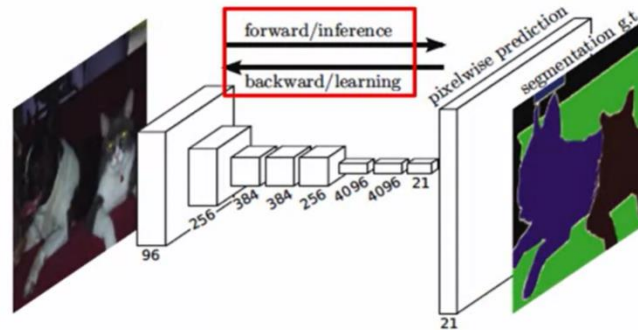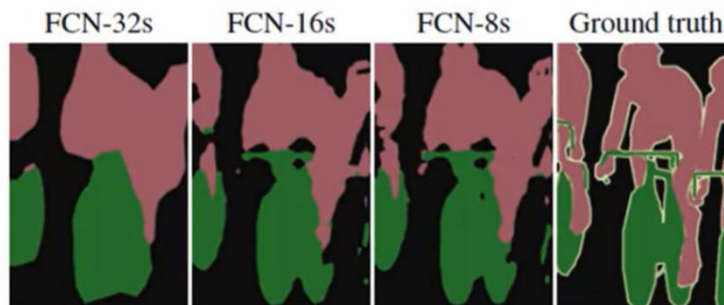# Fully Convolutional Neural Networks



https://people.eecs.berkeley.edu/~jonlong/long_shelhamer_fcn.pdf

# Comparison of Different FCNs



FCN-32s  FCN-16s  FCN-8s  Ground truth

# SegNet

# UNet



| | |
|---|---|
