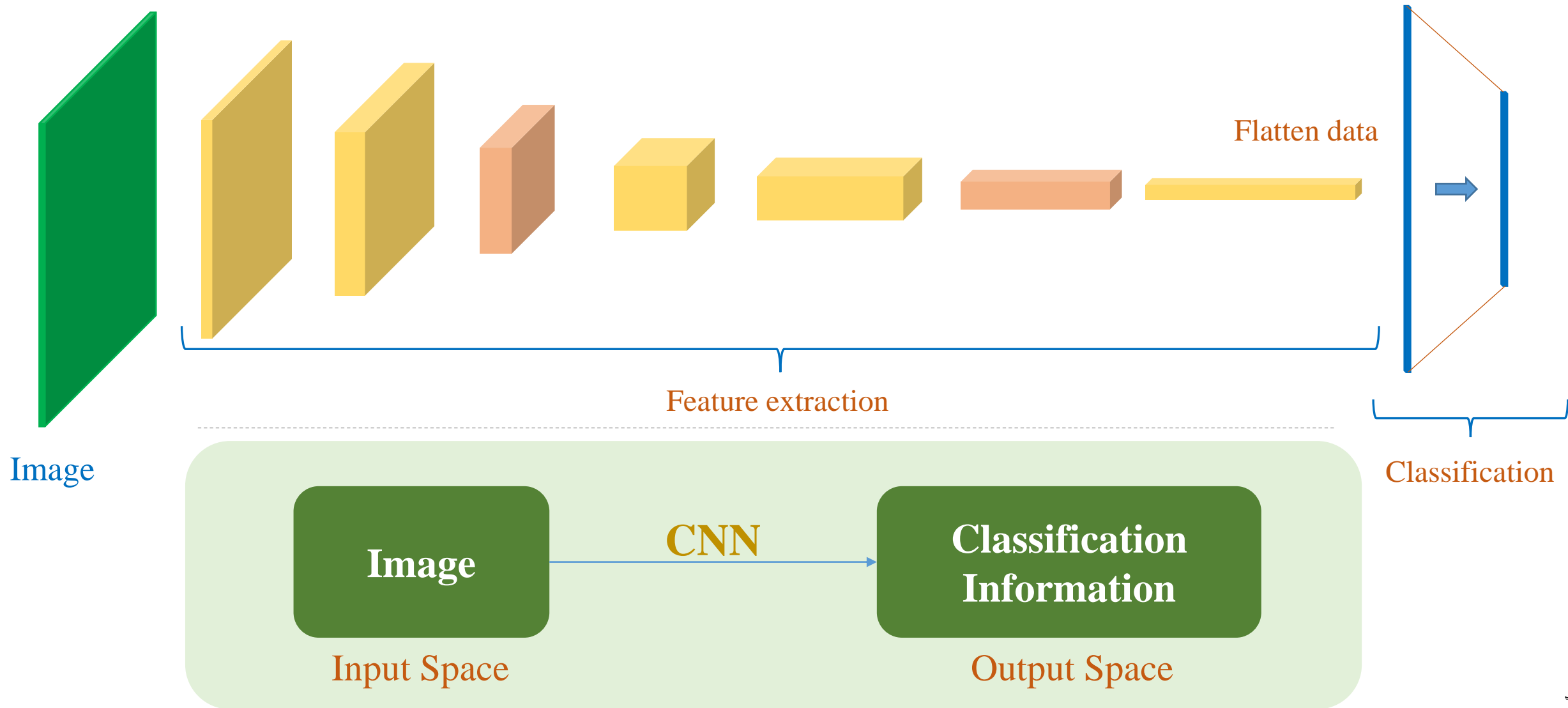


Image Domain Conversion (Model Construction)

Quang-Vinh Dinh
Ph.D. in Computer Science

Motivation



Motivation

Image
Classification

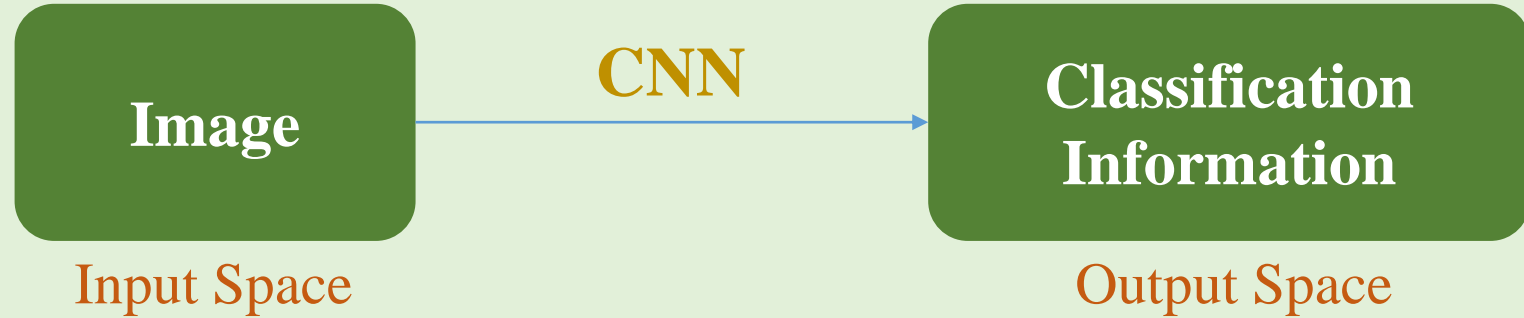


Image
Localization

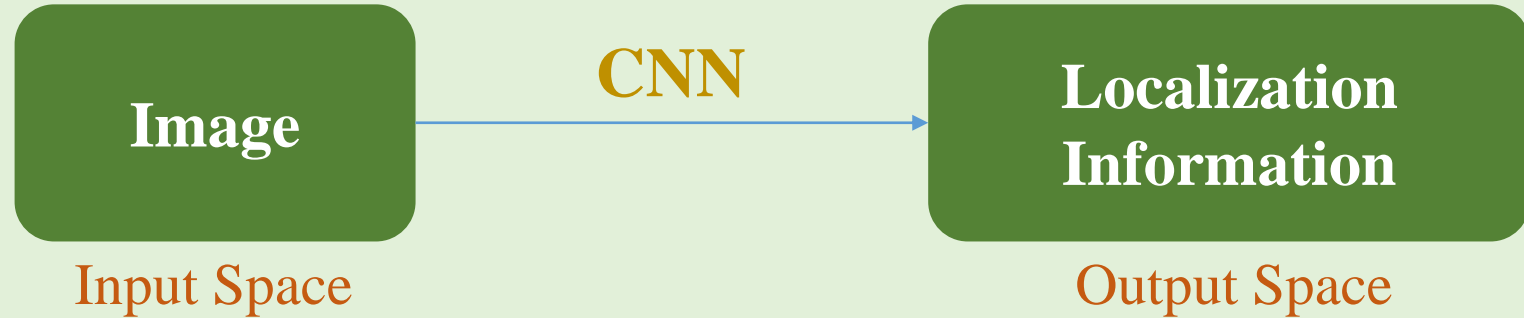
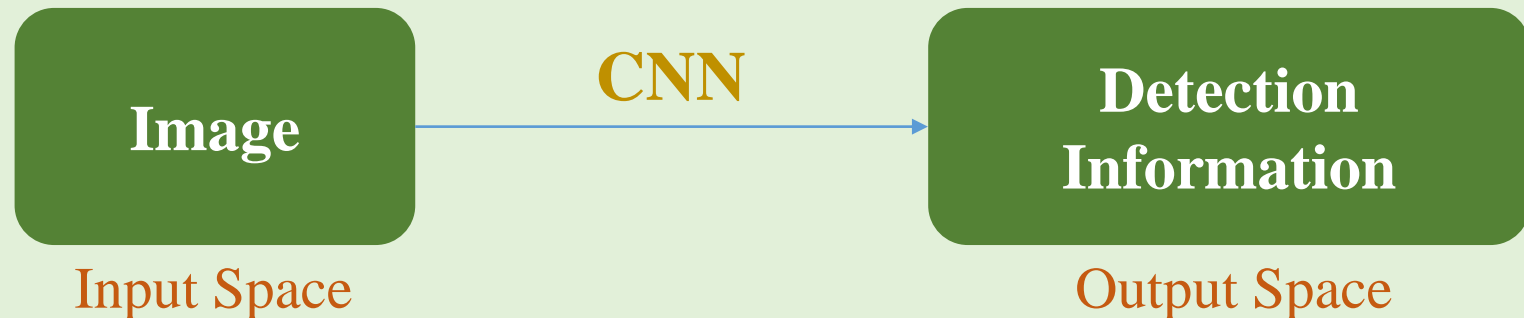
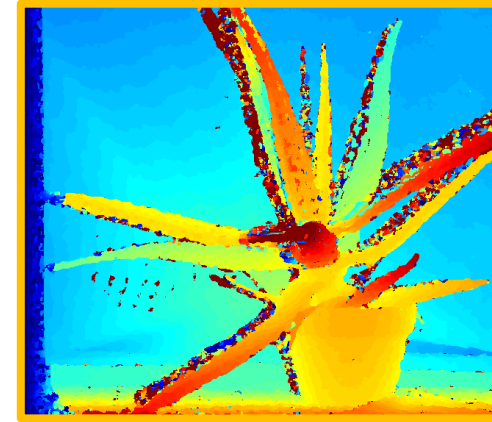


Image
Detection



Motivation



Image

Input Space

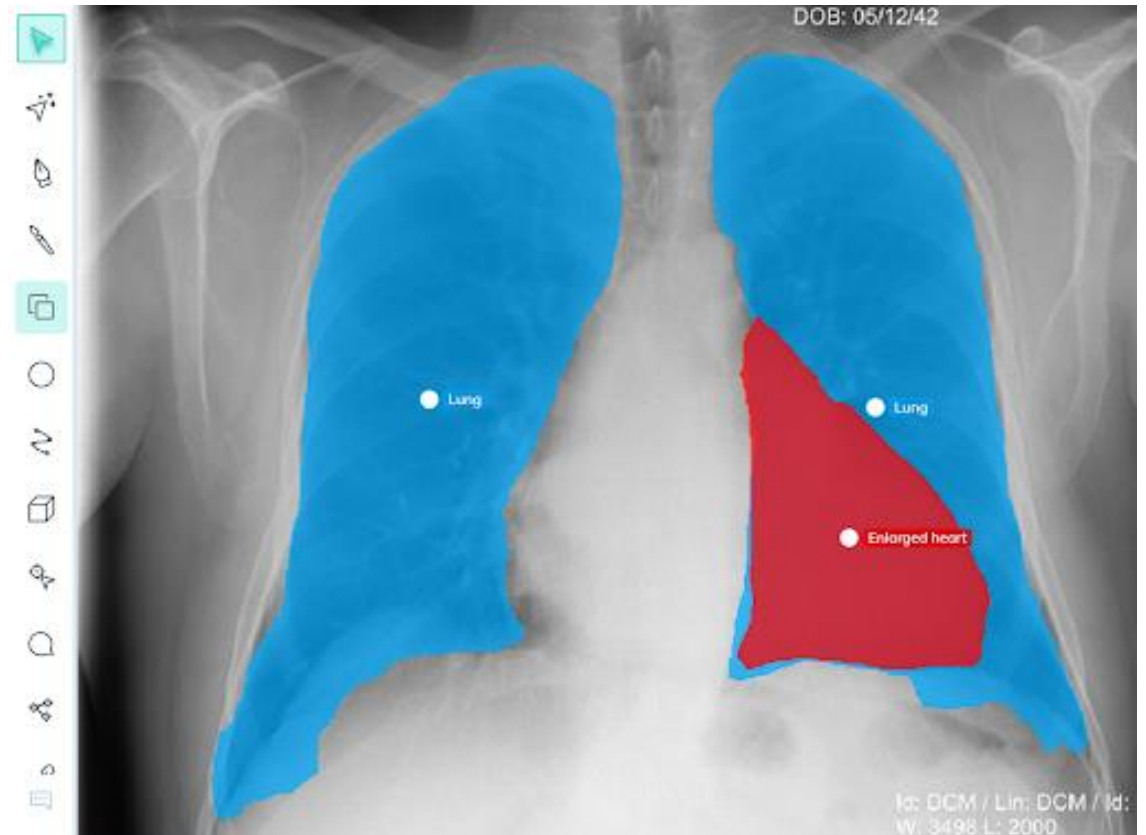
How?

