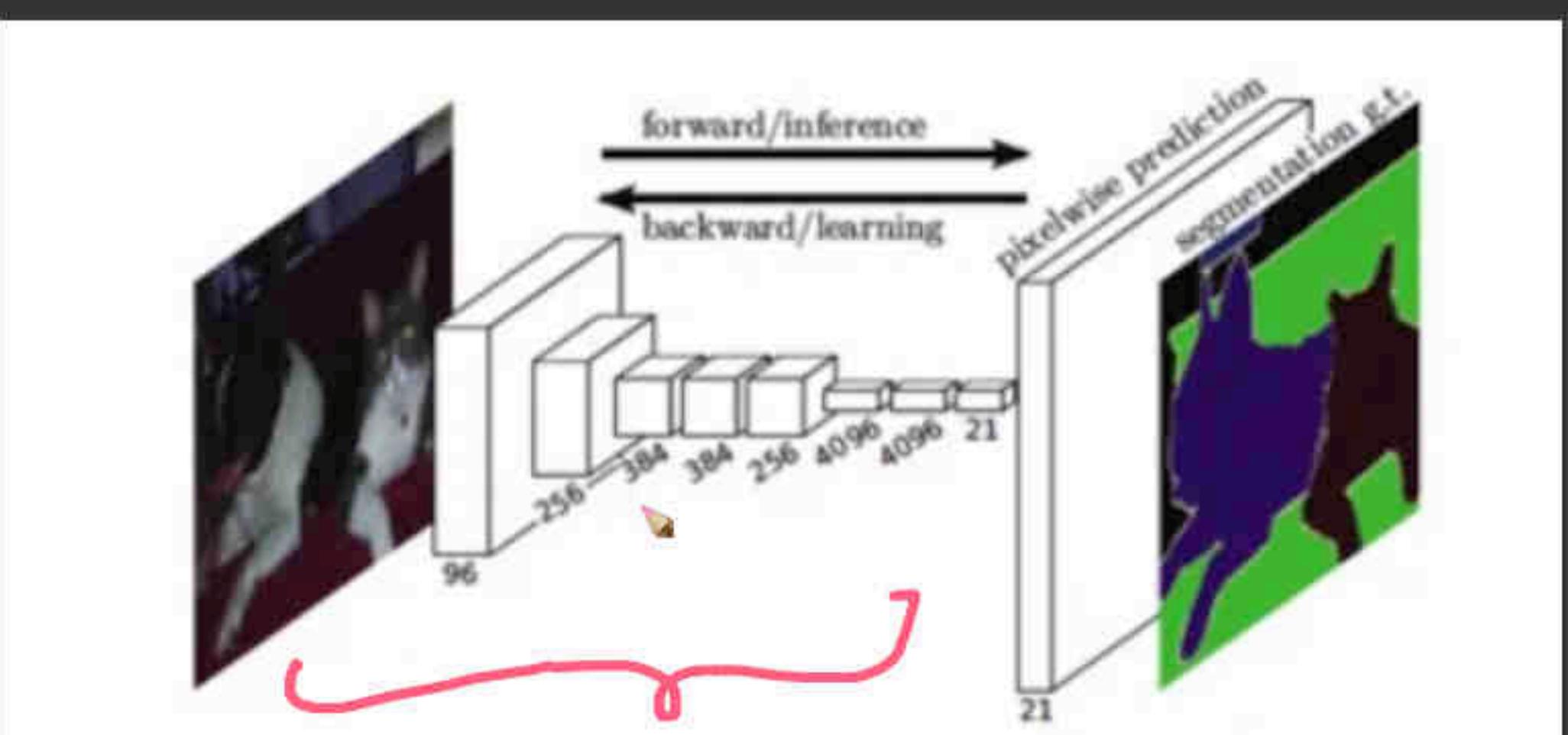
76 / 76

+ Code + Text

Connect ▾

FCN Coding Tutorial

- With This, we will dive deeper into the Convolutional Encoder-Decoder models by implementing a Fully Convolutional Neural Networks from scratch



How does the model understand the key features of the image ?

- The FCN model uses convolutional layers as feature extractor
 - When used in the Encoder part of the FCN model by downsampling the image

Now, the last layers contains all the key features we pass it to a decoder part of the model instead of Fully Connected Layers

- The decoder consists of Deconvolutional layers

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=-YT7bRAR05EJ

+ Code + Text Connect |  

```
!rm -rf ./logs/
```

[] def FCN8():
 ✓ { VGG16 = tf.keras.applications.VGG16(weights="imagenet", include_top=False, input_tensor=tf.keras.Input(shape=(128, 128, 3)))
 VGG16.trainable = False
 x = VGG16(VGG16.input)