| → | Conv (3 x3), RELU |
| → | Max Pool (2 x 2) |
| → | De conv (2 x 2) |
| ⇨ | Crop and concatenate |

# Mask R-CNN



Faster R-CNN w/ ResNet [19]
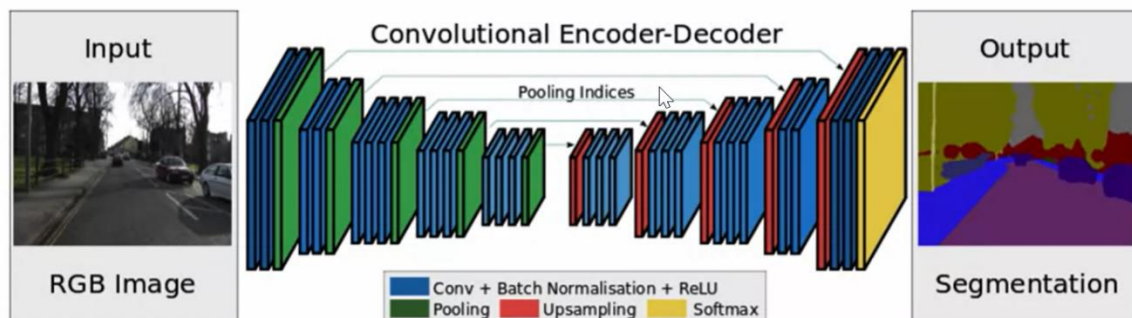
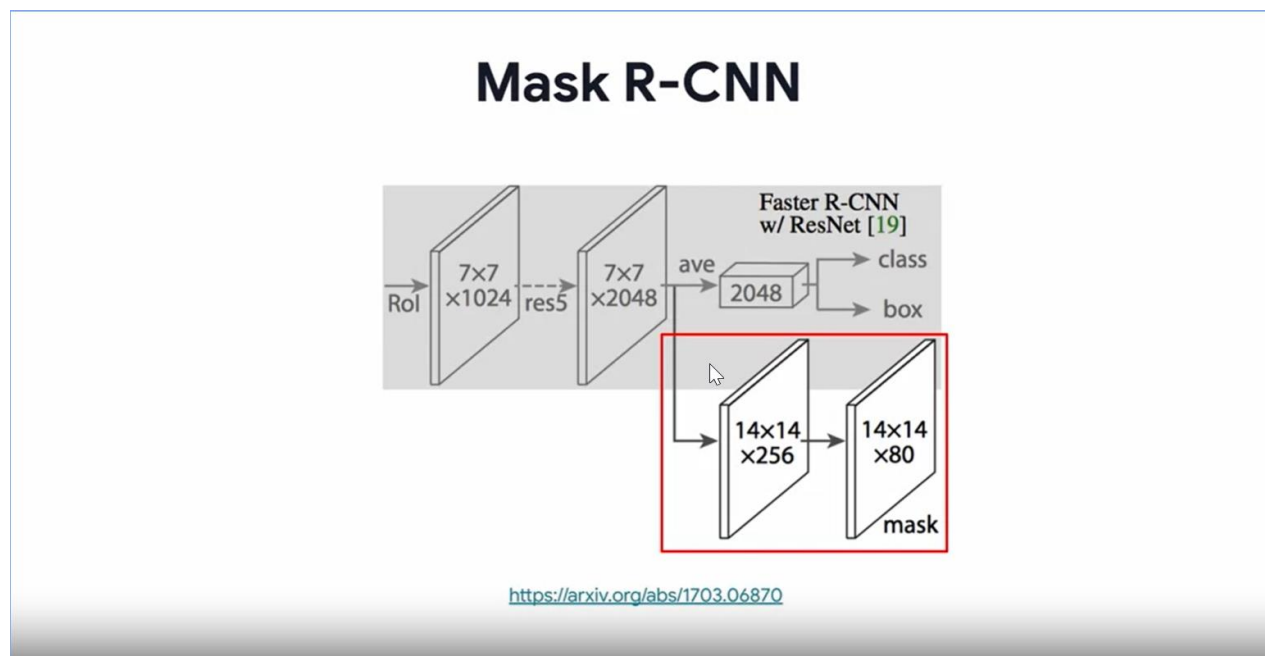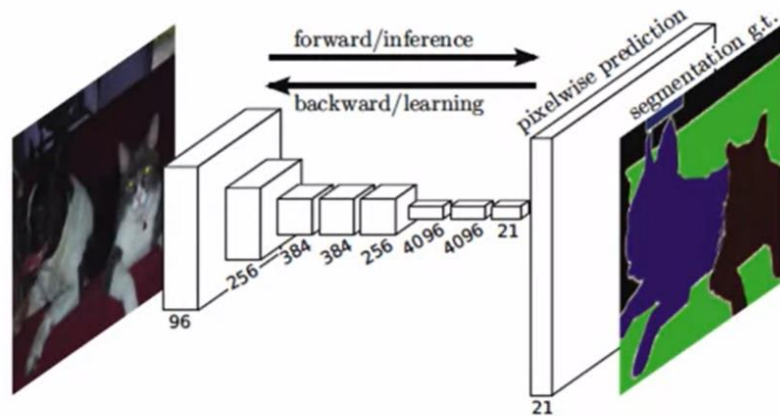https://arxiv.org/abs/1703.06870

