

江南大学

课程设计报告

题 目： 用于语义分割的全卷积网络

题 目（英文）： Fully Convolutional Networks for Semantic Segment

学 号： 1033190614、1033190502、1033190609

姓 名： 张梓琦，关乐平，刘鸿嘉

课 程： 人工智能

指 导 教 师： 徐天阳

江南大学

地址：无锡市蠡湖大道 1800 号

二〇二二 年 六 月

Catalogue

1	Introduction	2
1.1	Development History	2
1.2	Purpose of the study	2
1.3	Significance of the study	3
2	Approach	3
2.1	Fully convolutional networks	3
2.2	Translation invariance	5
2.3	Receptive Field.....	5
2.3.1	Receptive Field Introduction	5
2.3.2	Receptive Field Calculation	6
2.4	ResNet	6
2.5	Transposed Convolution.....	7
2.5.1	Upsample.....	7
2.5.2	Transposed Convolution.....	7
2.6	Jump structure	9
3	Experimental Evaluation	9
3.1	Experimental conditions.....	9
3.2	Experimental procedure	9
3.3	Visualisation analysis	20
3.4	Experimental experience	23
	References	25
	Group division of labour	26

1 Introduction

1.1 Development History

All developments are long periods of technical accumulation, coupled with qualitative changes that occur when some external conditions are met. We have briefly summarised several periods of image segmentation.

Before 2000, when digital image processing we used methods based on several categories: threshold segmentation, region segmentation, edge segmentation, texture features, clustering, etc.

Between 2000 and 2010, the main methods were based on four categories: graph theory, clustering, classification, and a combination of clustering and classification.

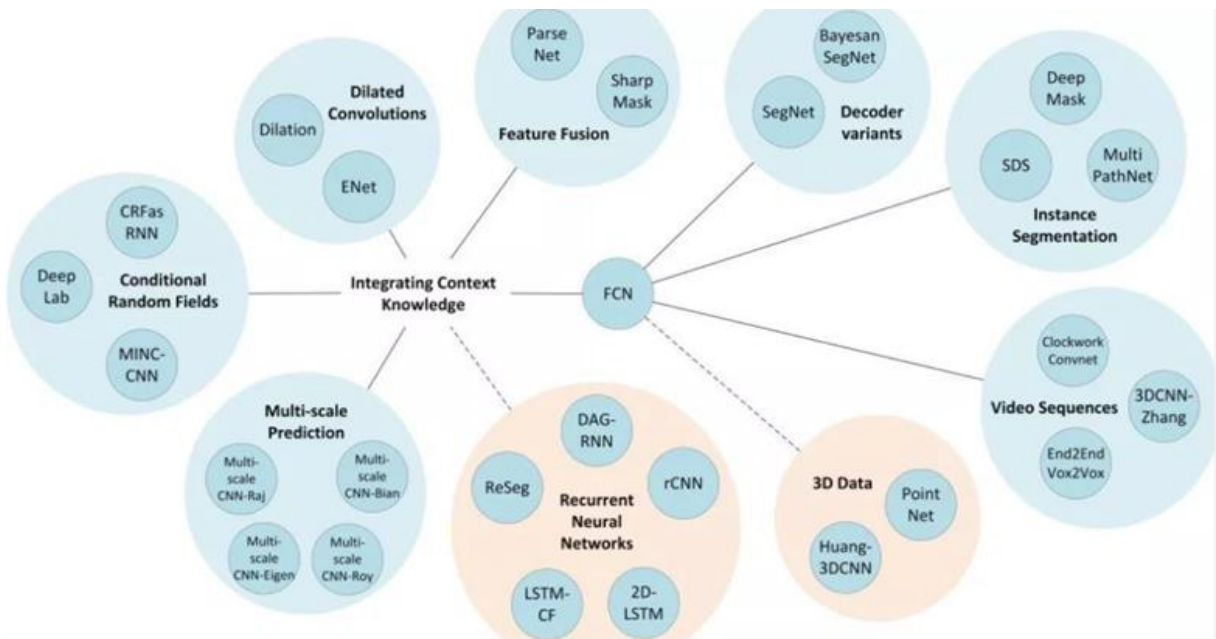
From 2010 to the present, the rise of neural network models and the development of deep learning has involved several main models

TABLE 2: Summary of semantic segmentation methods based on deep learning.

Name and Reference	Architecture	Accuracy	Efficiency	Training	Targets	Sequences	Multi-modal	3D	Source Code	Contribution(s)
Fully Convolutional Network [65]	VGG-16(FCN)	***	+	+	Instance	✓	✓	✓	✓	Forerunner
SegNet [66]	VGG-16 + Decoder	***	++	+	✓	✓	✓	✓	✓	Encoder-decoder
Bayesian SegNet [67]	SegNet	***	+	+	✓	✓	✓	✓	✓	Uncertainty modeling
DeepLab [68]	VGG-16/ResNet-101	***	+	+	✓	✓	✓	✓	✓	Standalone CRF, atrous convolutions
MINC-CNN [43]	GoogLeNet(FCN)	***	+	+	✓	✓	✓	✓	✓	Patchwise CNN, Standalone CRF
CRFasRNN [69]	FCN-8s	***	++	++	✓	✓	✓	✓	✓	CRF reformulated as RNN
Dilation [74]	VGG-16	***	+	+	✓	✓	✓	✓	✓	Dilated convolutions
ENet [72]	ENet bottleneck	***	++	+	✓	✓	✓	✓	✓	Bottleneck module for efficiency
Multi-scale-CNN-Raj [73]	VGG-16(FCN)	***	+	+	✓	✓	✓	✓	✓	Multi-scale architecture
Multi-scale-CNN-Eigen [74]	Custom	***	+	+	✓	✓	✓	✓	✓	Multi-scale sequential refinement
Multi-scale-CNN-Roy [75]	Multi-scale-CNN-Eigen	***	+	+	✓	✓	✓	✓	✓	Multi-scale coarse-to-fine refinement
Multi-scale-CNN-Bian [76]	FCN	***	++	++	✓	✓	✓	✓	✓	Independently trained multi-scale FCNs
ParseNet [77]	VGG-16	***	+	+	✓	✓	✓	✓	✓	Global context feature fusion
ReSeg [78]	VGG-16 + ReNet	***	+	+	✓	✓	✓	✓	✓	Extension of ReNet to semantic segmentation
LSTM-CF [79]	Fast R-CNN + DeepMask	***	+	+	✓	✓	✓	✓	✓	Fusion of contextual information from multiple sources
2D-LSTM [80]	MDRNN	***	++	+	✓	✓	✓	✓	✓	Image context modelling
rCNN [81]	MDRNN	***	++	+	✓	✓	✓	✓	✓	Different input sizes, image context
DAG-RNN [82]	Elman network	***	+	+	✓	✓	✓	✓	✓	Graph image structure for context modelling
SDS [10]	R-CNN + Box CNN	***	+	+	✓	✓	✓	✓	✓	Simultaneous detection and segmentation
DeepMask [83]	VGG-A	***	+	+	✓	✓	✓	✓	✓	Proposals generation for segmentation
SharpMask [84]	DeepMask	***	+	+	✓	✓	✓	✓	✓	Top-down refinement module
MultiPathNet [85]	Fast R-CNN + DeepMask	***	+	+	✓	✓	✓	✓	✓	Multi path information flow through network
Huang-3DCNN [86]	Own 3DCNN	+	+	+	✓	✓	✓	✓	✓	3DCNN for voxels
PointNet [87]	Own MLP-based	++	++	+	✓	✓	✓	✓	✓	Segmentation of uns. 3D point clouds
Clockwork Convnet [88]	FCN	++	++	+	✓	✓	✓	✓	✓	Clockwork scheduling for sequences
3DCNN-Zhang [89]	Own 3DCNN	++	+	+	✓	✓	✓	✓	✓	3D convolutions and graph cut for sequences
End2End Vox2Vox [90]	3DCNN	++	+	+	✓	✓	✓	✓	✓	3D convolutions/deconvolutions for sequences

1-A

By the end of 2017, we had differentiated hundreds of model structures. Of course, after examining the techniques and principles, we have discovered the characteristic that the most successful current deep learning techniques for image segmentation are based on a common pioneer: the FCN (Fully Convolutional Network).