 #Decoder
 o = tf.keras.layers.Conv2D(4096 , (7 , 7) , activation='relu' , padding='same' , name="conv6")(x)
 conv7 = tf.keras.layers.Conv2D(4096 , (1 , 1) , activation='relu' , padding='same' , name="conv7")(o)
 print(o.shape,conv7.shape)
 ## 4 times upsampling for pool4 layer
 conv7_4 = tf.keras.layers.Conv2DTranspose(2 , kernel_size=(4,4) , strides=(4,4) , use_bias=False)(conv7)
 ## 2 times upsampling for pool411

 pool4 = VGG16.get_layer('block4_pool').output
 pool411 = tf.keras.layers.Conv2D(2 , (1 , 1) , activation='relu' , padding='same' , name="pool4_11")(pool4)
 pool411_2 = tf.keras.layers.Conv2DTranspose(2 , kernel_size=(2,2) , strides=(2,2) , use_bias=False)(pool411)

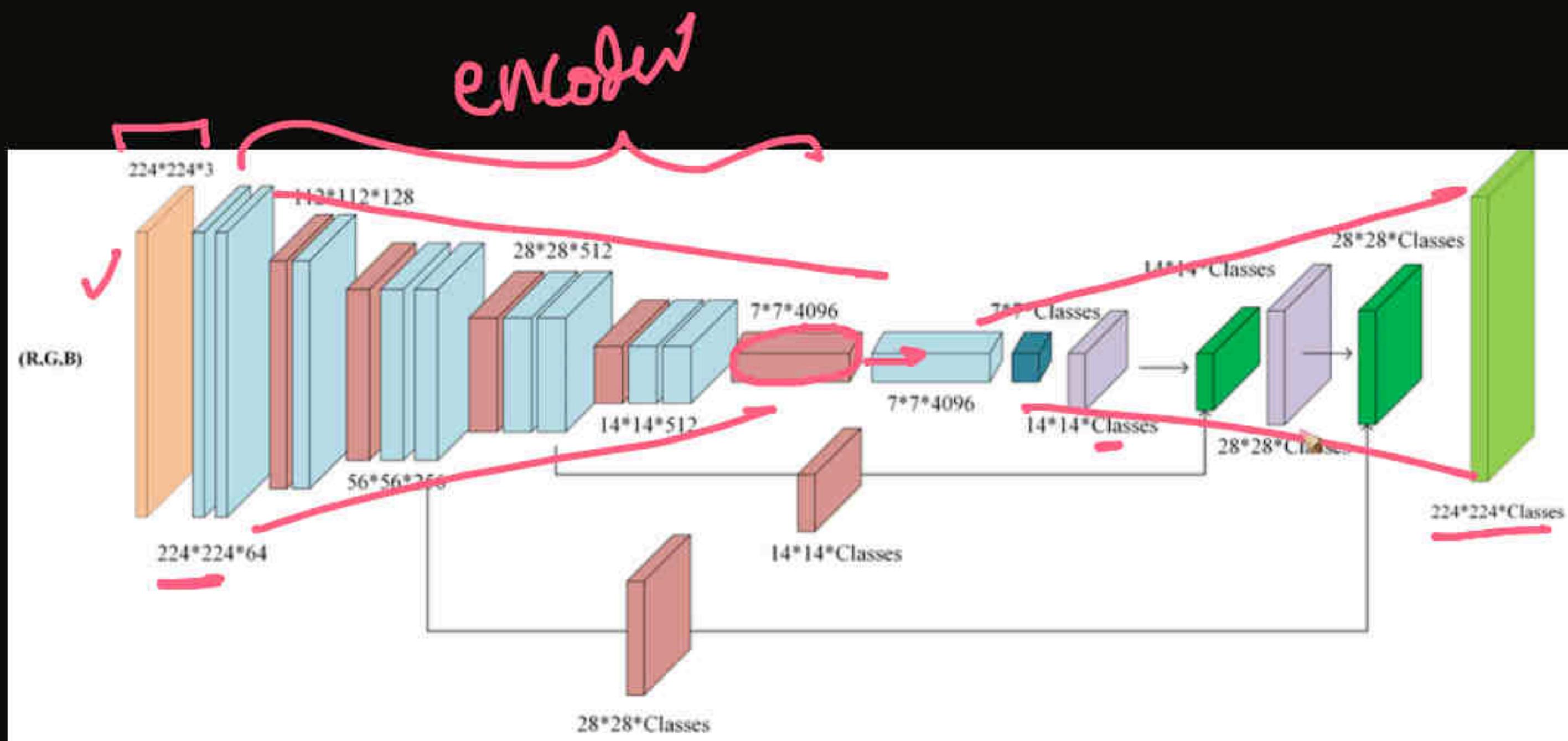
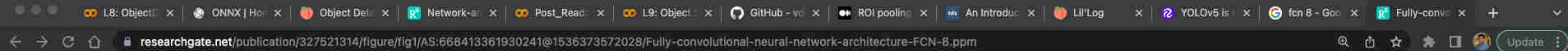
 pool3 = VGG16.get_layer('block3_pool').output
 pool311 = tf.keras.layers.Conv2D(2 , (1 , 1) , activation='relu' , padding='same' , name="pool3_11")(pool3)
 o = tf.keras.layers.Add(name="add")([pool411_2, pool311])