# FCN Architecture Details



## Fully Convolutional Neural Networks

forward/inference

backward/learning

pixelwise prediction

segmentation g.t.

96 256 384 384 256 4096 4096 21

21

Fully Convolutional Networks for Semantic Segmentation
By Jonathan Long, Evan Shelhamer, Trevor Darrell
https://arxiv.org/pdf/1411.4038.pdf

# Encoders

- Popular encoder architectures:
    - VGG-16
    - ResNet-50
    - MobileNet
- Reuse convolutional layers for feature extraction.
    - Do not reuse fully connected layers

# Decoders

# FCN-32



Image | Pool 1 | Pool 2 | Pool 3 | Pool 4 | Pool 5 | 32x Upsampled Prediction

Encoder (Down Sampling) | Decoder (Up Sampling)

# FCN-16



Image | Pool 1 | Pool 2 | Pool 3 | Pool 4 | Pool 5

2x Upsampled Prediction

1x1 conv layer

Pool 4 Prediction

Σ

16x Upsampled Prediction

Encoder (Down Sampling) | Decoder (Up Sampling)

# FCN-8



## Upsampling

- Upsampling is increasing height and width of the feature map.
- Two types of layers used in TensorFlow:
    - Simple Scaling - UpSampling2D
    - Transposed Convolution(Deconvolution) - Conv2DTranspose
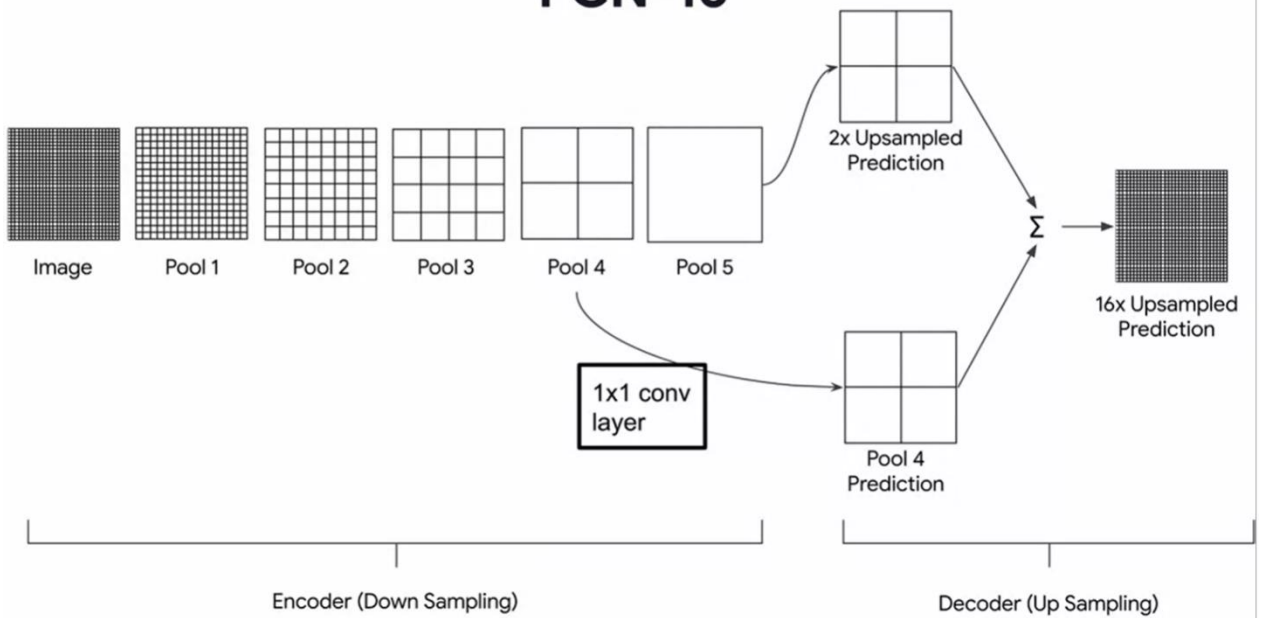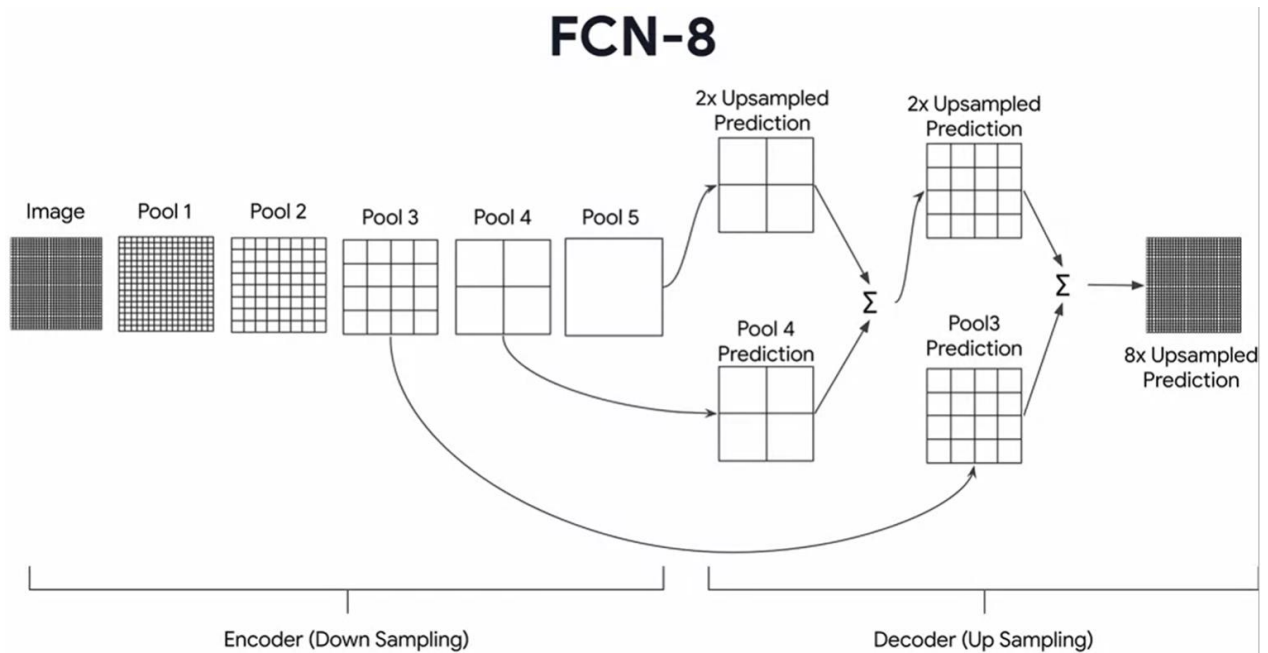
# Simple Scaling - UpSampling2D

- Upsampling2D scales up the image
- Two Types of scaling:
  - Nearest
    - Copies value from nearest pixel.
  - Bilinear
    - linear interpolation from nearby pixels.