1-B

1.2 Purpose of the study

For FCN, it means that all layers in the model network are convolutional layers, but the traditional CNN has several disadvantages: firstly, the storage overhead is large, the sliding window is large, each window needs storage space to hold features and discriminative categories, and using a fully connected structure, the last few layers are nearly exponentially storage incremental; secondly, the computational efficiency is low, a lot of repeated computations; finally, the sliding window size is relatively independent, and the use of full connectivity at the end can only

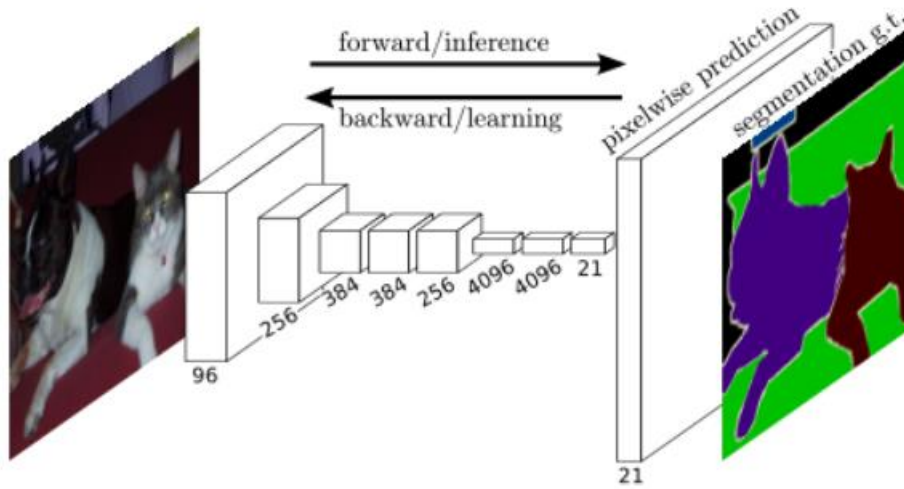
constrain local features. So to solve some of the problems of CNN, FCN converts the last three fully-connected layers in CNN into three convolutional layers.

1.3 Significance of the study

Convolutional networks are driving advances in recognition technology. Convolutional networks are not only improving overall image classification, but are also making advances in local tasks with structured outputs. These advances include bounding box target detection, partial and critical point prediction, and local correspondence. From rough inference to fine inference, it is natural that the next step is to make predictions for each pixel.

Fully convolutional versions of existing networks can predict dense outputs from inputs of arbitrary size. Both learning and inference are performed over the entire image by intensive feed-forward computation and back propagation. Pixel prediction and learning in the network is achieved by subsampling pooling in the upsampling layer within the network.

FCN classifies images at the pixel level, thus solving the problem of semantic segmentation. Unlike the classical CNN, which uses a fully-connected layer after the convolution layer to obtain a fixed length feature vector for classification (fully connected layer + softmax output), the FCN can accept an input image of arbitrary size and uses a deconvolution layer to upsample the feature map of the last convolution layer to restore it to the same size as the input image, thus generating a prediction for each pixel. A prediction is made for each pixel, while preserving the spatial information in the original input image, and the final pixel-by-pixel classification is performed on the upsampled feature map.



1-C

2 Approach

2.1 Fully convolutional networks

The only difference between the fully connected and convolutional layers is that the neurons in the convolutional layer are connected to only one local region in the input data and that the neurons in the convolutional column share parameters. However, in both types of layers, the neurons compute a dot product, so their functional form is the same. It is therefore possible to interconvert the two.

For either convolutional layer, there exists a fully connected layer that implements the same forward propagation function as it does. The weight matrix is a huge matrix, with all but some specific blocks being zero. And in most of these blocks, the elements are equal.

Conversely, any fully-connected layer can be transformed into a convolutional layer. For example, a fully connected layer with $K=4096$ and an input data body of size $7*7*512$ can be equivalently viewed as a convolutional layer with $F=7, P=0, S=1, K=4096$. In other words, the size of the filter is set to match the size of the input data body. Since there is only a single depth column covering and sliding over the input data body, the output will be $1*1*4096$, which is the same as using the initial fully-connected layer.

Fully-connected layers into convolutional layers: Of the two transformations, converting fully-connected layers into convolutional layers is more useful in practice. Suppose the input to a convolutional neural network is

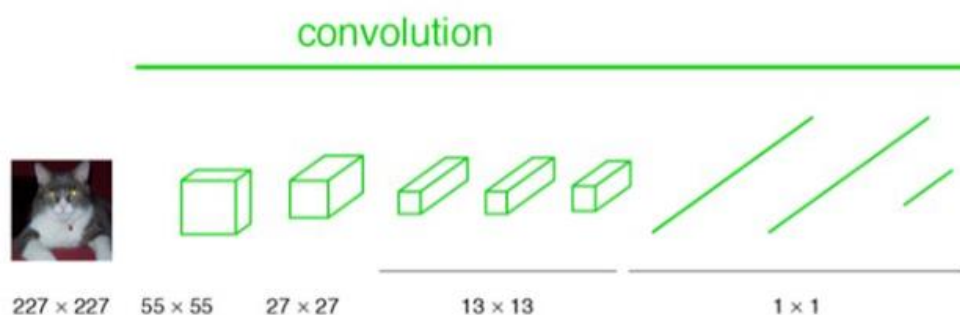
a $224 \times 224 \times 3$ image, and a series of convolutional and downsampling layers transform the image data into a body of activation data of size $7 \times 7 \times 512$. AlexNet uses two fully-connected layers of size 4096, and the last fully-connected layer with 1000 neurons is used to compute classification scores. We can transform any of these 3 fully-connected layers into a convolutional layer.

For the first fully-connected layer with a connection area of $[7 \times 7 \times 512]$, let its filter size be $F=7$, so that the output data body is $[1 \times 1 \times 4096]$.

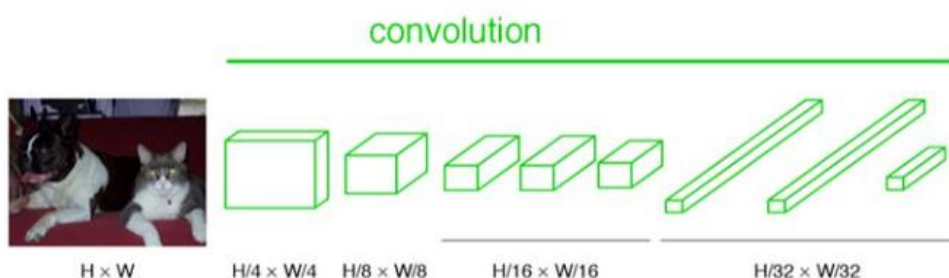
For the second fully connected layer, make the filter size $F=1$ so that the output data body is $[1 \times 1 \times 4096]$.

Do the same for the last fully-connected layer, making its $F=1$, so that the final output is $[1 \times 1 \times 1000]$.

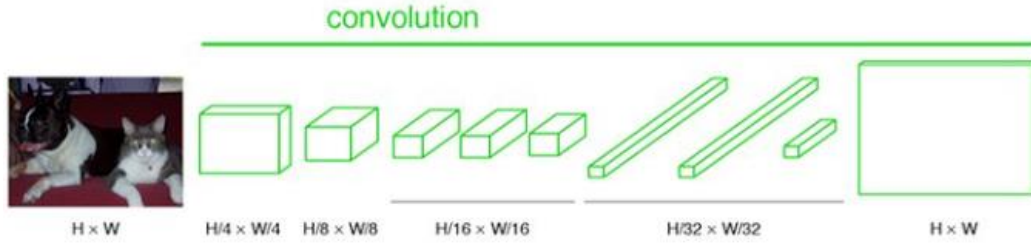
As shown below, the FCN converts the fully-connected layers in a traditional CNN into convolutional layers, corresponding to the CNN network FCN converts the last three fully-connected layers into three convolutional layers. In the traditional CNN structure, the first 5 layers are convolutional layers, layers 6 and 7 are each a 1D vector of length 4096, and layer 8 is a 1D vector of length 1000, corresponding to 1000 different classes of probabilities. fcnn represents these 3 layers as convolutional layers, and the sizes of the convolutional kernels (number of channels, width, height) are $(4096, 1, 1)$, $(4096, 1, 1)$, $(1000, 1, 1)$. There is no difference in the numbers, but convolution is a different concept and computation process from fully connected, using the weights and biases already trained in the previous CNN, but the difference is that the weights and biases have their own range and belong to a convolution kernel of their own. Therefore, all the layers in the FCN are convolutional layers, so it is called a full convolutional network.



The input image size in the CNN is agreed to be fixed resize to 227×227 , 55×55 after pooling in the first layer, 27×27 after pooling in the second layer, and 13×13 after pooling in the fifth layer. The FCN input image is $H \times W$, which becomes $1/4$ of the original image size after pooling in the first layer, $1/8$ of the original image size in the second layer, $1/16$ of the original image size in the fifth layer, and $1/32$ of the original image size in the eighth layer (erratum: in fact the real code is $1/2$ in the first layer, and so on).