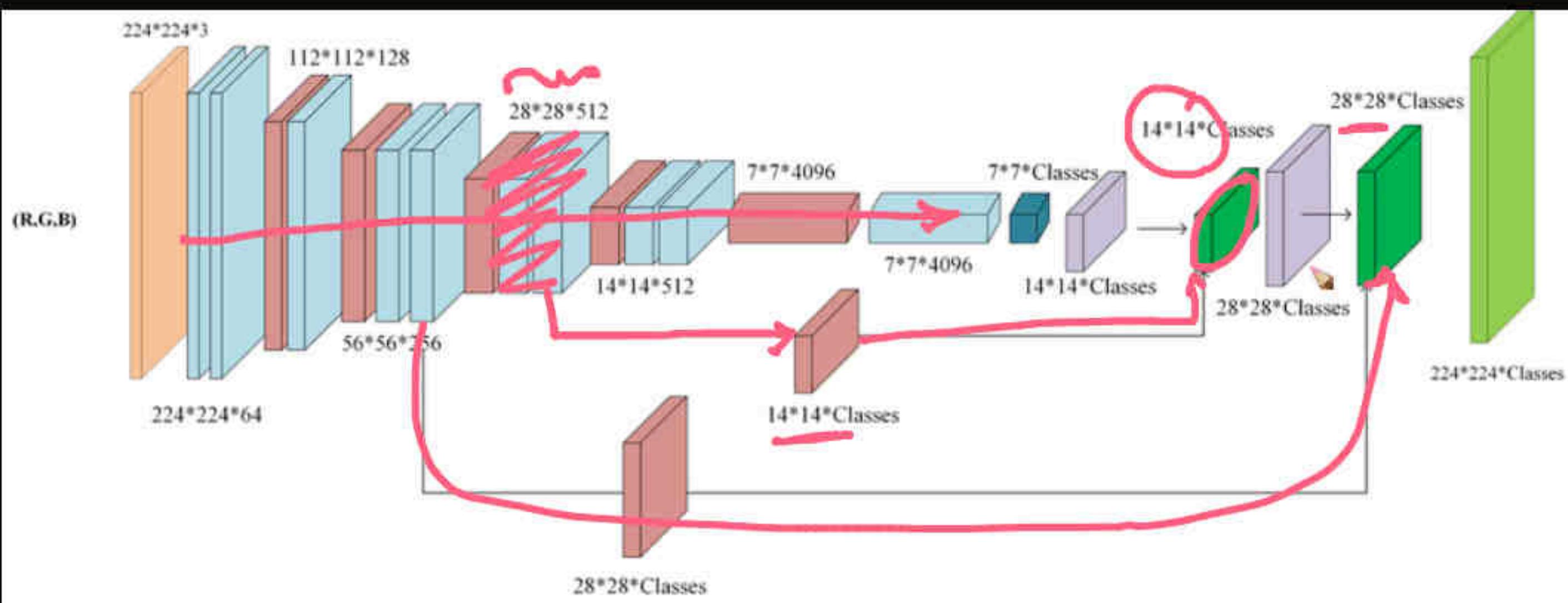
 print(o.shape,conv7_4.shape)
 o = tf.keras.layers.Add(name="add1")([o, conv7_4])
 o = tf.keras.layers.Conv2DTranspose(2 , kernel_size=(8,8) , strides=(8,8) , use_bias=False)(o)
 o = tf.keras.layers.Activation('softmax')(o)