Image

Output Space



Image Segmentation



Segmentation

Traditional method



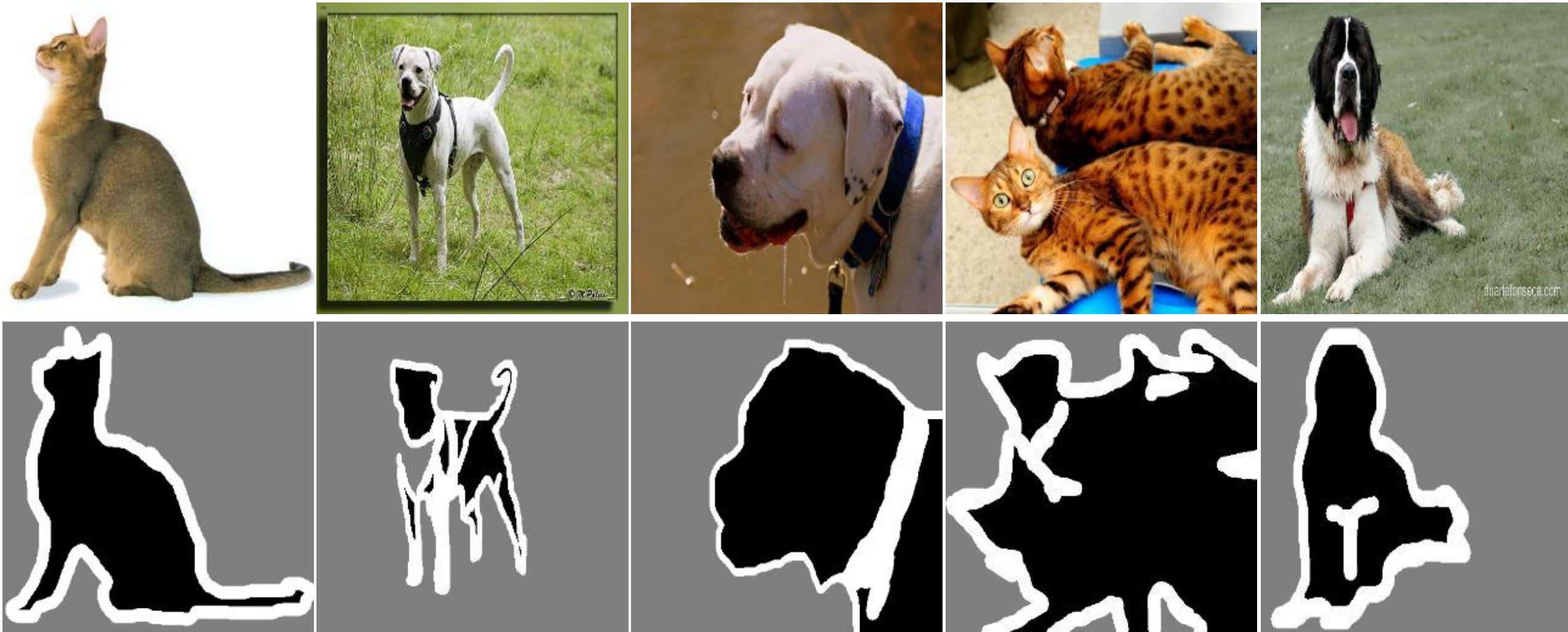
Original Image

K-Mean Segmentation

Segmentation

Oxford-IIT-Pet

<https://www.robots.ox.ac.uk/~vgg/data/pets/>

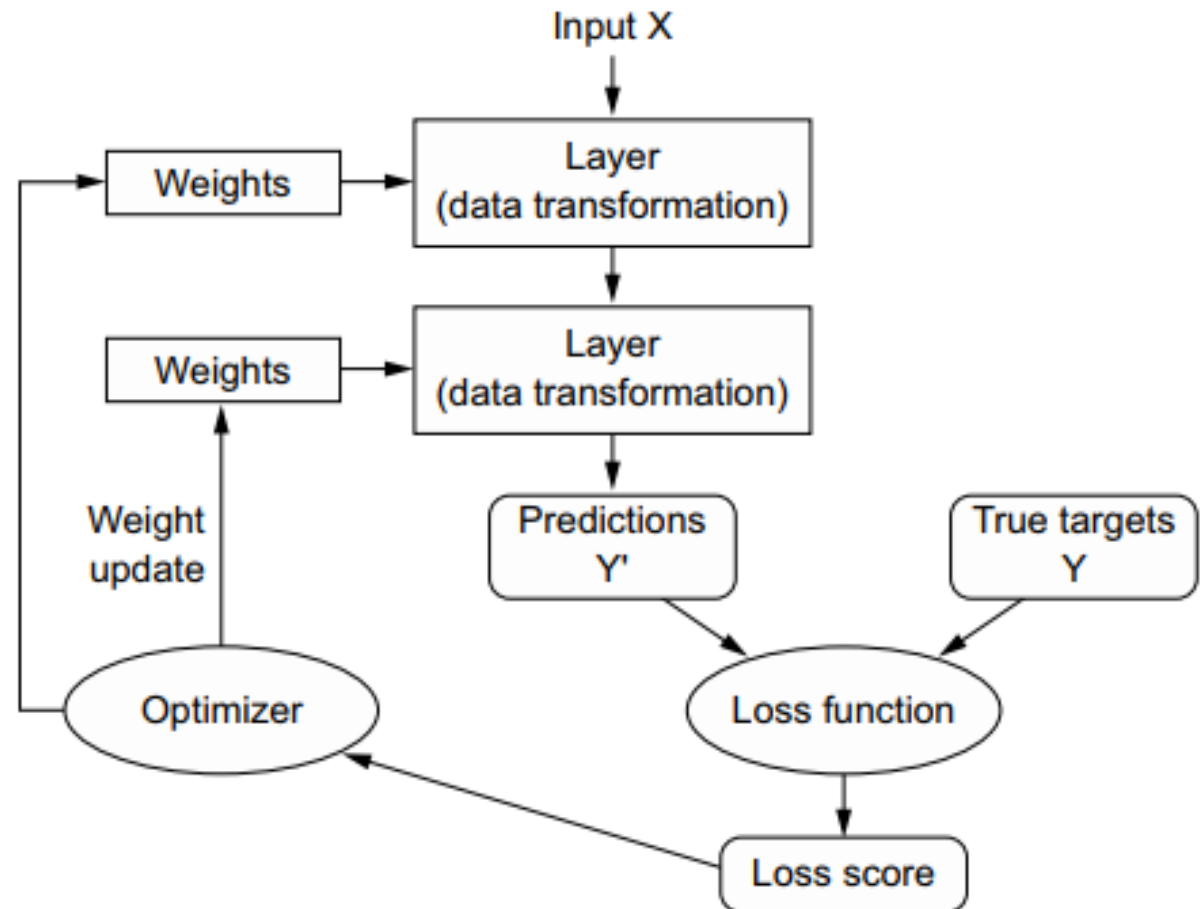


Discussion

Input/Output + Loss Function

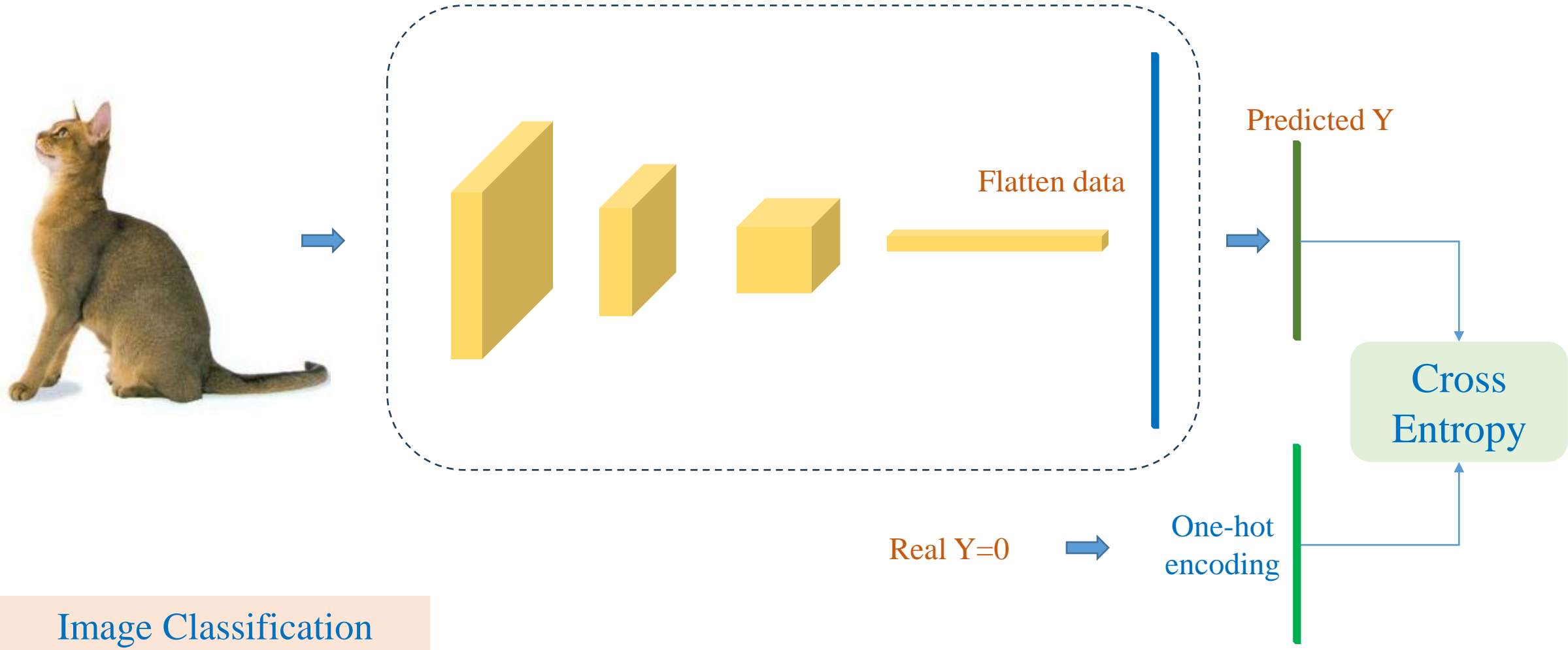
Image Classification

Temperature Forecasting



Discussion

Input/Output + Loss Function



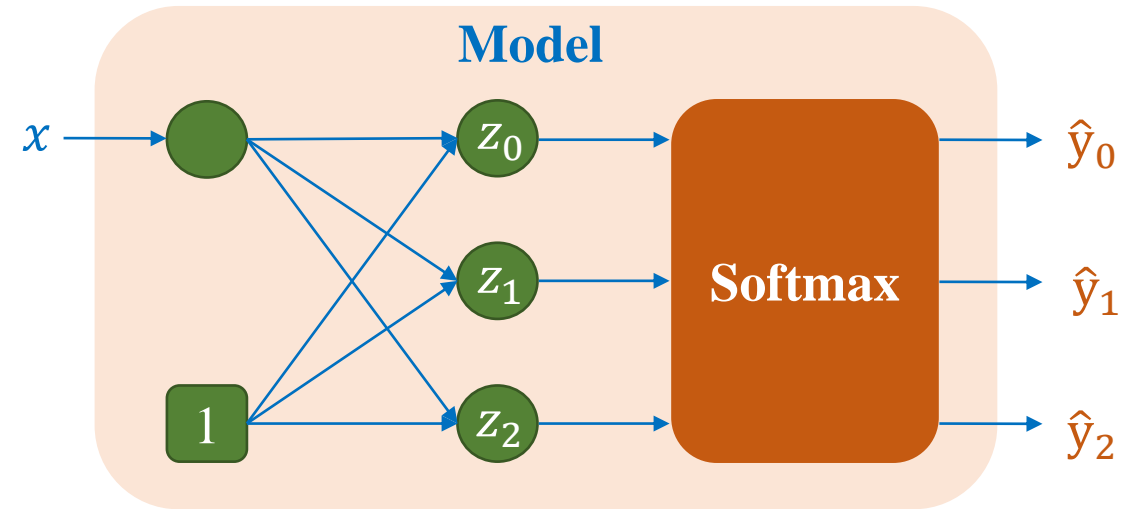
One-Hot Encoding and Cross-entropy

Feature	Label
Petal_Length	Label
1.4	0
1.3	0
1.5	0
4.5	1
4.1	1
4.6	1
5.2	2
5.6	2
5.9	2

#features = 1

#classes = 3

$y \in \{0,1,2\}$



One-hot encoding for label

$$\mathbf{y} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix} \quad y_i \in \{0,1\} \quad \sum_i y_i = 1$$

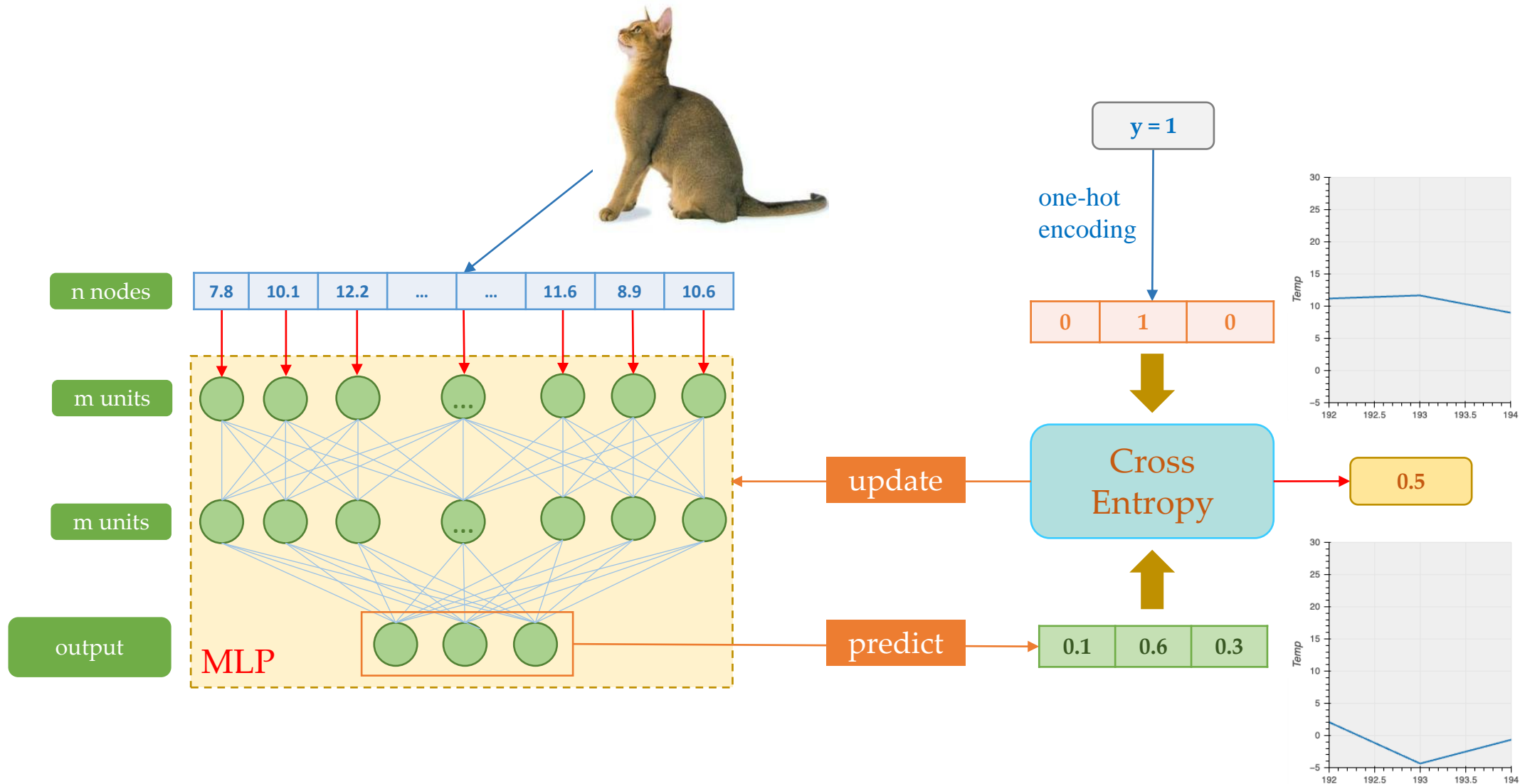
$$y = 0 \rightarrow \mathbf{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad y = 1 \rightarrow \mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad y = 2 \rightarrow \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Loss function

$$\begin{aligned} L(\boldsymbol{\theta}) &= -y_2 \log(\hat{y}_2) - y_1 \log(\hat{y}_1) - y_0 \log(\hat{y}_0) \\ &= -\sum_i y_i \log(\hat{y}_i) \end{aligned}$$