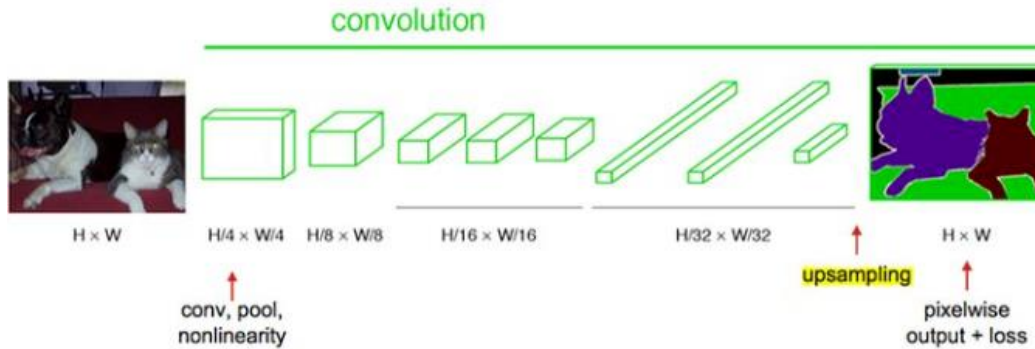
After multiple convolution and pooling, the image gets smaller and smaller, and the resolution gets lower and lower. When the image reaches $H/32 \times W/32$, the image is the smallest layer, and the resulting map is called heatmap. The heatmap is our most important high-dimensional diagnostic map, and after getting the heatmap of high-dimensional features is the most important step and the last step to upsampling the original image, enlarging, enlarging the image to the size of the original image.



2-C

Each layer of data in a convnet is a three-dimensional array of size $h \times w \times d$, where h and w are spatial dimensions, and d is the feature or channel dimension. The first layer is the image, with pixel size $h \times w$, and d color channels. Locations in higher layers correspond to the locations in the image they are path-connected to, which are called their receptive fields.

The final output is 1000 heatmap images that have been upsampled to the original size. In order to classify each pixel by predicting the label into the final image that has been semantically segmented, there is a small trick here, which is to finally classify the pixel by finding the maximum numerical description (probability) of its position in the 1000 images, pixel by pixel, as the pixel. This results in an image that has already been classified, as shown below with the dog and cat on the right.



2-D

Convnets are built on translation invariance. Their basic components (convolution, pooling, and activation functions) operate on local input regions, and depend only on relative spatial coordinates.

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

2-E

2.2 Translation invariance

In Euclidean geometry, translation is a geometric change that represents moving every point in an image or a space the same distance in the same direction. For example, for an image classification task, the target in the image should give the same result regardless of where it is moved in the image, this is the translation invariance in a convolutional neural network.

Translational invariance means that the system produces exactly the same response (output), regardless of how its input has been translated. Translational equivariance means that the system works the same at different locations, but its response changes as the target location changes.

2.3 Receptive Field

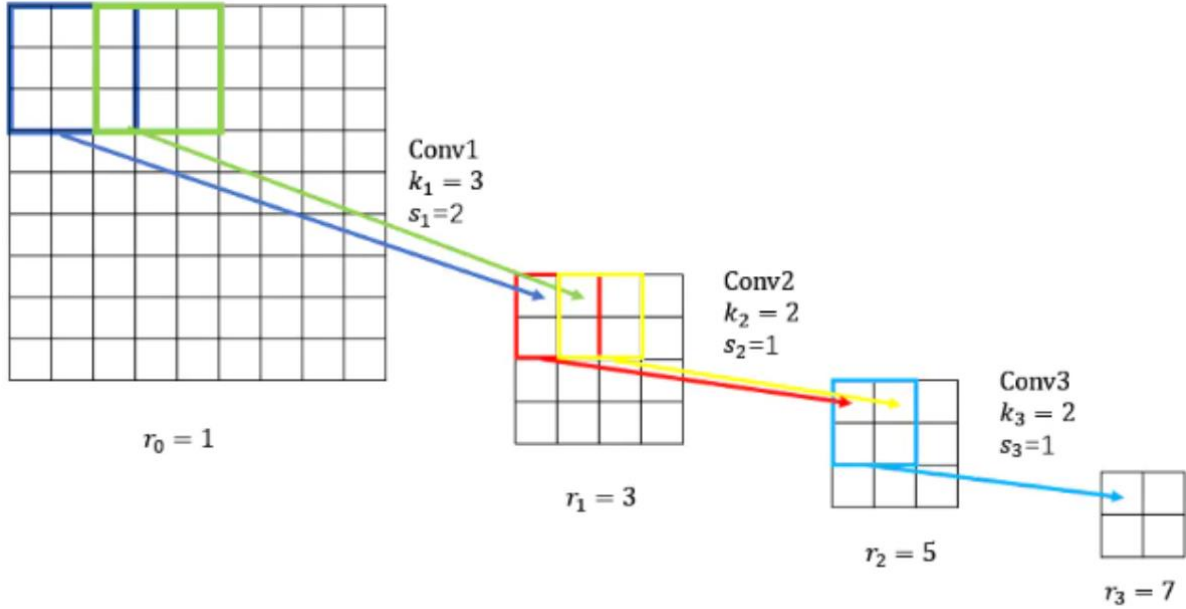
2.3.1 Receptive Field Introduction

In a convolutional neural network, Receptive Field is defined as the size of the area mapped on the original image for each pixel on the featuremap output from each layer of the convolutional neural network, where the original image is the input image to the network, which is pre-processed (e.g. resize, warp, crop).

The reason why neurons cannot perceive all the information of the original image is that convolutional and pooling

layers are commonly used in convolutional neural networks, and are locally connected between layers.

A larger value of the neuron's receptive field means that it has access to a larger range of the original image, which means that it may contain more global, higher-level semantic features; conversely, a smaller value means that the features it contains tend to be more local and detailed. The value of the perceptual field can therefore be used to roughly determine the level of abstraction of each layer.



2-F

2.3.2 Receptive Field Calculation

The receptive field size of a deep convolutional layer is related to the filter size and step size of all the layers before it, and these two parameters are involved in convolutional and pooling layers. We use k_n , s_n , and r_n to denote the kernel_size, stride, and receptive_field of the n th layer, respectively, so that for a convolutional neural network, the receptive field is calculated as follows:

$$r_0 = 1, r_1 = k_1$$

$$r_n = r_{n-1} * k_n - (k_n - 1) * (r_{n-1} - \prod_{i=1}^{n-1} s_i) \quad n \geq 2$$

2-G

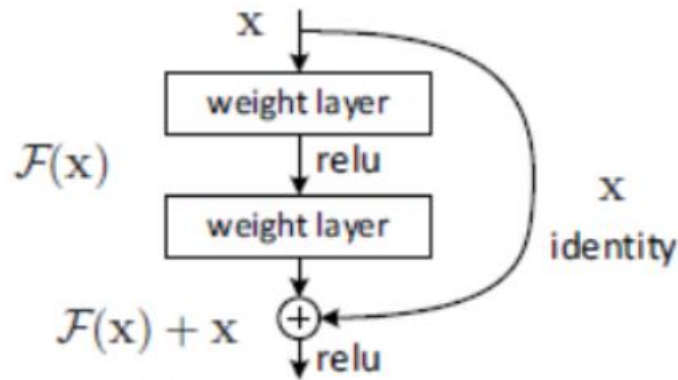
A larger value for the receptive field of a neuron means that it has access to a larger range of original images, which means that it may contain more global, higher semantic features; conversely, a smaller value means that the features it contains tend to be more local and detailed. The value of the perceptual field can therefore be used to roughly determine the level of abstraction of each layer

The effective receptive field is often smaller than the theoretical receptive field, because the input layer uses the edge points less often than the middle points, and therefore makes a different contribution, so that after multiple layers of convolution, the input layer's contribution to the feature map points is distributed in the shape of a Gaussian distribution

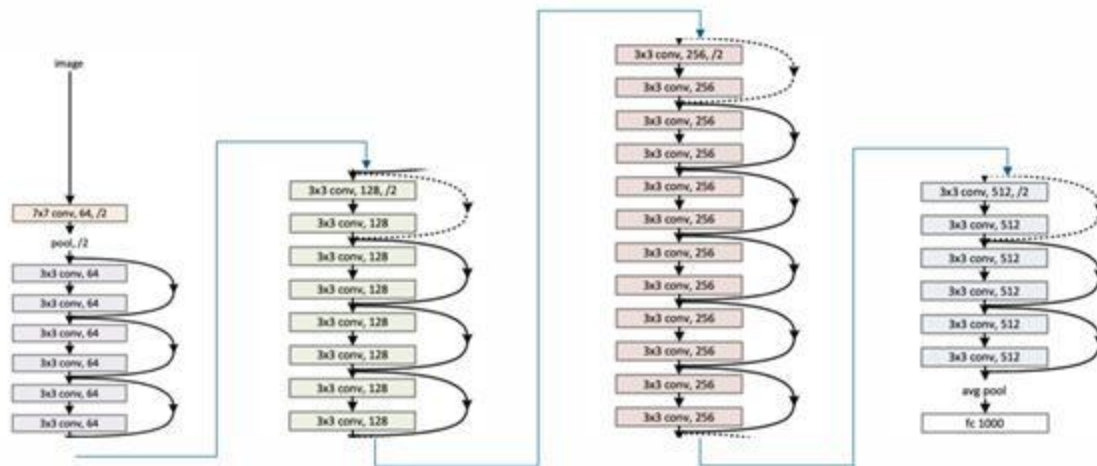
2.4 ResNet

He Kaiming proposed the residual neural network, or Resnet, in 2015 and won the classification competition at ILSVRC-2015.