 model = tf.keras.Model(VGG16.inpu

Input  VGG16

78 / 78





+ Code + Text

Connect | User Settings

```
VGG16.trainable = False
x = VGG16(VGG16.input)

#Decoder
o = tf.keras.layers.Conv2D( 4096 , ( 7 , 7 ) , activation='relu' , padding='same' , name="conv6")(x)
conv7 = tf.keras.layers.Conv2D( 4096 , ( 1 , 1 ) , activation='relu' , padding='same' , name="conv7")(o)
print(o.shape,conv7.shape)
## 4 times upsampling for pool4 layer
conv7_4 = tf.keras.layers.Conv2DTranspose( 2 , kernel_size=(4,4) , strides=(4,4) , use_bias=False)(conv7)

## 2 times upsampling for pool411
✓ pool4 = VGG16.get_layer('block4_pool').output
pool411 = tf.keras.layers.Conv2D( 2 , ( 1 , 1 ) , activation='relu' , padding='same' , name="pool4_11")(pool4)
pool411_2 = tf.keras.layers.Conv2DTranspose(2 , kernel_size=(2,2) , strides=(2,2) , use_bias=False)(pool411)

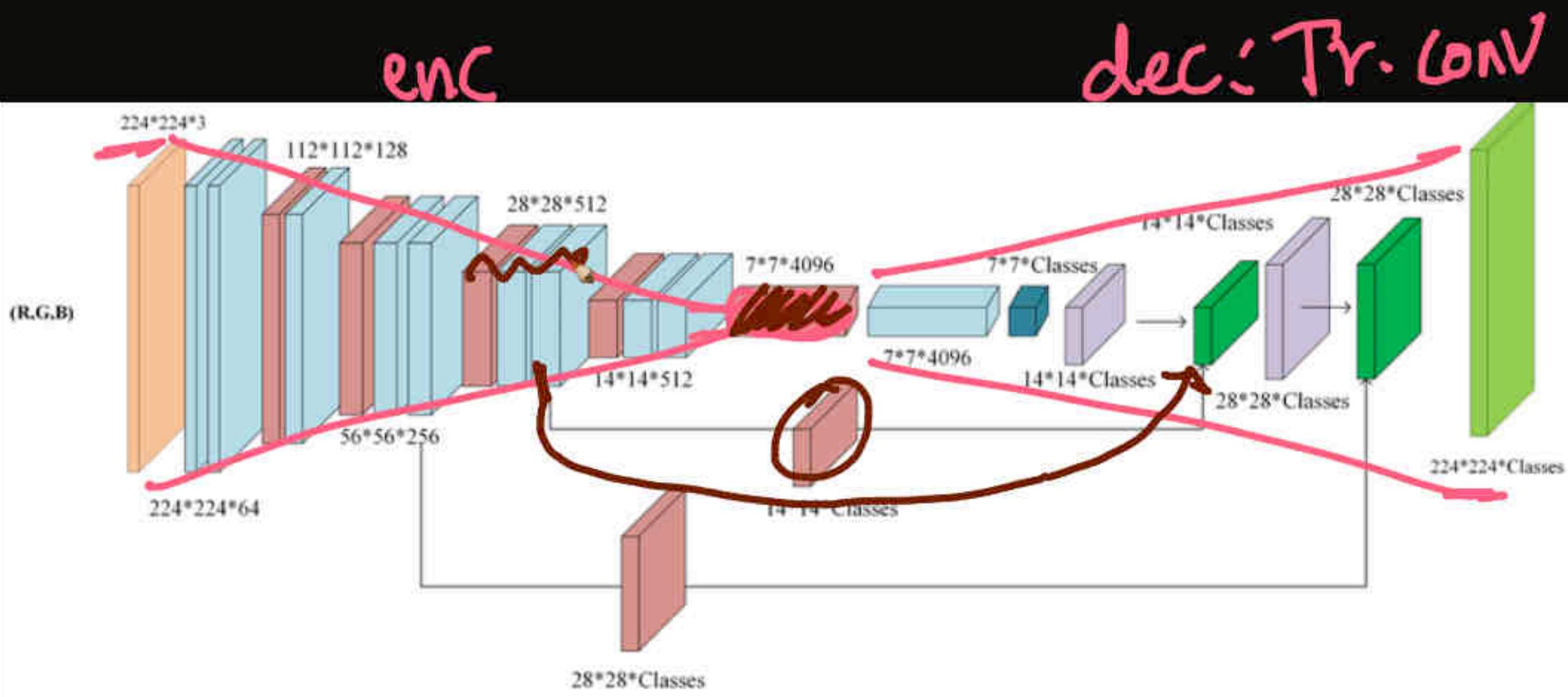
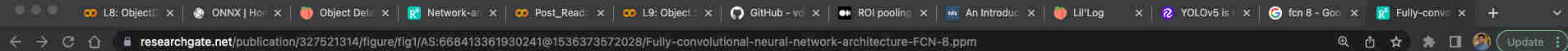
pool3 = VGG16.get_layer('block3_pool').output
pool311 = tf.keras.layers.Conv2D(2 , ( 1 , 1 ) , activation='relu' , padding='same' , name="pool3_11")(pool3)
o = tf.keras.layers.Add(name="add")([pool411_2, pool311])

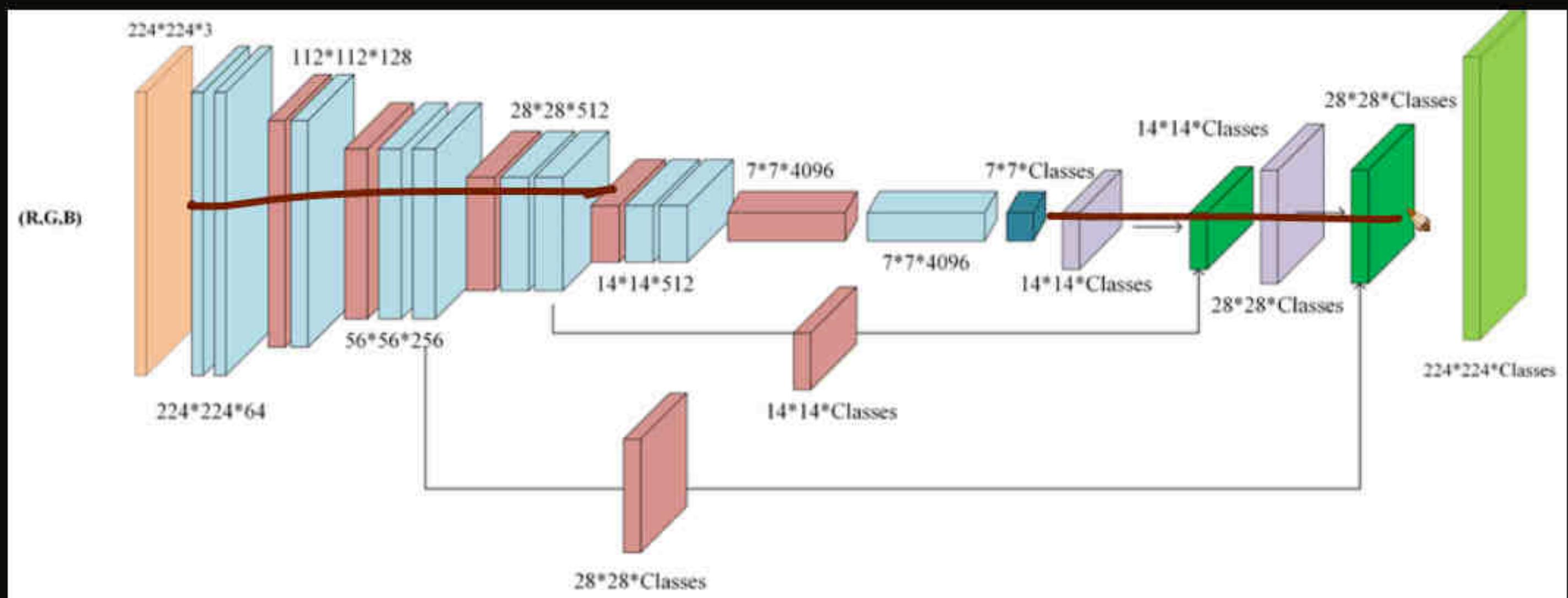
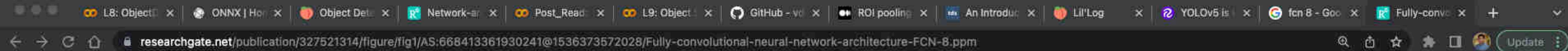
print(o.shape,conv7_4.shape)
o = tf.keras.layers.Add(name="add1")([o, conv7_4 ])
o = tf.keras.layers.Conv2DTranspose(2 , kernel_size=(8,8) , strides=(8,8) , use_bias=False)(o)
o = tf.keras.layers.Activation('softmax')(o)

model = tf.keras.Model(VGG16.input, o)

return model
```





