



# **Segmentations**

한양대학교

컴퓨터소프트웨어학부

박현준

# Contents

## [Segmentation]

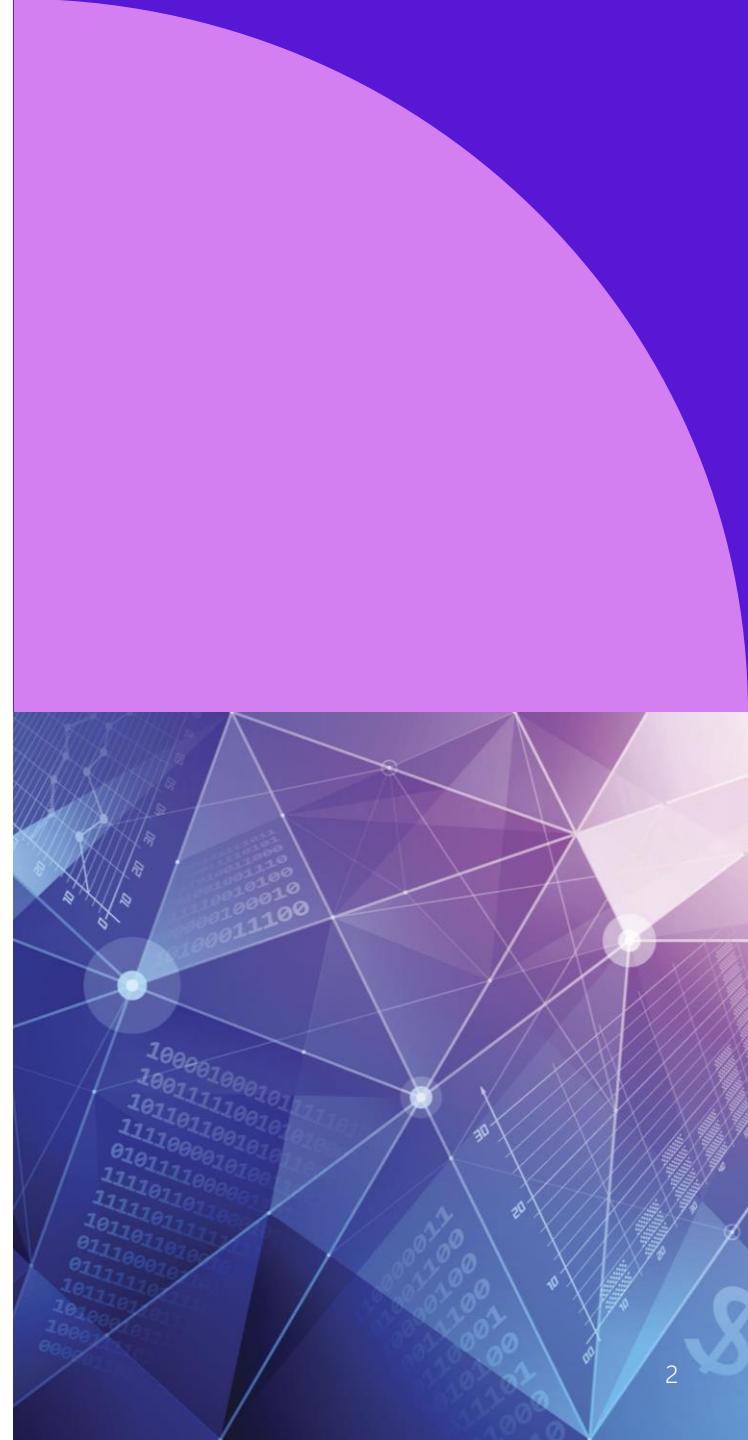
FCN : Semantic Segmentation

Mask R-CNN : Instance Segmentation

Panoptic Segmentation

## [Video Object Detection]

Deep Feature Flow for Video Recognition



# FCN

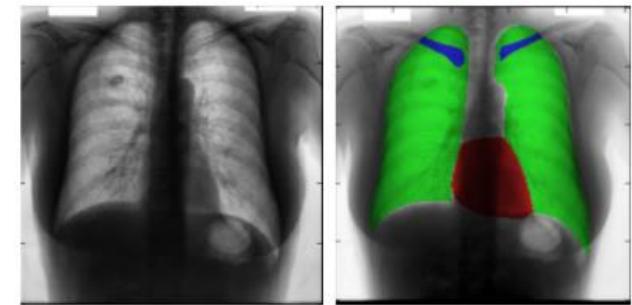
**Fully Convolutional Networks for Semantic Segmentation**

# Introduction

CNN-based Image Classification Network (**AlexNet**, **VGG16**, **GoogLeNet**)  
+  
Semantic Segmentation

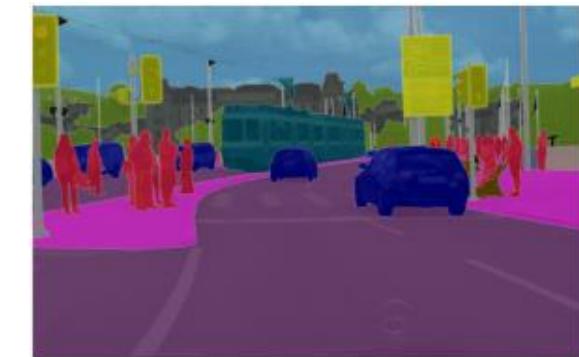
## Transfer-Learning

- ① Visual Recognition
- ② Detection
- ③ Instance & Semantic Segmentation



## Image Classification → Semantic Segmentation

- ① Convolutionization
- ② Deconvolution (Upsampling)
- ③ Skip Architecture



# Fully Convolutional Network

## ConvNets – Translation Invariant

↳ Operates on local input regions, and depend only on relative spatial coordinates

$x_{ij}$  : Data vector at location  $(i, j)$  in a particular layer

$y_{ij}$  : Following Layer

$k$  : kernel size

$s$  : stride or subsampling factor

$f_{ks}$  : layer type (e.g. matmul for conv or avg pool)

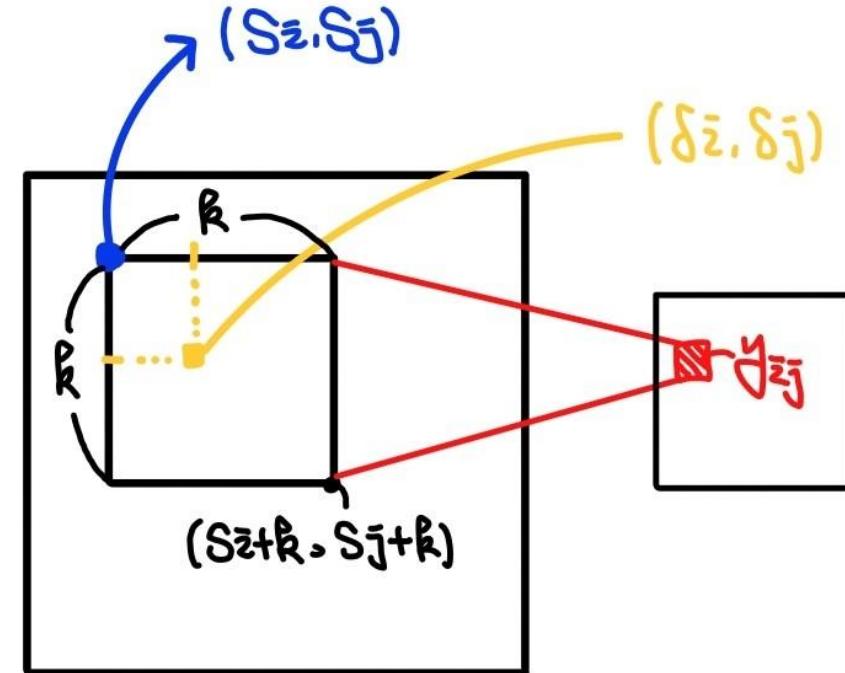
$$y_{ij} = f_{ks} \left( \{x_{si+\delta i, sj+\delta j}\}_{0 \leq \delta i, \delta j \leq k} \right)$$

- Compute Kernel Size, Stride

$$f_{ks} \circ g_{k's'} = (f \circ g)_{k'+(k-1)s', ss'}$$

- Loss Function

$$\ell(x; \theta) = \sum_{ij} \ell'(x_{ij}; \theta)$$



# Fully Convolutional Network

## Adapting classifiers for dense prediction

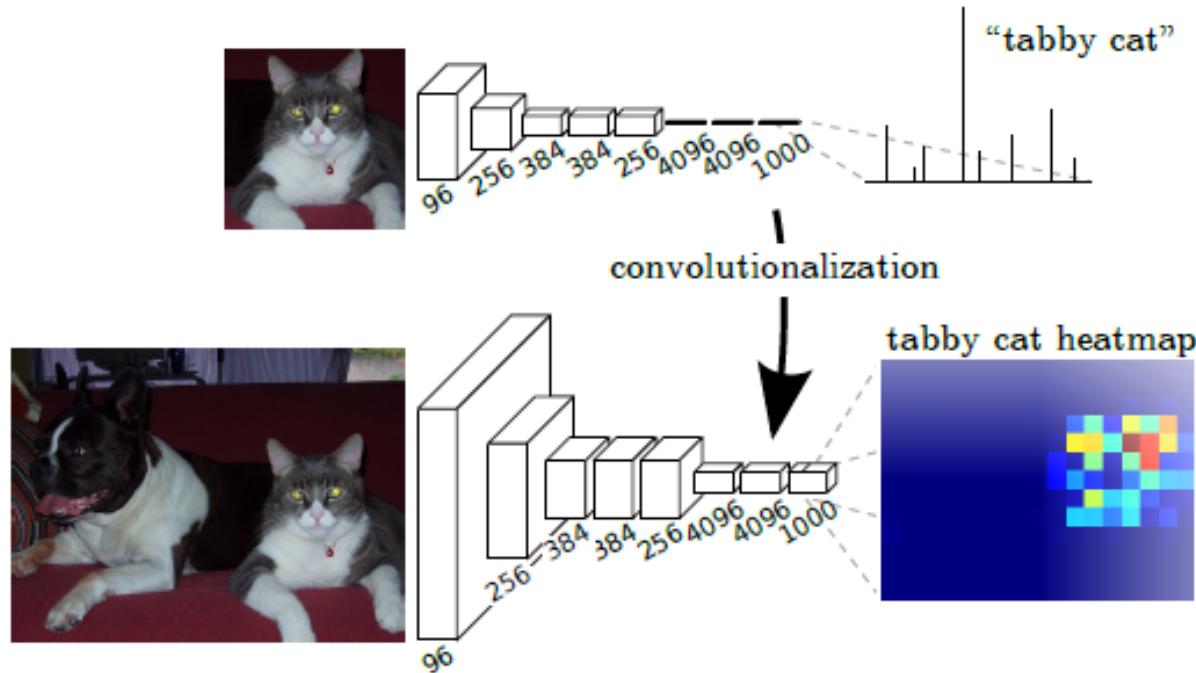
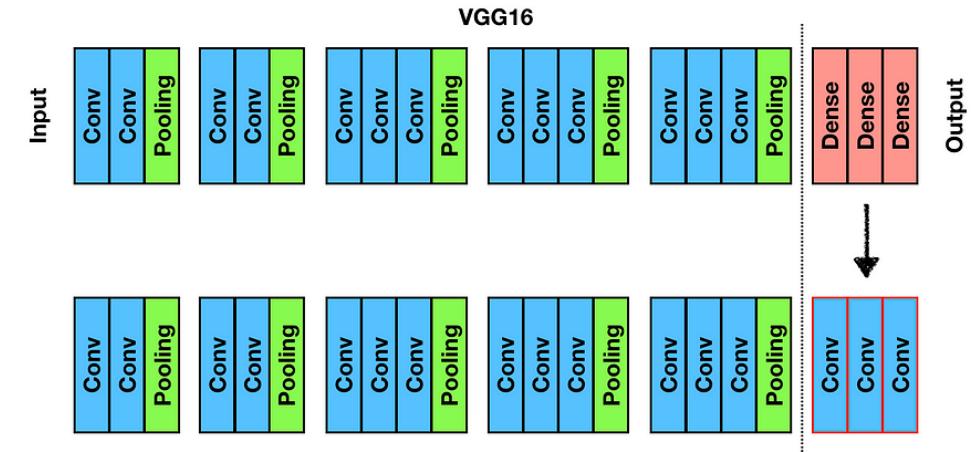


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

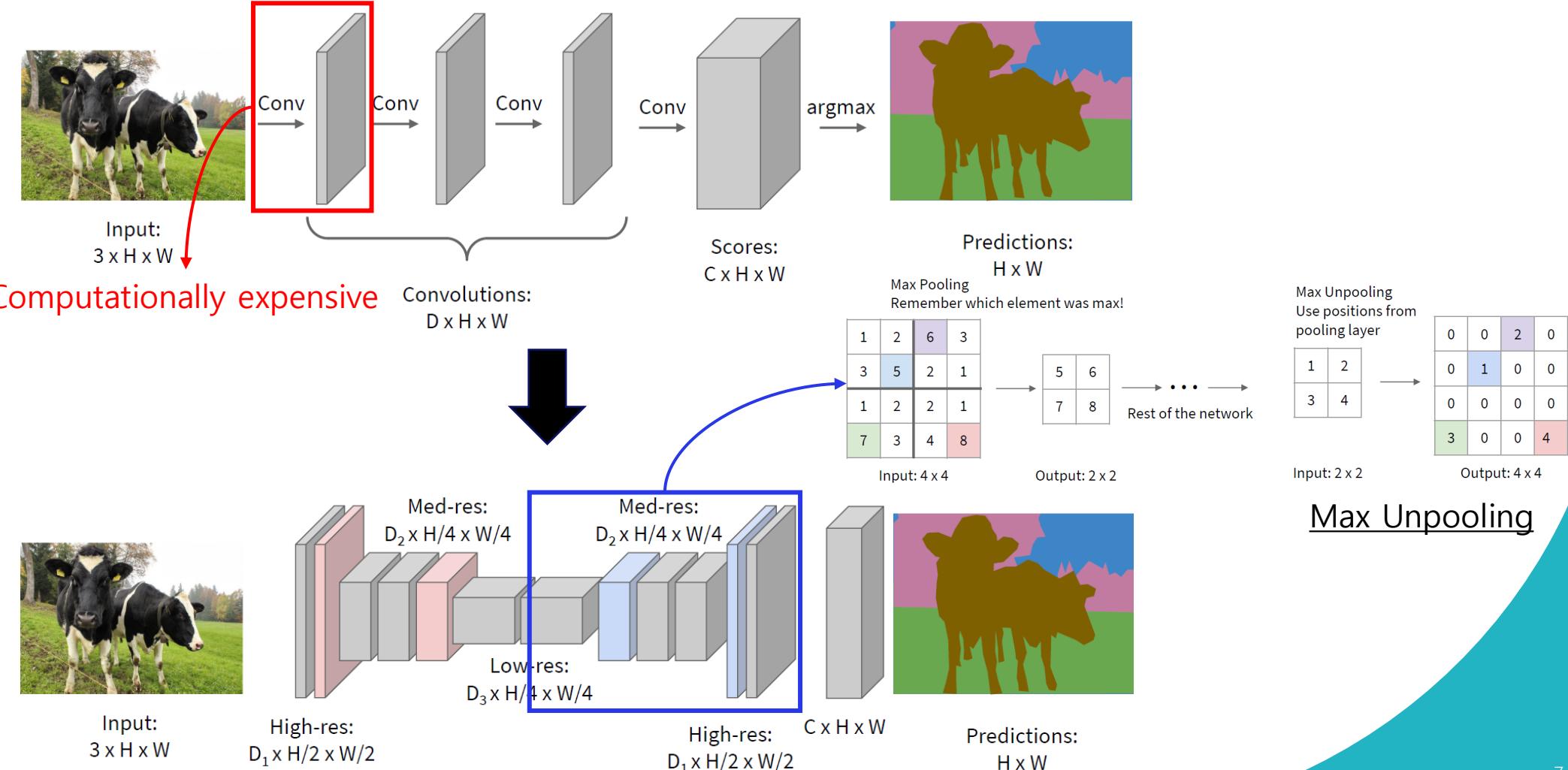


**"Convolutionize"** fc layers

⇒ Preserve spatial information

# Fully Convolutional Network

## Deconvolution (Upsampling)



# Fully Convolutional Network

## Deconvolution (Upsampling)

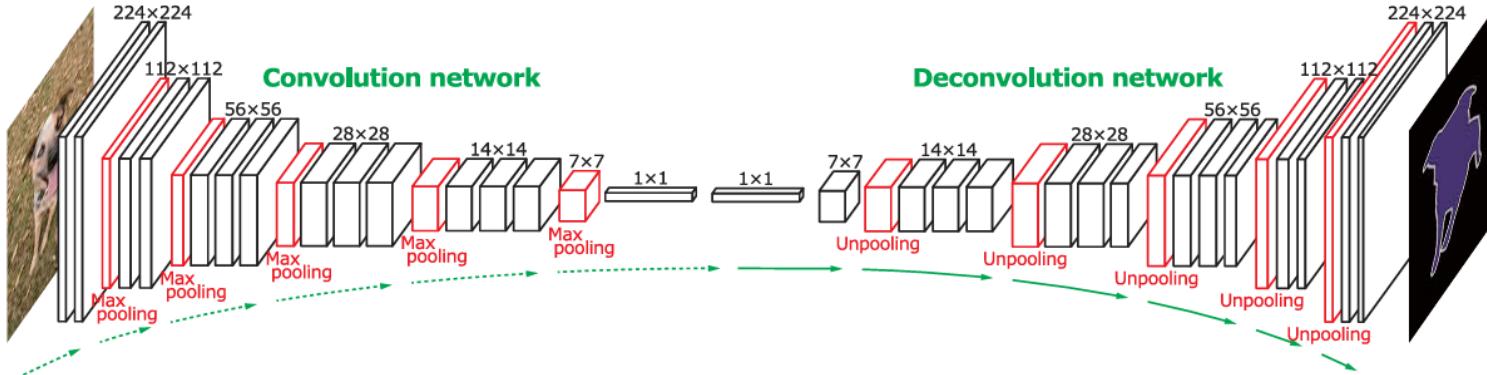


Figure 2. Overall architecture of the proposed network. On top of the convolution network based on VGG 16-layer net, we put a multi-layer deconvolution network to generate the accurate segmentation map of an input proposal. Given a feature representation obtained from the convolution network, dense pixel-wise class prediction map is constructed through multiple series of unpooling, deconvolution and rectification operations.