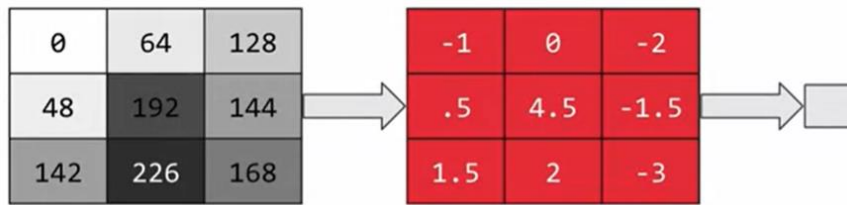# UpSampling2D - Usage

```python
x = UpSampling2D(
    size=(2, 2),
    data_format=None,
    interpolation='nearest')(x)

size: int or tuple of two ints
data_format: 'channels_first', 'channels_last' or None
interpolation: 'nearest' or 'bilinear'
```

# Transposed Convolution

| 0 | 64 | 128 |
|---|---|---|
| 48 | 192 | 144 |
| 142 | 226 | 168 |

| -1 | 0 | -2 |
|---|---|---|
| .5 | 4.5 | -1.5 |
| 1.5 | 2 | -3 |

Convolution->

| 32 | 64 | 64 |
|---|---|---|
| 48 | 192 | 128 |
| 128 | 192 | 168 |

| -1 | 0 | -2 |
|---|---|---|
| .5 | 4.5 | -1.5 |
| 1.5 | 2 | -3 |

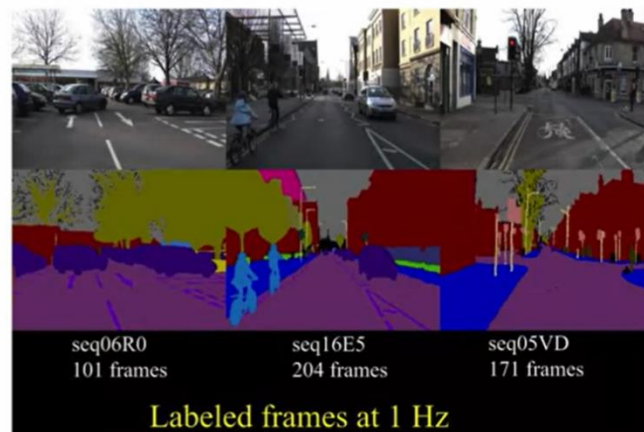<- Transposed Convolution

# Conv2DTranspose

- Reverse of Convolution.
- Applied to output of a convolution operation.
- Uses a kernel of a specified size and stride in order to recreate the original input before the convolution operation.

```
Conv2DTranspose(
    filters=32,
    kernel_size=(3, 3)
)
```

# Encoder in Code

- The Cambridge Driving, labelled video database (aka CamVid) contains 10 minutes of 30fps video, segmented and labelled with 32 classes
- GitHub account divamgupta has taken a subsample of the CamVid dataset to create a smaller dataset.



seq06R0     seq16E5     seq05VD
101 frames    204 frames    171 frames

Labeled frames at 1 Hz

- Divam Gupta's GitHub account containing a subsample of the CamVid dataset to create a smaller dataset.



```python
def block(x, n_convs, filters, kernel_size, activation, pool_size, pool_stride, block_name):
    for i in range(n_convs):
        x = tf.keras.layers.Conv2D(filters=filters,
                                   kernel_size=kernel_size, activation=activation,
                                   padding='same',
                                   name="{}_conv{}".format(block_name, i + 1))(x)

    x = tf.keras.layers.MaxPooling2D(pool_size=pool_size, strides=pool_stride,
                                     name="{}_pool{}".format(block_name, i+1 ))(x)

    return x
```

```python
def VGG_16(image_input):
    x = block(image_input, n_convs=2, filters=64, kernel_size=(3,3),
              activation='relu',pool_size=(2,2), pool_stride=(2,2),
              block_name='block1')
    p1= x

    x = block(x, n_convs=2, filters=128, kernel_size=(3,3),
              activation='relu',pool_size=(2,2), pool_stride=(2,2),
              block_name='block2')
    p2 = x

    ...
```

| conv1_1 |
| conv1_2 |
| pooling2 |

| conv2_1 |
| conv2_2 |
| pooling2 |

```python
    x = block(x, n_convs=3, filters=256, kernel_size=(3,3), activation='relu',pool_size=(2,2),
              pool_stride=(2,2), block_name='block3')
    p3 = x
```