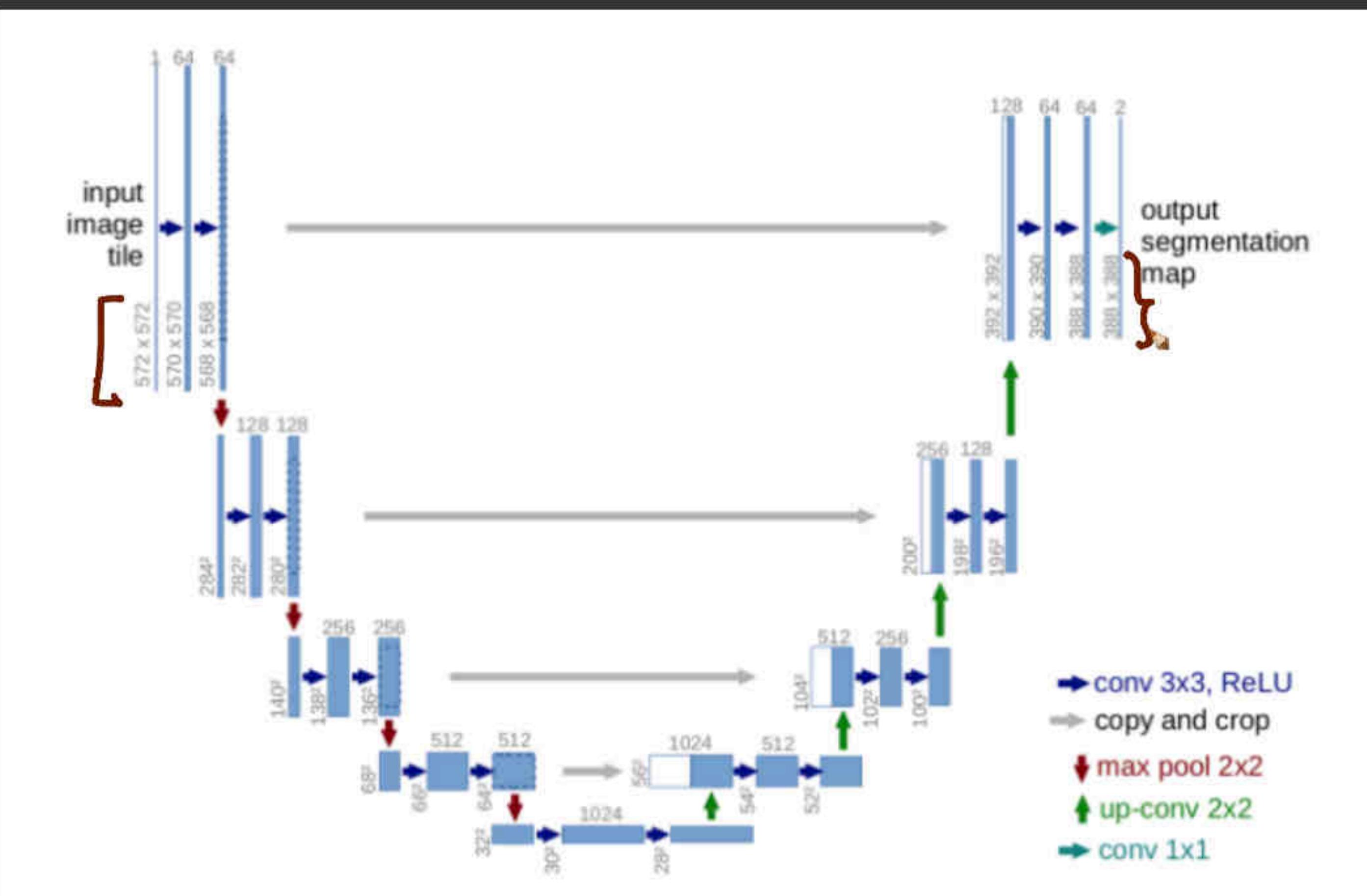


+ Code + Text

Connect



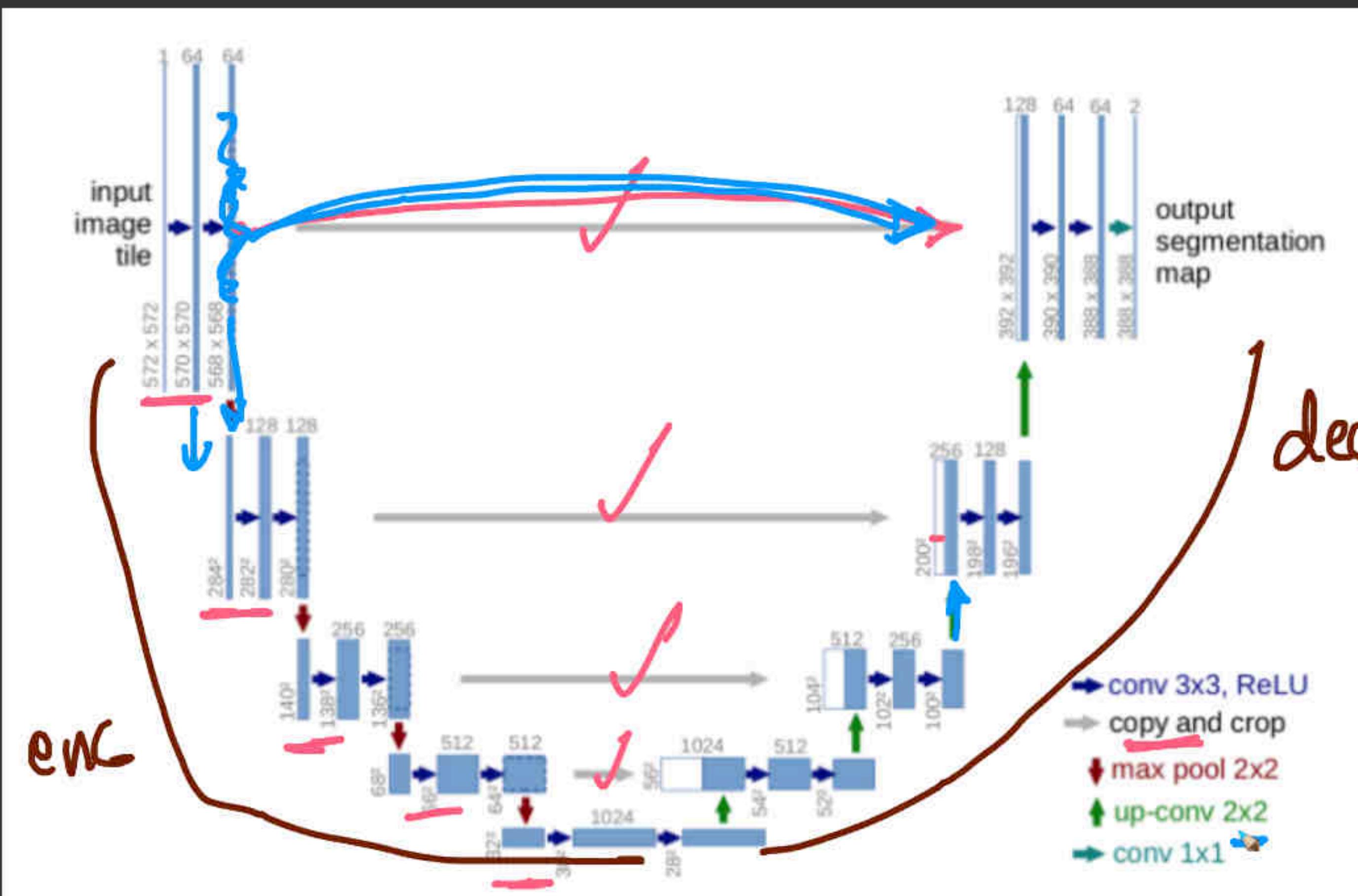
- Original paper: <https://arxiv.org/pdf/1505.04597.pdf>



+ Code + Text

FCN model.