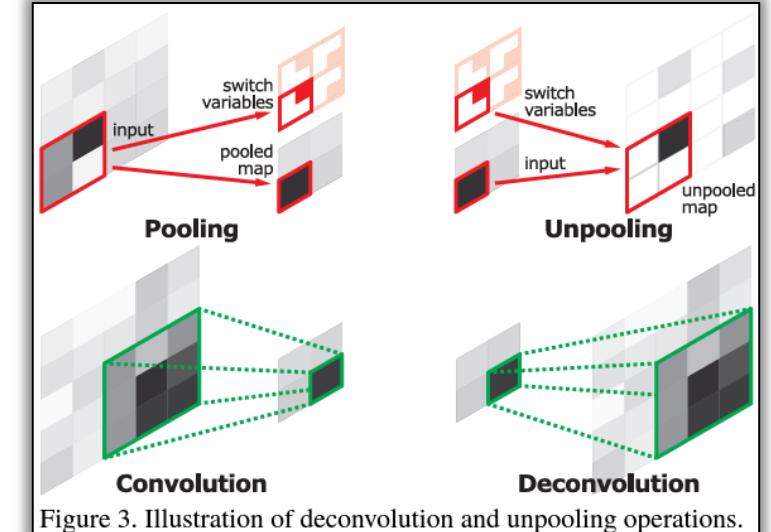
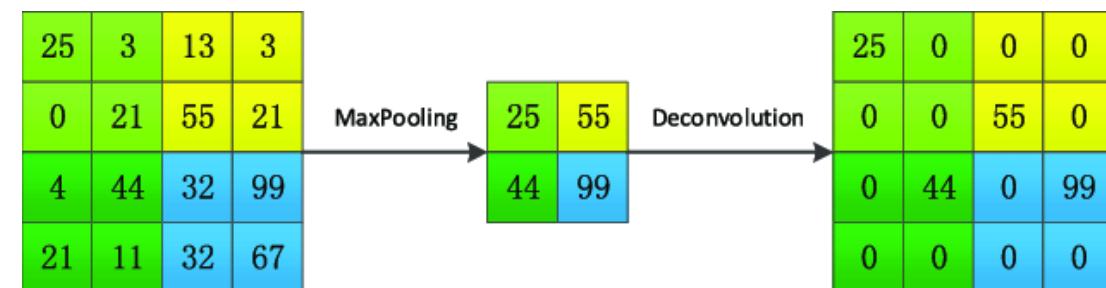
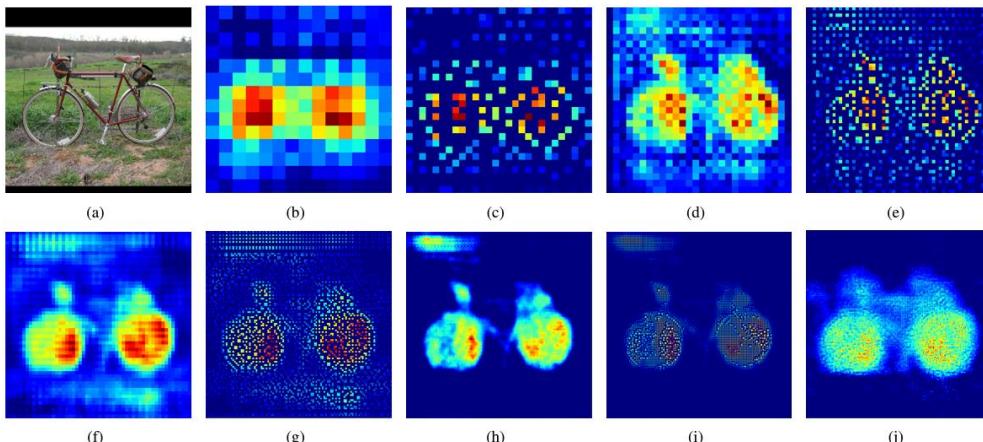
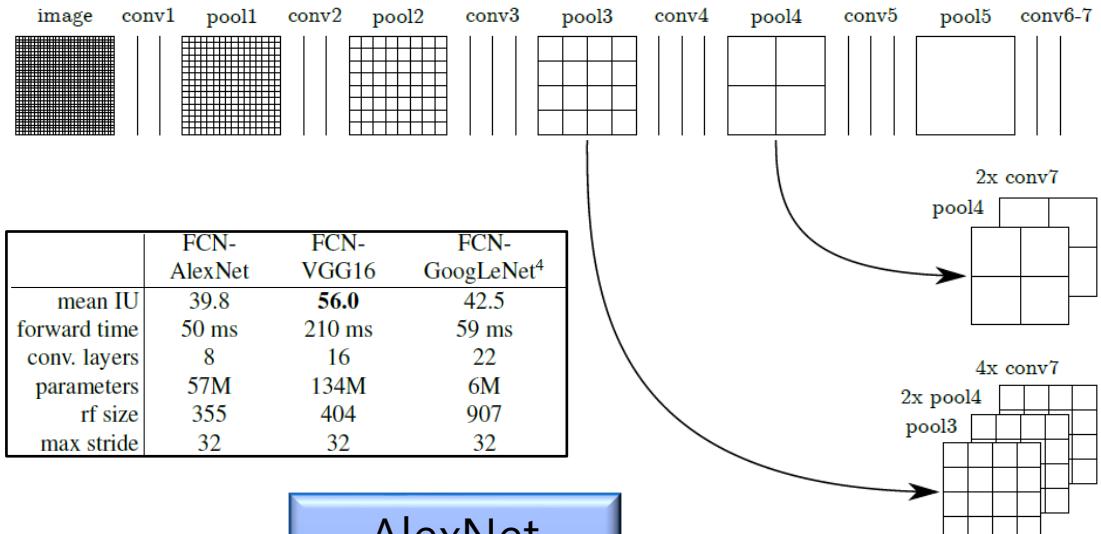


Figure 3. Illustration of deconvolution and unpooling operations.

# Fully Convolutional Network

## Skip Architecture



	FCN-AlexNet	FCN-VGG16	FCN-GoogLeNet <sup>4</sup>
mean IU	39.8	<b>56.0</b>	42.5
forward time	50 ms	210 ms	59 ms
conv. layers	8	16	22
parameters	57M	134M	6M
rf size	355	404	907
max stride	32	32	32

AlexNet  
VGG16  
GoogLeNet

- ILSVRC classifier + FCNs
- Fine-tuning for segmentation
- Skip connection between layers to fuse coarse, semantic and local, appearance info  
→ learned end-to-end to refine the semantics and spatial precision of the output

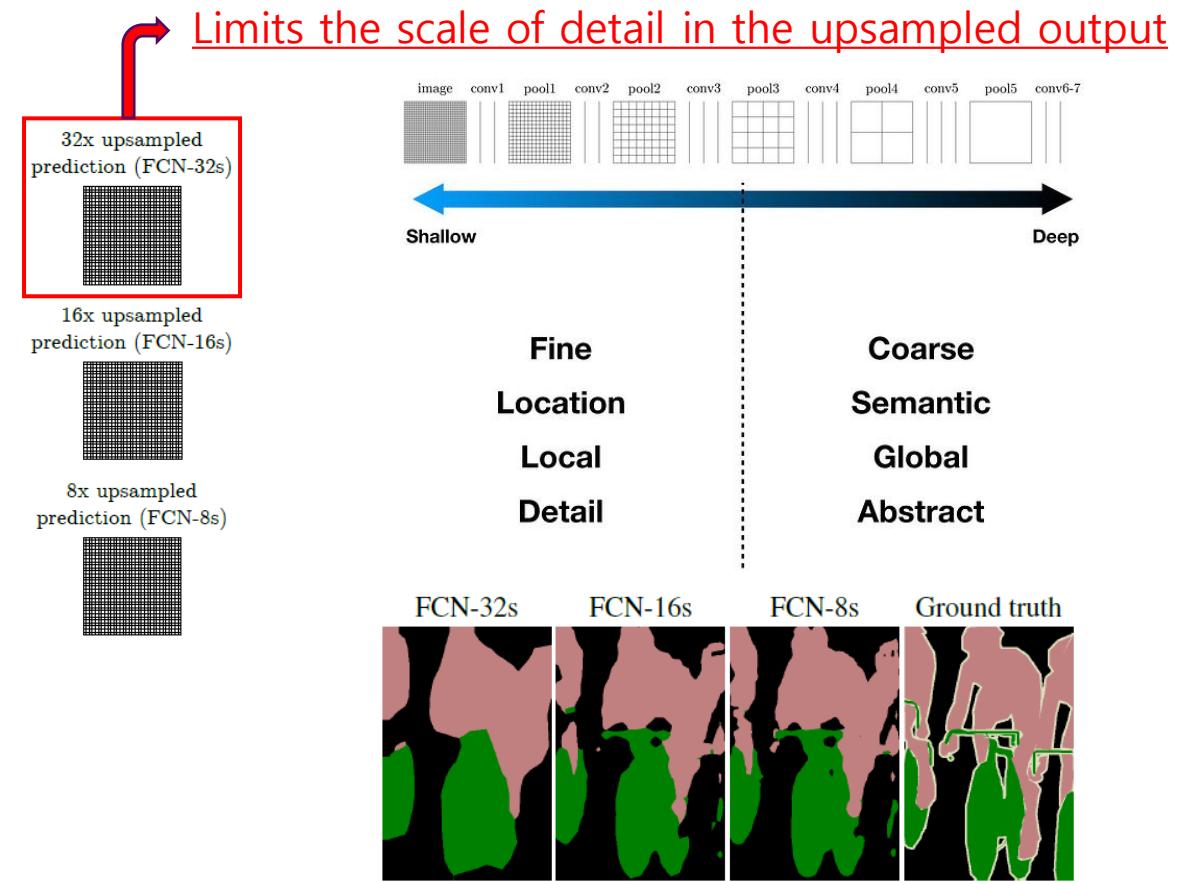
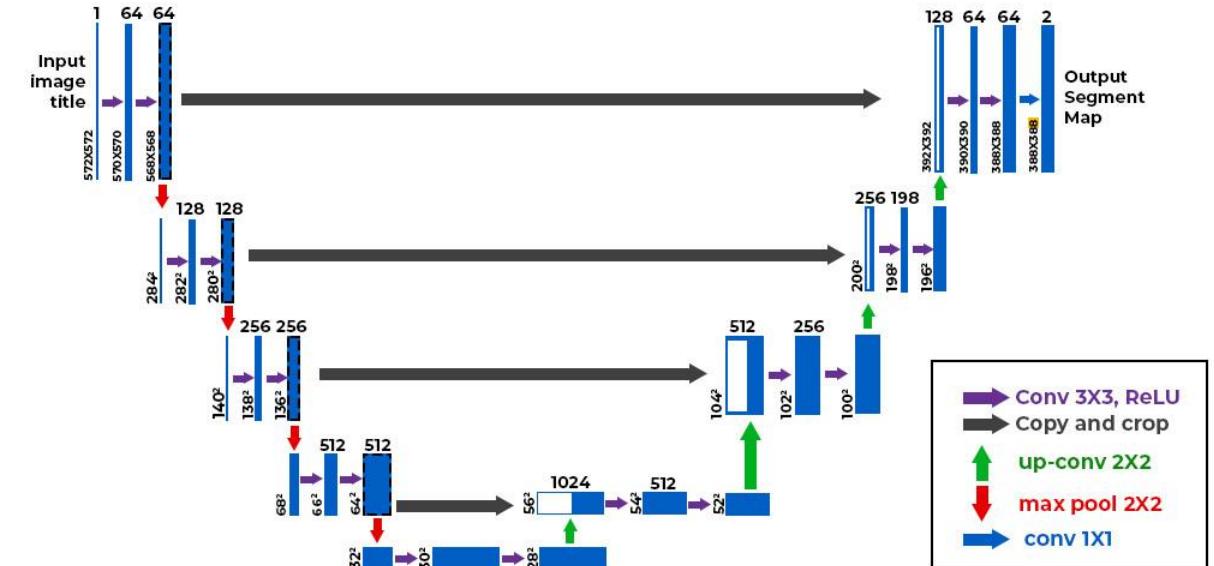
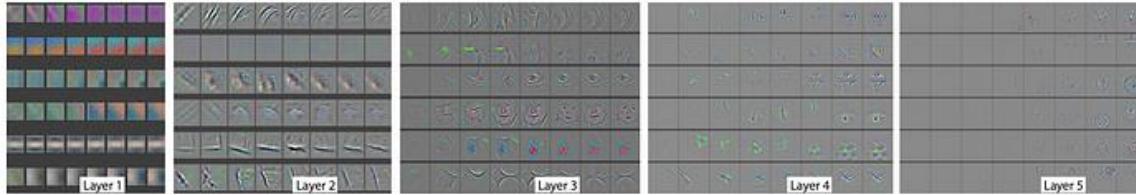


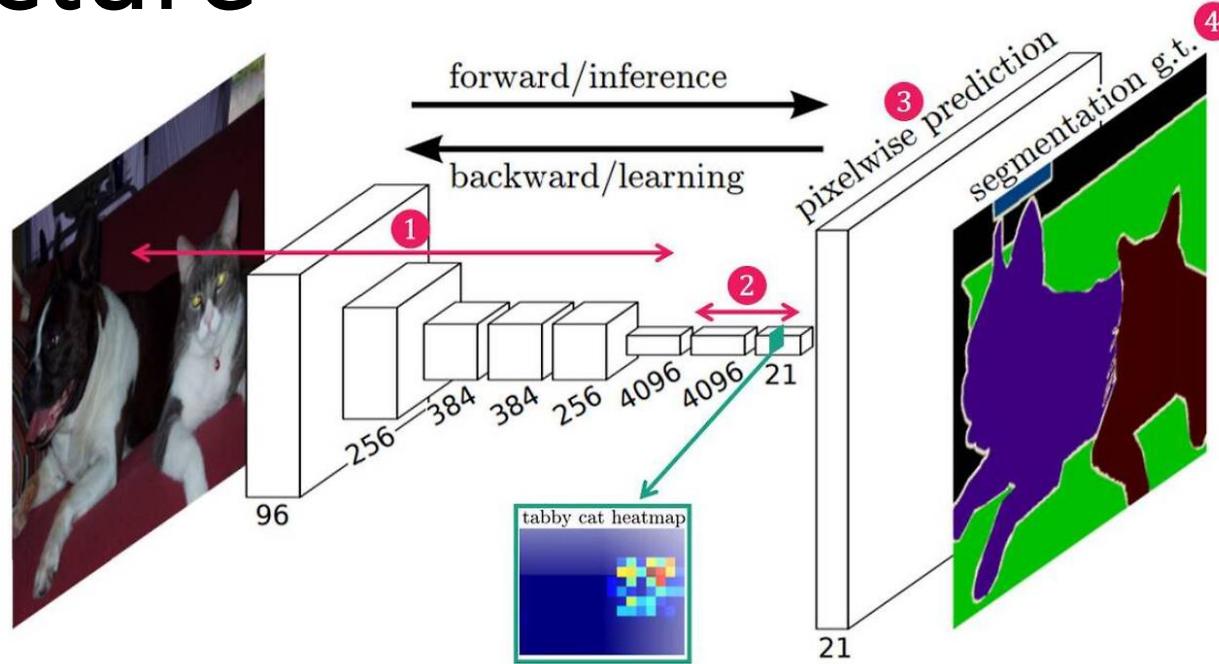
Figure 4. Refining fully convolutional nets by fusing information from layers with different strides improves segmentation detail. The first three images show the output from our 32, 16, and 8 pixel stride nets (see Figure 3).

# Fully Convolutional Network

## Skip Architecture

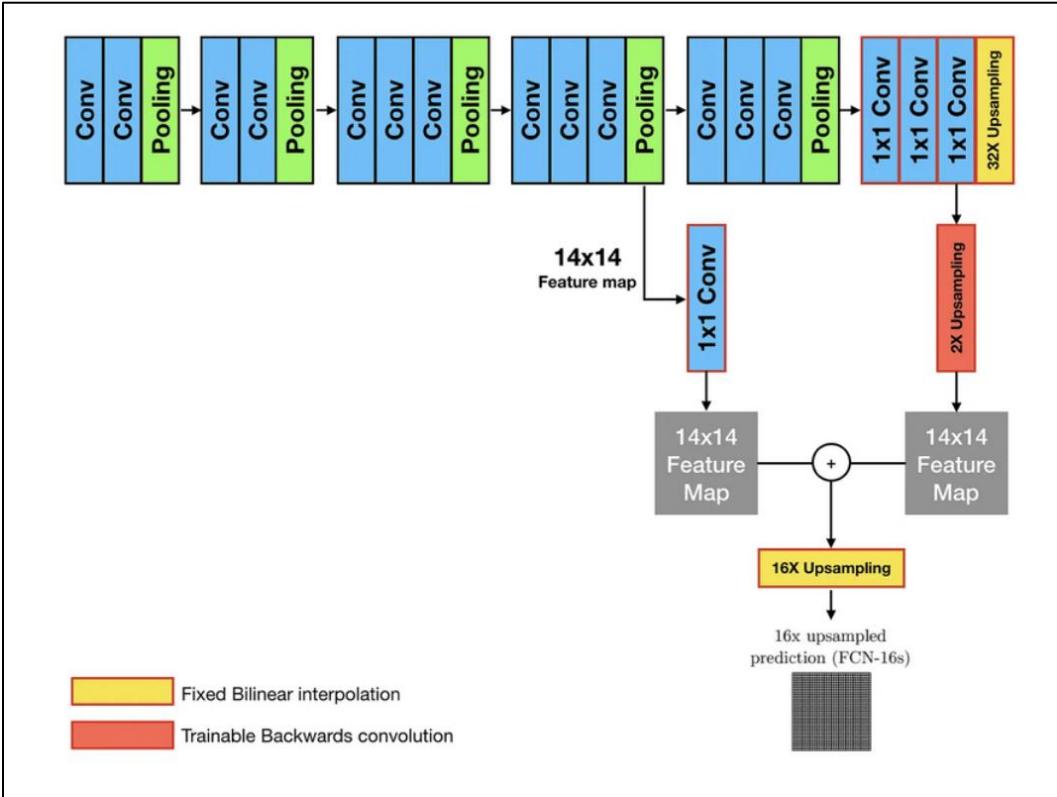


# Architecture

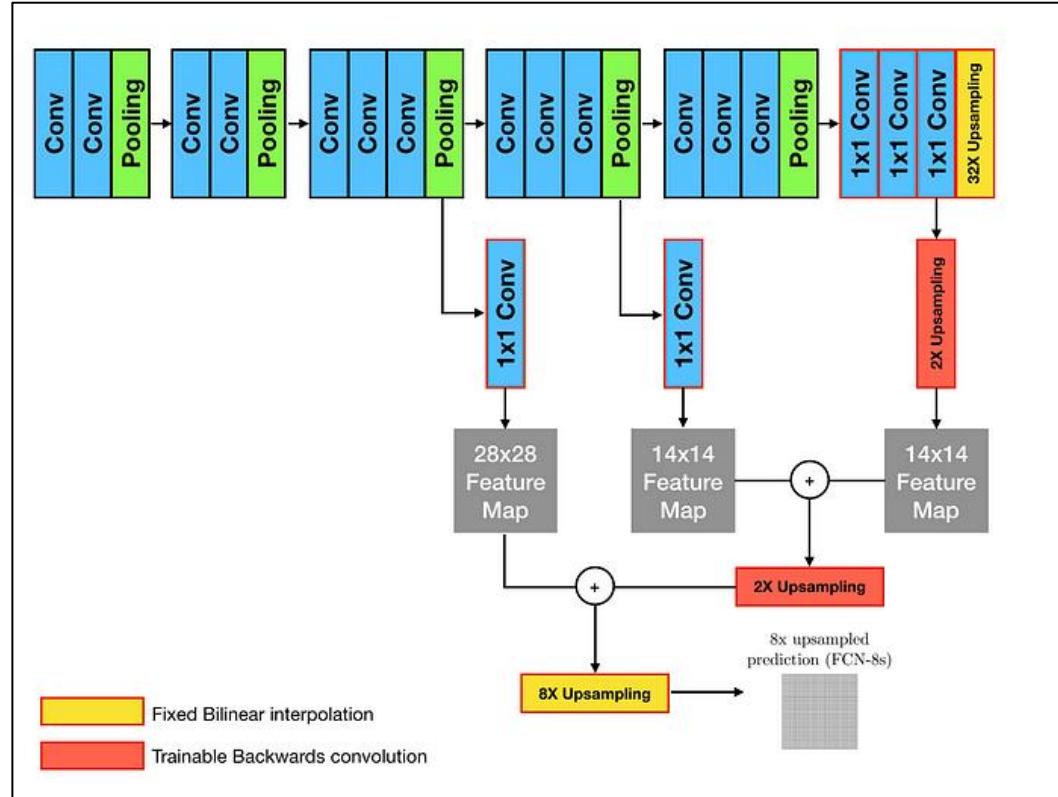


1. Feature map extraction by conv layers
2. Make feature map  $\#channel = \#object$  in dataset  
by using  $1 \times 1$  conv layer (extract class presence heat map)
3. Up-sampling (Transposed convolution)  
→ Create map which size is same as input image
4. Training network using difference between final feature map  
and labeled feature map