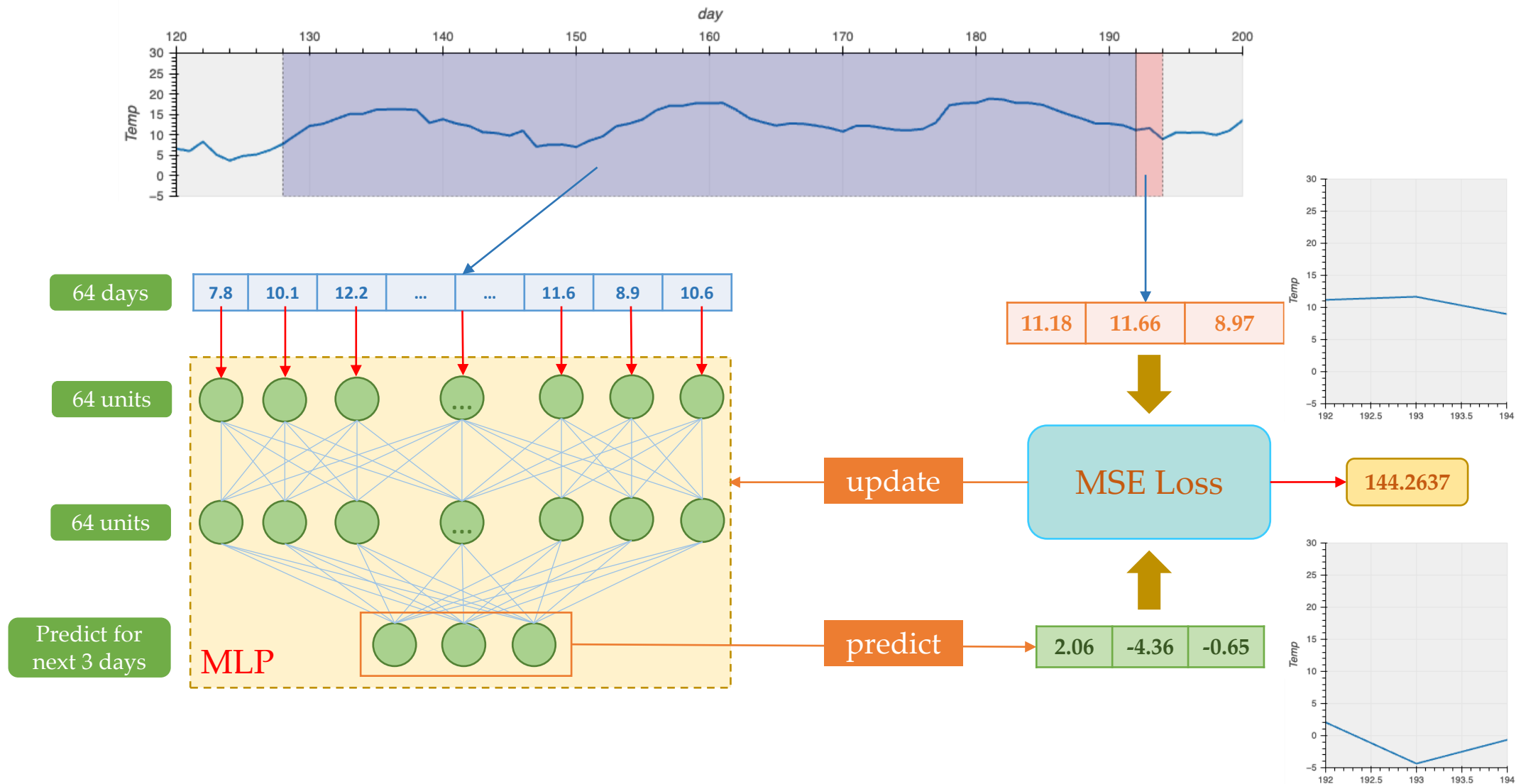
Discussion – Classification Problem

Input/Output + Loss Function



Discussion – Regression Problem

Input/Output + Loss Function



Discussion – Segmentation Problem

Input/Output + Loss Function



Model



Inference

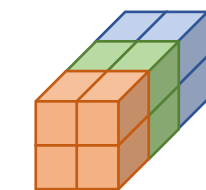
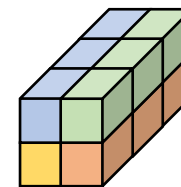


image
(2x2)



Model



Output
(3x2x2)



0	2
2	1

Predicted
Segmentation

Training

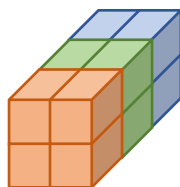
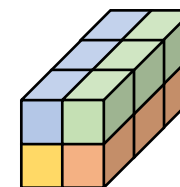


image
(2x2)



Model



Output
(3x2x2)

1	0
2	0

Ground-truth

one-hot encoding



0	1	0
---	---	---

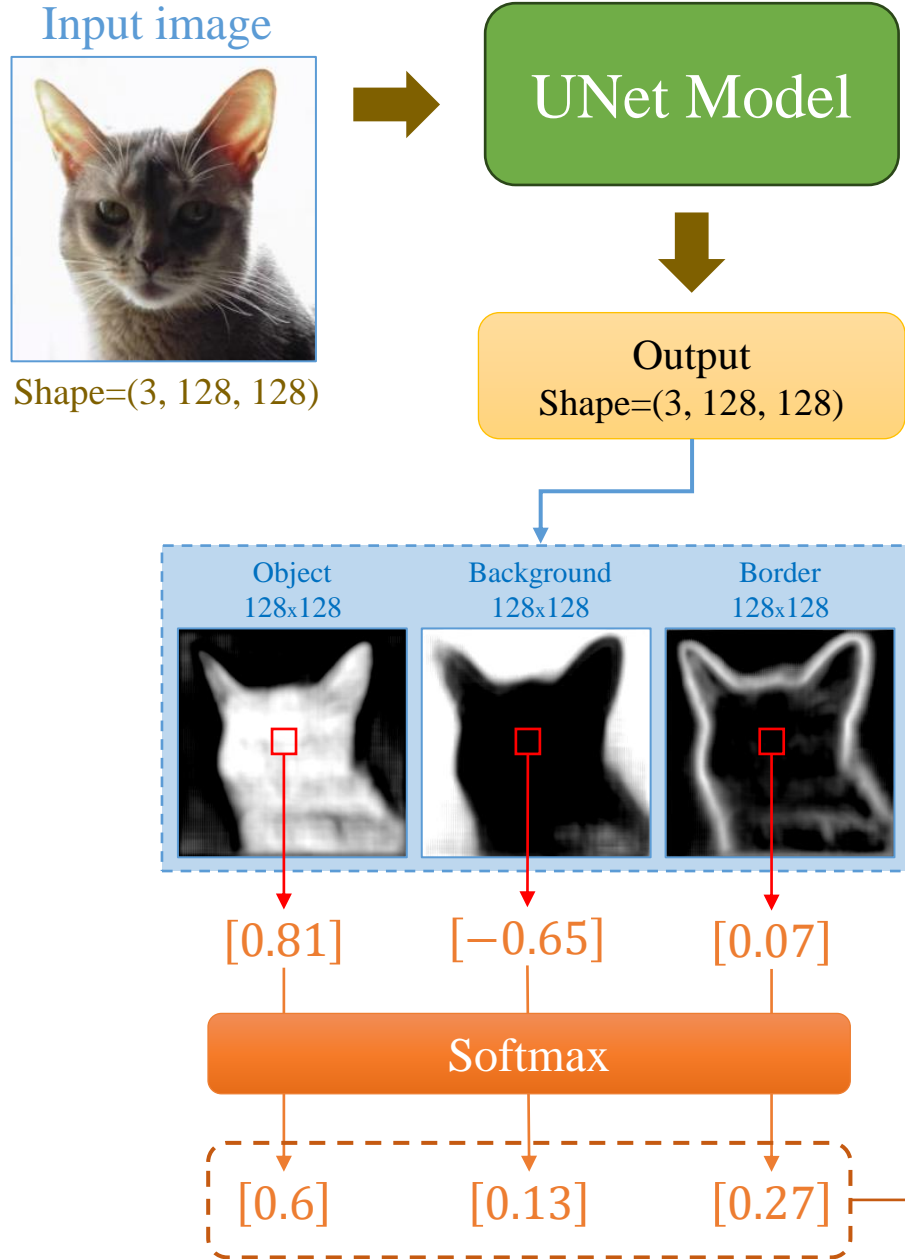
1	0	0
---	---	---

0	0	1
---	---	---

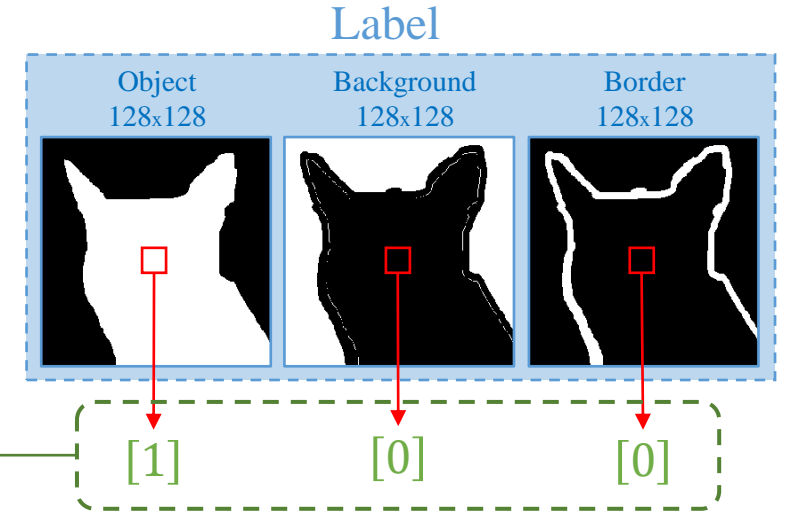
1	0	0
---	---	---

Cross
Entropy

Loss Function



$$CE = - \sum_i y_i \log(\hat{y}_i)$$



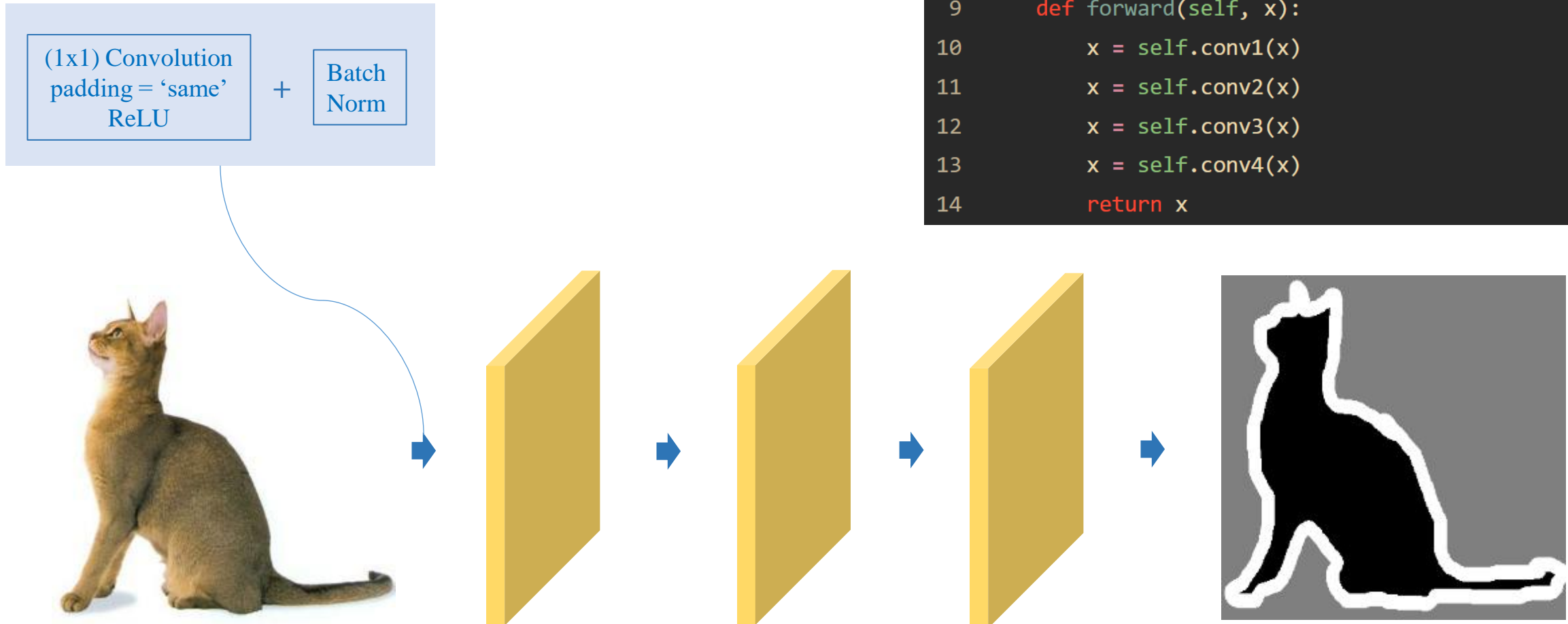
$$CE \left(\begin{bmatrix} 0.60 \\ 0.13 \\ 0.27 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right) = 0.5$$

$$\text{Total_loss} = \frac{1}{H * W} \sum_{j=1}^H \sum_{k=1}^W CE(\hat{y}_{j,k}, y_{j,k})$$