| conv3_1 |
| conv3_2 |
| conv3_3 |
| pooling3 |

```python
    x = block(x, n_convs=3, filters=512, kernel_size=(3,3), activation='relu',pool_size=(2,2),
              pool_stride=(2,2), block_name='block4')
    p4 = x
```

| conv4_1 |
| conv4_2 |
| conv4_3 |
| pooling4 |

```python
    x = block(x, n_convs=3, filters=512, kernel_size=(3,3), activation='relu',pool_size=(2,2),
              pool_stride=(2,2), block_name='block5')
    p5 = x
..
```
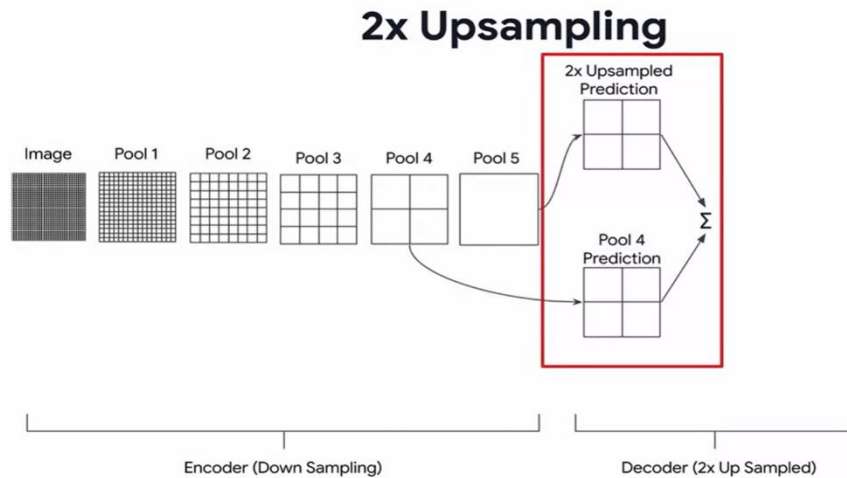
| conv5_1 |
| conv5_2 |
| conv5_3 |
| pooling5 |

# Decoder in Code (FCN8)

## 2x Upsampling



Encoder (Down Sampling)                  Decoder (2x Up Sampled)

## Define Decoder - 2x UpSampling

```python
def fcn8_decoder(convs, n_classes):
    f1, f2, f3, f4, f5 = convs

    o = tf.keras.layers.Conv2DTranspose(n_classes, kernel_size=(4,4),
                                        strides=(2,2), use_bias=False )(f5)

    o = tf.keras.layers.Cropping2D(cropping=(1,1))(o)

    o2 = f4
    o2 = ( tf.keras.layers.Conv2D(n_classes, (1,1),
                                  activation='relu', padding='same'))(o2)

    o = tf.keras.layers.Add()([o, o2])
    ...
```
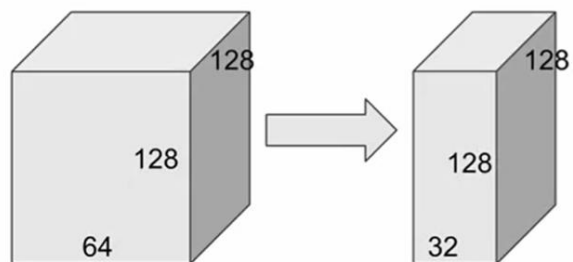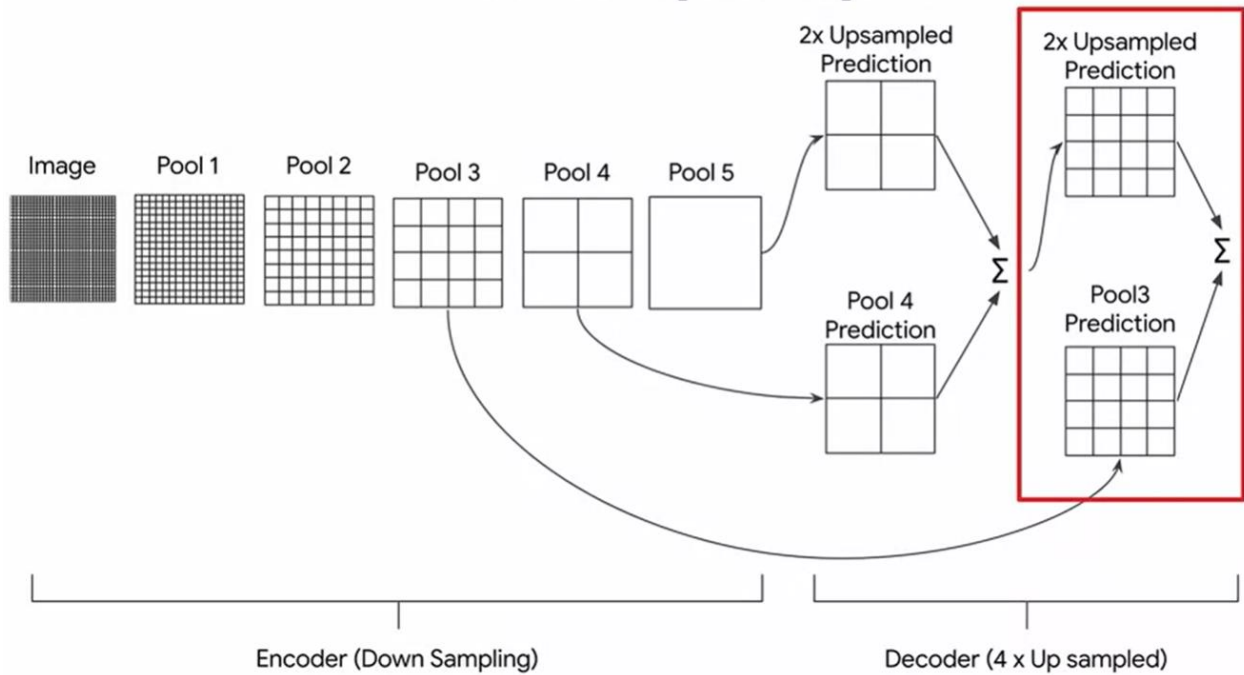
## 1 x 1 Convolutions