- Original paper: <https://arxiv.org/pdf/1505.04597.pdf>



U-net

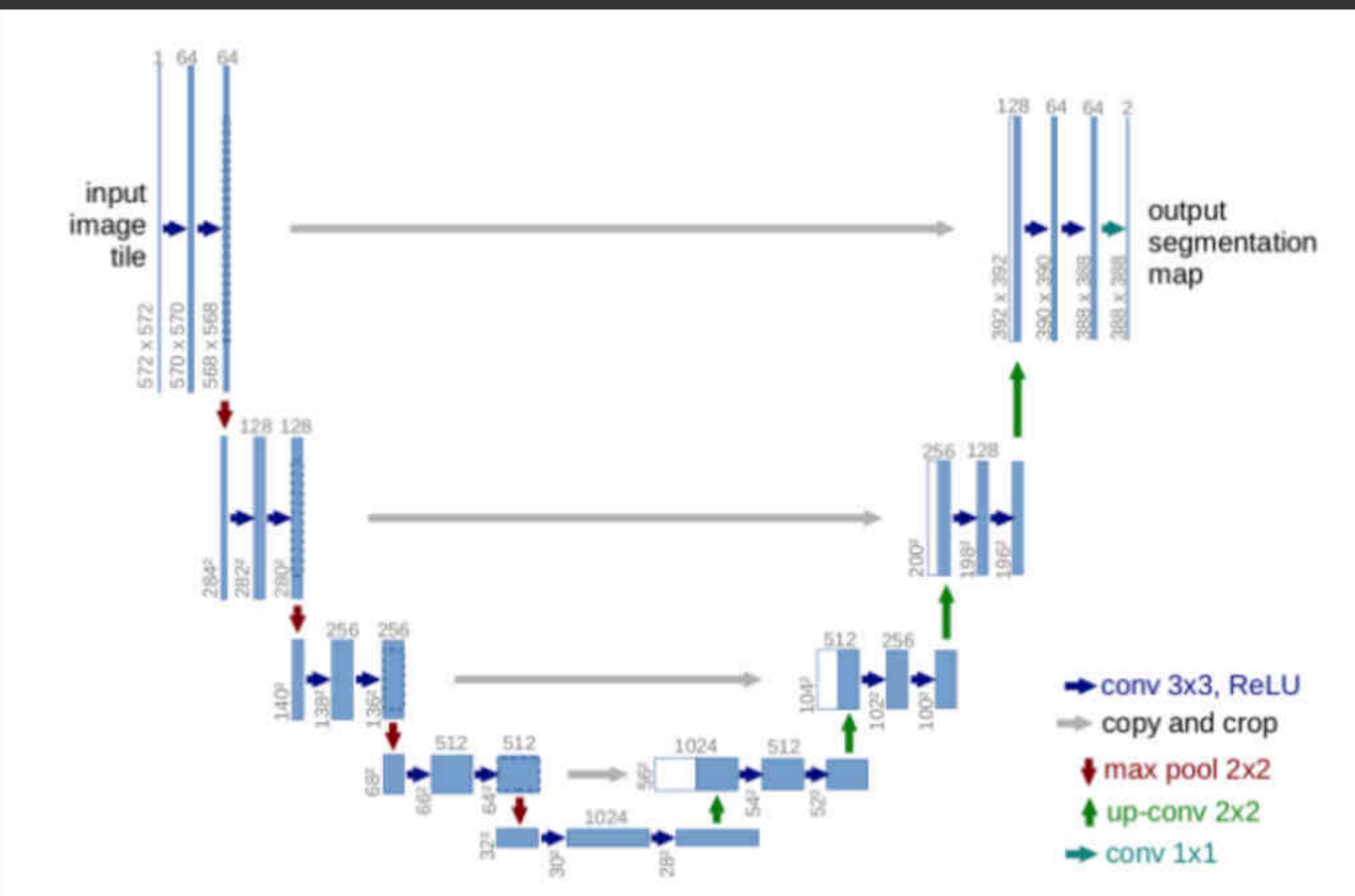
- ① ~~UpSampling~~ : cheaper
- ② many ~~skip~~-connections

+ Code + Text

Connect ▾

FCN model.