ResNet can effectively eliminate the gradient dispersion or gradient explosion problem caused by the increase in the number of convolutional layers. The core idea of ResNet is that the network output is divided into 2 parts: identity mapping, residual mapping, i.e. $y = x + F(x)$, which is illustrated as follows.



2-H



2-I

2.5 Transposed Convolution

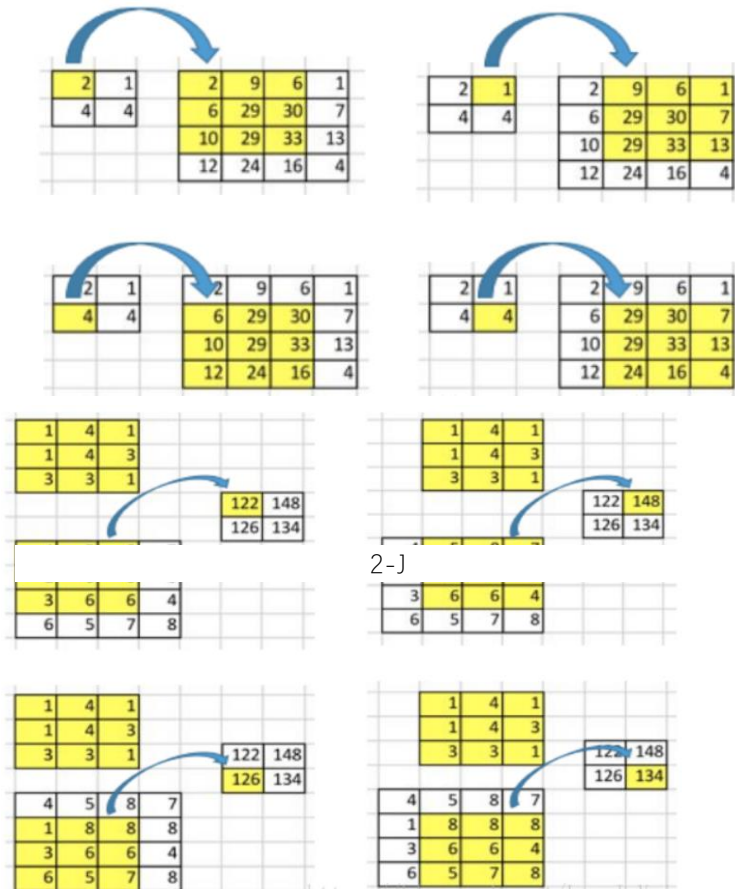
2.5.1 Upsample

In the field of deep learning applied to computer vision, as the input image is extracted by convolutional neural network (CNN), the size of the output often becomes smaller, and sometimes we need to restore the image to its original size for further computation, this operation of mapping the image from small resolution to large resolution by expanding the image size is called upsampling.

2.5.2 Transposed Convolution

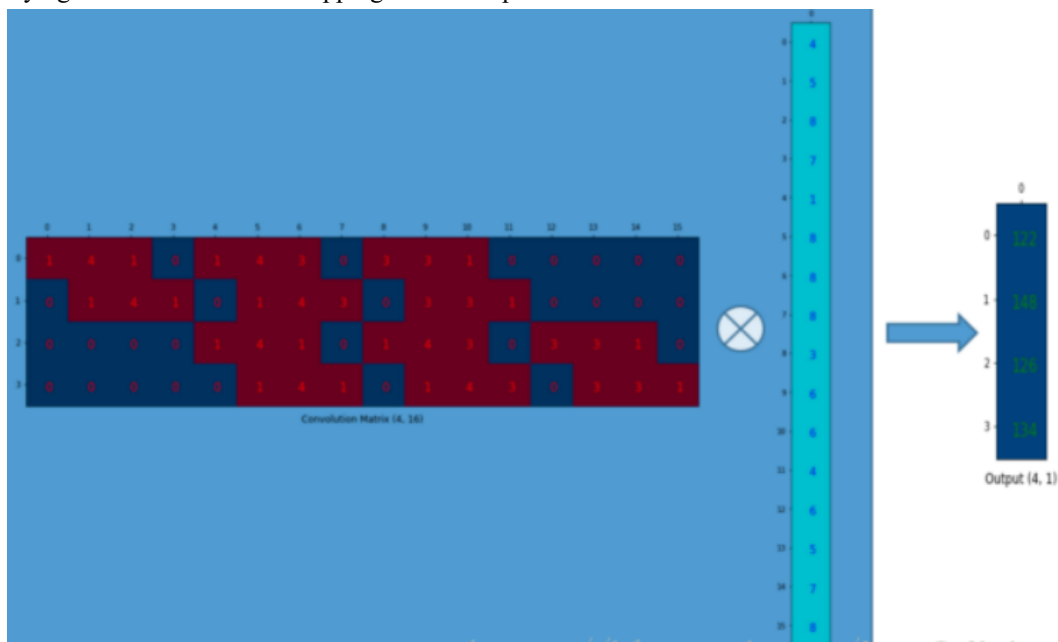
Deconvolution can be understood as the inverse of convolution. Deconvolution does not recover the loss of values caused by the convolution operation, it simply transforms the steps in the convolution process in reverse once, and can therefore be called transposed convolution. The parameters can be learned in end-to-end training using deconvolution. The formula for deconvolution is: $\text{output} = (\text{input} - 1) \times \text{stride} + k - 2p$

A convolution, for example, takes a large matrix and passes it through a convolution kernel to call it a small matrix; a convolution operation is a many-to-one mapping relationship. A deconvolution operation is a one-to-many mapping relationship. A deconvolution is the reversal of a convolution operation, converting a small matrix into a large one.



2-K

We want to go from $4(2 \times 2)$ to $16(4 \times 4)$, so we use a 16×4 matrix, but one more thing to note is that we are trying to maintain a 1 to 9 mapping relationship.



2-L

Suppose we transpose that this convolution matrix $C (4 \times 16)$ becomes $C_T (16 \times 4)$. We can matrix multiply the C^T and column vectors (4×1) to produce a 16×1 output matrix. This transpose matrix maps exactly one element to nine elements.

2.6 Jump structure

The 32-pixel step size of the final prediction layer limits the scale of detail in the upsampled output. FCN therefore solves this problem by adding a treaty structure, which fuses the final prediction layer with a lower layer with a finer step size. This turns the linear topology into a DAG, with edges jumping from lower to higher layers. Because they see fewer pixels, finer scale predictions should require fewer layers, so it makes sense to output from a shallower network. Combining fine and coarse layers allows the model to make local predictions that fit the global structure.

A jump structure is added between layers to fuse coarse, semantic and local appearance information. This jump structure is learned end-to-end to optimise the semantic and spatial accuracy of the output.

3 Experimental Evaluation

3.1 Experimental conditions

For the selection of the dataset, we first used the base dataset from VOC2012, which was divided into 21 categories, namely 'background', 'aeroplane', 'bicycle', 'bird', 'boat', 'bottle', 'bus', 'car', 'cat', 'chair', 'cow', 'diningtable', 'dog', 'horse', 'motorbike', 'person', 'potted plant', 'sheep', 'sofa', 'train', 'tv/monitor'. In addition to this, we tested this on our own photos around us and observed the segmentation effect.