( B, F, H, W) - B = # batches; F= # filters, H, W = Height/Width

Apply a layer with N 1x1 Convolutions with stride of 1:

( B, N, H, W) - B = # batches; N= # filters, H, W = Height/Width

# 2 x 2 Upsampled



```python
def fcn8_decoder(convs, n_classes):
  ...
  o = (tf.keras.layers.Conv2DTranspose( n_classes, kernel_size=(4,4),
                                        strides=(2,2))(o)

  o = tf.keras.layers.Cropping2D(cropping=(1, 1))(o)

  o2 = ( tf.keras.layers.Conv2D(n_classes,(1,1), activation='relu',
                                padding='same'))(f3)

  o = tf.keras.layers.Add()([o, o2])
```
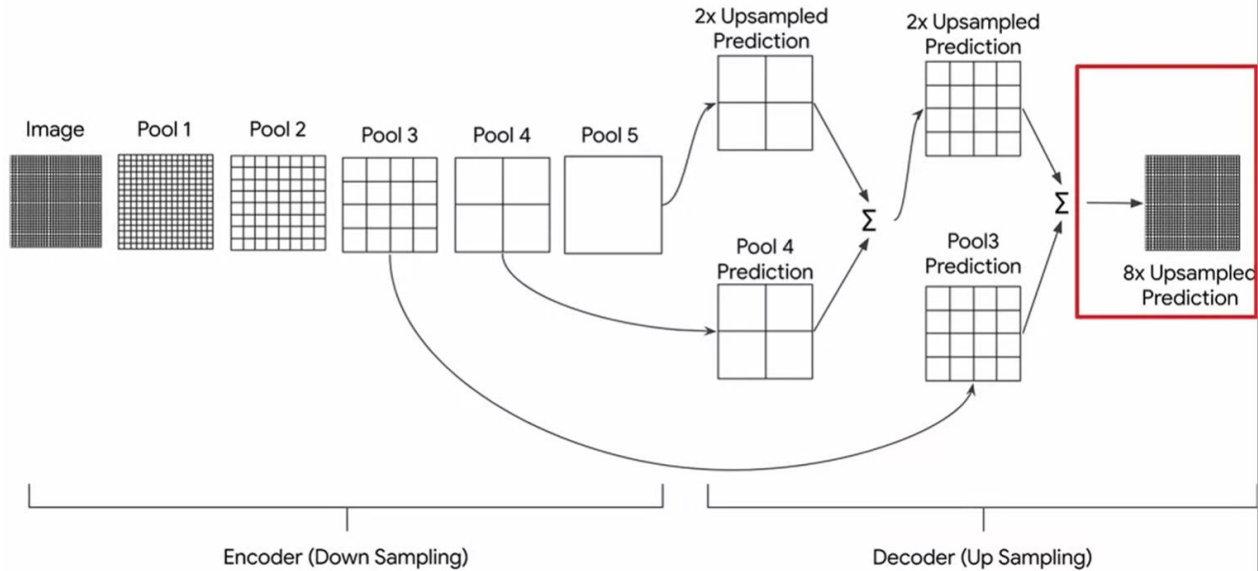
# 2 x 2 x 8 Up Sampled