- Original paper: <https://arxiv.org/pdf/1505.04597.pdf>



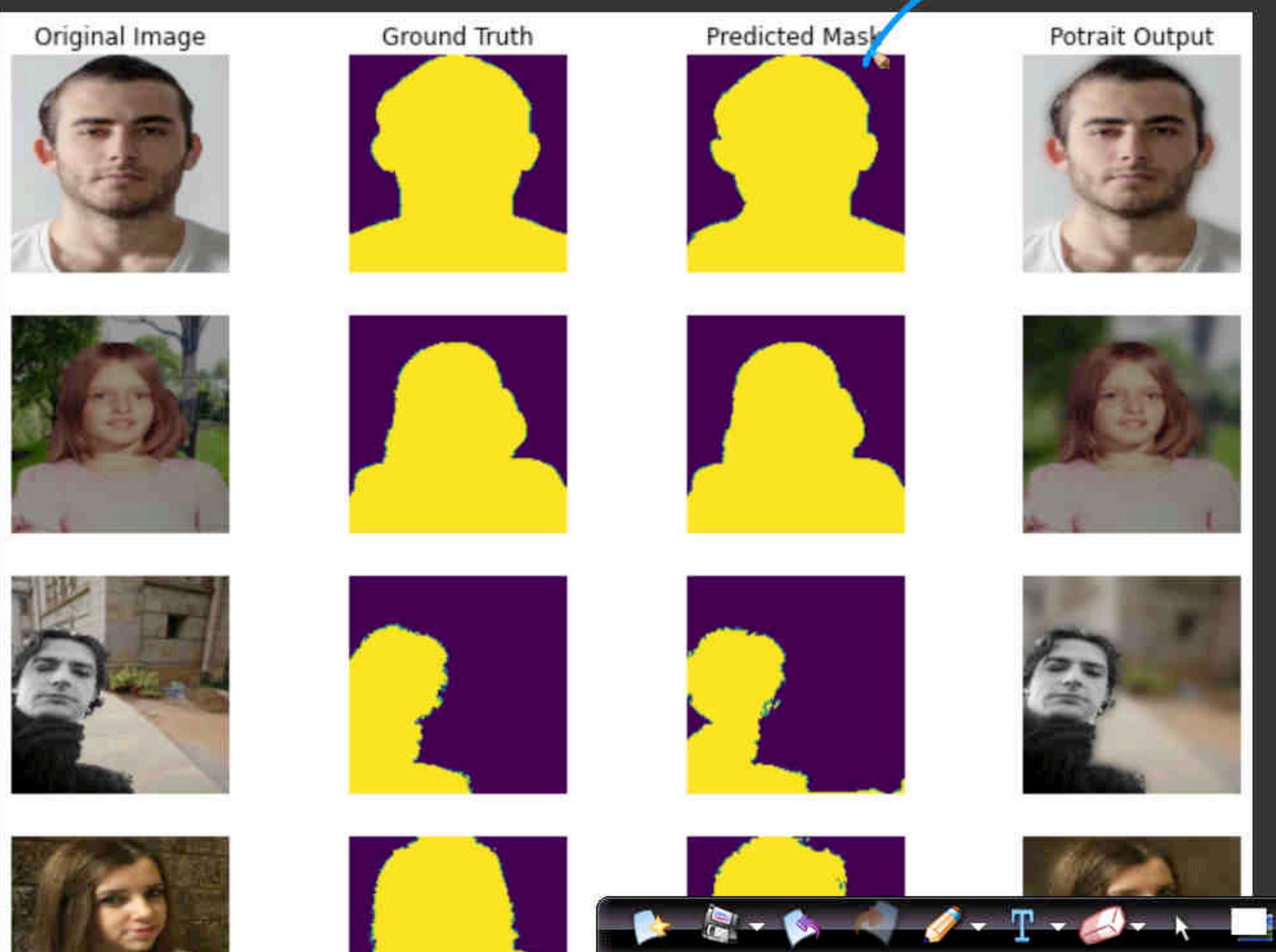
Initiatively

Conv = Trconv
Maxpooling = UpSampling

+ Code + Text

Connect | User Settings

```
    ax[idx1][2].axis('off')
    ax[idx1][3].axis('off')
plt.show()
```



colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=6kE9uTGs5mYz

+ Code + Text

```
ax[idx1][2].axis('off')
ax[idx1][3].axis('off')
plt.show()
```

Original Image Ground Truth Predicted Mask Potrait Output

colab.research.google.com/drive/1J0CD_eybAUmKmr0aR-ISBaWeQk9mtSIP#scrollTo=6kE9uTGs5mYz

+ Code + Text Connect |

iii {x} □

Yolo (recap)

YOLO

YOLOV3

RetinaNet
(FPN)

Seg: ...

FCN8

U-net

Mask RCNN

DeepLabV3

dilated conv

89 / 89

The image shows a grid of 20 images arranged in five rows. The first four rows each contain four images, while the fifth row contains a single image. Each image consists of a headshot on the left and its corresponding yellow segmentation mask on the right. The segmentation masks are used to identify the human subjects in the headshots.

[+ Code](#) [+ Text](#)

Connect



- Dataset : https://drive.google.com/drive/folders/1NKEogbBaPsKUw_6qEkcB_nCB4P84Alwl?usp=sharing

{x} ▾ Outline

- Motivation
- Business Case
- EDA
- image segmentation over classification / detection
- Semantic vs Instance Segmentation
- The idea of Encoder-Decoder network
- Dice loss, Dice coeff
- FCN coding tutorial
- Unet coding tutorial
- Intuition for Mask R-CNN
- Intuition for DeeplabV3

◀ Motivation

- We started with image classification, which is determining if an image contains an object or not
- Then in Object Detection we learnt about object localization : where exactly the object is present in the image and we drew a rectangular

[+ Code](#)[+ Text](#)

90 / 90