Segmentation (1)

Pixel-wise and color-based method



```
1 class CNN_Segmentation(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.conv1 = ConvBlock(3, 64, kernel_size=1)
5         self.conv2 = ConvBlock(64, 64, kernel_size=1)
6         self.conv3 = ConvBlock(64, 64, kernel_size=1)
7         self.conv4 = nn.Conv2d(64, 3, kernel_size=1)
8
9     def forward(self, x):
10        x = self.conv1(x)
11        x = self.conv2(x)
12        x = self.conv3(x)
13        x = self.conv4(x)
14        return x
```

Segmentation (1)

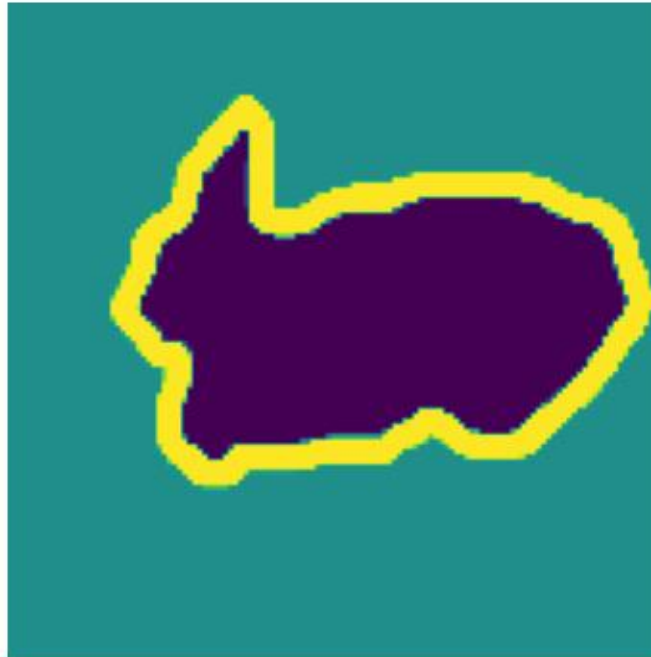
Result

```
1 for epoch in range(max_epoch):
2     model.train()
3     for inputs, targets in train_loader:
4         inputs, targets = inputs.to(device), targets.to(device)
5
6         optimizer.zero_grad()
7         outputs = model(inputs)
8         loss = criterion(outputs, targets)
9
10        loss.backward()
11        optimizer.step()
```

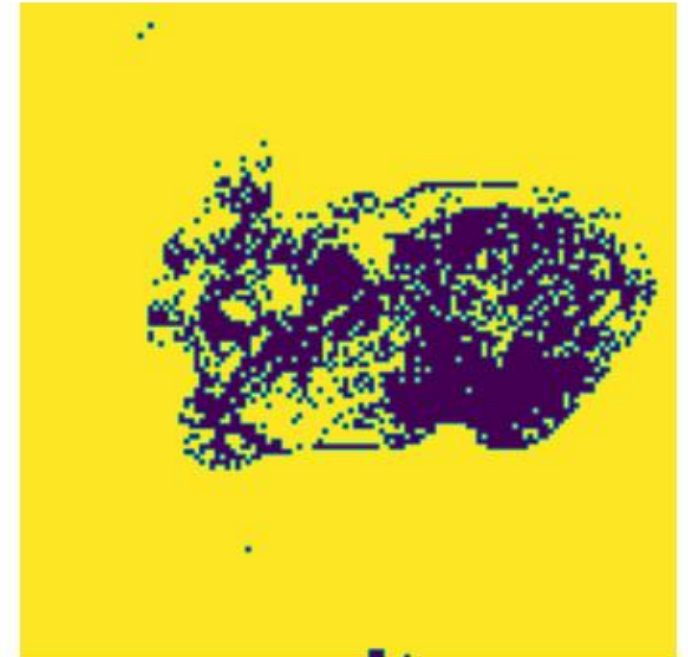
Input Image



True Mask

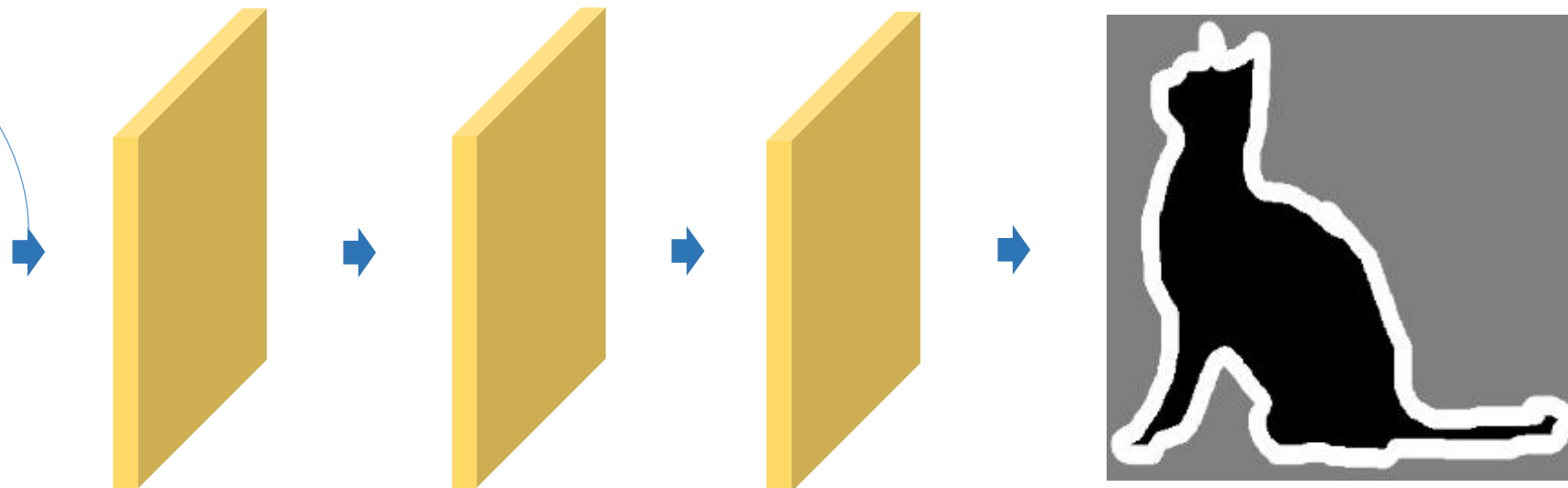
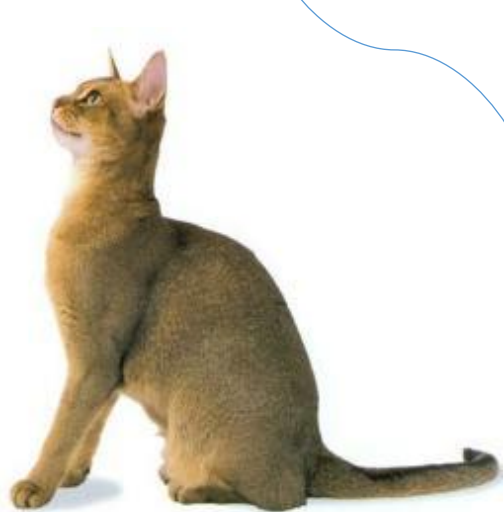
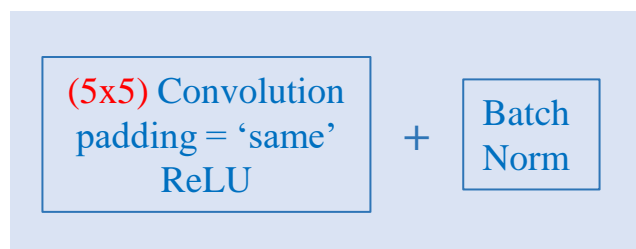


Predicted Mask



Segmentation (2)

Increase kernel sizes



```
1 class CNN_Segmentation(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.conv1 = ConvBlock(3, 64, kernel_size=5)
5         self.conv2 = ConvBlock(64, 64, kernel_size=5)
6         self.conv3 = ConvBlock(64, 64, kernel_size=5)
7         self.conv4 = nn.Conv2d(64, 3, kernel_size=5)
8
9     def forward(self, x):
10         x = self.conv1(x)
11         x = self.conv2(x)
12         x = self.conv3(x)
13         x = self.conv4(x)
14         return x
```

Segmentation (2)

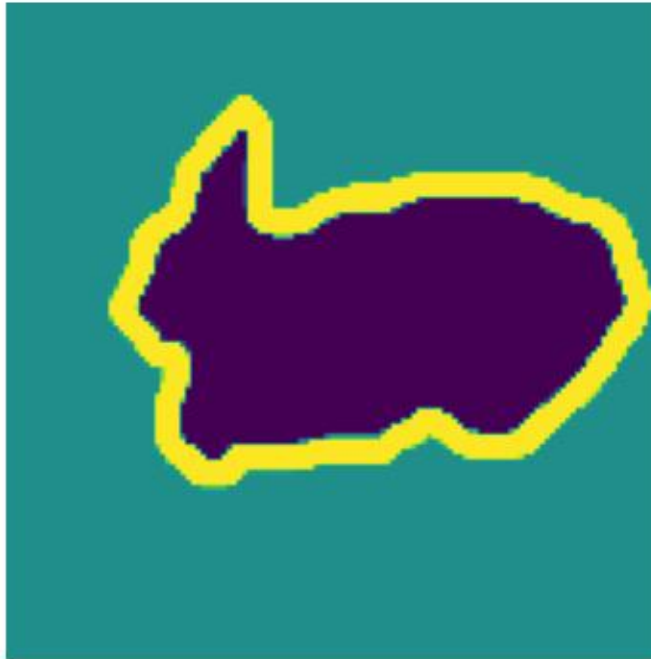
Result

```
1 for epoch in range(max_epoch):
2     model.train()
3     for inputs, targets in train_loader:
4         inputs, targets = inputs.to(device), targets.to(device)
5
6         optimizer.zero_grad()
7         outputs = model(inputs)
8         loss = criterion(outputs, targets)
9
10        loss.backward()
11        optimizer.step()
```

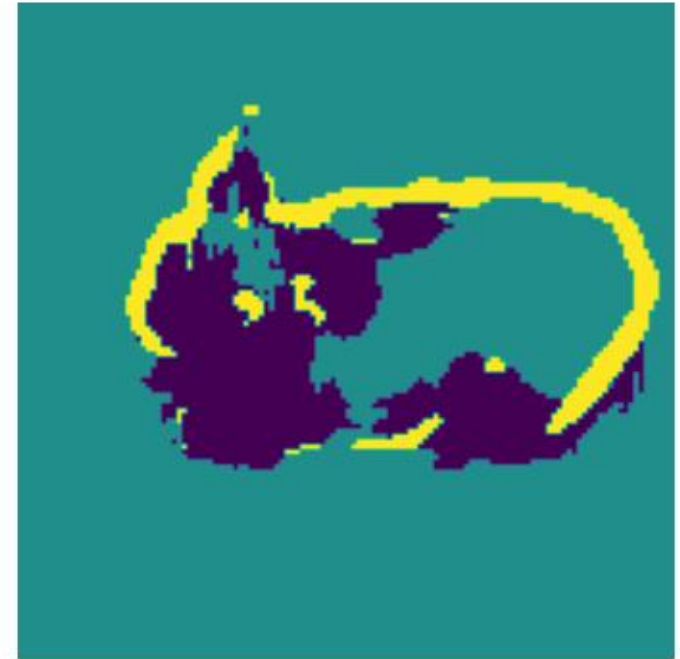
Input Image



True Mask

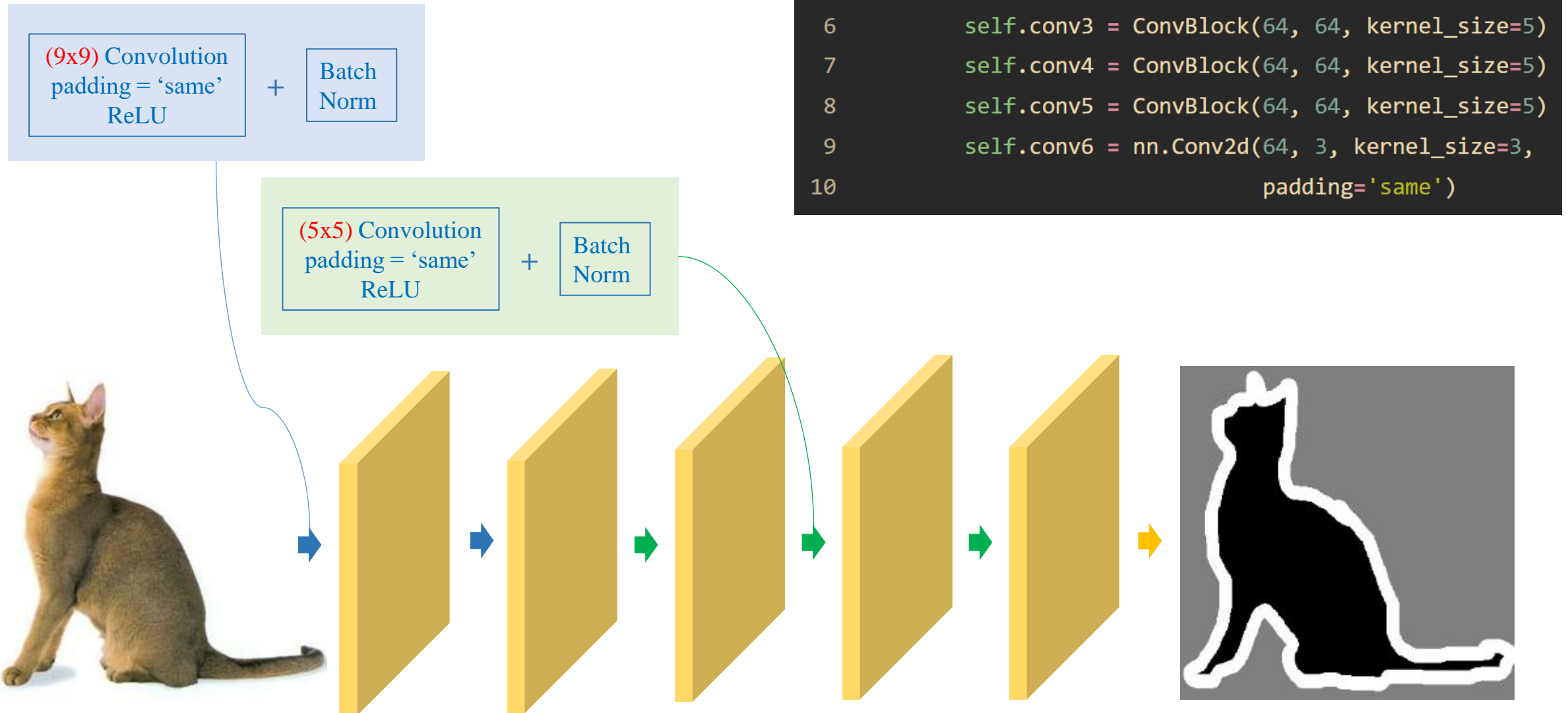


Predicted Mask



Segmentation (3)