3.2 Experimental procedure

Reading data from a dataset:

```
1. voc_root = 'data\\VOCdevkit\\VOC2012'
2. print("voc_root: \n", os.listdir(voc_root), "\n")
3. print("JPEGImages: \n", os.listdir(voc_root+'/JPEGImages')[:10], "\n")
4. print("SegmentationClass: \n", os.listdir(voc_root+'/SegmentationClass')[:10], "\n")

5. def read_images(root = voc_root, train = True):
6.     # 从数据集中读取数据
7.     # train == True 读取训练集
8.     # train == False 读取测试集
9.     txt_fname = root + '/ImageSets/Segmentation/' + ('train.txt' if train else 'val.txt')
10.    with open(txt_fname, 'r') as f:
11.        images = f.read().split()
12.        data = [os.path.join(root, 'JPEGImages', i + '.jpg') for i in images]
13.        label = [os.path.join(root, 'SegmentationClass', i + '.png') for i in images]
14.    return data, label
```

Cutting image:

```
1. def crop_image(data, label, height, width):
2.     # 切割图像
```

2-G

```
3.     .....
4.     data is PIL.Image object
5.     label is PIL.Image object
6.     ...
7.     box = (0, 0, width, height)
8.     data = data.crop(box)
9.     label = label.crop(box)
10.    return data, label
```

Define labels and colours for each category

```
1. # 定义数据集每个类别的标签
2. classes = ['background', 'aeroplane', 'bicycle', 'bird', 'boat',
3.            'bottle', 'bus', 'car', 'cat', 'chair', 'cow', 'diningtable',
4.            'dog', 'horse', 'motorbike', 'person', 'potted plant',
5.            'sheep', 'sofa', 'train', 'tv/monitor']
6.
7. # 定义数据集每个类别的显示颜色(RGB)
8. colormap = [[0, 0, 0], [128, 0, 0], [0, 128, 0], [128, 128, 0], [0, 0, 128],
9.             [128, 0, 128], [0, 128, 128], [128, 128, 128], [64, 0, 0], [192, 0, 0],
10.            [64, 128, 0], [192, 128, 0], [64, 0, 128], [192, 0, 128],
11.            [64, 128, 128], [192, 128, 128], [0, 64, 0], [128, 64, 0],
12.            [0, 192, 0], [128, 192, 0], [0, 64, 128]]
13.
14. cm2lbl = np.zeros(256 ** 3)
15. for i, cm in enumerate(colormap):
16.     cm2lbl[(cm[0] * 256 + cm[1]) * 256 + cm[2]] = i
```

Fill the labels with the subscript information of the corresponding category according to the RGB values

```
1. def image2label(im):
2.     # 将标签按照 RGB 值填入对应类别的下标信息
3.     data = np.array(im, dtype="int32")
4.     idx = (data[:, :, 0] * 256 + data[:, :, 1]) * 256 + data[:, :, 2]
5.     # print(np.shape(cm2lbl), np.shape(idx))
6.     # print(cm2lbl[idx])
7.     return np.array(cm2lbl[idx], dtype="int32")
```

Cropping data to h*w size

```
1. def image_transforms(data, label, height, width):
2.     # 将数据裁切为 h*w 大小
3.     data, label = crop_image(data, label, height, width)
4.     # 将数据转换成 tensor, 并且做标准化处理
5.     im_tfs = tfs.Compose([
6.         # 将 PIL Image 或 numpy.ndarray 转换为 tensor, 并除 255 归一化到[0,1]之间
7.         tfs.ToTensor(),
8.         # 标准化处理-->转换为标准正太分布, 使模型更容易收敛
9.         tfs.Normalize(
10.            mean=[0.485, 0.456, 0.406],
11.            std=[0.229, 0.224, 0.225])
12.     ])
13.     data = im_tfs(data)
14.     label = image2label(label)
15.     label = torch.from_numpy(label)
16.     return data, label
```

Data pre-processing, filtered images

```
1. class VOCSegDataset(Dataset):
2.     # 数据预处理
3.
4.     # 构造函数
5.     def __init__(self, train, height, width, transforms):
6.         self.height = height
7.         self.width = width
8.         self.fnum = 0 # 用来记录被过滤的图片数
9.         self.transforms = transforms
10.        data_list, label_list = read_images(train=train)
11.        # 过滤不符合规则的图片
```

```

12.         self.data_list = self._filter(data_list)
13.         self.label_list = self._filter(label_list)
14.         if (train == True):
15.             print("训练集: 加载了 " + str(len(self.data_list)) + " 张图片和标签" + ",
过滤了" + str(self.fnum) + "张图片")
16.         else:
17.             print("测试集: 加载了 " + str(len(self.data_list)) + " 张图片和标签" + ",
过滤了" + str(self.fnum) + "张图片")
18.
19.     # 过滤掉长小于 height 和宽小于 width 的图片
20.     def _filter(self, images):
21.         img = []
22.         for im in images:
23.             if (Image.open(im).size[1] >= self.height and
24.                 Image.open(im).size[0] >= self.width):
25.                 img.append(im)
26.             else:
27.                 self.fnum += 1
28.         return img
29.
30.     # 重载 getitem 函数, 使类可以迭代
31.     def __getitem__(self, idx):
32.         data = self.data_list[idx]
33.         label = self.label_list[idx]
34.         data = Image.open(data)
35.         label = Image.open(label).convert('RGB')
36.         data, label = self.transforms(data, label, self.height, self.width)
37.         return data, label
38.
39.     def __len__(self):
40.         return len(self.data_list)

```

Load the dataset and output the number of filtered images and the number loaded

```

1. height = 224
2. width = 224
3.
4. voc_train = VOCSegDataset(True, height, width, image_transforms)
5. voc_test = VOCSegDataset(False, height, width, image_transforms)
6. train_data = DataLoader(voc_train, batch_size=8, shuffle=True)
7. valid_data = DataLoader(voc_test, batch_size=8)

```

训练集: 加载了 1456 张图片和标签, 过滤了16张图片

测试集: 加载了 1436 张图片和标签, 过滤了26张图片

Display of data sets

```

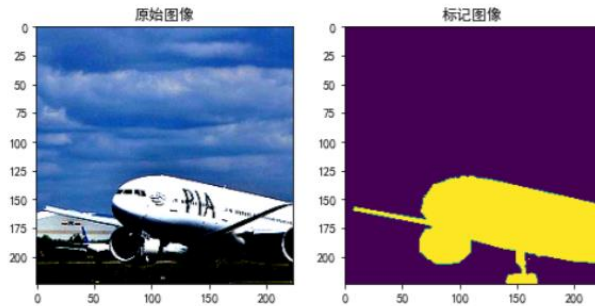
1. x = np.array(voc_train[9][0]).transpose(1,2,0)
2. y = np.array(voc_train[9][1])
3. print(np.shape(y))
4. fig, axs = plt.subplots(1, 2, figsize = (8, 8))
5. axs[0].imshow(x)
6. axs[0].set_title("原始图像")
7. axs[1].imshow(y)
8. axs[1].set_title("标记图像")

```

9. plt.show()

Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).

(224, 224)



3-A

Constructing the kernel: Here we use a bilinear interpolation method to construct the kernel, calculating each output y_{ij} from the four nearest inputs through a linear mapping that depends only on the relative positions of the input and output units

```
1. def bilinear_kernel(in_channels, out_channels, kernel_size):
2.     '''
3.     return a bilinear filter tensor
4.     '''
5.     factor = (kernel_size + 1) // 2
6.     if kernel_size % 2 == 1:
7.         center = factor - 1
8.     else:
9.         center = factor - 0.5
10.    og = np.ogrid[:kernel_size, :kernel_size]
11.    # 生成横向和纵向的 0-kernel_size、步长为一的两个二维数组
12.    filt = (1 - abs(og[0] - center) / factor) * (1 - abs(og[1] - center) / factor)
13.    weight = np.zeros((in_channels, out_channels, kernel_size, kernel_size), dtype='
float32')
14.    weight[range(in_channels), range(out_channels), :, :] = filt
15.    return torch.from_numpy(weight)
```

When constructing fcn, we use the pre-trained resnet34 instead of vgg

```
1. model_urls = {
2.     'resnet18': 'https://download.pytorch.org/models/resnet18-5c106cde.pth',
3.     'resnet34': 'https://download.pytorch.org/models/resnet34-333f7ec4.pth',
4.     'resnet50': 'https://download.pytorch.org/models/resnet50-19c8e357.pth',
5.     'resnet101': 'https://download.pytorch.org/models/resnet101-5d3b4d8f.pth',
6.     'resnet152': 'https://download.pytorch.org/models/resnet152-b121ed2d.pth',
7. }
8.
9. # 使用预训练的 resnet 34 (34 层残差网络) 代替论文中的 vgg 实现 fcn
10. model_root = "./model/resnet34-333f7ec4.pth"
11. pretrained_net = models.resnet34(pretrained=False)
12. pre = torch.load(model_root)
13. pretrained_net.load_state_dict(pre)
14. # 分类的总数
15. num_classes = len(classes)
```