# Architecture



FCN-16s



FCN-8s

# Result

## Metrics

- Pixel Accuracy :  $\frac{\sum_i n_{ii}}{\sum_i t_i}$
- Mean Accuracy :  $\frac{1}{n_{cl}} * \frac{\sum_i n_{ii}}{\sum_i t_i}$
- Mean IoU :  $\frac{1}{n_{cl}} * \frac{\sum_i n_{ii}}{t_i + \sum_j n_{ji} - n_{ii}}$
- Frequency weighted IoU :  $(\sum_k t_k)^{-1} \frac{\sum_i t_i n_{ii}}{t_i + \sum_j n_{ji} - n_{ii}}$

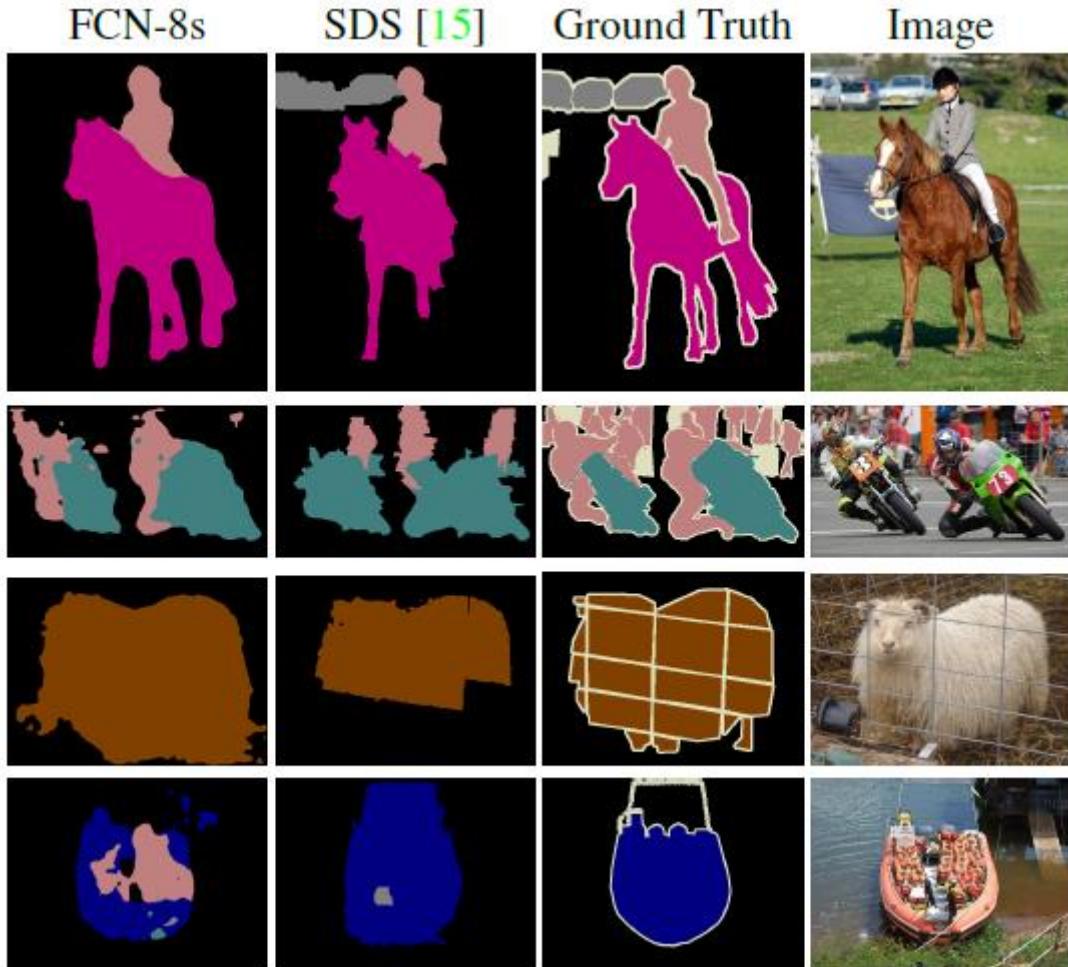
$n_{ij}$  : #pixels  $\in$  *class<sub>i</sub>* predicted to belong to *class<sub>j</sub>*

$n_{cl}$  : #class

$t_i = \sum_j n_{ij}$  : #pixels  $\in$  *class<sub>i</sub>*

# Result

	pixel acc.	mean acc.	mean IU	f.w. IU	
Gupta <i>et al.</i> [13]	60.3	-	28.6	47.0	
FCN-32s RGB	60.0	42.2	29.2	43.9	
FCN-32s RGBD	61.5	42.4	30.5	45.5	
FCN-32s HHA	57.1	35.2	24.2	40.4	
FCN-32s RGB-HHA	64.3	44.9	32.8	48.0	
FCN-16s RGB-HHA	<b>65.4</b>	<b>46.1</b>	<b>34.0</b>	<b>49.5</b>	
	pixel acc.	mean acc.	mean IU	f.w. IU	geom. acc.
Liu <i>et al.</i> [23]	76.7	-	-	-	-
Tighe <i>et al.</i> [33]	-	-	-	-	90.8
Tighe <i>et al.</i> [34] 1	75.6	41.1	-	-	-
Tighe <i>et al.</i> [34] 2	78.6	39.2	-	-	-
Farabet <i>et al.</i> [7] 1	72.3	50.8	-	-	-
Farabet <i>et al.</i> [7] 2	78.5	29.6	-	-	-
Pinheiro <i>et al.</i> [28]	77.7	29.8	-	-	-
FCN-16s	<b>85.2</b>	<b>51.7</b>	39.5	76.1	<b>94.3</b>



# Conclusion

## FCN



## Strength

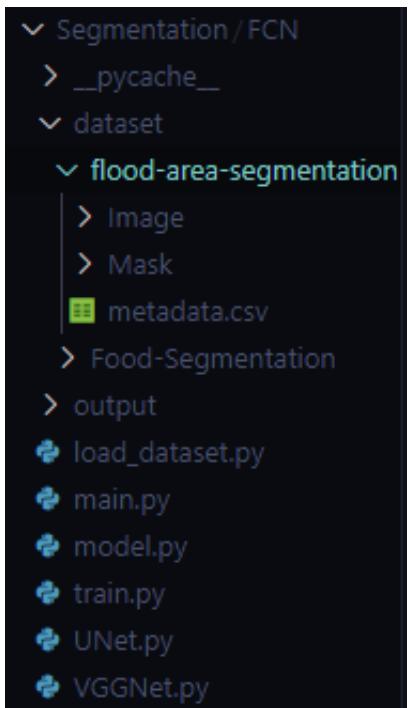
- ▶ “Skip architecture” ⇒ Performance ↑ dramatically
- ▶ Apply image classification network to image segmentation task
- ▶ Convolutionize fc layer  
(∵ lose spatial info, fixed input image size, too much #parameters)

# Implementation

## <Dataset>

: Flood Area  
Segmentation  
(from Kaggle)

- Images : 290
- Labels : 290



```
15:46:03 ➜ jaehak(HJooon) ➜ ...Hyunjoon/Segmentation/FCN ➤
$ python model.py
-----

| Layer (type)       | Output Shape       | Param #   |
|--------------------|--------------------|-----------|
| Conv2d-1           | [1, 64, 224, 224]  | 1,792     |
| ReLU-2             | [1, 64, 224, 224]  | 0         |
| Conv2d-3           | [1, 64, 224, 224]  | 36,928    |
| ReLU-4             | [1, 64, 224, 224]  | 0         |
| MaxPool2d-5        | [1, 64, 112, 112]  | 0         |
| Conv2d-6           | [1, 128, 112, 112] | 73,856    |
| ReLU-7             | [1, 128, 112, 112] | 0         |
| Conv2d-8           | [1, 128, 112, 112] | 147,584   |
| ReLU-9             | [1, 128, 112, 112] | 0         |
| MaxPool2d-10       | [1, 128, 56, 56]   | 0         |
| Conv2d-11          | [1, 256, 56, 56]   | 295,168   |
| ReLU-12            | [1, 256, 56, 56]   | 0         |
| Conv2d-13          | [1, 256, 56, 56]   | 590,080   |
| ReLU-14            | [1, 256, 56, 56]   | 0         |
| MaxPool2d-15       | [1, 256, 28, 28]   | 0         |
| Conv2d-16          | [1, 512, 28, 28]   | 1,180,160 |
| ReLU-17            | [1, 512, 28, 28]   | 0         |
| Conv2d-18          | [1, 512, 28, 28]   | 2,359,808 |
| ReLU-19            | [1, 512, 28, 28]   | 0         |
| MaxPool2d-20       | [1, 512, 14, 14]   | 0         |
| Conv2d-21          | [1, 512, 14, 14]   | 2,359,808 |
| ReLU-22            | [1, 512, 14, 14]   | 0         |
| Conv2d-23          | [1, 512, 14, 14]   | 2,359,808 |
| ReLU-24            | [1, 512, 14, 14]   | 0         |
| MaxPool2d-25       | [1, 512, 7, 7]     | 0         |
| ConvTranspose2d-26 | [1, 512, 14, 14]   | 2,359,808 |
| ReLU-27            | [1, 512, 14, 14]   | 0         |
| BatchNorm2d-28     | [1, 512, 14, 14]   | 1,024     |
| ConvTranspose2d-29 | [1, 256, 28, 28]   | 1,179,904 |
| ReLU-30            | [1, 256, 28, 28]   | 0         |
| BatchNorm2d-31     | [1, 256, 28, 28]   | 512       |
| ConvTranspose2d-32 | [1, 128, 56, 56]   | 295,040   |
| ReLU-33            | [1, 128, 56, 56]   | 0         |
| BatchNorm2d-34     | [1, 128, 56, 56]   | 256       |
| ConvTranspose2d-35 | [1, 64, 112, 112]  | 73,792    |
| ReLU-36            | [1, 64, 112, 112]  | 0         |
| BatchNorm2d-37     | [1, 64, 112, 112]  | 128       |
| ConvTranspose2d-38 | [1, 32, 224, 224]  | 18,464    |
| ReLU-39            | [1, 32, 224, 224]  | 0         |
| BatchNorm2d-40     | [1, 32, 224, 224]  | 64        |
| Conv2d-41          | [1, 2, 224, 224]   | 66        |


Total params: 13,334,050  
Trainable params: 13,334,050  
Non-trainable params: 0


Input size (MB): 0.57  
Forward/backward pass size (MB): 270.46


```

dataset > flood-area-segmentation > metadata.csv > data		
1	Image, Mask	
2	0.jpg,0.png	
3	1.jpg,1.png	
4	2.jpg,2.png	
5	3.jpg,3.png	
6	4.jpg,4.png	
7	5.jpg,5.png	
8	6.jpg,6.png	
9	7.jpg,7.png	
10	8.jpg,8.png	
11	9.jpg,9.png	
12	10.jpg,10.png	
13	11.jpg,11.png	
14	12.jpg,12.png	
15	13.jpg,13.png	
16	14.jpg,14.png	
17	15.jpg,15.png	
18	16.jpg,16.png	
19	17.jpg,17.png	
20	18.jpg,18.png	
21	19.jpg,19.png	
22	20.jpg,20.png	

# Implementation

**FCN**

```

import torch
import torch.nn as nn
from torchsummary import summary as model_summary

class FCN(nn.Module):
    def __init__(self):
        super(FCN, self).__init__()

        self.block1 = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv2d(64, 64, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2)
        )
        self.block2 = nn.Sequential(
            nn.Conv2d(64, 128, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv2d(128, 128, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2)
        )
        self.block3 = nn.Sequential(
            nn.Conv2d(128, 256, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2)
        )
        self.block4 = nn.Sequential(
            nn.Conv2d(256, 512, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv2d(512, 512, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2)
        )
        self.block5 = nn.Sequential(
            nn.Conv2d(512, 512, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv2d(512, 512, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2)
        )

        self.deconv1 = nn.ConvTranspose2d(512, 512, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn1 = nn.BatchNorm2d(512)
        self.deconv2 = nn.ConvTranspose2d(512, 256, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn2 = nn.BatchNorm2d(256)
        self.deconv3 = nn.ConvTranspose2d(256, 128, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn3 = nn.BatchNorm2d(128)
        self.deconv4 = nn.ConvTranspose2d(128, 64, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn4 = nn.BatchNorm2d(64)
        self.deconv5 = nn.ConvTranspose2d(64, 32, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn5 = nn.BatchNorm2d(32)

        self.classifier = nn.Conv2d(32, 2, kernel_size=1)

    def forward(self, x):
        x = self.block1(x)
        x1 = x
        x = self.block2(x)
        x2 = x
        x = self.block3(x)
        x3 = x
        x = self.block4(x)
        x4 = x
        x = self.block5(x)
        x5 = x

        score = self.bn1(self.relu(self.deconv1(x5)))
        score = score + x4
        score = self.bn2(self.relu(self.deconv2(score)))
        score = score + x3
        score = self.bn3(self.relu(self.deconv3(score)))
        score = score + x2
        score = self.bn4(self.relu(self.deconv4(score)))
        score = score + x1
        score = self.bn5(self.relu(self.deconv5(score)))
        score = self.classifier(score)

        return score

    return score

```

**Pre-trained VGG16**

```

import torch
import torch.nn as nn
from torchvision import models
from torchsummary import summary as model_summary

ranges = {'vgg16': ((0, 5), (5, 10), (10, 17), (17, 24), (24, 31))}

class VGGNet(nn.Module):
    def __init__(self, pretrained=True):
        super(VGGNet, self).__init__()
        self.ranges = ranges['vgg16']
        self.features = models.vgg16(weights=pretrained).features

    def forward(self, x):
        output = []
        for idx in range(len(self.ranges)):
            for layer in range(self.ranges[idx][0], self.ranges[idx][1]):
                x = self.features[layer](x)
                output.append(x)
            return output

class FCNs(nn.Module):
    def __init__(self, pretrained_net, n_class):
        super().__init__()
        self.n_class = n_class
        self.pretrained_net = pretrained_net
        self.relu = nn.ReLU(inplace=True)
        self.deconv1 = nn.ConvTranspose2d(512, 512, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn1 = nn.BatchNorm2d(512)
        self.deconv2 = nn.ConvTranspose2d(512, 256, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn2 = nn.BatchNorm2d(256)
        self.deconv3 = nn.ConvTranspose2d(256, 128, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn3 = nn.BatchNorm2d(128)
        self.deconv4 = nn.ConvTranspose2d(128, 64, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn4 = nn.BatchNorm2d(64)
        self.deconv5 = nn.ConvTranspose2d(64, 32, kernel_size=3, stride=2, padding=1, dilation=1, output_padding=1)
        self.bn5 = nn.BatchNorm2d(32)

        self.classifier = nn.Conv2d(32, n_class, kernel_size=1)

    def forward(self, x):
        output = self.pretrained_net(x)

        x5 = output['x5']
        x4 = output['x4']
        x3 = output['x3']
        x2 = output['x2']
        x1 = output['x1']

        score = self.bn1(self.relu(self.deconv1(x5)))
        score = score + x4
        score = self.bn2(self.relu(self.deconv2(score)))
        score = score + x3
        score = self.bn3(self.relu(self.deconv3(score)))
        score = score + x2
        score = self.bn4(self.relu(self.deconv4(score)))
        score = score + x1
        score = self.bn5(self.relu(self.deconv5(score)))
        score = self.classifier(score)

        return score

    vgg16 = VGGNet(pretrained=True)
    model = FCNs(vgg16, 2)
    model

```

**UNet**

```

import torch
import torch.nn as nn
def dual_conv(in_channel, out_channel):
    conv = nn.Sequential(
        nn.Conv2d(in_channel, out_channel, kernel_size=3, padding=1),
        nn.ReLU(inplace=True),
        nn.Conv2d(out_channel, out_channel, kernel_size=3, padding=1),
        nn.ReLU(inplace=True),
    )
    return conv

def crop_tensor(target_tensor, tensor):
    target_size = target_tensor.size()[2]
    tensor_size = tensor.size()[2]
    delta = target_size - tensor_size
    delta = delta // 2
    return tensor[:, :, delta:target_size - delta, delta:target_size - delta]

class UNet(nn.Module):
    def __init__(self):
        super(UNet, self).__init__()

        # Left side (contracting path)
        self.down_conv1 = dual_conv(3, 64)
        self.down_conv2 = dual_conv(64, 128)
        self.down_conv3 = dual_conv(128, 256)
        self.down_conv4 = dual_conv(256, 512)
        self.down_conv5 = dual_conv(512, 1024)
        self.maxpool = nn.MaxPool2d(kernel_size=2, stride=2)

        # Right side (expansion path)
        self.trans1 = nn.ConvTranspose2d(1024, 512, kernel_size=2, stride=2)
        self.up_conv1 = dual_conv(1024, 512)
        self.trans2 = nn.ConvTranspose2d(512, 256, kernel_size=2, stride=2)
        self.up_conv2 = dual_conv(512, 256)
        self.trans3 = nn.ConvTranspose2d(256, 128, kernel_size=2, stride=2)
        self.up_conv3 = dual_conv(256, 128)
        self.trans4 = nn.ConvTranspose2d(128, 64, kernel_size=2, stride=2)
        self.up_conv4 = dual_conv(128, 64)

        # Output layer
        self.out = nn.Conv2d(64, 2, kernel_size=1)

    def forward(self, image):
        # Left side
        x1 = self.down_conv1(image)
        x2 = self.maxpool(x1)
        x3 = self.down_conv2(x2)
        x4 = self.maxpool(x3)
        x5 = self.down_conv3(x4)
        x6 = self.maxpool(x5)
        x7 = self.down_conv4(x6)
        x8 = self.maxpool(x7)
        x9 = self.down_conv5(x8)

        # Right side
        x = self.trans1(x9)
        y = crop_tensor(x, x7)
        x = self.up_conv1(torch.cat((x, y), 1))
        x = self.trans2(x)
        y = crop_tensor(x, x5)
        x = self.up_conv2(torch.cat((x, y), 1))
        x = self.trans3(x)
        y = crop_tensor(x, x3)
        x = self.up_conv3(torch.cat((x, y), 1))
        x = self.trans4(x)
        y = crop_tensor(x, x1)
        x = self.up_conv4(torch.cat((x, y), 1))

        x = self.out(x)

        return x