Keep increasing



Segmentation (3)

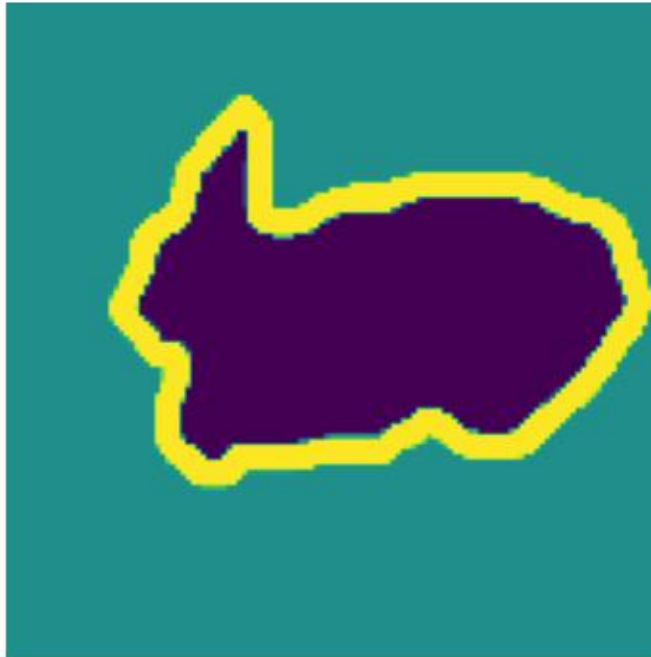
Result

```
1 for epoch in range(max_epoch):
2     model.train()
3     for inputs, targets in train_loader:
4         inputs, targets = inputs.to(device), targets.to(device)
5
6         optimizer.zero_grad()
7         outputs = model(inputs)
8         loss = criterion(outputs, targets)
9
10        loss.backward()
11        optimizer.step()
```

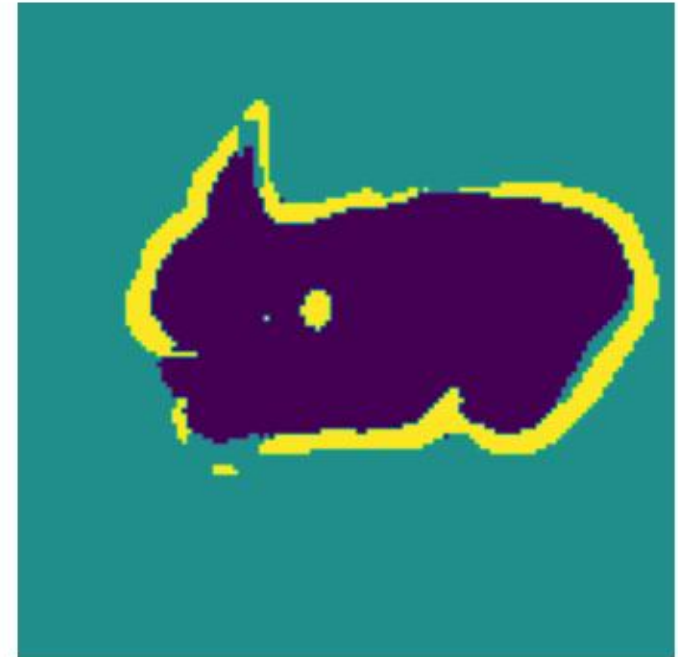
Input Image



True Mask

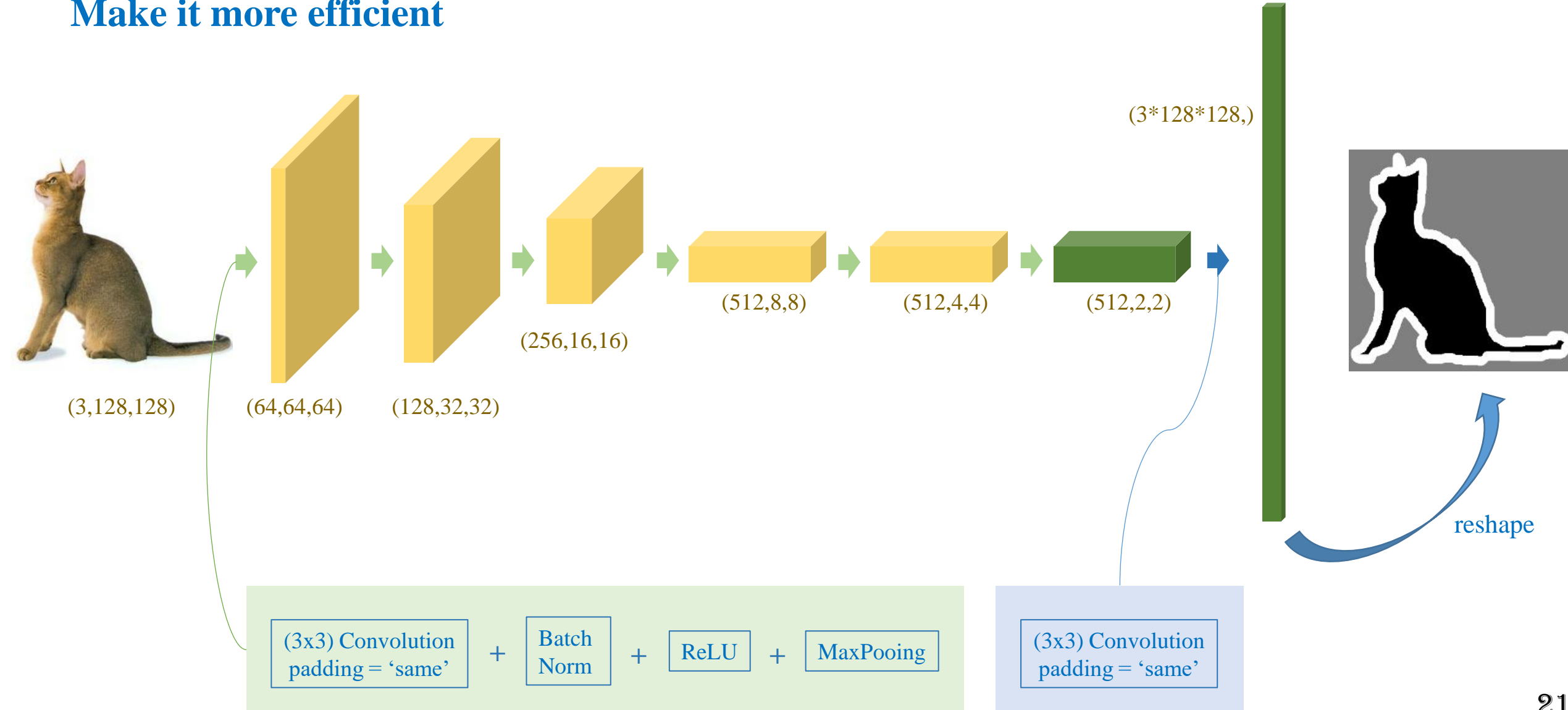


Predicted Mask



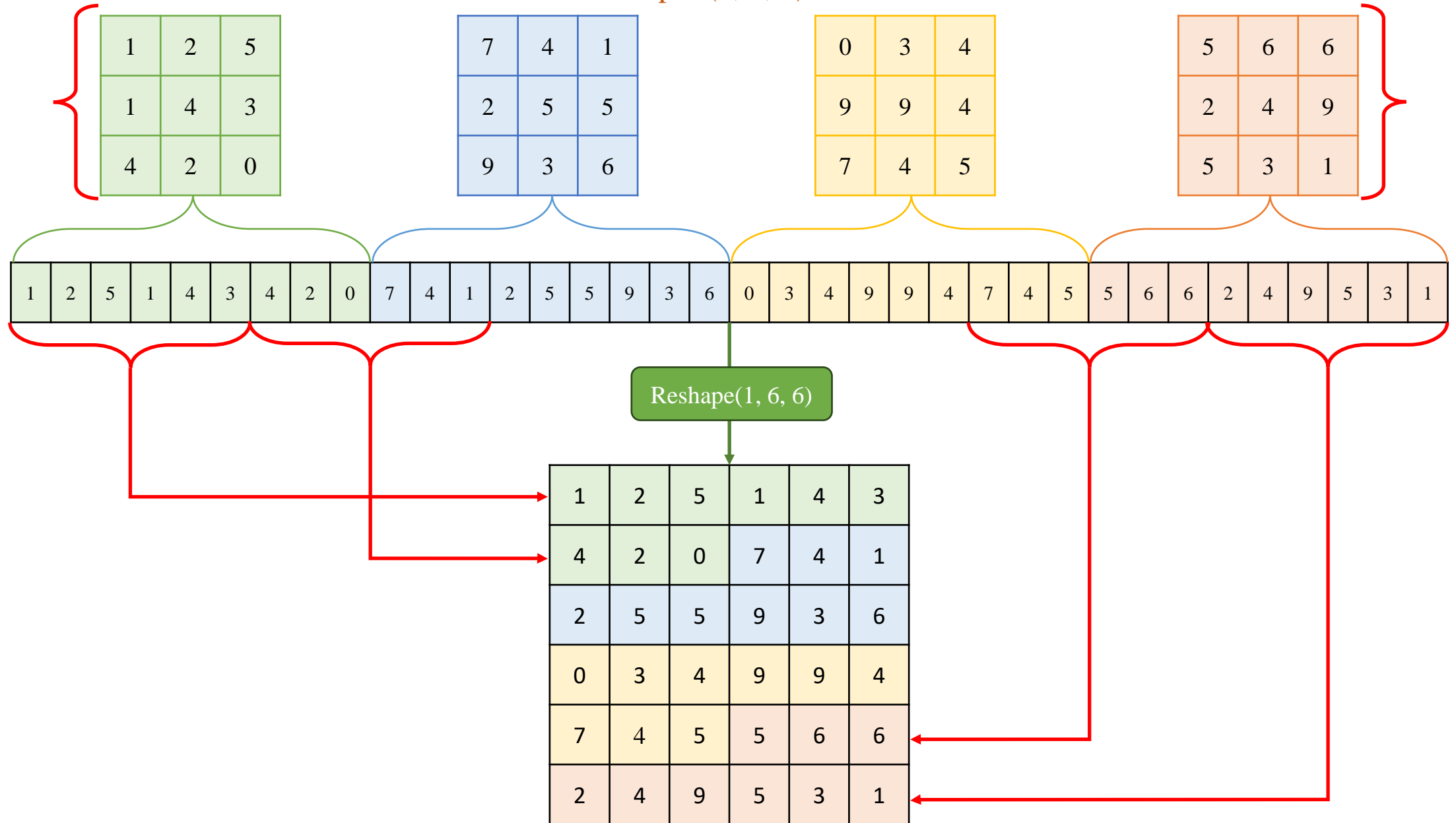
Segmentation (4)

Make it more efficient



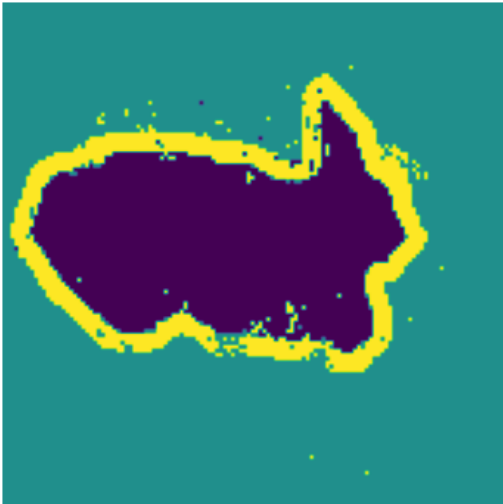
ReShape

Shape=(4, 3, 3)



Segmentation (4)

Make it more efficient



Problem?

```
1  class FCN_Segmentation(nn.Module):
2      def __init__(self):
3          super().__init__()
4          self.enc_1 = Encoder(3, 64)
5          self.enc_2 = Encoder(64, 128)
6          self.enc_3 = Encoder(128, 256)
7          self.enc_4 = Encoder(256, 512)
8          self.enc_5 = Encoder(512, 512)
9          self.enc_6 = Encoder(512, 512)
10         self.out_conv = nn.Conv2d(512, 128*128*3,
11                                     kernel_size=2)
12
13     def forward(self, x):
14         x = self.enc_1(x)
15         x = self.enc_2(x)
16         x = self.enc_3(x)
17         x = self.enc_4(x)
18         x = self.enc_5(x)
19         x = self.enc_6(x)
20         x = self.out_conv(x)
21         return torch.reshape(x, (-1, 3, 128, 128))
```