Construction of fully convolutional neural network (FCN) models

```
1. class fcn(nn.Module):
2.     def __init__(self, num_classes):
3.         super(fcn, self).__init__()
4.
```

```

5.         # 第一段, 通道数为 128, 输出特征图尺寸为 28*28
6.         # conv1, conv2_x, conv3_x
7.         self.stage1 = nn.Sequential(*list(pretrained_net.children())[:-4])
8.         # 第二段, 通道数为 256, 输出特征图尺寸为 14*14
9.         # conv4_x
10.        self.stage2 = list(pretrained_net.children())[-4]
11.        # 第三段, 通道数为 512, 输出特征图尺寸为 7*7
12.        # conv5_x
13.        self.stage3 = list(pretrained_net.children())[-3]
14.
15.        # 三个 kernel 为 1*1 的卷积操作, 各个通道信息融合
16.        self.scores1 = nn.Conv2d(512, num_classes, 1)
17.        self.scores2 = nn.Conv2d(256, num_classes, 1)
18.        self.scores3 = nn.Conv2d(128, num_classes, 1)
19.
20.        # 反卷积, 将特征图尺寸放大八倍
21.        self.upsample_8x = nn.ConvTranspose2d(num_classes, num_classes, kernel_size=
16, stride=8, padding=4, bias=False)
22.        self.upsample_8x.weight.data = bilinear_kernel(num_classes, num_classes, 16)
23.
24.        # 反卷积, 将特征图尺寸放大两倍
25.        self.upsample_2x_1 = nn.ConvTranspose2d(num_classes, num_classes, kernel_size=
4, stride=2, padding=1, bias=False)
26.        self.upsample_2x_1.weight.data = bilinear_kernel(num_classes, num_classes, 4
)
27.        self.upsample_2x_2 = nn.ConvTranspose2d(num_classes, num_classes, kernel_size=
4, stride=2, padding=1, bias=False)
28.        self.upsample_2x_2.weight.data = bilinear_kernel(num_classes, num_classes, 4
)
29.
30.        def forward(self, x):
31.            x = self.stage1(x)
32.            s1 = x # 224/8 = 28
33.
34.            x = self.stage2(x)
35.            s2 = x # 224/16 = 14
36.
37.            x = self.stage3(x)
38.            s3 = x # 224/32 = 7
39.
40.            # 将各通道信息融合
41.            s3 = self.scores1(s3)
42.            # 上采样 放大二倍
43.            s3 = self.upsample_2x_1(s3)
44.
45.            s2 = self.scores2(s2)
46.            # 将二三层训练特征合成 14*14
47.            s2 = s2 + s3
48.            s2 = self.upsample_2x_2(s2)
49.
50.            # 28*28
51.            s1 = self.scores3(s1)
52.            # 将一二三层训练特征合成 28*28
53.            s = s1 + s2
54.            # 将 s 放大八倍, 变为原图像尺寸 224*224
55.            s = self.upsample_8x(s)
56.
57.            # 返回特征图
58.            return s

```

FCN network structure.

print(net)

```
(stage1): Sequential(
  (0): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)
  (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (2): ReLU(inplace=True)
  (3): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
  (4): Sequential(
    (0): BasicBlock(
      (conv1): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
    (1): BasicBlock(
      (conv1): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
    (2): BasicBlock(
      (conv1): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)
```

3-B


```

(stage3): Sequential(
  (0): BasicBlock(
    (conv1): Conv2d(256, 512, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (downsample): Sequential(
      (0): Conv2d(256, 512, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (1): BasicBlock(
    (conv1): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  )
  (2): BasicBlock(
    (conv1): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  )
)
(scores1): Conv2d(512, 21, kernel_size=(1, 1), stride=(1, 1))
(scores2): Conv2d(256, 21, kernel_size=(1, 1), stride=(1, 1))
(scores3): Conv2d(128, 21, kernel_size=(1, 1), stride=(1, 1))
(upsample_8x): ConvTranspose2d(21, 21, kernel_size=(16, 16), stride=(8, 8), padding=(4, 4), bias=False)
(upsample_2x_1): ConvTranspose2d(21, 21, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
(upsample_2x_2): ConvTranspose2d(21, 21, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)

```

3-E

Model evaluation indicators, calculation of confusion matrix and determination of accuracy.

```

1. # 计算混淆矩阵
2. def _fast_hist(label_true, label_pred, n_class):
3.     # mask 在和 label_true 相对应的索引的位置上填入 true 或者 false
4.     # label_true[mask] 会把 mask 中索引为 true 的元素输出
5.     mask = (label_true >= 0) & (label_true < n_class)
6.     hist = np.bincount(
7.         n_class * label_true[mask].astype(int) +
8.         label_pred[mask], minlength=n_class ** 2).reshape(n_class, n_class)
9.     return hist
10.
11.
12. """
13. label_trues 正确的标签值
14. label_preds 模型输出的标签值
15. n_class 数据集中的分类数
16. """
17. def label_accuracy_score(label_trues, label_preds, n_class):
18.     """Returns accuracy score evaluation result.
19.         - overall accuracy
20.         - mean accuracy
21.         - mean IU
22.         - fwavacc
23.     """
24.     hist = np.zeros((n_class, n_class))
25.     # 通过迭代器将一个个数据进行计算
26.     for lt, lp in zip(label_trues, label_preds):
27.         hist += _fast_hist(lt.flatten(), lp.flatten(), n_class)
28.
29.     # np.diag(a) 假如 a 是一个二维矩阵，那么会输出矩阵的对角线元素
30.     # np.sum() 可以计算出所有元素的和。如果 axis=1，则表示按行相加
31.     acc = np.diag(hist).sum() / hist.sum()
32.     # np.diag 以一维数组的形式返回方阵的对角线
33.     acc_cls = np.diag(hist) / hist.sum(axis=1)
34.     # nanmean 会自动忽略 nan 的元素求平均
35.     acc_cls = np.nanmean(acc_cls)
36.     iu = np.diag(hist) / (hist.sum(axis=1) + hist.sum(axis=0) - np.diag(hist))
37.     mean_iu = np.nanmean(iu)
38.     freq = hist.sum(axis=1) / hist.sum()
39.     fwavacc = (freq[freq > 0] * iu[freq > 0]).sum()
40.

```

```
41.     return acc, acc_cls, mean_iu, fwavacc
```

Defining hyperparameters and saving training data.

```
1. net = fcn(num_classes)
2. PATH = "./model/fcn-resnet34.pth"
3. net.load_state_dict(torch.load(PATH))
4. if torch.cuda.is_available():
5.     net = net.cuda()
6.
7. # 训练时的数据
8. train_loss = []
9. train_acc = []
10. train_acc_cls = []
11. train_mean_iu = []
12. train_fwavacc = []
13.
14. # 验证时的数据
15. eval_loss = []
16. eval_acc = []
17. eval_acc_cls = []
18. eval_mean_iu = []
19. eval_fwavacc = []
```

Conducting network training.

```
1. # 损失
2. criterion = nn.NLLLoss()
3.
4. # 加速器 SGD
5. Eta = 1e-2
6. basic_optim = torch.optim.SGD(net.parameters(), lr=Eta, weight_decay=1e-4)
7. optimizer = basic_optim
8.
9. # 网络训练
10. EPOCHES = 60
11.
12. exp_lr_scheduler = lr_scheduler.StepLR(optimizer, step_size=7, gamma=0.1)
```

Model training.

```
1. for e in range(EPOCHES):
2.
3.     _train_loss = 0 # 记录一轮训练总损失
4.     _train_acc = 0
5.     _train_acc_cls = 0
6.     _train_mean_iu = 0
7.     _train_fwavacc = 0
8.     exp_lr_scheduler.step()
9.     prev_time = datetime.now()
10.    net = net.train()
11.    for img_data, img_label in train_data:
12.        if torch.cuda.is_available():
13.            im = Variable(img_data).cuda()
14.            label = torch.tensor(img_label, dtype=torch.int64)
15.            label = Variable(label).cuda()
16.        else:
17.            im = Variable(img_data)
18.            label = torch.tensor(img_label, dtype=torch.int64)
19.            label = Variable(label)
20.
21.        # 前向传播
22.        out = net(im)
```

```

23.         out = F.log_softmax(out, dim=1)
24.         loss = criterion(out, label)
25.
26.         # 反向传播
27.         # 梯度清零
28.         optimizer.zero_grad()
29.         loss.backward()
30.         # 更新
31.         optimizer.step()
32.         _train_loss += loss.item()
33.
34.         # label_pred 输出的是 21*224*224 的向量，对于每一个点都有 21 个分类的概率
35.         # 我们取概率值最大的那个下标作为模型预测的标签，然后计算各种评价指标
36.         label_pred = out.max(dim=1)[1].data.cpu().numpy()
37.         label_true = label.data.cpu().numpy()
38.         # label_pred: (8, 224, 224) label_true: (8, 224, 224)
39.         for lbt, lbp in zip(label_true, label_pred):
40.             # lbt: (224, 224) lbp: (224, 224)
41.             acc, acc_cls, mean_iu, fwavacc = label_accuracy_score(lbt, lbp, num_classes)
42.             _train_acc += acc
43.             _train_acc_cls += acc_cls
44.             _train_mean_iu += mean_iu
45.             _train_fwavacc += fwavacc
46.
47.         # 记录当前轮的数据
48.         train_loss.append(_train_loss / len(train_data))
49.         train_acc.append(_train_acc / len(voc_train))
50.         train_acc_cls.append(_train_acc_cls)
51.         train_mean_iu.append(_train_mean_iu / len(voc_train))
52.         train_fwavacc.append(_train_fwavacc)
53.
54.         net = net.eval()
55.
56.         _eval_loss = 0
57.         _eval_acc = 0
58.         _eval_acc_cls = 0
59.         _eval_mean_iu = 0
60.         _eval_fwavacc = 0
61.
62.         for img_data, img_label in valid_data:
63.             if torch.cuda.is_available():
64.                 im = Variable(img_data).cuda()
65.                 label = torch.tensor(img_label, dtype=torch.int64)
66.                 label = Variable(label).cuda()
67.             else:
68.                 im = Variable(img_data)
69.                 label = torch.tensor(img_label, dtype=torch.int64)
70.                 label = Variable(label)
71.
72.             # forward
73.             out = net(im)
74.             # 对结果进行归一化
75.             out = F.log_softmax(out, dim=1)
76.             loss = criterion(out, label)
77.             _eval_loss += loss.item()
78.
79.             label_pred = out.max(dim=1)[1].data.cpu().numpy()
80.             label_true = label.data.cpu().numpy()
81.             for lbt, lbp in zip(label_true, label_pred):
82.                 acc, acc_cls, mean_iu, fwavacc = label_accuracy_score(lbt, lbp, num_classes)
83.                 _eval_acc += acc
84.                 _eval_acc_cls += acc_cls
85.                 _eval_mean_iu += mean_iu