```

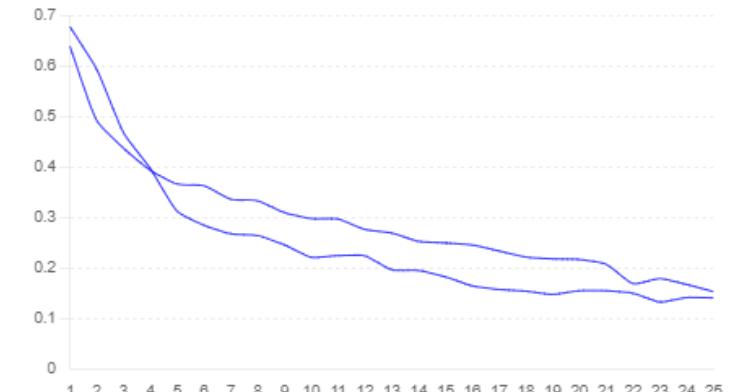
# Implementation

VGG16 (Pre-trained O)

Epoch=25	Epoch=100
Epoch 1/25, Loss: 0.6961	Epoch 1/100, Loss: 0.7409
Epoch 2/25, Loss: 0.6707	Epoch 2/100, Loss: 0.7003
Epoch 3/25, Loss: 0.6345	Epoch 3/100, Loss: 0.6663
Epoch 4/25, Loss: 0.5762	Epoch 4/100, Loss: 0.6038
Epoch 5/25, Loss: 0.5114	Epoch 5/100, Loss: 0.5306
Epoch 6/25, Loss: 0.4536	Epoch 6/100, Loss: 0.5136
Epoch 7/25, Loss: 0.4294	Epoch 7/100, Loss: 0.4545
Epoch 8/25, Loss: 0.4227	Epoch 8/100, Loss: 0.4434
Epoch 9/25, Loss: 0.3976	Epoch 9/100, Loss: 0.4258
Epoch 10/25, Loss: 0.3728	Epoch 10/100, Loss: 0.4258
Epoch 11/25, Loss: 0.3540	Epoch 11/100, Loss: 0.3944
Epoch 12/25, Loss: 0.3486	Epoch 12/100, Loss: 0.3656
Epoch 13/25, Loss: 0.3501	Epoch 13/100, Loss: 0.3694
Epoch 14/25, Loss: 0.3627	Epoch 14/100, Loss: 0.3663
Epoch 15/25, Loss: 0.3440	Epoch 15/100, Loss: 0.3616
Epoch 16/25, Loss: 0.3471	Epoch 16/100, Loss: 0.3534
Epoch 17/25, Loss: 0.3410	Epoch 17/100, Loss: 0.3527
Epoch 18/25, Loss: 0.3186	Epoch 18/100, Loss: 0.3334
Epoch 19/25, Loss: 0.3102	Epoch 19/100, Loss: 0.3198
Epoch 20/25, Loss: 0.3048	Epoch 20/100, Loss: 0.3257
Epoch 21/25, Loss: 0.3012	Epoch 21/100, Loss: 0.3103
Epoch 22/25, Loss: 0.2929	Epoch 22/100, Loss: 0.3126
Epoch 23/25, Loss: 0.2925	Epoch 23/100, Loss: 0.3082
Epoch 24/25, Loss: 0.2979	Epoch 24/100, Loss: 0.3038
Epoch 25/25, Loss: 0.2950	Epoch 25/100, Loss: 0.3136
Training complete.	Training complete.
Epoch 1/25, Loss: 0.6776	Epoch 35/100, Loss: 0.1195
Epoch 2/25, Loss: 0.5996	Epoch 36/100, Loss: 0.1156
Epoch 3/25, Loss: 0.4081	Epoch 37/100, Loss: 0.1184
Epoch 4/25, Loss: 0.3467	Epoch 38/100, Loss: 0.1256
Epoch 5/25, Loss: 0.3132	Epoch 39/100, Loss: 0.1211
Epoch 6/25, Loss: 0.2825	Epoch 40/100, Loss: 0.1192
Epoch 7/25, Loss: 0.2617	Epoch 41/100, Loss: 0.1157
Epoch 8/25, Loss: 0.2435	Epoch 42/100, Loss: 0.1116
Epoch 9/25, Loss: 0.2307	Epoch 43/100, Loss: 0.1086
Epoch 10/25, Loss: 0.2193	Epoch 44/100, Loss: 0.1041
Epoch 11/25, Loss: 0.2091	Epoch 45/100, Loss: 0.0944
Epoch 12/25, Loss: 0.2007	Epoch 46/100, Loss: 0.0864
Epoch 13/25, Loss: 0.1963	Epoch 47/100, Loss: 0.0809
Epoch 14/25, Loss: 0.1914	Epoch 48/100, Loss: 0.0867
Epoch 15/25, Loss: 0.1892	Epoch 49/100, Loss: 0.0810
Epoch 16/25, Loss: 0.1937	Epoch 50/100, Loss: 0.0872
Epoch 17/25, Loss: 0.1905	Epoch 51/100, Loss: 0.0830
Epoch 18/25, Loss: 0.1872	Epoch 52/100, Loss: 0.0850
Epoch 19/25, Loss: 0.1835	Epoch 53/100, Loss: 0.0824
Epoch 20/25, Loss: 0.1763	Epoch 54/100, Loss: 0.0862
Epoch 21/25, Loss: 0.1738	Epoch 55/100, Loss: 0.0881
Epoch 22/25, Loss: 0.1693	Epoch 56/100, Loss: 0.0853
Epoch 23/25, Loss: 0.1648	Epoch 57/100, Loss: 0.0877
Epoch 24/25, Loss: 0.1605	Epoch 58/100, Loss: 0.0857
Epoch 25/25, Loss: 0.1531	Epoch 59/100, Loss: 0.0893
Training complete.	Training complete.

VGG16 (Pre-trained X)

Epoch=25	Epoch=100
Epoch 1/25, Loss: 0.7234	Epoch 35/100, Loss: 0.2248
Epoch 2/25, Loss: 0.6561	Epoch 36/100, Loss: 0.2892
Epoch 3/25, Loss: 0.5858	Epoch 37/100, Loss: 0.1987
Epoch 4/25, Loss: 0.5328	Epoch 38/100, Loss: 0.1893
Epoch 5/25, Loss: 0.4786	Epoch 39/100, Loss: 0.1224
Epoch 6/25, Loss: 0.4627	Epoch 40/100, Loss: 0.1799
Epoch 7/25, Loss: 0.4348	Epoch 41/100, Loss: 0.1085
Epoch 8/25, Loss: 0.4087	Epoch 42/100, Loss: 0.1692
Epoch 9/25, Loss: 0.3899	Epoch 43/100, Loss: 0.1915
Epoch 10/25, Loss: 0.3769	Epoch 44/100, Loss: 0.1789
Epoch 11/25, Loss: 0.3817	Epoch 45/100, Loss: 0.1893
Epoch 12/25, Loss: 0.3641	Epoch 46/100, Loss: 0.1542
Epoch 13/25, Loss: 0.3669	Epoch 47/100, Loss: 0.1774
Epoch 14/25, Loss: 0.3298	Epoch 48/100, Loss: 0.1052
Epoch 15/25, Loss: 0.3346	Epoch 49/100, Loss: 0.1095
Epoch 16/25, Loss: 0.3409	Epoch 50/100, Loss: 0.1533
Epoch 17/25, Loss: 0.3588	Epoch 51/100, Loss: 0.1108
Epoch 18/25, Loss: 0.3589	Epoch 52/100, Loss: 0.1481
Epoch 19/25, Loss: 0.3587	Epoch 53/100, Loss: 0.1087
Epoch 20/25, Loss: 0.3542	Epoch 54/100, Loss: 0.1431
Epoch 21/25, Loss: 0.3385	Epoch 55/100, Loss: 0.1315
Epoch 22/25, Loss: 0.3328	Epoch 56/100, Loss: 0.0978
Epoch 23/25, Loss: 0.3367	Epoch 57/100, Loss: 0.139
Epoch 24/25, Loss: 0.3407	Epoch 58/100, Loss: 0.0939
Epoch 25/25, Loss: 0.3098	Epoch 59/100, Loss: 0.0892
Training complete.	Training complete.
Epoch 1/25, Loss: 0.7695	Epoch 69/100, Loss: 0.1542
Epoch 2/25, Loss: 0.6292	Epoch 70/100, Loss: 0.1869
Epoch 3/25, Loss: 0.5497	Epoch 71/100, Loss: 0.1774
Epoch 4/25, Loss: 0.4949	Epoch 72/100, Loss: 0.1436
Epoch 5/25, Loss: 0.4233	Epoch 73/100, Loss: 0.1224
Epoch 6/25, Loss: 0.4184	Epoch 74/100, Loss: 0.1121
Epoch 7/25, Loss: 0.3980	Epoch 75/100, Loss: 0.1085
Epoch 8/25, Loss: 0.3699	Epoch 76/100, Loss: 0.1283
Epoch 9/25, Loss: 0.3571	Epoch 77/100, Loss: 0.1157
Epoch 10/25, Loss: 0.3585	Epoch 78/100, Loss: 0.1072
Epoch 11/25, Loss: 0.3588	Epoch 79/100, Loss: 0.1154
Epoch 12/25, Loss: 0.3589	Epoch 80/100, Loss: 0.1151
Epoch 13/25, Loss: 0.3574	Epoch 81/100, Loss: 0.1052
Epoch 14/25, Loss: 0.3346	Epoch 82/100, Loss: 0.1293
Epoch 15/25, Loss: 0.3409	Epoch 83/100, Loss: 0.1108
Epoch 16/25, Loss: 0.3587	Epoch 84/100, Loss: 0.1087
Epoch 17/25, Loss: 0.3588	Epoch 85/100, Loss: 0.1032
Epoch 18/25, Loss: 0.3589	Epoch 86/100, Loss: 0.1064
Epoch 19/25, Loss: 0.3587	Epoch 87/100, Loss: 0.0978
Epoch 20/25, Loss: 0.3588	Epoch 88/100, Loss: 0.0939
Epoch 21/25, Loss: 0.3573	Epoch 89/100, Loss: 0.0892
Epoch 22/25, Loss: 0.3588	Epoch 90/100, Loss: 0.0863
Epoch 23/25, Loss: 0.3589	Epoch 91/100, Loss: 0.0829
Epoch 24/25, Loss: 0.3587	Epoch 92/100, Loss: 0.0804
Epoch 25/25, Loss: 0.3588	Epoch 93/100, Loss: 0.0767
Training complete.	Training complete.



UNet

Epoch=25	Epoch=100
Epoch 1/25, Loss: 0.6906	Epoch 35/100, Loss: 0.3421
Epoch 2/25, Loss: 0.6832	Epoch 36/100, Loss: 0.3344
Epoch 3/25, Loss: 0.6787	Epoch 37/100, Loss: 0.3446
Epoch 4/25, Loss: 0.6483	Epoch 38/100, Loss: 0.3322
Epoch 5/25, Loss: 0.5734	Epoch 39/100, Loss: 0.3226
Epoch 6/25, Loss: 0.5198	Epoch 40/100, Loss: 0.3197
Epoch 7/25, Loss: 0.5267	Epoch 41/100, Loss: 0.3197
Epoch 8/25, Loss: 0.4898	Epoch 42/100, Loss: 0.3025
Epoch 9/25, Loss: 0.3372	Epoch 43/100, Loss: 0.3130
Epoch 10/25, Loss: 0.4567	Epoch 44/100, Loss: 0.3224
Epoch 11/25, Loss: 0.4123	Epoch 45/100, Loss: 0.3224
Epoch 12/25, Loss: 0.3742	Epoch 46/100, Loss: 0.3043
Epoch 13/25, Loss: 0.3284	Epoch 47/100, Loss: 0.3043
Epoch 14/25, Loss: 0.4126	Epoch 48/100, Loss: 0.2961
Epoch 15/25, Loss: 0.3582	Epoch 49/100, Loss: 0.3043
Epoch 16/25, Loss: 0.3407	Epoch 50/100, Loss: 0.3097
Epoch 17/25, Loss: 0.3745	Epoch 51/100, Loss: 0.3049
Epoch 18/25, Loss: 0.3745	Epoch 52/100, Loss: 0.3026
Epoch 19/25, Loss: 0.3825	Epoch 53/100, Loss: 0.3093
Epoch 20/25, Loss: 0.3768	Epoch 54/100, Loss: 0.3093
Epoch 21/25, Loss: 0.3892	Epoch 55/100, Loss: 0.3093
Epoch 22/25, Loss: 0.3745	Epoch 56/100, Loss: 0.3093
Epoch 23/25, Loss: 0.3745	Epoch 57/100, Loss: 0.3093
Epoch 24/25, Loss: 0.3688	Epoch 58/100, Loss: 0.3093
Epoch 25/25, Loss: 0.3677	Epoch 59/100, Loss: 0.3093
Training complete.	Training complete.

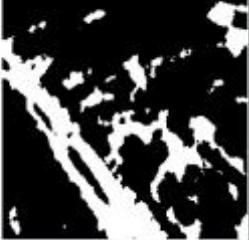
# Implementation

## FCN

Label	Epoch=25			Epoch=100		
Image	  			  		
Ground Truth	  			  		
Prediction	  			  		

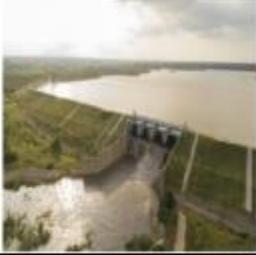
# Implementation

## VGG16 (Pre-trained O)

Label	Epoch=25			Epoch=100		
Image						
Ground Truth						
Prediction						

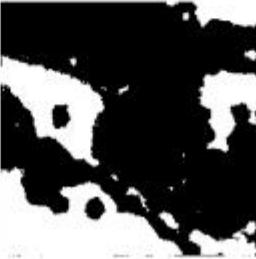
# Implementation

## VGG16 (Pre-trained X)

Label	Epoch=25			Epoch=100		
Image						
Ground Truth						
Prediction						

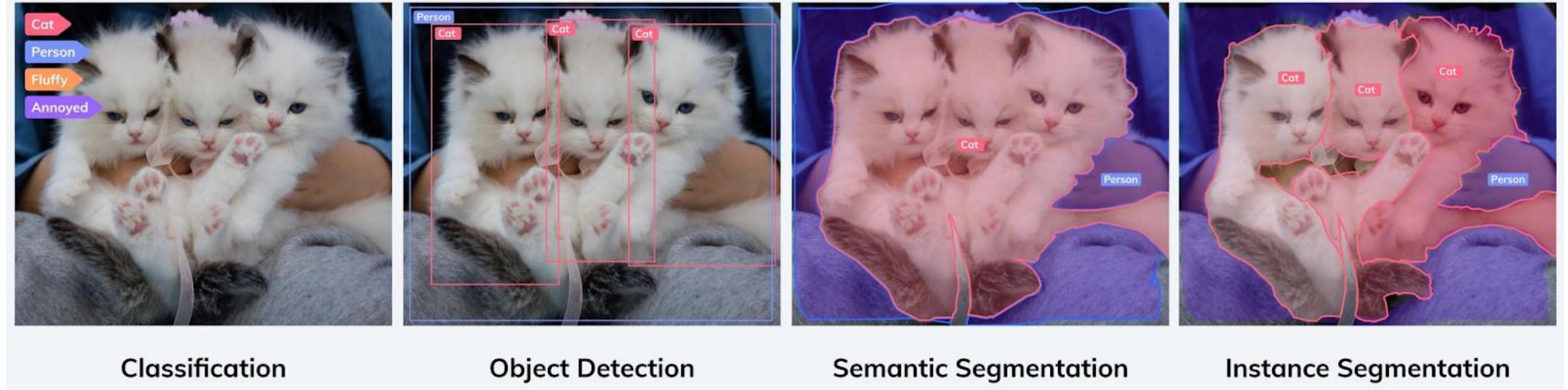
# Implementation

## UNet

Label	Epoch=25			Epoch=100		
Image						
Ground Truth						
Prediction						

# **Mask R-CNN**

# Introduction



Classification

Object Detection

Semantic Segmentation

Instance Segmentation

# Introduction

## Instance Segmentation

Correct detection of all objects + Precisely segmenting each instance

## Faster R-CNN + Mask branch

Adding a branch for predicting segmentation masks on each RoI in parallel with existing branch for classification and bbox regression

## RoIAlign

Simple, quantization-free layer that preserves exact spatial locations  
( $\because$  Faster R-CNN [RPN] was not designed for pixel-to-pixel alignment)

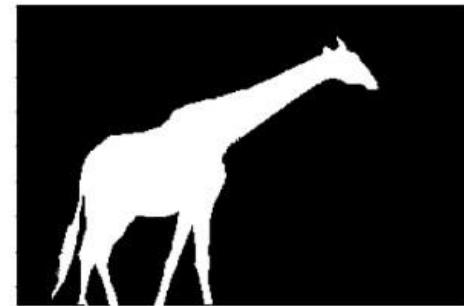
→ Improves mask accuracy 10% to 50%

# Introduction

Base



Binary Mask



"Without bells and whistles"



Mask R-CNN surpasses all previous SOTA single-model results on COCO instance segmentation task without adding additional supplementary functions.