## Define Decoder

```python
def fcn8_decoder(convs, n_classes):

    ...

    o = tf.keras.layers.Conv2DTranspose(n_classes , kernel_size=(8,8),
                                        strides=(8,8))(o)

    o = (tf.keras.layers.Activation('softmax'))(o)

    return o
```

## Define Final Model

```python
def segmentation_model():
    inputs = tf.keras.layers.Input(shape=(224,224,3,))
    convs = VGG_16(image_input=inputs)
    outputs = fcn8_decoder(convs, 12)
    model = tf.keras.Model(inputs=inputs, outputs=outputs)
    return model


model = segmentation_model()
```

# Evaluation with IoU and Dice Score

## Sample Visualization of Predicted Segments



Original Image

Predicted Segments

Ground Truth Segments

# Area of Overlap

Area of Overlap = sum(True Positives)



Ground Truth

Predictions

**Combined Area** = Total Pixels in predicted segmentation mask + Total Pixels in True Segmentation mask



+

**Area of Union** = Total Pixels in predicted segmentation mask + Total Pixels in True Segmentation mask - Area of Overlap
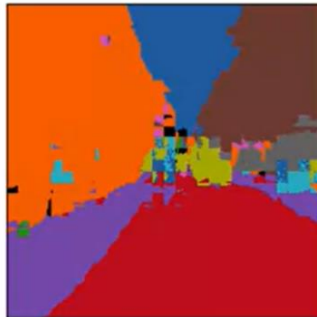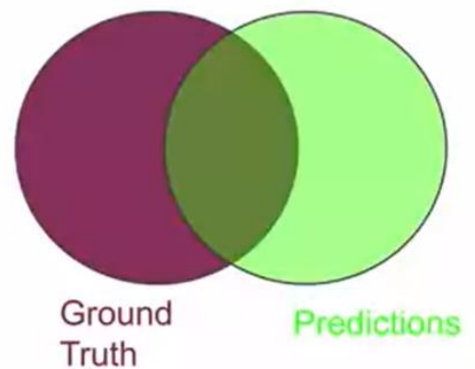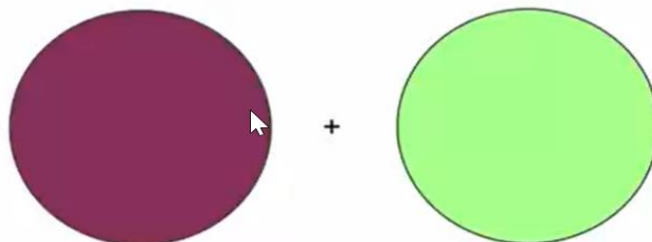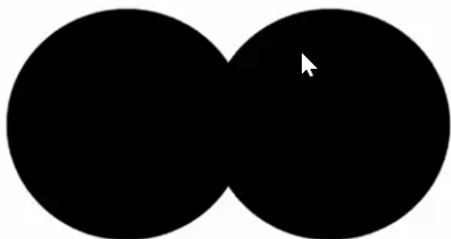
## Calculate Areas

```python
def class_wise_metrics(y_true, y_pred):
    ...
    smoothening_factor = 0.00001

    for i in range(n_classes):
        intersection = np.sum((y_pred == i) * (y_true == i))
        y_true_area = np.sum((y_true == i))
        y_pred_area = np.sum((y_pred == i))
        combined_area = y_true_area + y_pred_area
        union_area = combined_area - intersection
```

## Intersection Over Union

$$IOU = \frac{Area\ of\ Overlap}{Area\ of\ Union}$$

## Calculate IOU

```python
def class_wise_metrics(y_true, y_pred):
    class_wise_iou = []

    ...

    for i in range(n_classes):
        ...
        iou = (intersection) / (union_area)
        class_wise_iou.append(iou)

    return class_wise_iou
```

## IOU Results

| Class | IOU |
|---|---|
| sky | 0.8779669959482955 |
| building | 0.7570989578412737 |
| column/pole | 4.57875457665808e-10 |
| road | 0.915543155822588 |
| side walk | 0.7235628237658467 |
| vegetation | 0.7664541807647628 |
| traffic light | 3.0202657798187055e-05 |
| fence | 0.006380242448568188 |
| vehicle | 0.2950299461448835 |
| pedestrian | 0.0001264333276608086 |
| byciclist | 0.023621930993270864 |
| void | 0.16456276759816527 |

You'll also use the dice score as one of the evaluation metrics in the Colab. The dice score is twice the area of overlap divided by the combined area. It can be used in similar circumstances to the intersection over union score, and they're often both used. The subtle difference between them is that the dice score tends to veer towards the average performance. Whereas the IOU helps you understand worst case performance. Code for this calculated class wise is also pretty straightforward. Iterate through the classes, calculate the intersection and the combined area for that class, and then calculate the score according to the formula. Again, we can get a metric for our classification using it. As you can see here, we did very well in some classes and we did quite poorly in others. Particularly alarming is that the pedestrian and cyclist scores are very poor.