```

```

86.         _eval_fwavacc += fwavacc
87.
88.     # 记录当前轮的数据
89.     eval_loss.append(_eval_loss / len(valid_data))
90.     eval_acc.append(_eval_acc / len(voc_test))
91.     eval_acc_cls.append(_eval_acc_cls)
92.     eval_mean_iu.append(_eval_mean_iu / len(voc_test))
93.     eval_fwavacc.append(_eval_fwavacc)
94.
95.     # 打印当前轮训练的结果
96.     cur_time = datetime.now()
97.     h, remainder = divmod((cur_time - prev_time).seconds, 3600)
98.     # divmod() 函数返回当参数 1 除以参数 2 时包含商和余数的元组。
99.     m, s = divmod(remainder, 60)
100.    epoch_str = ('Epoch: {}, Train Loss: {:.5f}, Train Acc: {:.5f}, Train Mean IU:
        {:.5f}, '
101.                'Valid Loss: {:.5f}, Valid Acc: {:.5f}, Valid Mean IU: {:.5f} '.fo
        rmat(
102.            e, _train_loss / len(train_data), _train_acc / len(voc_train), _train_mean_
        iu / len(voc_train),
103.            _eval_loss / len(valid_data), _eval_acc / len(voc_test), _eval_mean_iu /
        len(voc_test)))
104.    time_str = 'Time: {:.0f}:{:.0f}:{:.0f}'.format(h, m, s)
105.    print(epoch_str + time_str)

```

Training process.

```

s: 0.49830, Valid Acc: 0.89420, Valid Mean IU: 0.57202 Time: 0:0:42
Epoch: 27, Train Loss: 0.04224, Train Acc: 0.98345, Train Mean IU: 0.89445, Valid Los
s: 0.49892, Valid Acc: 0.89381, Valid Mean IU: 0.57287 Time: 0:0:42
Epoch: 28, Train Loss: 0.04185, Train Acc: 0.98349, Train Mean IU: 0.89513, Valid Los
s: 0.50622, Valid Acc: 0.89322, Valid Mean IU: 0.57142 Time: 0:0:49
Epoch: 29, Train Loss: 0.04206, Train Acc: 0.98348, Train Mean IU: 0.89375, Valid Los
s: 0.50191, Valid Acc: 0.89426, Valid Mean IU: 0.57371 Time: 0:0:48
Epoch: 30, Train Loss: 0.04165, Train Acc: 0.98360, Train Mean IU: 0.89621, Valid Los
s: 0.49704, Valid Acc: 0.89346, Valid Mean IU: 0.56481 Time: 0:0:49
Epoch: 31, Train Loss: 0.04201, Train Acc: 0.98344, Train Mean IU: 0.89478, Valid Los
s: 0.49330, Valid Acc: 0.89418, Valid Mean IU: 0.57040 Time: 0:0:48
Epoch: 32, Train Loss: 0.04198, Train Acc: 0.98345, Train Mean IU: 0.89484, Valid Los
s: 0.49907, Valid Acc: 0.89426, Valid Mean IU: 0.57120 Time: 0:0:48
Epoch: 33, Train Loss: 0.04201, Train Acc: 0.98351, Train Mean IU: 0.89314, Valid Los
s: 0.49317, Valid Acc: 0.89445, Valid Mean IU: 0.57292 Time: 0:0:47
Epoch: 34, Train Loss: 0.04163, Train Acc: 0.98364, Train Mean IU: 0.89520, Valid Los
s: 0.49903, Valid Acc: 0.89384, Valid Mean IU: 0.57186 Time: 0:0:48
Epoch: 35, Train Loss: 0.04166, Train Acc: 0.98360, Train Mean IU: 0.89532, Valid Los
s: 0.49857, Valid Acc: 0.89341, Valid Mean IU: 0.56861 Time: 0:0:48

```

3-F

Saving model training results and plotting curves.

```

1. PATH = "./model/fcn-resnet34.pth"
2. torch.save(net.state_dict(), PATH)
3. plt.plot(np.array(train_loss))
4. plt.title("Train Loss")
5. plt.figure()
6. plt.plot(np.array(train_acc))
7. plt.title("Pixel Accuracy")
8. plt.figure()
9. plt.plot(np.array(train_mean_iu))
10. plt.title("Mean IU")

```

Test training results.

```

1. print("Eval Pixel Accuracy: ", np.array(eval_acc)[-1])
2. print("Eval Mean IU: ", np.array(eval_mean_iu)[-1])

```

```

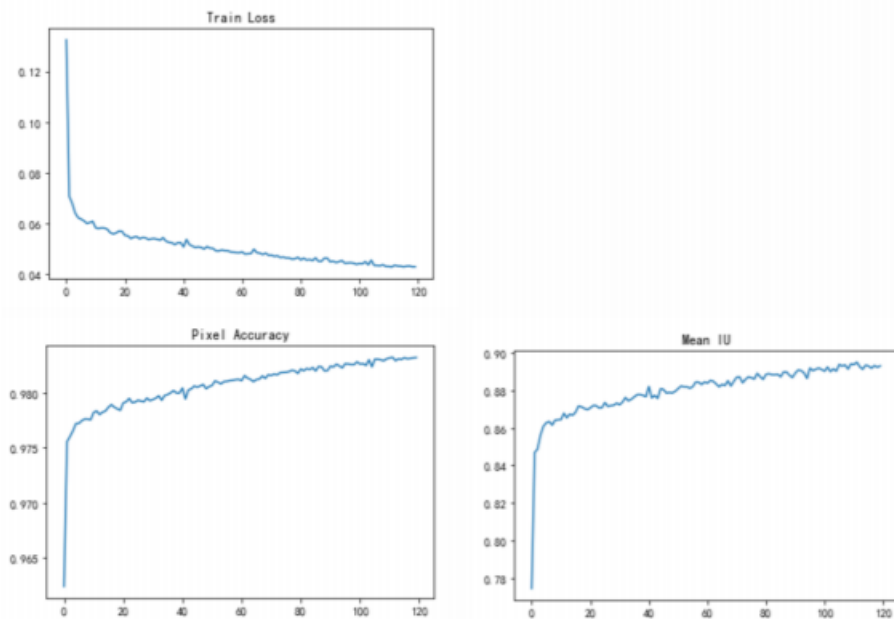
Eval Pixel Accuracy: 0.8944751938358038
Eval Mean IU: 0.5766601620395866

```

3-G

3.3 Visualisation analysis

Plotting training data evaluation curves.