	FCN	Mask R-CNN
Perform	Per-pixel multi-class categorization	Binary mask
Segmentation & Classification	Couple	<b>Decouple</b> (Predict binary mask for each class) → Rely on network's ROI classification

# Related Work

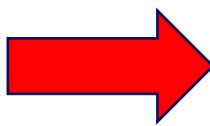
## R-CNN

### RoI-Pooling

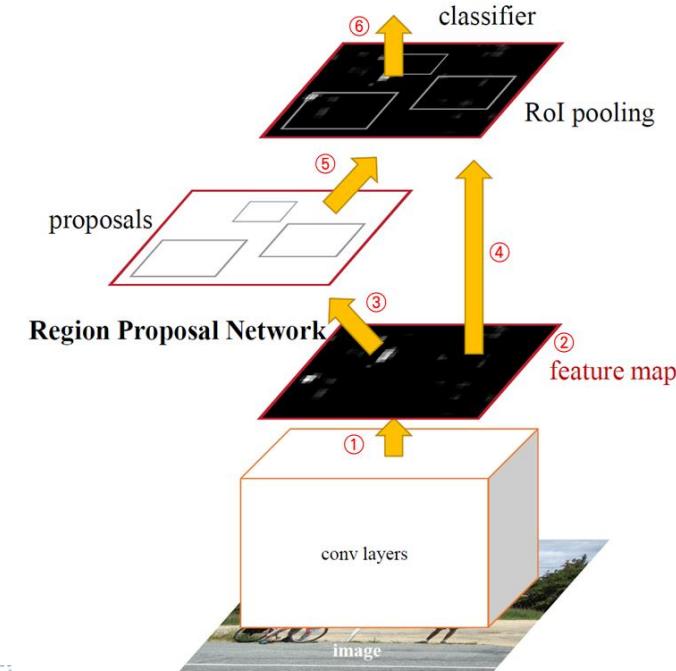
$$\frac{\partial L}{\partial x_i} = \sum_r \sum_j [i = i^*(r, j)] \frac{\partial L}{\partial y_{rj}}$$

$$y_{rj} = x_{i^*(r, j)}$$

$$i^*(r, j) = \operatorname{argmax}_{i' \in \mathcal{R}(r, j)} x_{i'}$$

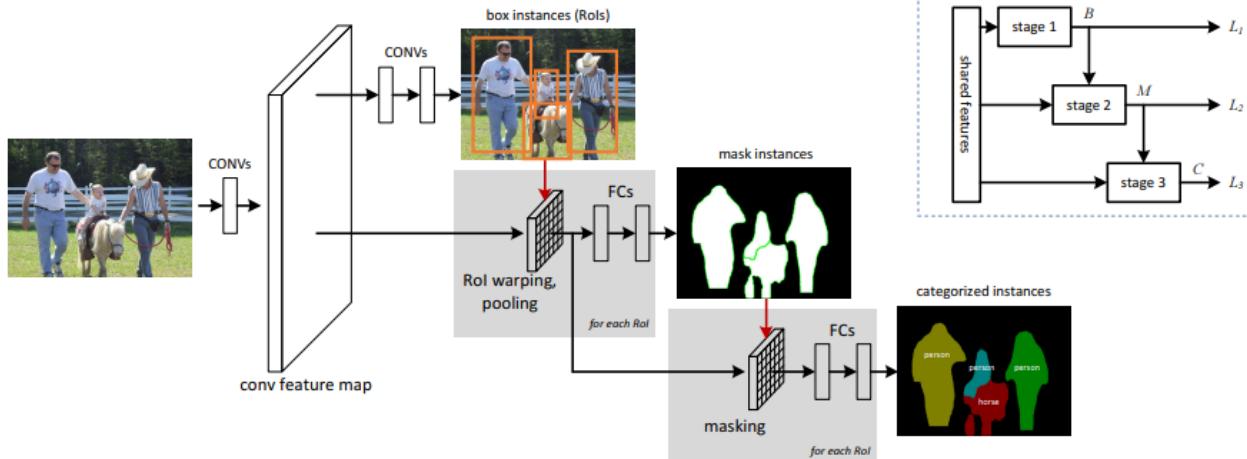


## RPN



Robust  
Flexible

## Instance Segmentation



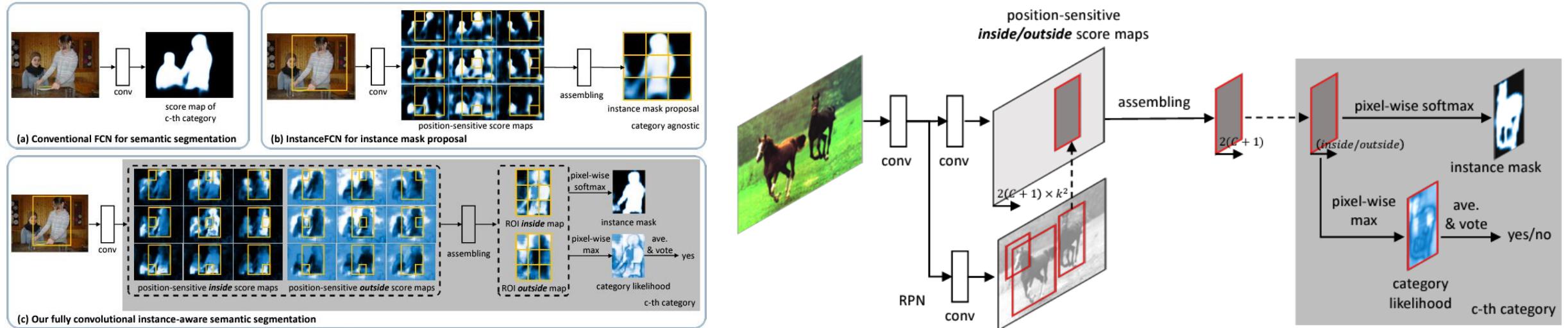
Multiple-Stage  $\Rightarrow$  parallel  
Simple  
Flexible

Figure 2. Multi-task Network Cascades for instance-aware semantic segmentation. At the top right corner is a simplified illustration.

# Related Work

## Instance Segmentation

### Instance-aware Semantic Segmentation



## Segment proposal system + Object detection system

In order to predict a set of position sensitive output channels fully convolutionally

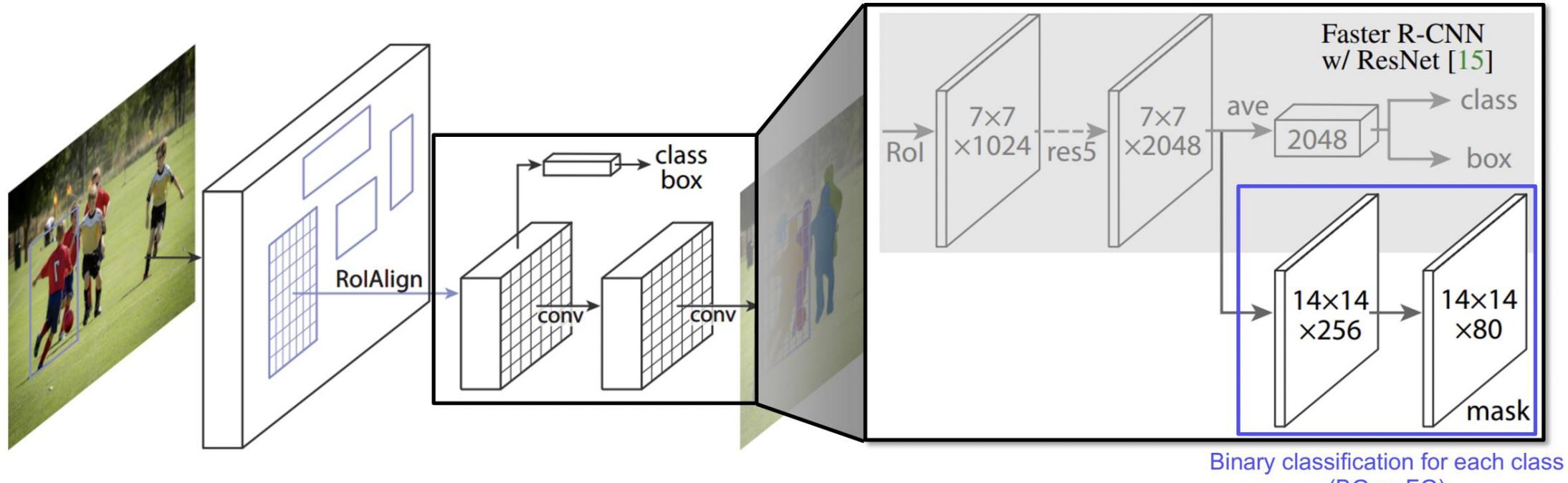
Systematic errors on overlapping instances and creates spurious edges

# Mask R-CNN

## Faster R-CNN + Mask branch

Mask R-CNN

[He et al., ICCV 2017]



Mask R-CNN = Faster R-CNN + Mask branch

# Mask R-CNN

**RoIAlign** – remove harsh quantization of RoIPool

0.8	1.1	2.1	3.9	0.4
2.2	4.3	4.2	1.8	6.3
0.2	5.5	4.8	7.7	9.5
1.1	0.8	2.2	0.8	5.2
5.0	1.8	1.9	2.0	4.1

Quantization

0.8	1.1	2.1	3.9	0.4
2.2	4.3	4.2	1.8	6.3
0.2	5.5	4.8	7.7	9.5
1.1	0.8	2.2	0.8	5.2
5.0	1.8	1.9	2.0	4.1

Cause **misalignment**  
between RoI and extracted features

→ Negative effect on  
predicting pixel-accurate mask

**"Bilinear Interpolation"**

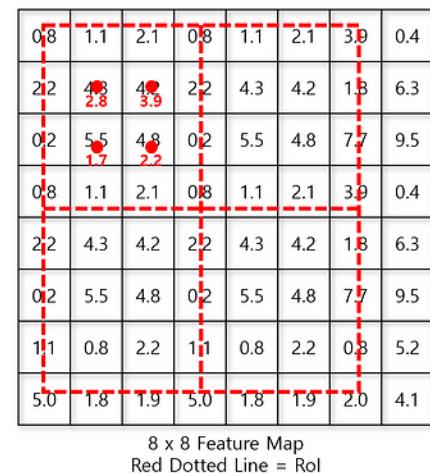
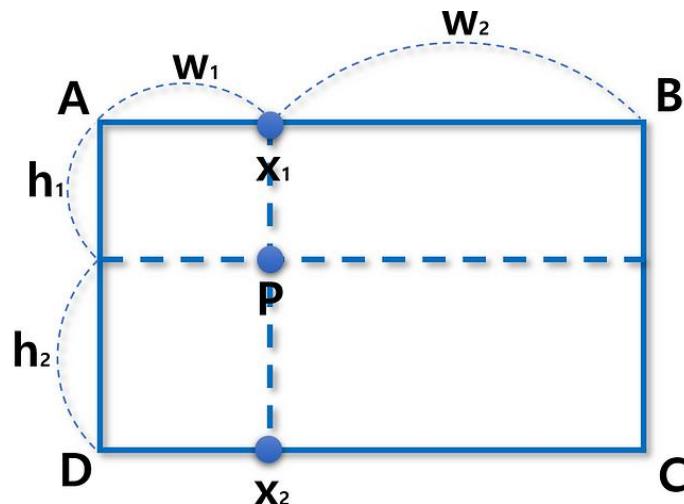


Table 6. RoIAlign vs. RoIPool for keypoint detection on minival.

3.9	4.2
3.6	3.7

2 x 2 RoI  
Fixed Size

# Mask R-CNN

## Loss Function

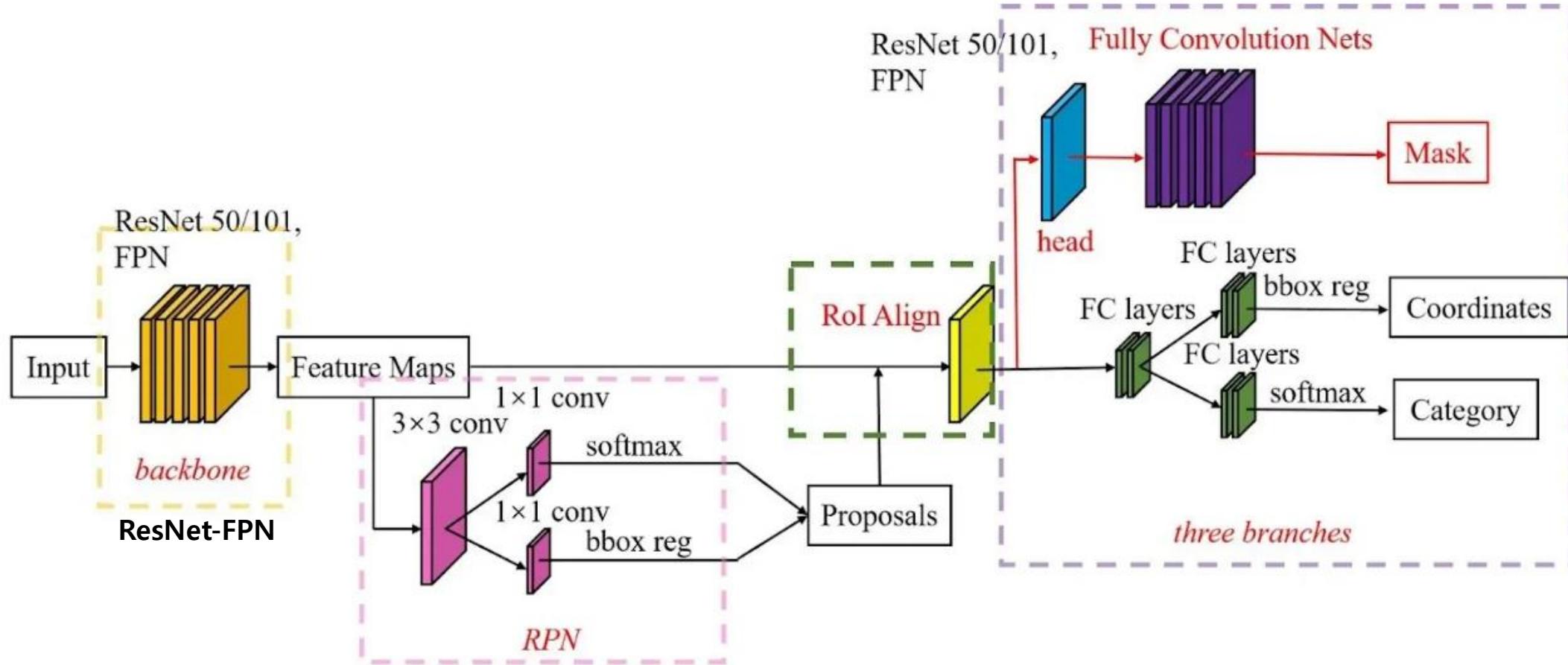
$$L = L_{cls} + L_{box} + L_{mask}$$

- $L_{mask}$
- $km^2$ -dimension ( $m \times m$  images \*  $k$  binary mask class)
  - Apply per-pixel sigmoid
  - Average binary cross-entropy loss

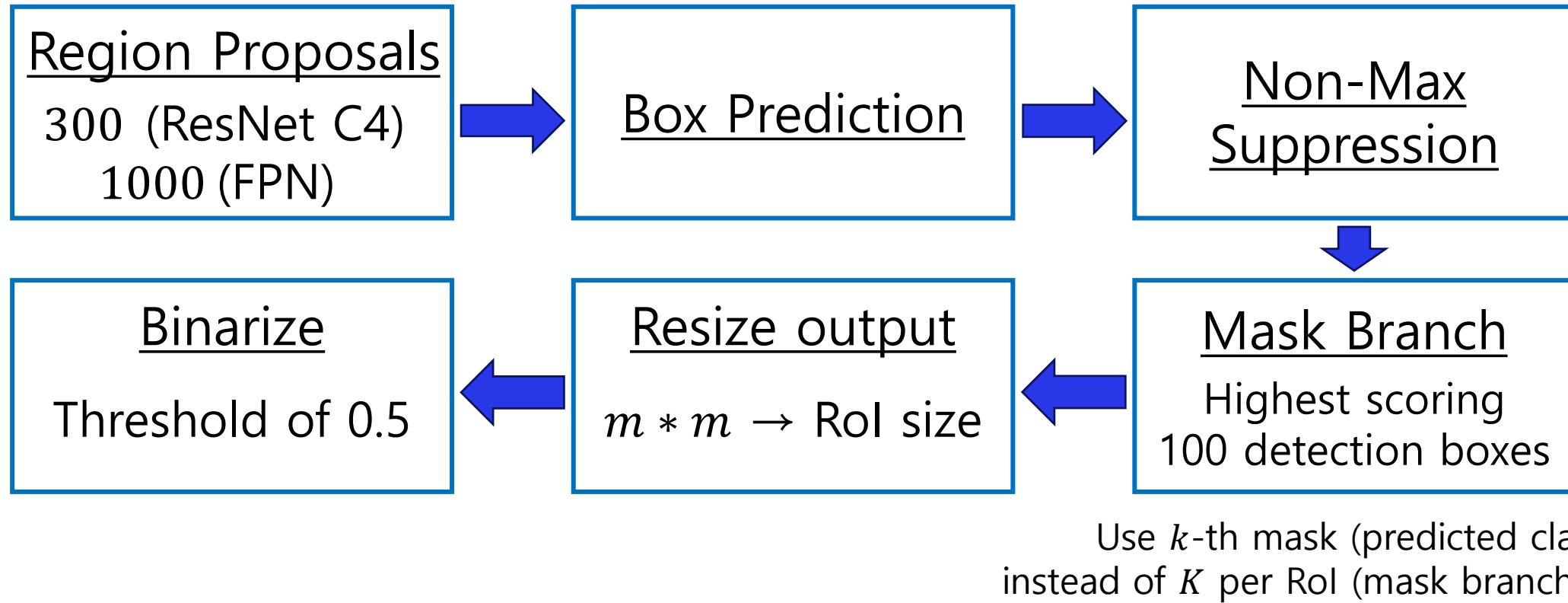
Predict  $m * m$  mask from each RoI using FCN  
⇒ Allow each layer maintain explicit  $m * m$  object spatial layout

Use convolution layer instead of fc layer

# Architecture



# Inference



# Ablations

<i>net-depth-features</i>	AP	AP <sub>50</sub>	AP <sub>75</sub>
ResNet-50-C4	30.3	51.2	31.5
ResNet-101-C4	32.7	54.2	34.3
ResNet-50-FPN	33.6	55.2	35.3
ResNet-101-FPN	35.4	57.3	37.5
ResNeXt-101-FPN	<b>36.7</b>	<b>59.5</b>	<b>38.9</b>

(a) **Backbone Architecture:** Better backbones bring expected gains: deeper networks do better, FPN outperforms C4 features, and ResNeXt improves on ResNet.

	AP	AP <sub>50</sub>	AP <sub>75</sub>
<i>softmax</i>	24.8	44.1	25.1
<i>sigmoid</i>	<b>30.3</b>	<b>51.2</b>	<b>31.5</b>
	+5.5	+7.1	+6.4

(b) **Multinomial vs. Independent Masks (ResNet-50-C4):** *Decoupling* via per-class binary masks (*sigmoid*) gives large gains over multinomial masks (*softmax*).