You would not want to use a model like this in production.
Hopefully, that was a useful introduction into how the fully convolutional layers architecture can be used for image segmentation.
Now there's a lot to cover and there's no better way to do that than to dig into the code.
Next, you're going to practice everything that you've seen so far.
Building an image segmentation model to segment dashcam views for self-driving car tasks.
You're going to use VGG16 for the encoder and FCN8 for the decoder, and you'll also use things like IOU and the dice score to evaluate your model's performance. In the next lesson, we'll switch gears and we'll look at the unit architecture for image segmentation

## Dice Score

$$\text{Dice Score} = 2 \times \frac{\text{Area of Overlap}}{\text{Combined Area}}$$



## Calculate Dice Score

```python
def class_wise_metrics(y_true, y_pred):
  class_wise_dice_score = []

  ...

  for i in range(n_classes):
    ...
    dice_score = 2 * (intersection) / (combined_area)
    class_wise_dice_score.append(dice_score)

  return class_wise_dice_score
```

## Dice Score Results

| | |
|---|---|
| sky | 0.9350185576821789 |
| building | 0.861760180856767 |
| column/pole | 9.15750915331616e-10 |
| road | 0.9559097147402678 |
| side walk | 0.8396129387346395 |
| vegetation | 0.8677883515092748 |
| traffic light | 6.040349126864078e-05 |
| fence | 0.012679586065167121 |
| vehicle | 0.45563416820437186 |
| pedestrian | 0.0002528346886971326 |
| byciclist | 0.04615362426586659 |
| void | 0.2826172572009191 |

First is the encoder path, on the left of the U.

Here, similar to the FCNs, we have an image fed into convolution layers and then down sampled using max pooling.

So, for example at the 1st level, if the image is fed in as a 128 by 128 and pass through two convolutional layers that have 64 filters each, the subsequent images when pooled, will be 64 by 64.

***Note that this isn't because of the number of filters, but because a max pooling layer with a 2 by 2 window and a stride of 2,*** can reduce the dimensionality by half. So 128 by 128 becomes 64 by 64.

In the 2nd level, the images are passed through two layers of convolutions with 128 filters each and they are then pulled from 64 by 64 to 32 by 32. In the 3rd level, the 32 by 32 matrices are passed through two layers of 256 filters each and then pulled to 16 by 16. And in the 4th level, the 16 by 16 images are fed through two layers of 512 filters each, before being pulled into an 8 by 8 at the fifth level.

Thus from the 1st level to the 5th, a 128 by 128 image is filtered and down sampled into 8 by 8 blocks.

UNet

An additional element in the unit architecture, is a simple convolutional layer, which can further extract features, but it doesn't have a pooling layer to follow it. This convolutional layer is named the bottleneck. Data flows through that, before we get to the decoder side of the ear.



UNet

And that's here and this is the right hand side of the U shape. Starting at the bottom, the 5th level with the 8 by 8 blocks gets up samples to 16 by 16. So, now were moving up to the 4th level of the U shape, but notice what happens next.

You'll take the 512 filters from the layer of the encoder that's at the same level as this decoder layer, in this case it's the 4th level. Notice that since the encoder layer and decoder layer are the same level in the unit, they also have the same height and width of 16 by 16 and they also have the same number of filters at 512 each.

So, you'll then concatenate the filters from the encoder, with the filters of the decoder, for a total of 1024 filters. You'll then pass this concatenated set of 1024 filters, through 2 convolutional layers.

And this pattern continues through the decoder, you'll up sample the blocks to 32 by 32 and move up to level 3. You'll take the filters from the encoder on the same level, concatenate them to the blocks from the decoder and pass the entire thing through 2 convolutional layers. And of course this works as well as up sampling the image, it will get fed through two layers of convolutional filters,
that match the one on the previous level. But with the additional ones that are the skip over, from the equivalent filters on the encoder side.

So up sampled to 64 by 64, move up to level 2, combine the filters from the encoder with a decoder and pass them through two convolution layers. And finally, up sample to 128 by 128 and move up to level 1, concatenate the filters from the encoder and decoder and pass them through the two convolution layers.