Further Reading

Fully Convolutional Networks

<https://arxiv.org/pdf/1411.4038.pdf>

Fully Convolutional Networks for Semantic Segmentation

Jonathan Long* Evan Shelhamer* Trevor Darrell
UC Berkeley

{jonlong, shelhamer, trevor}@cs.berkeley.edu

Abstract

Convolutional networks are powerful visual models that yield hierarchies of features. We show that convolutional networks by themselves, trained end-to-end, pixels-to-pixels, exceed the state-of-the-art in semantic segmentation. Our key insight is to build “fully convolutional” networks that take input of arbitrary size and produce correspondingly-sized output with efficient inference and learning. We define and detail the space of fully convolutional networks, explain their application to spatially dense prediction tasks, and draw connections to prior models. We adapt contemporary classification networks (AlexNet [19],

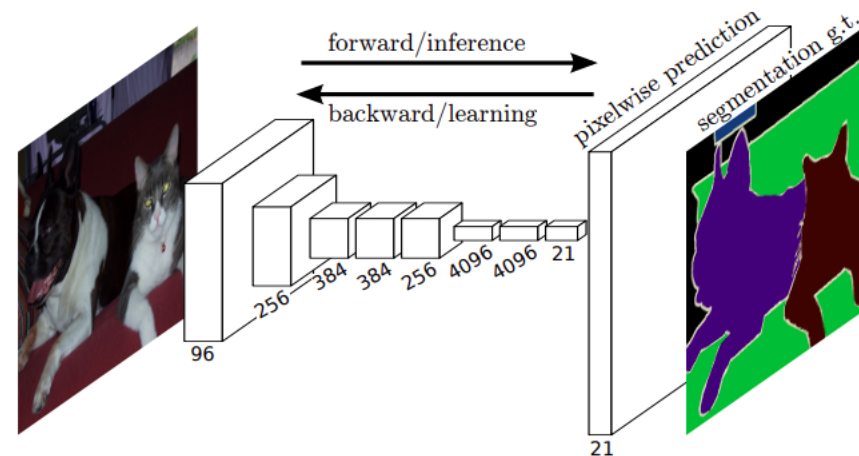
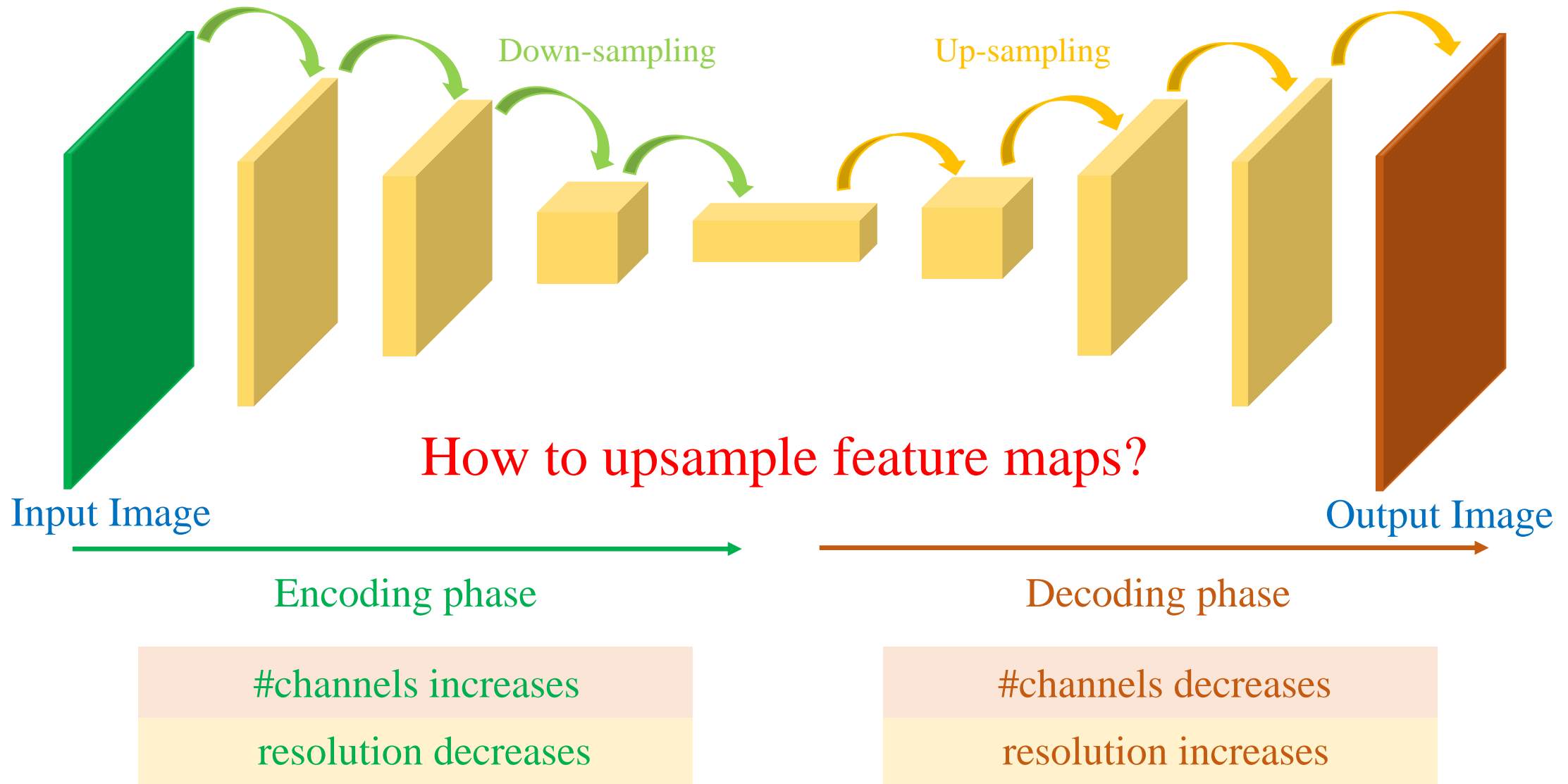


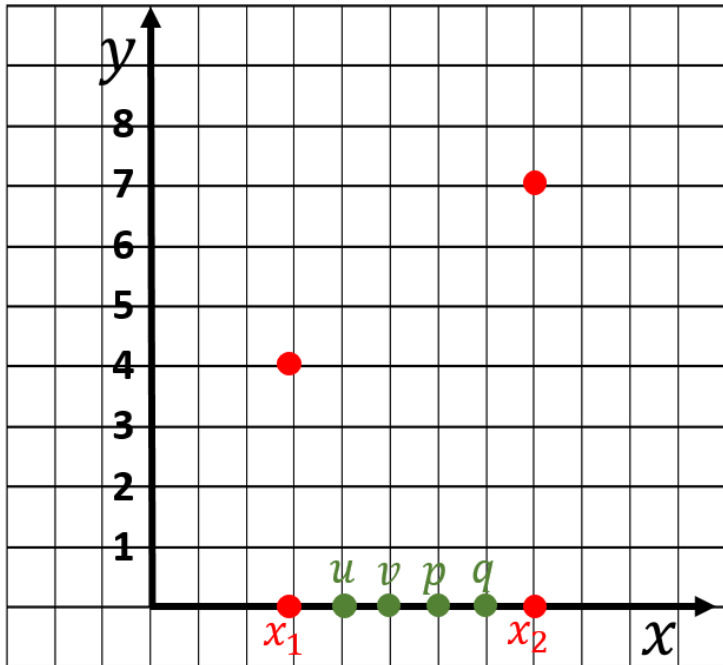
Figure 1. Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.

Segmentation (5)

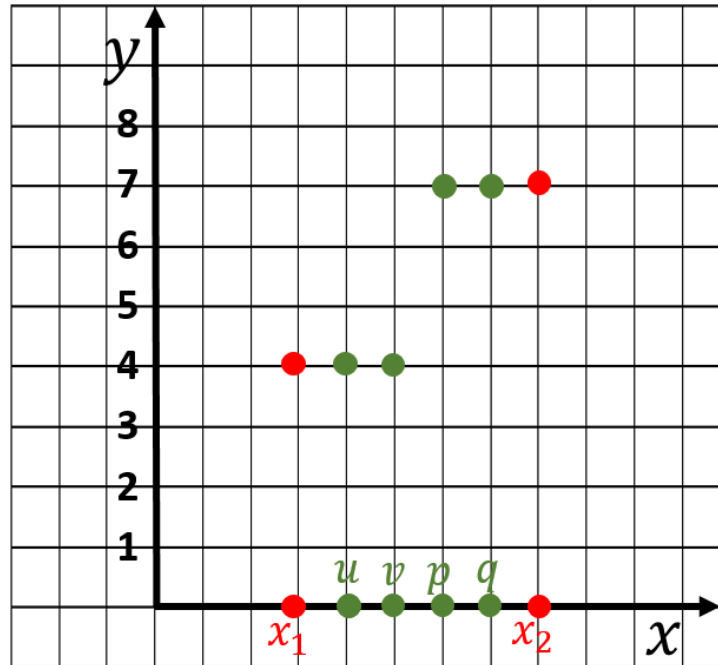


Segmentation (5)

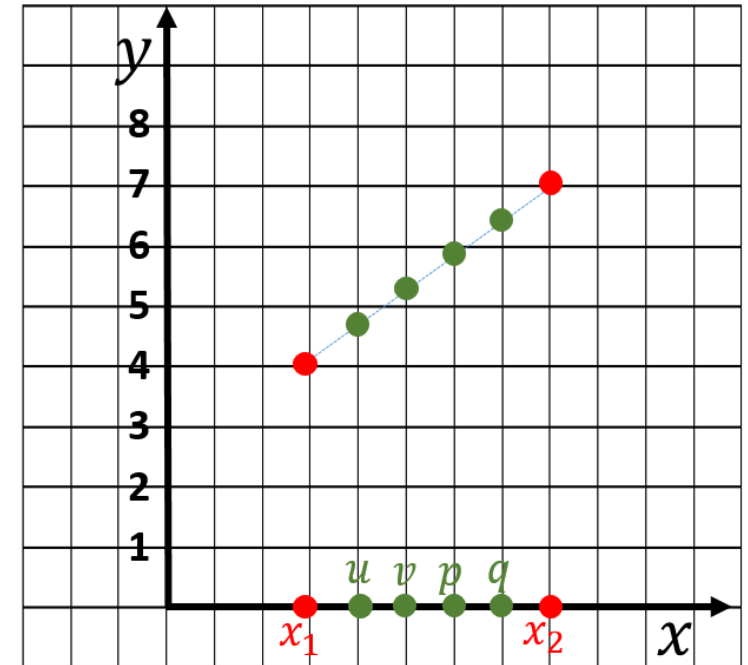
- ❖ Image upsampling
- ❖ Data interpolation



Tìm giá trị cho các vị trí u , v , p và q



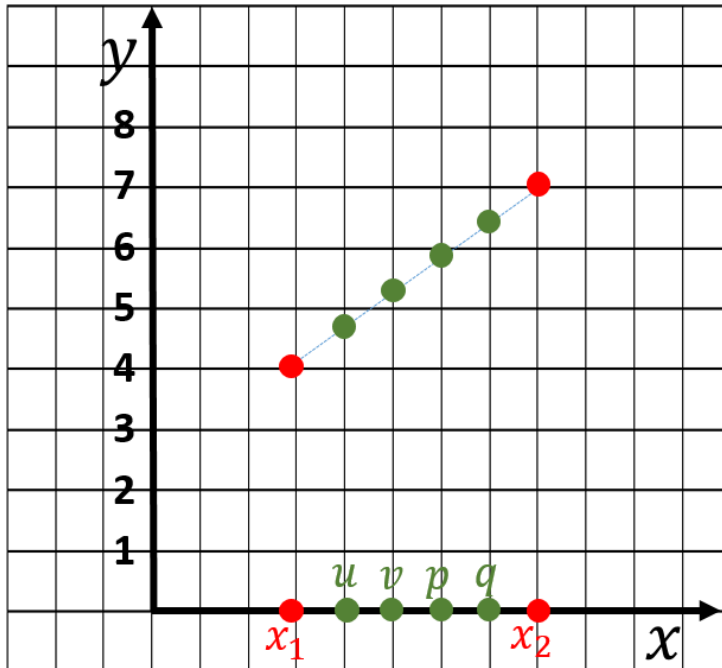
Nearest neighbor: Tính khoảng cách đến x_1 và x_2 , và lấy giá trị của x gần hơn



Nội suy theo hàm tuyến tính

Segmentation (5)

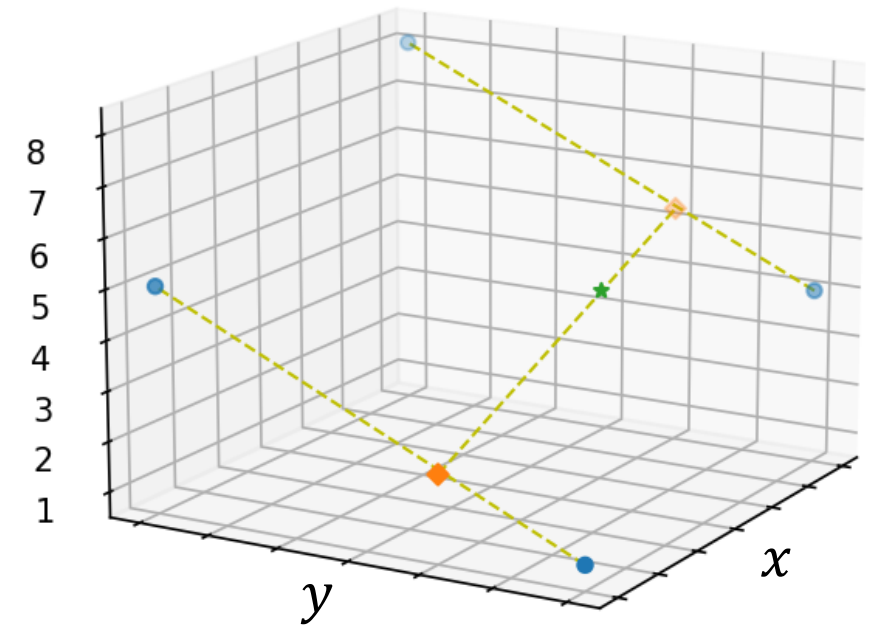
- ❖ Image upsampling
 - ❖ Data interpolation