	<th>bilinear?</th> <th>agg.</th> <th>AP</th> <th>AP<sub>50</sub></th> <th>AP<sub>75</sub></th>	bilinear?	agg.	AP	AP <sub>50</sub>	AP <sub>75</sub>
<i>RoIPool</i> [9]			max	26.9	48.8	26.4
<i>RoIWarp</i> [7]		✓	max	27.2	49.2	27.1
<i>RoIAlign</i>	✓	✓	max	<b>30.2</b>	<b>51.0</b>	<b>31.8</b>
	✓	✓	ave	<b>30.3</b>	<b>51.2</b>	<b>31.5</b>

(c) **RoIAlign (ResNet-50-C4):** Mask results with various RoI layers. Our RoIAlign layer improves AP by ~3 points and AP<sub>75</sub> by ~5 points. Using proper alignment is the only factor that contributes to the large gap between RoI layers.

	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>	AP <sub>75</sub> <sup>bb</sup>
<i>RoIPool</i>	23.6	46.5	21.6	28.2	52.7	26.9
<i>RoIAlign</i>	<b>30.9</b>	<b>51.8</b>	<b>32.1</b>	<b>34.0</b>	<b>55.3</b>	<b>36.4</b>
	+7.3	+5.3	+10.5	+5.8	+2.6	+9.5

(d) **RoIAlign (ResNet-50-C5, stride 32):** Mask-level and box-level AP using *large-stride* features. Misalignments are more severe than with stride-16 features (Table 2c), resulting in massive accuracy gaps.

	mask branch	AP	AP <sub>50</sub>	AP <sub>75</sub>
MLP	fc: 1024 → 1024 → 80 · 28 <sup>2</sup>	31.5	53.7	32.8
MLP	fc: 1024 → 1024 → 1024 → 80 · 28 <sup>2</sup>	31.5	54.0	32.6
FCN	conv: 256 → 256 → 256 → 256 → 256 → 80	<b>33.6</b>	<b>55.2</b>	<b>35.3</b>

(e) **Mask Branch (ResNet-50-FPN):** Fully convolutional networks (FCN) vs. multi-layer perceptrons (MLP, fully-connected) for mask prediction. FCNs improve results as they take advantage of explicitly encoding spatial layout.

Table 2. **Ablations for Mask R-CNN.** We train on `trainval35k`, test on `minival`, and report *mask* AP unless otherwise noted.

# Mask R-CNN for Human Pose Estimation



Figure 6. Keypoint detection results on COCO test using Mask R-CNN (ResNet-50-FPN), with person segmentation masks predicted from the same model. This model has a keypoint AP of 63.1 and runs at 5 fps.

## Keypoint detection

	AP <sup>kp</sup>	AP <sup>kp</sup> <sub>50</sub>	AP <sup>kp</sup> <sub>75</sub>	AP <sup>kp</sup> <sub>M</sub>	AP <sup>kp</sup> <sub>L</sub>
CMU-Pose+++ [4]	61.8	84.9	67.5	57.1	68.2
G-RMI [26] <sup>†</sup>	62.4	84.0	68.5	<b>59.1</b>	68.1
Mask R-CNN, keypoint-only	62.7	87.0	68.4	57.4	71.1
Mask R-CNN, keypoint & mask	<b>63.1</b>	<b>87.3</b>	<b>68.7</b>	57.8	71.4

## Multi-task learning

	AP <sup>bb</sup> <sub>person</sub>	AP <sup>mask</sup> <sub>person</sub>	AP <sup>kp</sup>
Faster R-CNN	52.5	-	-
Mask R-CNN, mask-only	<b>53.6</b>	<b>45.8</b>	-
Mask R-CNN, keypoint-only	50.7	-	64.2
Mask R-CNN, keypoint & mask	52.0	45.1	<b>64.7</b>

# Experiment

## MS COCO

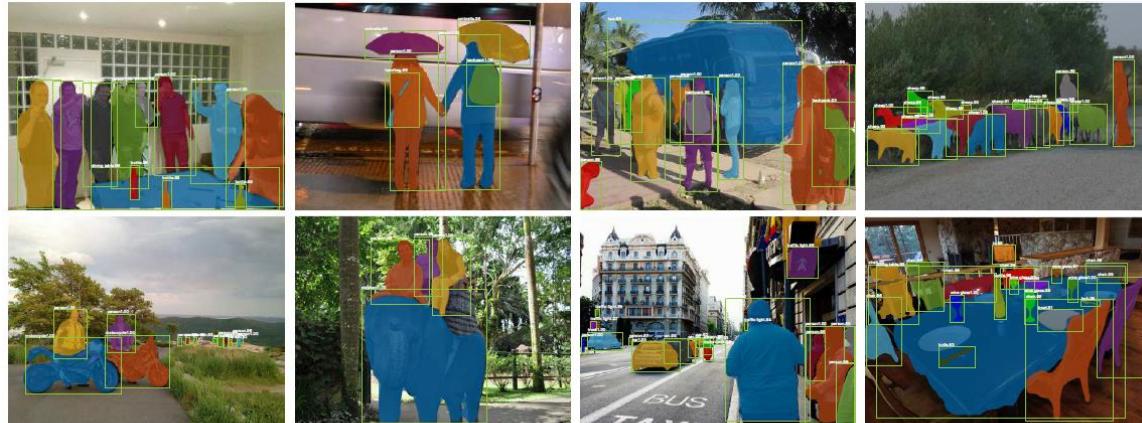


Figure 2. Mask R-CNN results on the COCO test set. These results are based on ResNet-101 [15], achieving a *mask AP* of 35.7 and running at 5 fps. Masks are shown in color, and bounding box, category, and confidences are also shown.



Figure 4. More results of Mask R-CNN on COCO test images, using ResNet-101-FPN and running at 5 fps, with 35.7 mask AP (Table 1).

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
MNC [7]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [21] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [21] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
<b>Mask R-CNN</b>	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
<b>Mask R-CNN</b>	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
<b>Mask R-CNN</b>	ResNeXt-101-FPN	<b>37.1</b>	<b>60.0</b>	<b>39.4</b>	<b>16.9</b>	<b>39.9</b>	<b>53.5</b>

# Experiment

## FCIS vs Mask R-CNN

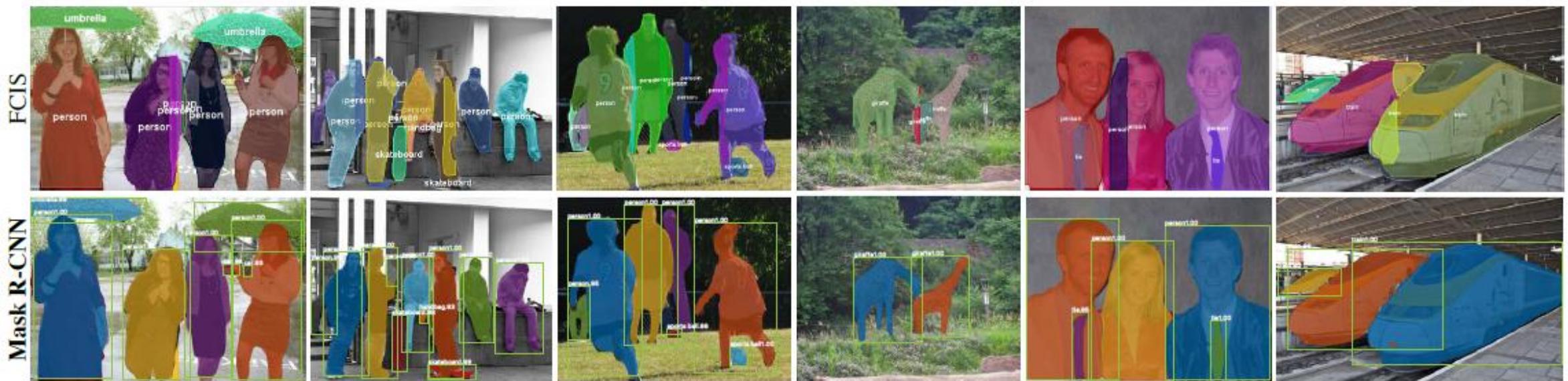
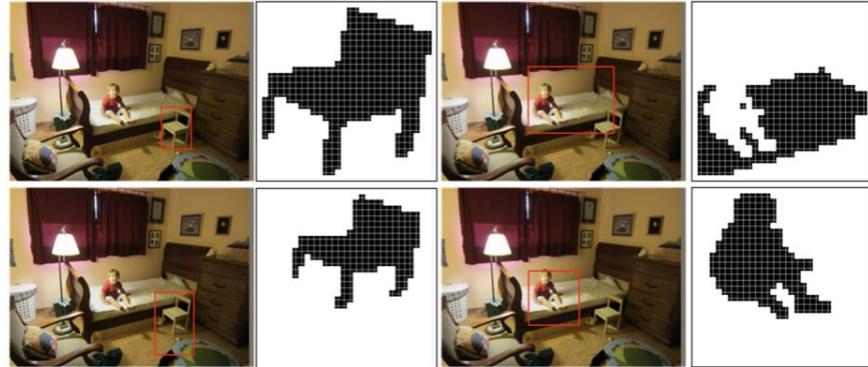


Figure 5. FCIS+++ [21] (top) vs. Mask R-CNN (bottom, ResNet-101-FPN). FCIS exhibits systematic artifacts on overlapping objects.

# Result



	training data	AP [val]	AP	AP <sub>50</sub>	person	rider	car	truck	bus	train	mcycle	bicycle
InstanceCut [23]	fine + coarse	15.8	13.0	27.9	10.0	8.0	23.7	14.0	19.5	15.2	9.3	4.7
DWT [4]	fine	19.8	15.6	30.0	15.1	11.7	32.9	17.1	20.4	15.0	7.9	4.9
SAIS [17]	fine	-	17.4	36.7	14.6	12.9	35.7	16.0	23.2	19.0	10.3	7.8
DIN [3]	fine + coarse	-	20.0	38.8	16.5	16.7	25.7	20.6	30.0	23.4	17.1	10.1
SGN [29]	fine + coarse	29.2	25.0	44.9	21.8	20.1	39.4	24.8	33.2	30.8	17.7	12.4
Mask R-CNN	fine	31.5	26.2	49.9	30.5	23.7	46.9	22.8	32.2	18.6	19.1	16.0
Mask R-CNN	fine + COCO	36.4	32.0	58.1	34.8	27.0	49.1	30.1	40.9	30.9	24.1	18.7

Table 7: Results on Cityscapes val ('AP [val]' column) and test (remaining columns) sets. Our method uses ResNet-50-FPN.



# Conclusion

**Mask R-CNN = Object Detection + Instance Segmentation**



## < Limitations >

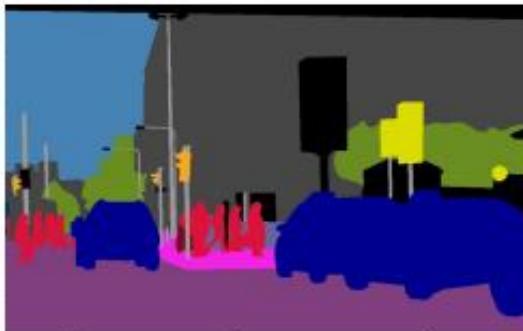
- Computational Complexity
- Small-Object Segmentation
- Data Requirements
- Limited Generalization to Unseen Categories

# Panoptic Segmentation

# Introduction



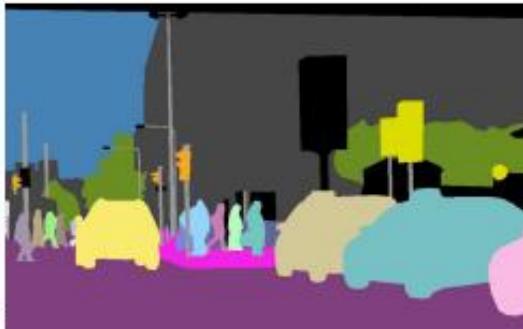
(a) image



(b) semantic segmentation



(c) instance segmentation



(d) panoptic segmentation

## Semantic Segmentation

(Assign class label to each pixel)

+

## Instance Segmentation

(Detect and Segment each object instance)

=

## Panoptic Segmentation

(Including everything visible in one view)

Lack of appropriate metrics or associated recognition challenges  
⇒ “**Panoptic Quality (PQ)**”

# Format

## Task format

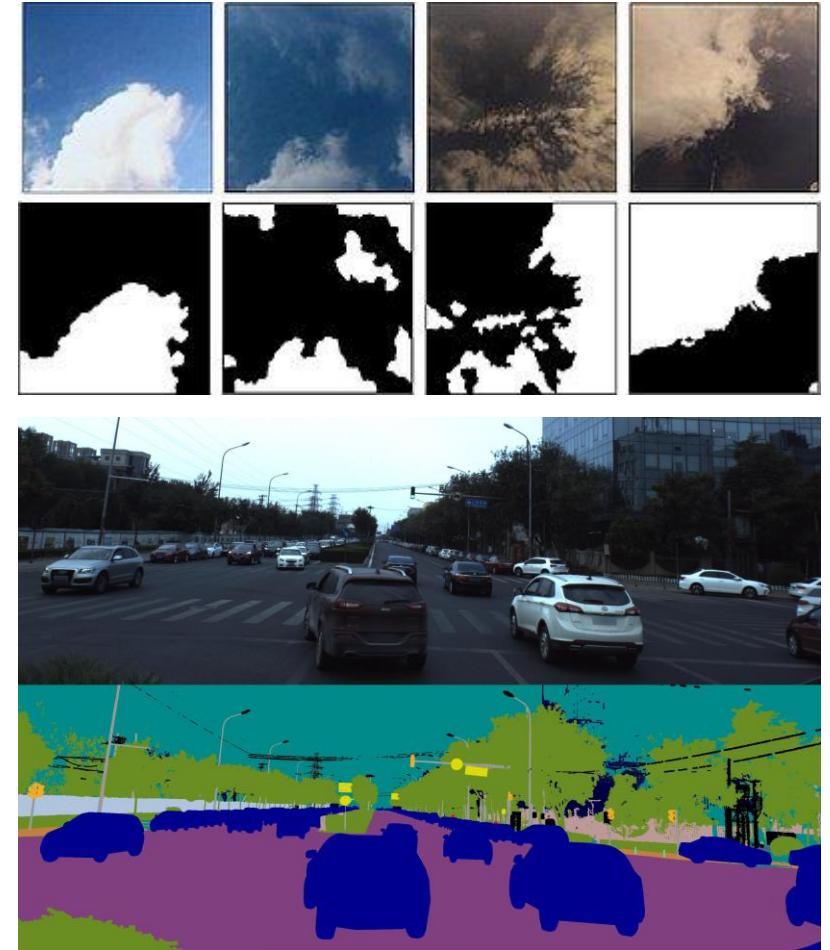
$$\mathcal{L} := \{0, \dots, L - 1\}, \quad (l_i, z_i) \in \mathcal{L} \times \mathbb{N}$$

- $l_i$  : semantic class of pixel  $i$
- $z_i$  : instance id

## Stuff and Thing labels

$$\mathcal{L} = \mathcal{L}^{St} \cup \mathcal{L}^{Th}, \quad \mathcal{L}^{St} \cap \mathcal{L}^{Th} = \emptyset$$

- $l_i \in \mathcal{L}^{St}$  : Stuff classes all pixels belong to same instance  
 $\Leftrightarrow$  its corresponding instance id  $z_i$  is irrelevant
- $l_i \in \mathcal{L}^{Th}$  : All pixels belonging to a single instance must have the same  $(l_i, z_i)$



# Metric

- Completeness
- Interpretability
- Simplicity

## Segment Matching

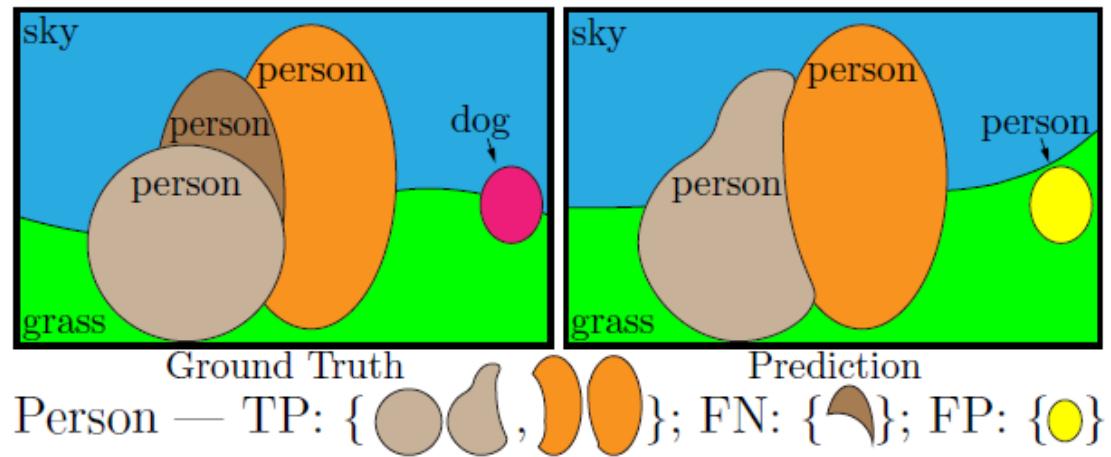
$$IoU(p_i, g) = \frac{|p_i \cap g|}{|p_i \cup g|} \leq \frac{|p_i \cap g|}{|g|} \quad \text{for } i \in \{1,2\}$$

$$|p_1 \cap g| + |p_2 \cap g| \leq |g| \quad (\because p_1 \cap p_2 = \emptyset)$$

$$\Rightarrow IoU(p_1, g) + IoU(p_2, g) \leq \frac{|p_1 \cap g| + |p_2 \cap g|}{|g|} \leq 1$$

$$\therefore \underline{IoU(p_1, g) > 0.5 \leftrightarrow IoU(p_2, g) < 0.5}$$

Only one ground truth segment can have IoU



## PQ Computation

$$PQ = \frac{\sum_{(p,g) \in TP} IoU(p, g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$$

Penalize segment without matches

$$\frac{\sum_{(p,g) \in TP} IoU(p, g)}{|TP|} : \text{avg IoU of matched segments}$$

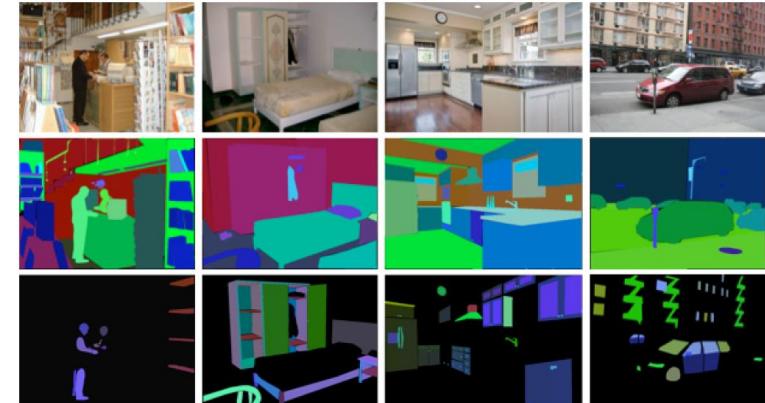
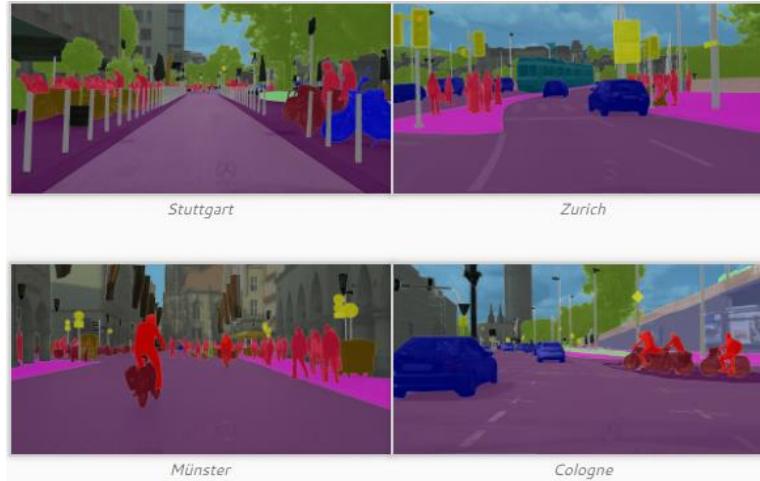
$$PQ = \boxed{\frac{\sum_{(p,g) \in TP} IoU(p, g)}{|TP|}} \times \boxed{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}$$

Segmentation Quality      Recognition Quality

# Result

## Datasets

- Cityscapes
  - 5000 images (2975 train, 500 val, 1525 test)
  - Dense pixel annotations (97% coverage) of 19 classes
- ADE20k
  - 25k ↑ images (20k train, 2k val, 3k test)
  - Densely annotated with an open-dictionary label set
- Mapillary Vistas
  - 25k street-view images (18k train, 2k val, 5k test)
  - Densely annotated (98% pixel coverage)  
with 28 stuff and 37 thing classes



# Human Consistency

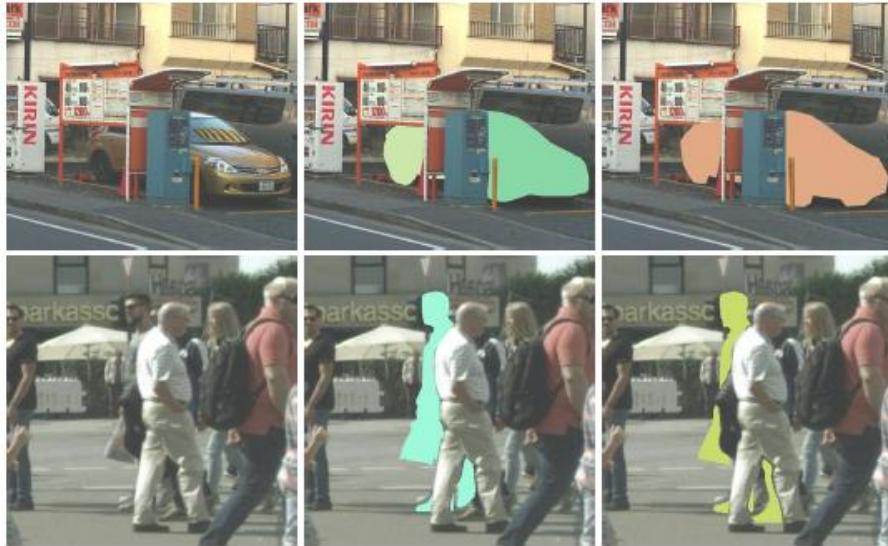


Figure 3: **Segmentation flaws.** Images are zoomed and cropped. Top row (Vistas image): both annotators identify the object as a car, however, one splits the car into two cars. Bottom row (Cityscapes image): the segmentation is genuinely ambiguous.

	PQ	$PQ^{St}$	$PQ^{Th}$	SQ	$SQ^{St}$	$SQ^{Th}$	RQ	$RQ^{St}$	$RQ^{Th}$
Cityscapes	69.7	71.3	67.4	84.2	84.4	83.9	82.1	83.4	80.2
ADE20k	67.1	70.3	65.9	85.8	85.5	85.9	78.0	82.4	76.4
Vistas	57.5	62.6	53.4	79.5	81.6	77.9	71.4	76.0	67.7

Table 1: **Human consistency for stuff vs. things.** Panoptic, segmentation, and recognition quality (PQ, SQ, RQ) averaged over classes ( $PQ = SQ \times RQ$  per class) are reported as percentages. Perhaps surprisingly, we find that **human consistency on each dataset** is relatively similar for both stuff and things.

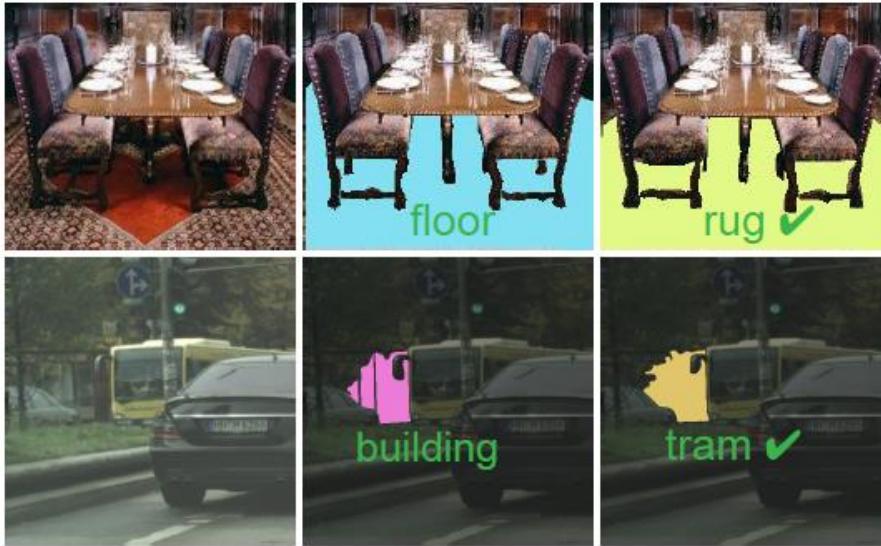


Figure 4: **Classification flaws.** Images are zoomed and cropped. Top row (ADE20k image): simple misclassification. Bottom row (Cityscapes image): the scene is extremely difficult, tram is the correct class for the segment. Many errors are difficult to resolve.

	$PQ^S$	$PQ^M$	$PQ^L$	$SQ^S$	$SQ^M$	$SQ^L$	$RQ^S$	$RQ^M$	$RQ^L$
Cityscapes	35.1	62.3	84.8	67.8	81.0	89.9	51.5	76.5	94.1
ADE20k	49.9	69.4	79.0	78.0	84.0	87.8	64.2	82.5	89.8
Vistas	35.6	47.7	69.4	70.1	76.6	83.1	51.5	62.3	82.6

Table 2: **Human consistency vs. scale**, for small (S), medium (M) and large (L) objects. **Scale plays a large role in determining human consistency for panoptic segmentation.** On large objects both SQ and RQ are above 80 on all datasets, while for small objects RQ drops precipitously. SQ for small objects is quite reasonable.

# Human Consistency

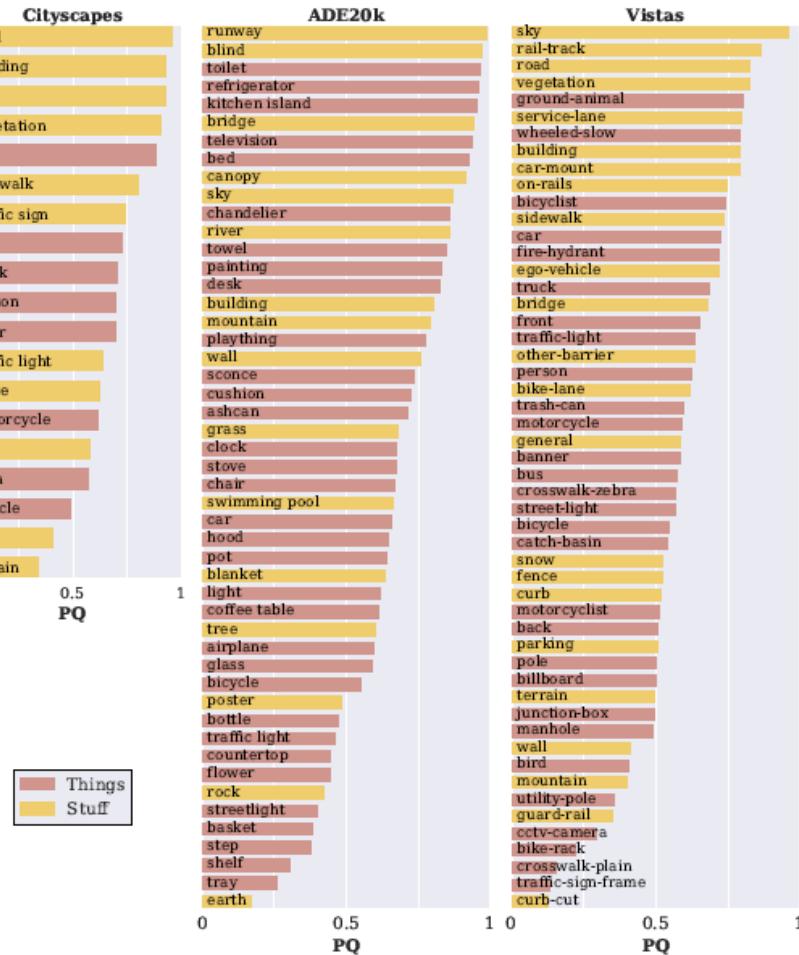


Figure 5: **Per-Class Human consistency, sorted by PQ.** Thing classes are shown in red, stuff classes in orange (for ADE20k every other class is shown, classes without matches in the dual-annotated test sets are omitted). Things and stuff are distributed fairly evenly, implying PQ balances their performance.

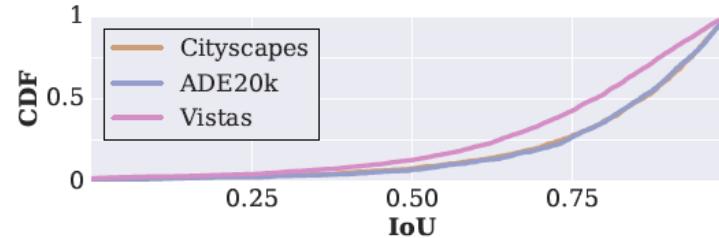


Figure 6: **Cumulative density functions of overlaps** for matched segments in three datasets when matches are computed by solving a maximum weighted bipartite matching problem [47]. After matching, less than 16% of matched objects have IoU below 0.5.

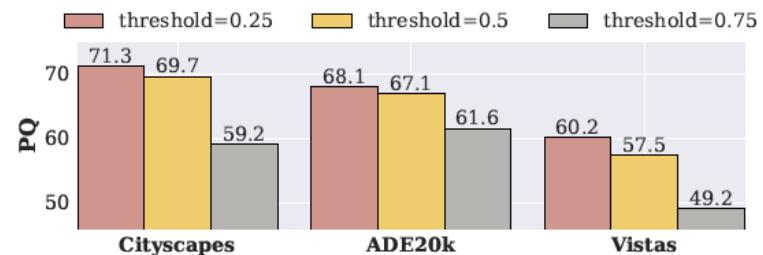


Figure 7: **Human consistency for different IoU thresholds.** The difference in PQ using a matching threshold of 0.25 vs. 0.5 is relatively small. For IoU of 0.25 matching is obtained by solving a maximum weighted bipartite matching problem. For a threshold greater than 0.5 the matching is unique and much easier to obtain.

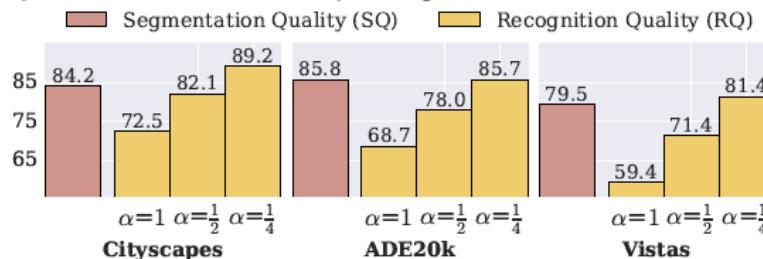


Figure 8: **SQ vs. RQ** for different  $\alpha$ , see (3). Lowering  $\alpha$  reduces the penalty of unmatched segments and thus increases the reported RQ (SQ is not affected). We use  $\alpha$  of 0.5 throughout but by tuning  $\alpha$  one can balance the influence of SQ and RQ in the final metric.

# Result

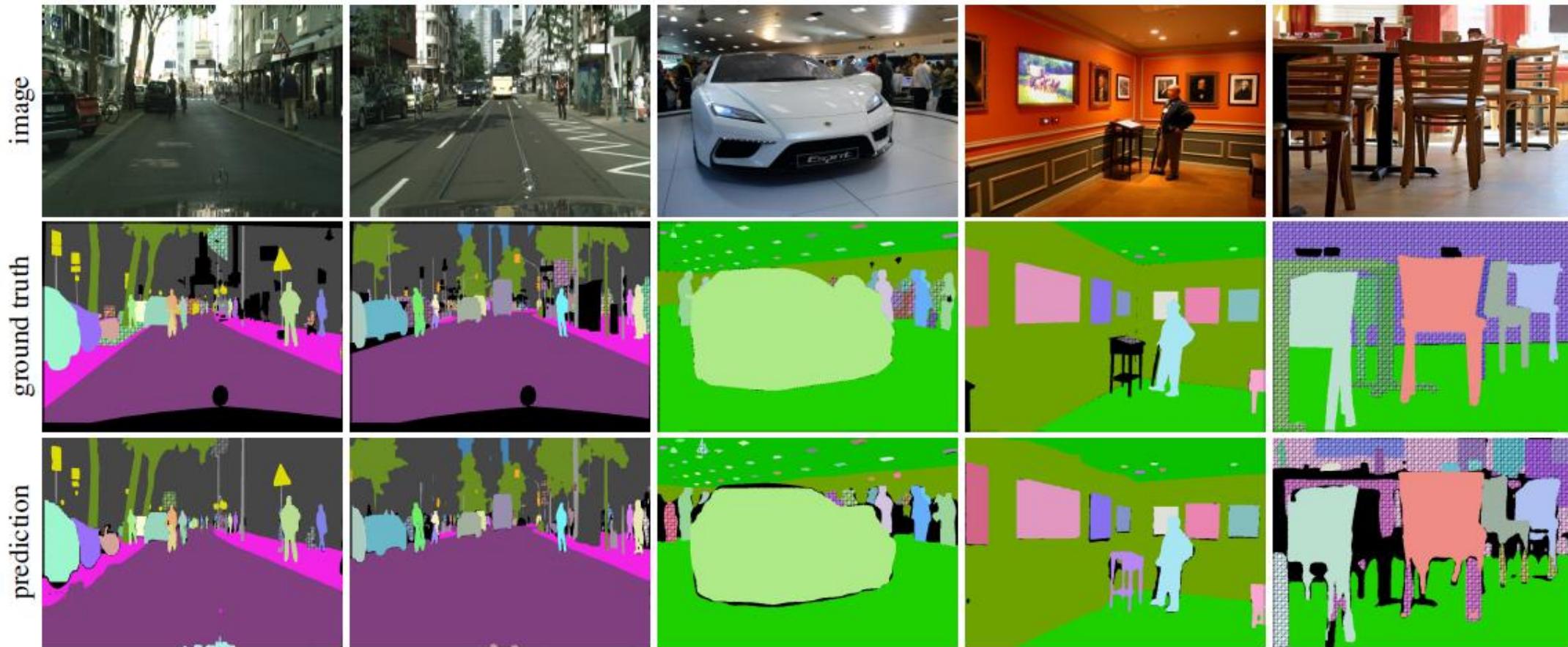


Figure 9: **Panoptic segmentation results** on Cityscapes (left two) and ADE20k (right three). Predictions are based on the merged outputs of state-of-the-art instance and semantic segmentation algorithms (see Tables 3 and 4). Colors for matched segments ( $\text{IoU} > 0.5$ ) match (crosshatch pattern indicates unmatched regions and black indicates unlabeled regions). Best viewed in color and with zoom.