Nội suy theo hàm tuyến tính

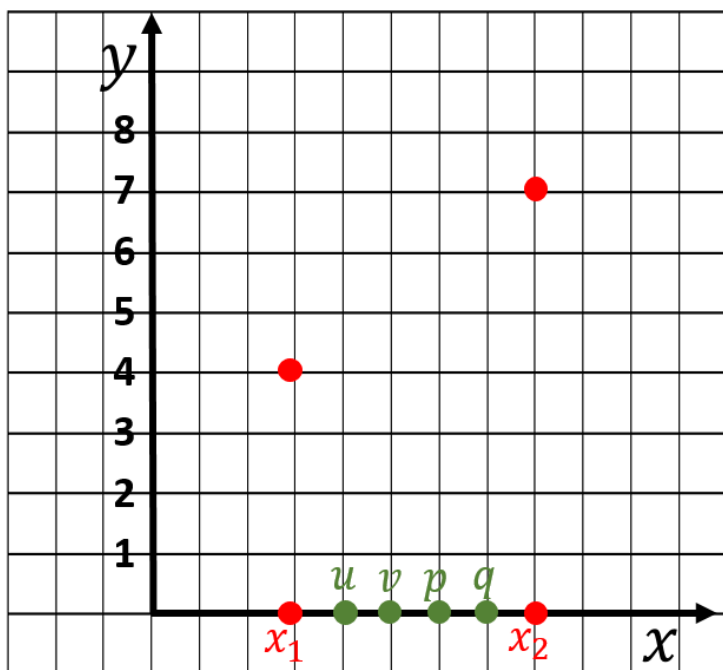
1.0			5.0
3.2	4.2		7.2
4.0			8.0

Bilinear interpolation

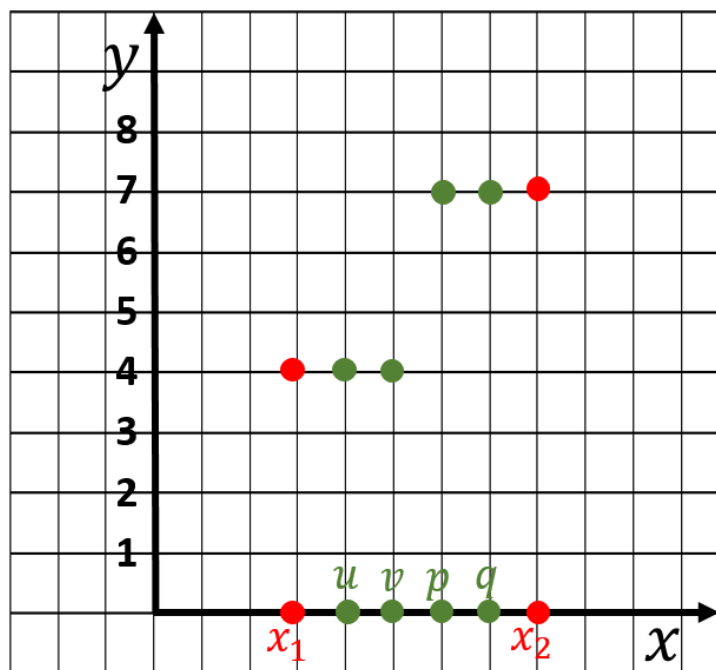


Segmentation (5)

❖ Nearest neighbor interpolation



Tìm giá trị cho các vị trí u , v , p và q



Nearest neighbor: Tính khoảng cách đến x_1 và x_2 , và lấy giá trị của x gần hơn

data =

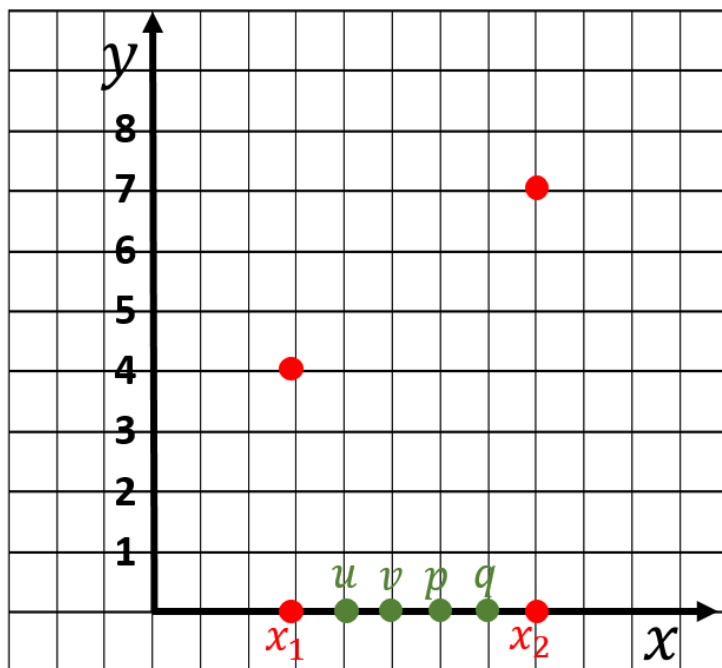
1	5
4	8

1	1	5	5
1	1	5	5
4	4	8	8
4	4	8	8

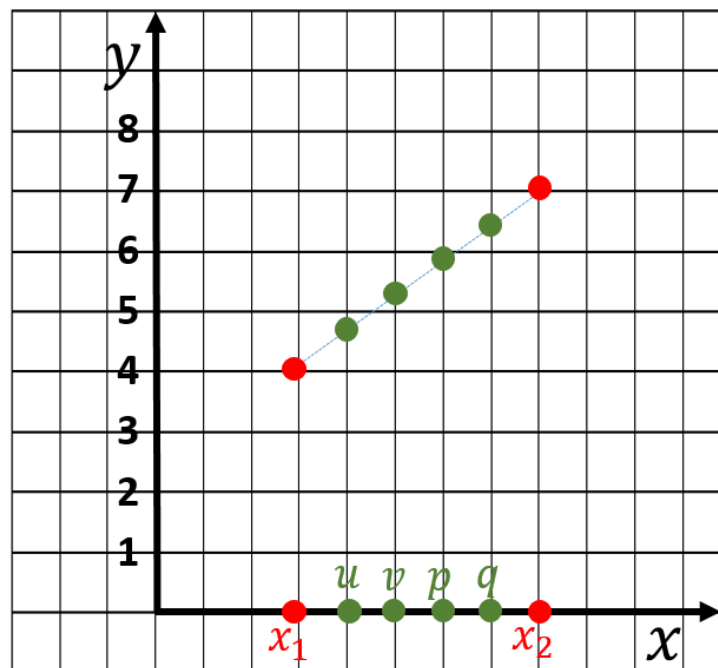
Nearest neighbor
interpolation

Segmentation (5)

❖ Bilinear interpolation



Tìm giá trị cho các vị trí u , v , p và q



Nội suy theo hàm tuyến tính

data =

1	5
4	8

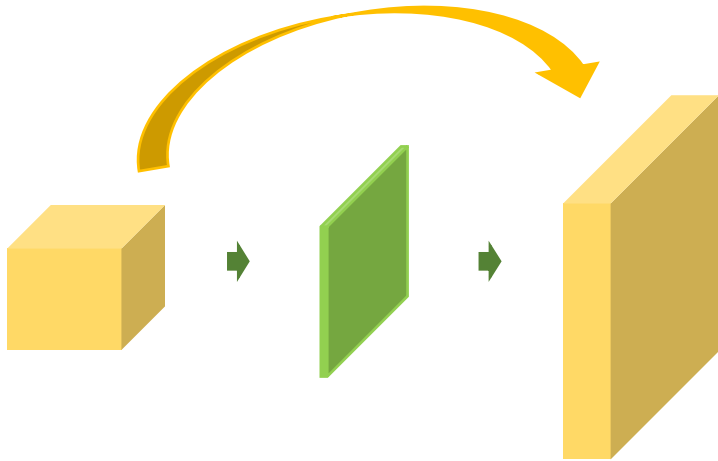
1.0	2.0	4.0	5.0
1.7	2.7	4.7	5.7
3.2	4.2	6.2	7.2
4.0	5.0	7.0	8.0

Bilinear
interpolation

How to Upsample Feature Maps

❖ Interpolation

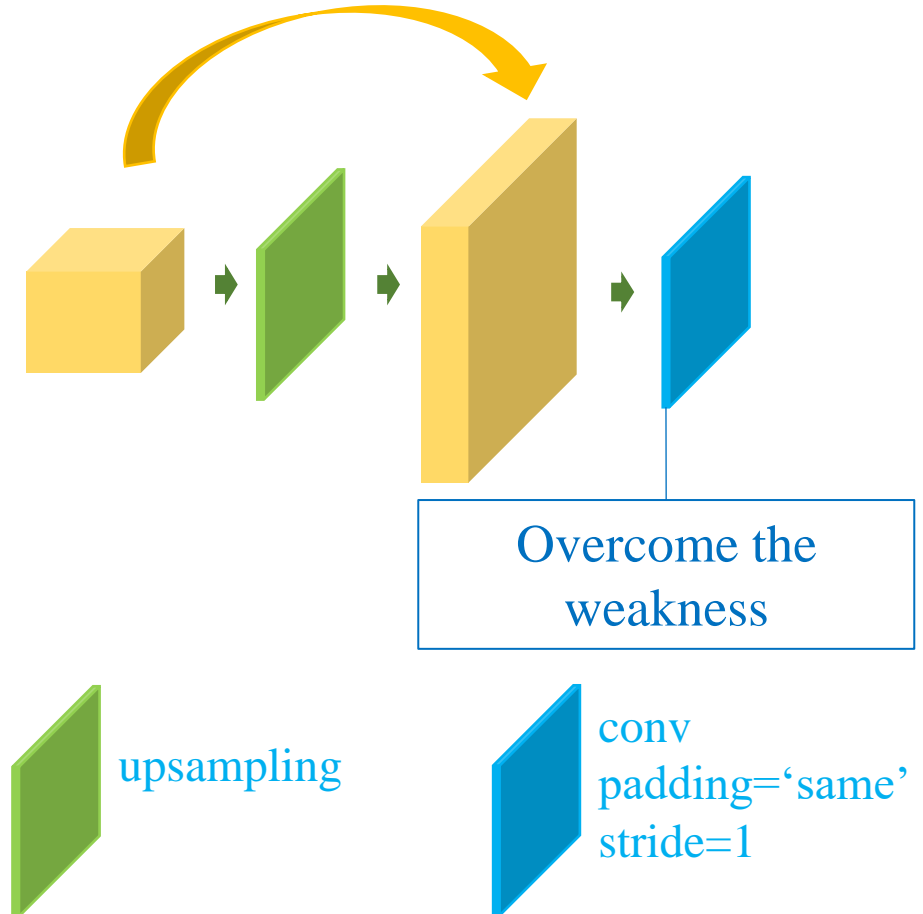
Naïve approach: Only use 'image upsampling'



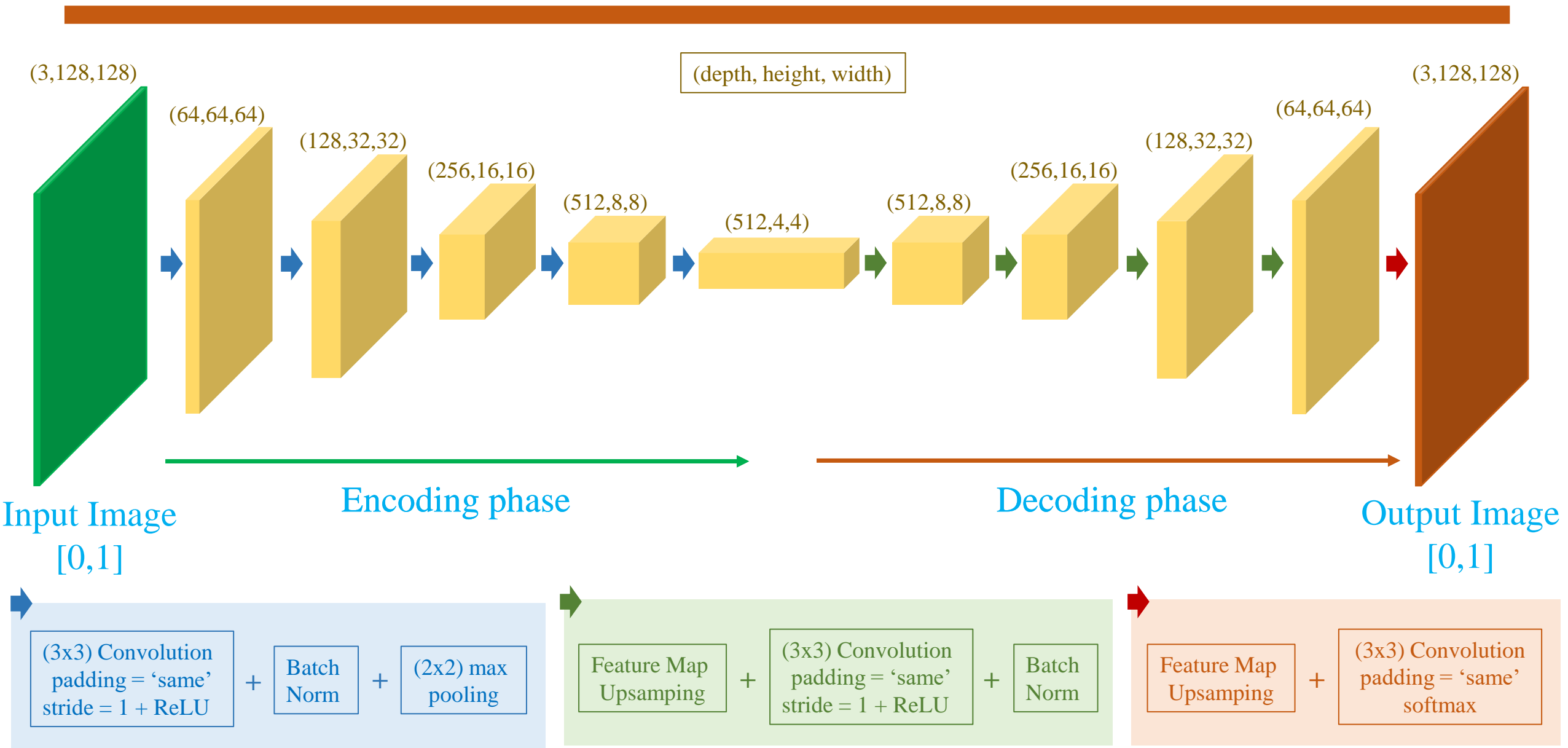
Output feature maps are lack of details

```
nn.Upsample(scale_factor=2,  
            mode='bilinear', align_corners=True)  
nn.Conv2d(in_channels, out_channels,  
          kernel_size=3, padding=1)
```