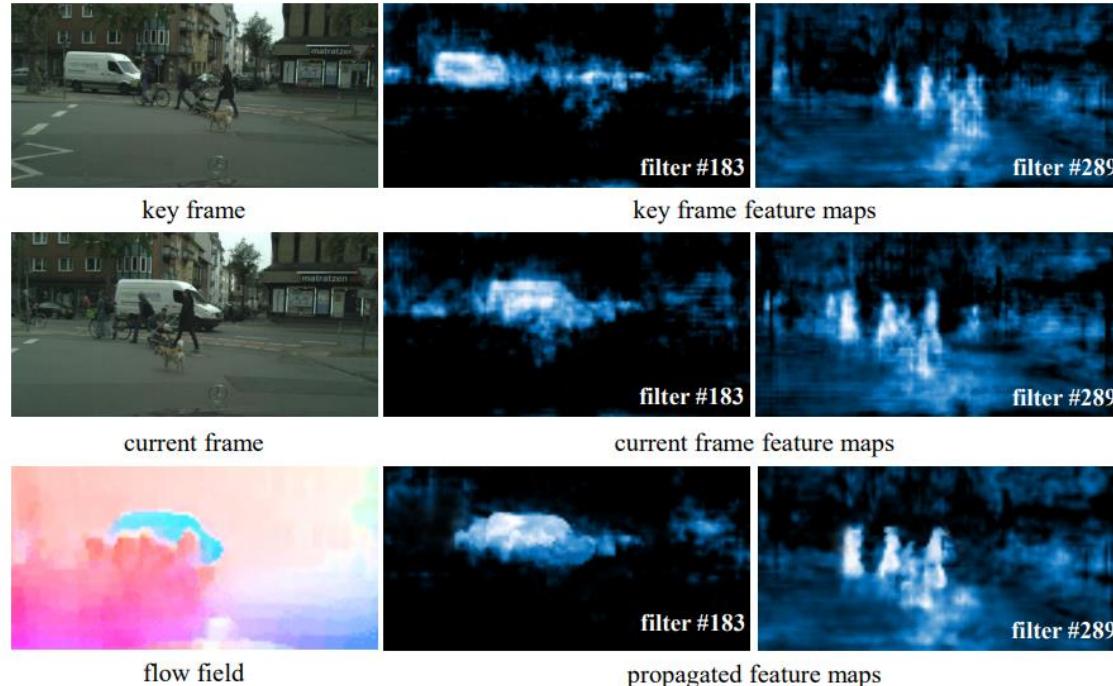
# **Video Object Detection**

## **: Deep Feature Flow for Video Recognition**

# Introduction

## “Fast and Accurate Video Recognition”

For Autonomous Driving, Video Surveillance



- **Flow estimation**  
: Apply image recognition network on sparse key frames
- **Feature propagation**  
Propagates the deep feature maps from key frames to other frames via a flow field

⇒ **fast, accurate, general, end-to-end framework**

# Deep Feature Flow

## Deep Feature Flow Inference

- Bilinear-interpolation  
:  $f_i^c(p) = \sum_q G(q, p + \delta p) f_k^c(q)$
- $G(q, p + \delta p) = g(q_x, p_x + \delta p_x) \cdot g(q_y, p_y + \delta p_y)$   
(where  $g(a, b) = \max(0, 1 - |a - b|)$ )
- Scale field :  $S_{i \rightarrow k} = S(I_k, I_i)$
- Feature propagation  
:  $f_i = \mathcal{W}(f_k, M_{i \rightarrow k}, S_{i \rightarrow k})$

## Feature Propagation

$$\frac{\partial f_i^c(p)}{\partial M_{i \rightarrow k}(p)} = S_{i \rightarrow k}^c(p) \sum_q \frac{\partial G(q, p + \delta p)}{\partial \delta p} f_k^c(q) \rightarrow \text{parameter free \& fully differentiable}$$

**Algorithm 1** Deep feature flow inference algorithm for video recognition.

---

```

1: input: video frames  $\{\mathbf{I}_i\}$ 
2:  $k = 0$ ; ▷ initialize key frame
3:  $\mathbf{f}_0 = \mathcal{N}_{feat}(\mathbf{I}_0)$ 
4:  $\mathbf{y}_0 = \mathcal{N}_{task}(\mathbf{f}_0)$ 
5: for  $i = 1$  to  $\infty$  do
6:   if is_key_frame( $i$ ) then ▷ key frame scheduler
7:      $k = i$  ▷ update the key frame
8:      $\mathbf{f}_k = \mathcal{N}_{feat}(\mathbf{I}_k)$ 
9:      $\mathbf{y}_k = \mathcal{N}_{task}(\mathbf{f}_k)$ 
10:    else ▷ use feature flow
11:       $\mathbf{f}_i = \mathcal{W}(\mathbf{f}_k, \mathcal{F}(\mathbf{I}_k, \mathbf{I}_i), \mathcal{S}(\mathbf{I}_k, \mathbf{I}_i))$  ▷ propagation
12:       $\mathbf{y}_i = \mathcal{N}_{task}(\mathbf{f}_i)$ 
13:    end if
14:  end for
15: output: recognition results  $\{\mathbf{y}_i\}$ 

```

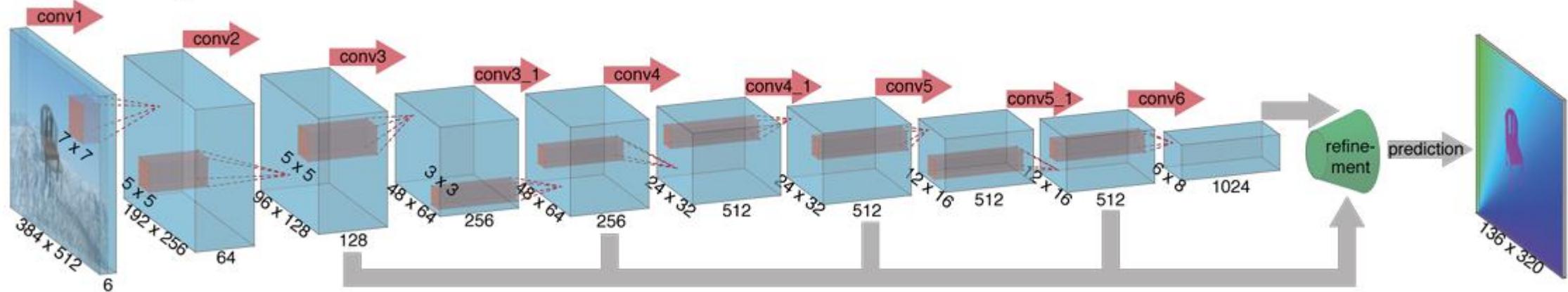
---

$k$	key frame index
$i$	current frame index
$r$	per-frame computation cost ratio, Eq. (5)
$l$	key frame duration length
$s$	overall speedup ratio, Eq. (7)
$\mathbf{I}_i, \mathbf{I}_k$	video frames
$\mathbf{y}_i, \mathbf{y}_k$	recognition results
$\mathbf{f}_k$	convolutional feature maps on key frame
$\mathbf{f}_i$	propagated feature maps on current frame
$\mathbf{M}_{i \rightarrow k}$	2D flow field
$\mathbf{p}, \mathbf{q}$	2D location
$S_{i \rightarrow k}$	scale field
$\mathcal{N}$	image recognition network
$\mathcal{N}_{feat}$	sub-network for feature extraction
$\mathcal{N}_{task}$	sub-network for recognition result
$\mathcal{F}$	flow estimation function
$\mathcal{W}$	feature propagation function, Eq. (3)

Table 1. Notations.

# FlowNet

FlowNetSimple



① FlowNet Half

: reduce convolutional kernels in each layer  $\frac{1}{2}$  and complexity to  $\frac{1}{4}$

② FlowNet Inception

: Adopt Inception structure and reduces complexity to  $\frac{1}{8}$



Pre-trained on synthetic flying chair dataset

# Structure

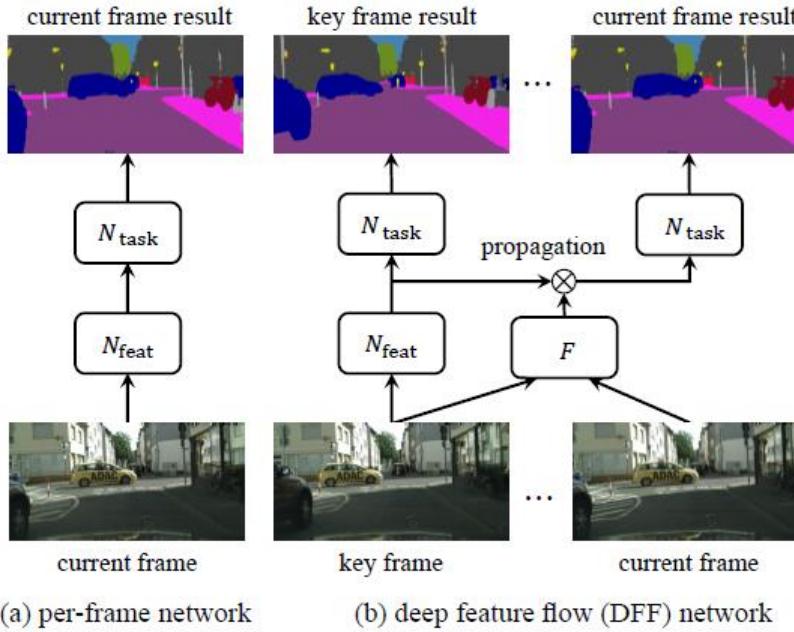
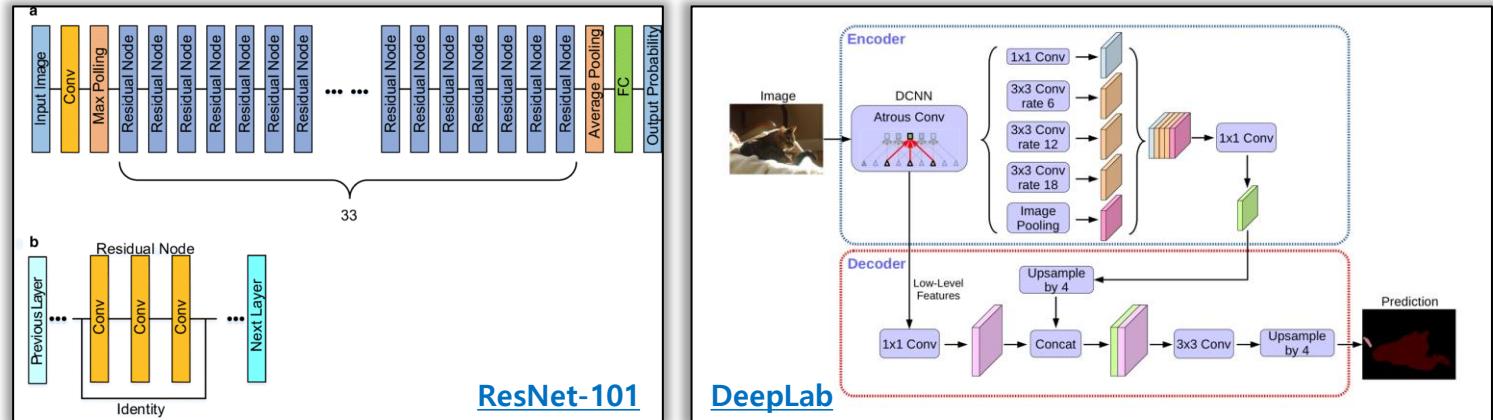


Figure 2. Illustration of video recognition using per-frame network evaluation (a) and the proposed deep feature flow (b).

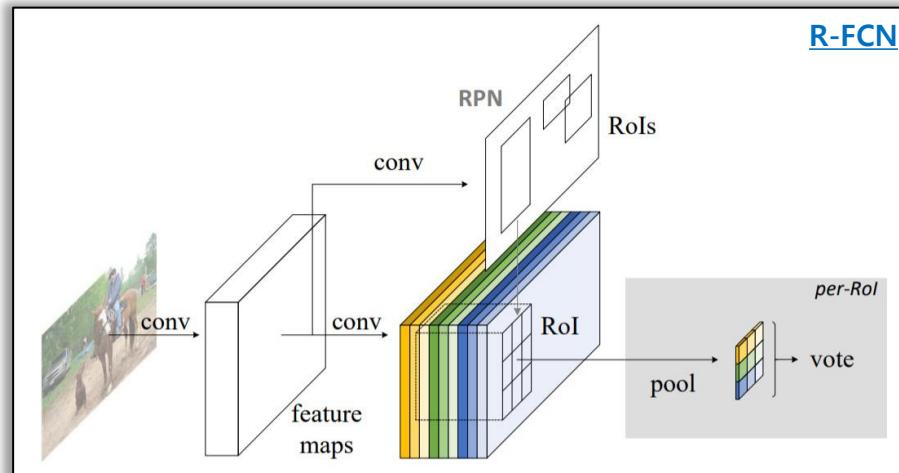
layer	type	stride	# output	Inception/Reduction			
				#1x1	#1x1#3x3	#1x1#3x3#3x3	#pool
conv1	7x7 conv	2	32				
pool1	3x3 max pool	2	32				
conv2	Inception		64	24-32	24-32-32		
conv3_1	3x3 conv	2	128				
conv3_2	Inception		128	48	32-64	8-16-16	
conv3_3	Inception		128	48	48-64	12-16-16	
conv4_1	Reduction	2	256	32	112-128	28-32-32	64
conv4_2	Inception		256	96	112-128	28-32-32	
conv5_1	Reduction	2	384	48	96-192	36-48-48	96
conv5_2	Inception		384	144	96-192	36-48-48	
conv6_1	Reduction	2	512	64	192-256	48-64-64	128
conv6_2	Inception		512	192	192-256	48-64-64	



ResNet-50 or ResNet-101 pre-trained for ImageNet

- Semantic Segmentation : DeepLab
- Object Detection : R-FCN

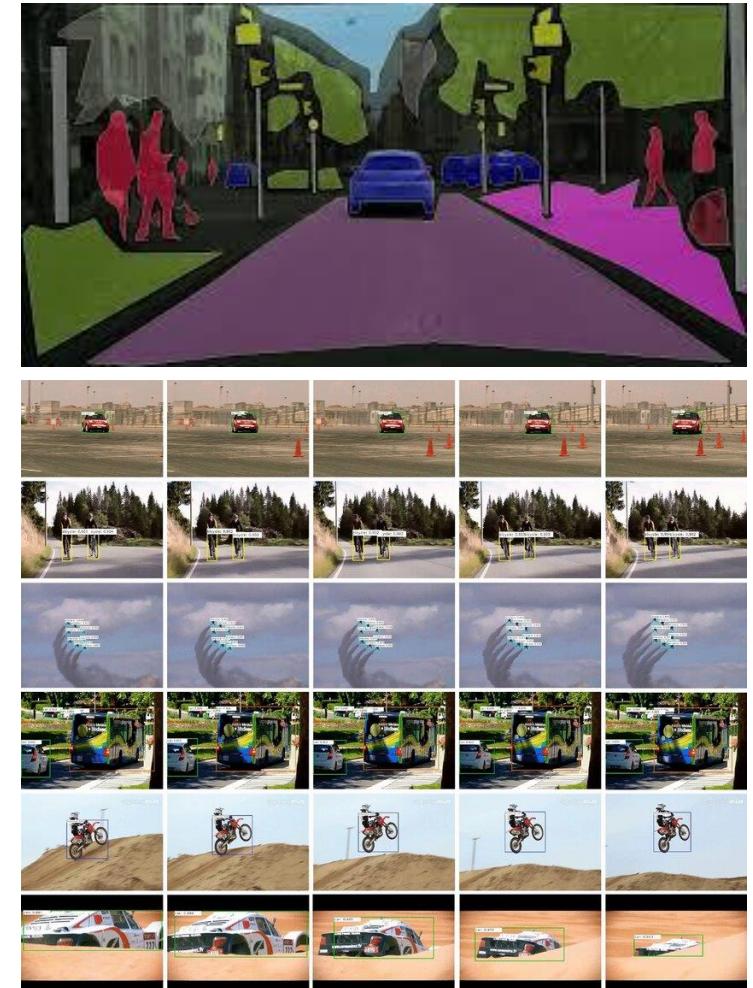
	FlowNet	FlowNet Half	FlowNet Inception
ResNet-50	9.20	33.56	68.97
ResNet-101	12.71	46.30	95.24



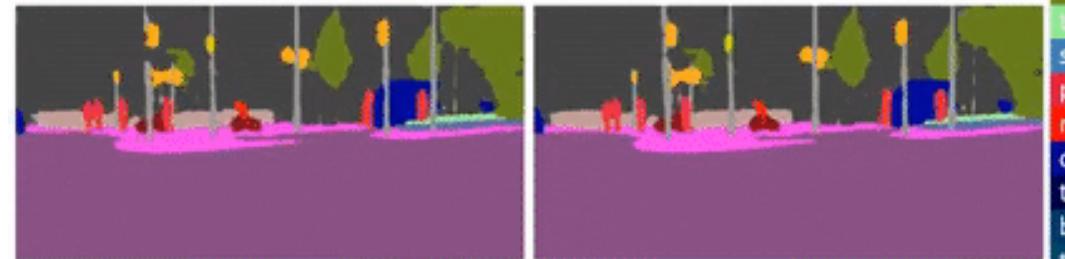
# Experiment

## Datasets

- Cityscapes
  - For urban scene understanding and autonomous driving
  - 2975,500 dataset, and 1525 snippets respectively
  - Feature network : 1024 pixels / Flow network : 512 pixels
- ImageNet VID
  - For object detection in Videos
  - 3,862,555 dataset, and 937 fully-annotated video snippets respectively
  - Feature network : 600 pixels / Flow network : 300 pixels



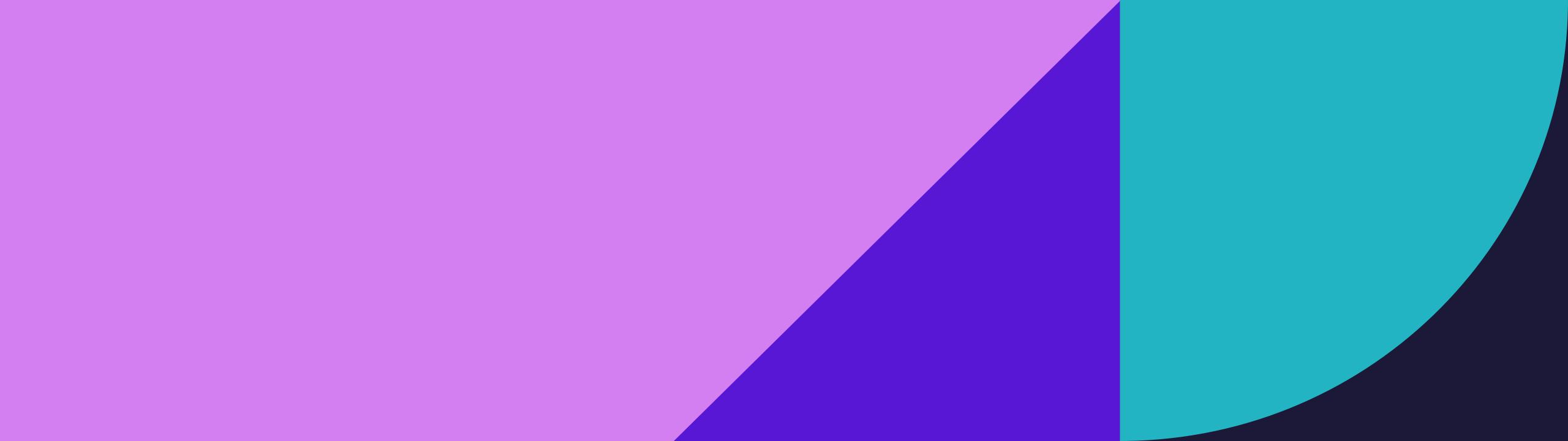
# Result



road
sidewalk
building
wall
fence
pole
traffic light
traffic sign
vegetation
terrain
sky
person
rider
car
truck
bus
train
motorcycle
bicycle

# Reference

- Kaiming He, Georgia Gkioxari, Piotr Dollar, Ross Girshick. 「**Mask R-CNN**」, 2017
- Jonathan Long, Evan Shelhamer, Trevor Darrell. 「**Fully Convolutional Networks for Semantic Segmentation**」, 2014
- Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, Piotr Dollar. 「**Panoptic Segmentation**」, 2018
- Jifeng Dai, Kaiming He, Jain Sun. 「**Instance-aware Semantic Segmentation via Multi-task Network Cascades**」, 2016
- Jifeng Dai, Yi Li, Haozhi Qi, Xiangyang Ji, Yichen Wei. 「**Fully Convolutional Instance-aware Semantic Segmentation**」, 2017
- Matthew D.Zeiler, Rob Fergus. 「**Visualizing and Understanding Convolutional Networks**」, 2013
- Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, Yichen Wei. 「**Deep Feature Flow for Video Recognition**」, 2017
- Philipp Fischer, Alexey Dosovitskiy, Eddy Ilg, Philip Häusser, Caner Hazirbaş, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, Thomas Brox. 「**FlowNet: Learning Optical Flow with Convolutional Networks**」, 2015
- Hyeonwoo Noh, Seunghoon Hong, Bohyung Han. 「**Learning Deconvolution Network for Semantic Segmentation**」, 2015
- [Stanford. CS231n](#)
- <https://www.kaggle.com/datasets/faizalkarim/flood-area-segmentation>
- <https://image-net.org/challenges/LSVRC/2017/>



THANK YOU