Use 'image upsampling'+Conv

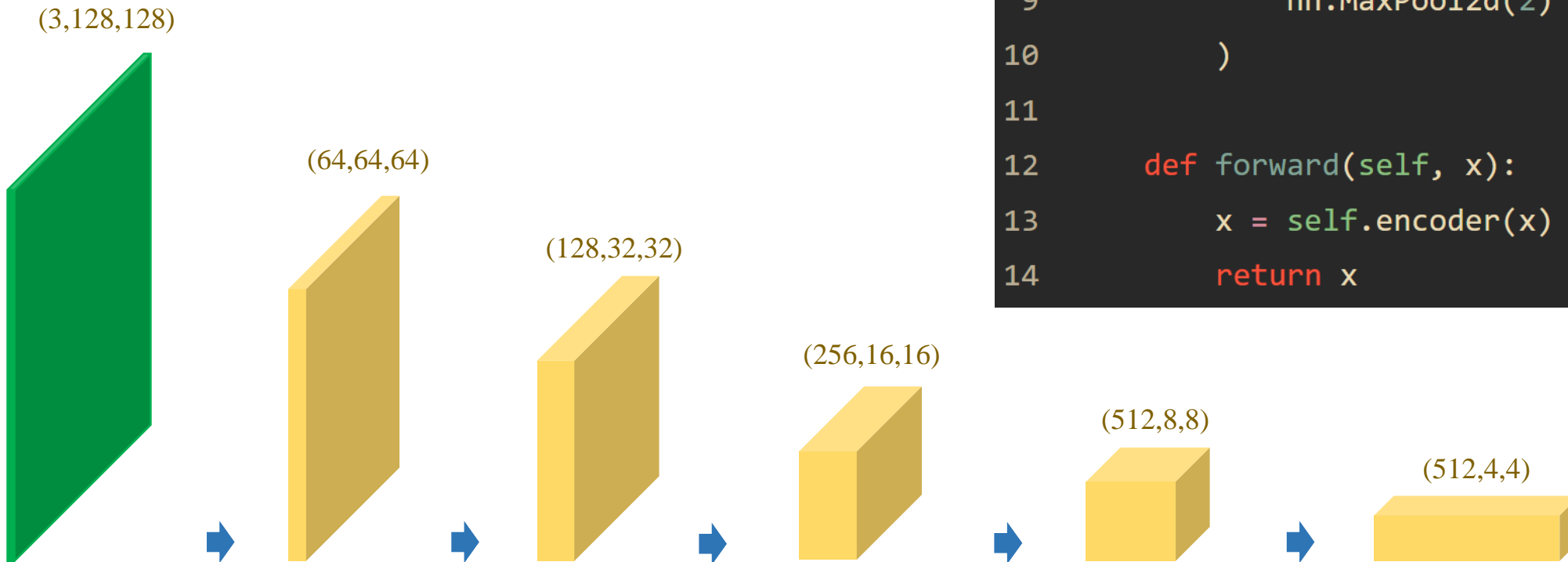


Segmentation (5)



Segmentation (5)

❖ Encoding



```
1 class Encoder(nn.Module):
2     def __init__(self, in_channels, out_channels):
3         super().__init__()
4         self.encoder = nn.Sequential(
5             nn.Conv2d(in_channels, out_channels,
6                       kernel_size=3, padding=1),
7             nn.BatchNorm2d(out_channels),
8             nn.ReLU(inplace=True),
9             nn.MaxPool2d(2)
10        )
11
12    def forward(self, x):
13        x = self.encoder(x)
14        return x
```

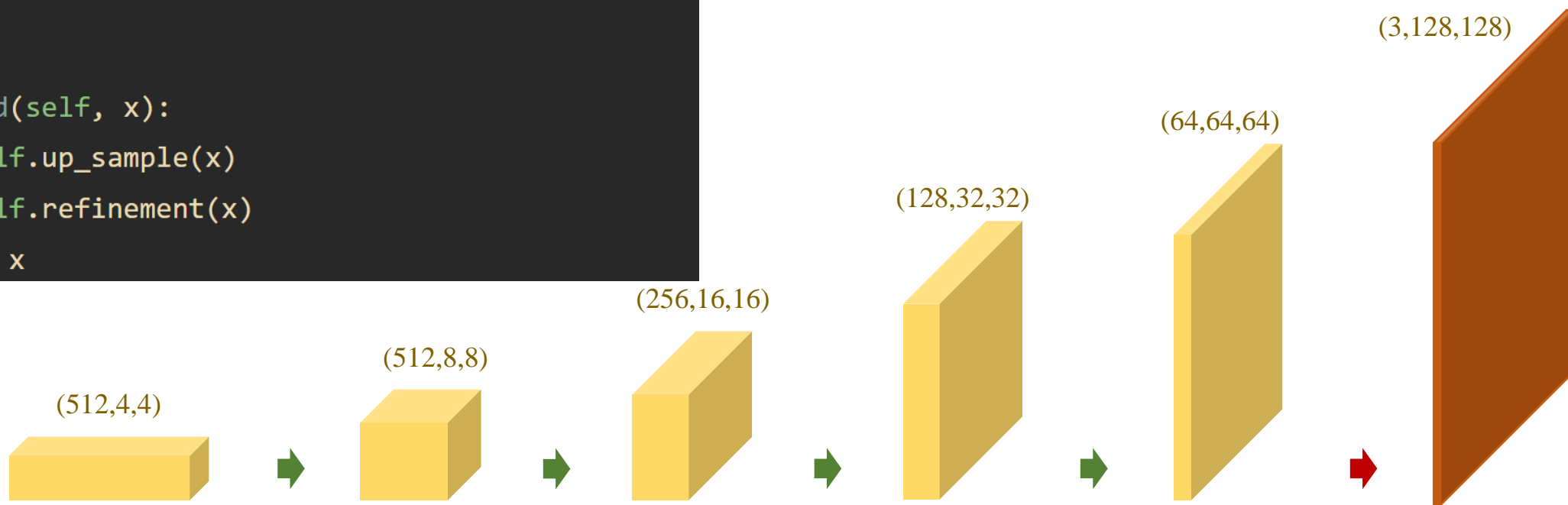
Segmentation (5)

❖ Decoding

ESRGAN: Enhanced Super-Resolution Generative Adversarial Networks, 2018

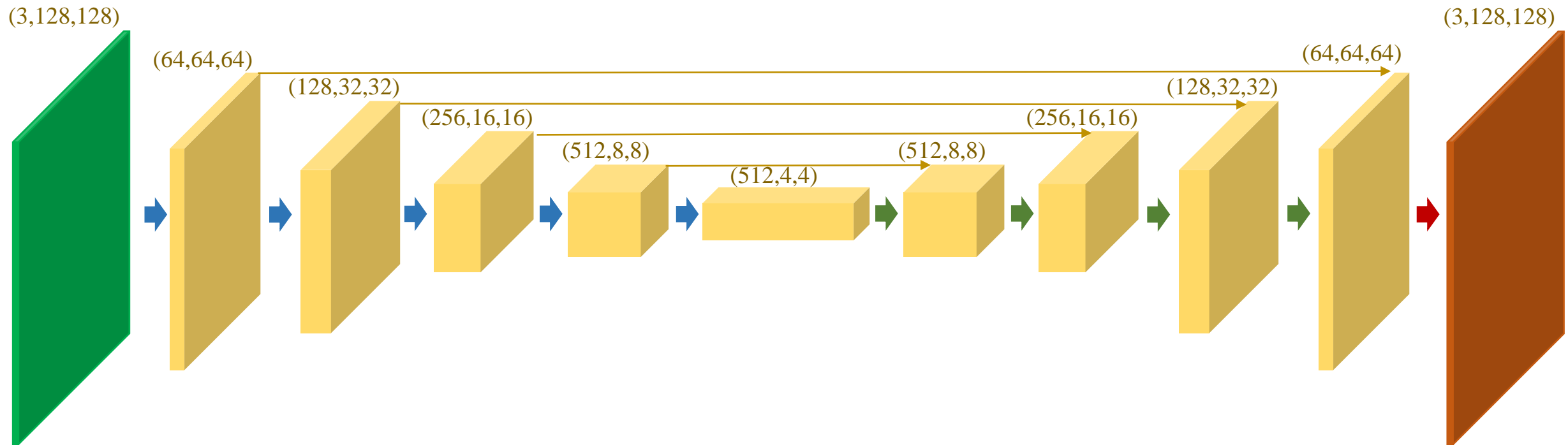
Perception-Oriented Single Image Super-Resolution using Optimal Objective Estimation, CVPR 2023

```
1 class Decoder(nn.Module):
2     def __init__(self, in_channels, out_channels):
3         super().__init__()
4         self.up_sample = nn.Upsample(scale_factor=2,
5                                     mode='bilinear',
6                                     align_corners=True)
7         self.refinement = nn.Sequential(
8             nn.Conv2d(in_channels, out_channels,
9                     kernel_size=3, padding=1),
10            nn.BatchNorm2d(out_channels),
11            nn.ReLU(inplace=True)
12        )
13
14    def forward(self, x):
15        x = self.up_sample(x)
16        x = self.refinement(x)
17        return x
```



Segmentation (6)

Using skip connections



(3x3) Convolution
padding = 'same'
stride = 1 + ReLU

+

Batch
Norm

+

(2x2) max
pooling

Feature Map
Upsampling

+

(3x3) Convolution
padding = 'same'
stride = 1 + ReLU

+

Batch
Norm

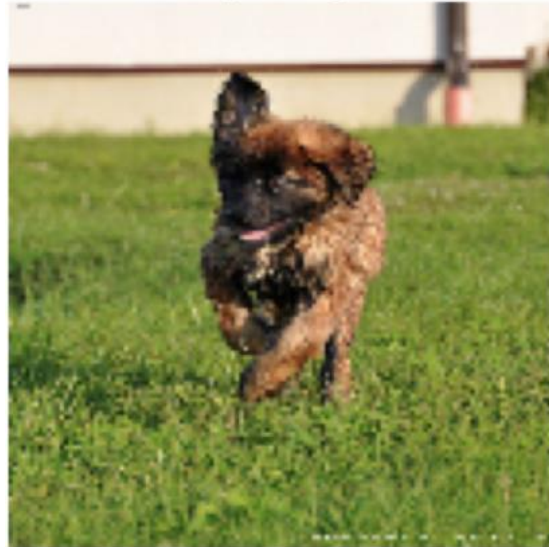
Feature Map
Upsampling

+

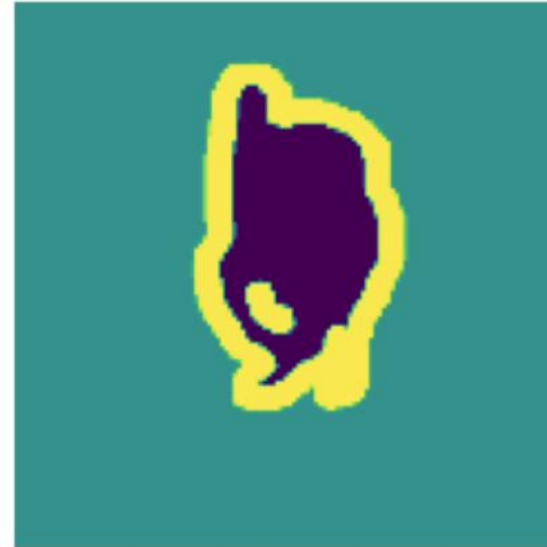
(3x3) Convolution
padding = 'same'
softmax

Segmentation (6)

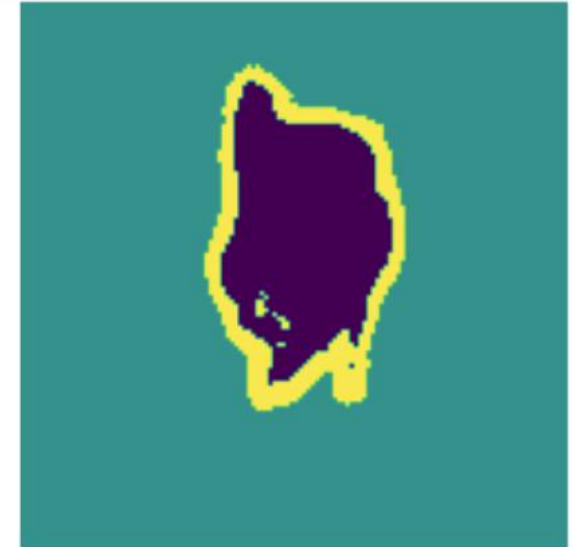
**With Skip
Connections**



Input Image



True Mask



Predicted Mask

**Without Skip
Connections**