3-H

The accuracy rate shows that the 89.4% accuracy of our model is closer to the data in the paper.

	pixel acc.	mean acc.	mean IU	f.w. IU
FCN-32s-fixed	83.0	59.7	45.4	72.0
FCN-32s	89.1	73.3	59.4	81.4
FCN-16s	90.0	75.7	62.4	83.0
FCN-8s	90.3	75.8	62.7	83.3

3-I

Visualisation of the test set prediction results.

```

1. # 加载模型
2. net = fcn(num_classes)
3. PATH = "./model/fcn-resnet34.pth"
4. net.load_state_dict(torch.load(PATH))
5. if torch.cuda.is_available():
6.     net = net.cuda()
7. cm = np.array(colormap).astype('uint8')
8. size = 224
9. num_image = 10
10. def predict(img, label): # 预测结果
11.     img = Variable(img.unsqueeze(0)).cuda()
12.     out = net(img)
13.     pred = out.max(1)[1].squeeze().cpu().data.numpy()
14.     # 将 pred 的分类值, 转换成各个分类对应的 RGB 值

```

```

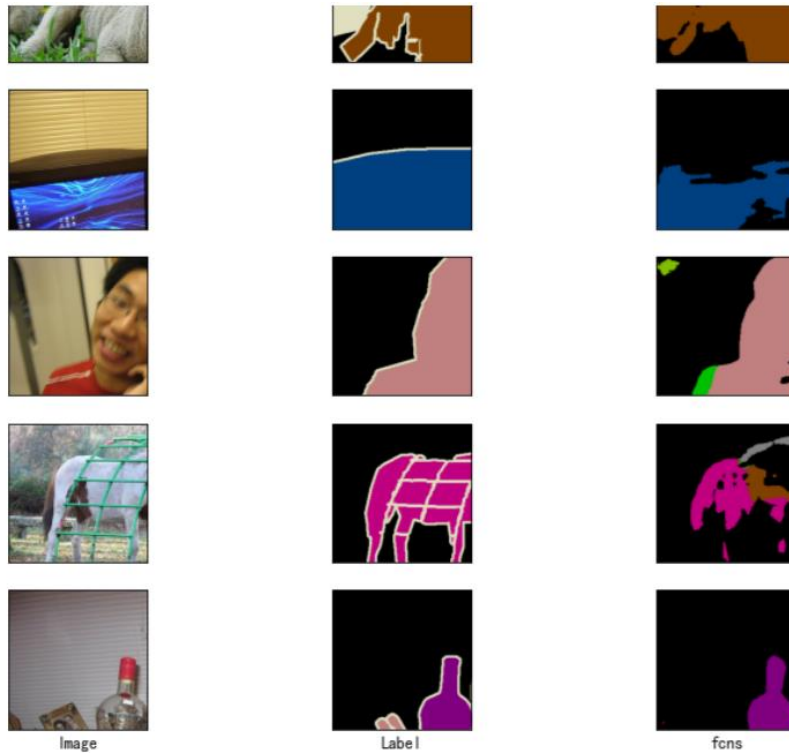
15.     pred = cm[pred]
16.     # 将 numpy 转换成 PIL 对象
17.     pred = Image.fromarray(pred)
18.     label = cm[label.numpy()]
19.     return pred, label
20. _, figs = plt.subplots(num_image, 3, figsize=(12, 22))
21. for i in range(num_image):
22.     img_data, img_label = voc_test[i]
23.     pred, label = predict(img_data, img_label)
24.     img_data = Image.open(voc_test.data_list[i])
25.     img_label = Image.open(voc_test.label_list[i]).convert("RGB")
26.     img_data, img_label = crop_image(img_data, img_label, 224, 224)
27.     figs[i, 0].imshow(img_data) # 原始图片
28.     figs[i, 0].axes.get_xaxis().set_visible(False) # 去掉 x 轴
29.     figs[i, 0].axes.get_yaxis().set_visible(False) # 去掉 y 轴
30.     figs[i, 1].imshow(img_label) # 标签
31.     figs[i, 1].axes.get_xaxis().set_visible(False) # 去掉 x 轴
32.     figs[i, 1].axes.get_yaxis().set_visible(False) # 去掉 y 轴
33.     figs[i, 2].imshow(pred) # 模型输出结果
34.     figs[i, 2].axes.get_xaxis().set_visible(False) # 去掉 x 轴
35.     figs[i, 2].axes.get_yaxis().set_visible(False) # 去掉 y 轴
36.
37. # 在最后一行图片下面添加标题
38. figs[num_image - 1, 0].set_title("Image", y=-0.2)
39. figs[num_image - 1, 1].set_title("Label", y=-0.2)
40. figs[num_image - 1, 2].set_title("fcns", y=-0.2)

```

Text(0.5, -0.2, 'fcns')



3-J



3-K

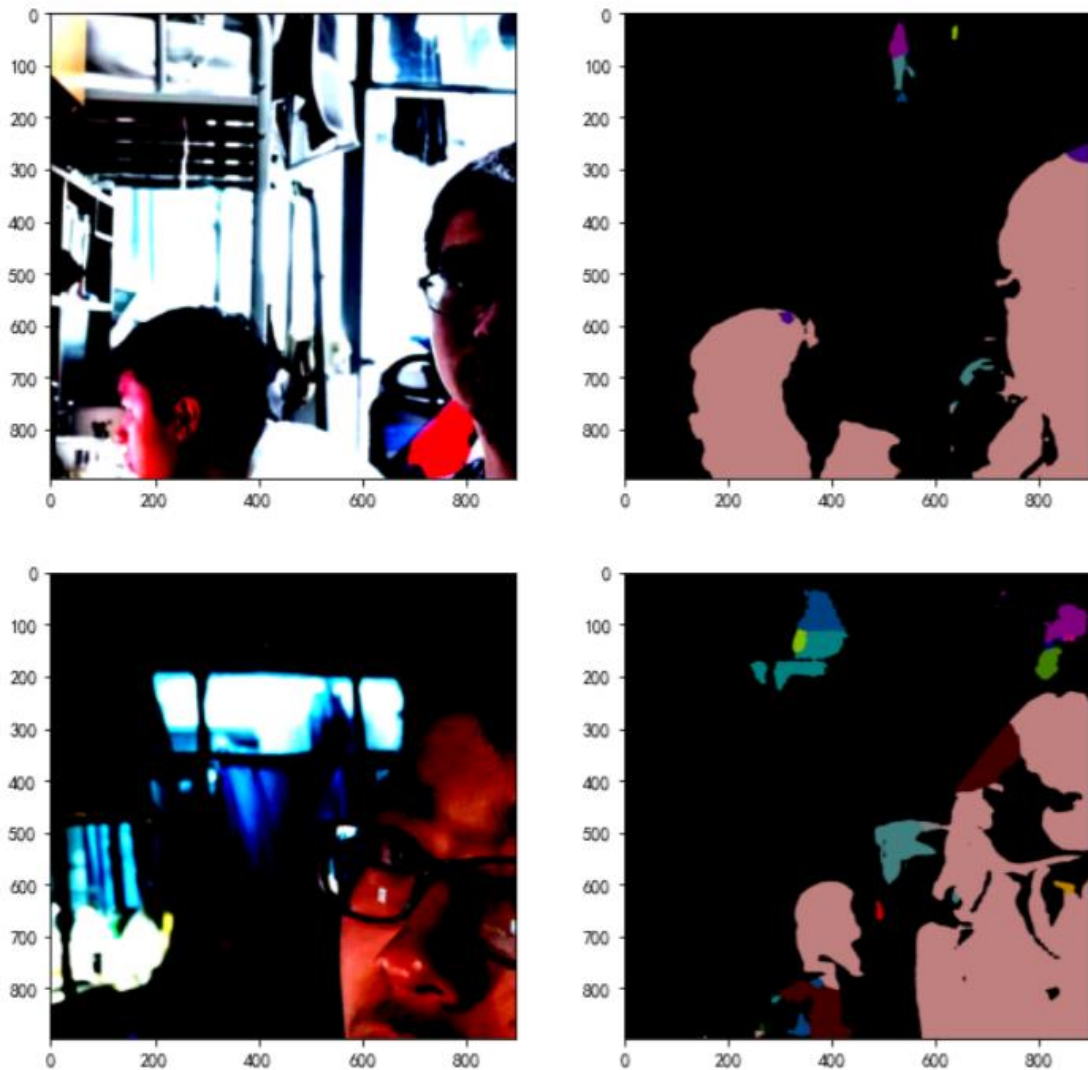
Testing of model segmentation using everyday photographs.

```

1. MY_DATA_DIR = './data'
2. print(os.listdir(MY_DATA_DIR+'\\Mydata')[:10])
3. Is = os.listdir(MY_DATA_DIR+'\\Mydata')
4. data = [os.path.join(MY_DATA_DIR, 'Mydata', i) for i in Is]
5. height = 896
6. width = 896
7. def My_Image(data, height, width):
8.     img = []
9.     for im in data:
10.         if (Image.open(im).size[1] >= height and Image.open(im).size[0] >= width):
11.             img.append(im)
12.     img = [Image.open(i) for i in img]
13.     box = (0, 0, width, height)
14.     img = [i.crop(box) for i in img]
15.     im_tfs = tfs.Compose([
16.         tfs.ToTensor(),
17.         tfs.Normalize(
18.             mean=[0.485, 0.456, 0.406],
19.             std=[0.229, 0.224, 0.225])
20.     ])
21.     img = [im_tfs(i) for i in img]
22.     return img
23. def predict_img():
24.     img = Variable(img.unsqueeze(0)).cuda()
25.     out = net(img)
26.     pred = out.max(1)[1].squeeze().cpu().data.numpy()
27.     pred = cm[pred]
28.     pred = Image.fromarray(pred)
29.     return pred
30. _, figs = plt.subplots(len(data), 2, figsize=(10, 10))
31. for i in range(len(data)):
32.     pred = predict_img[i]
33.     figs[i, 0].imshow(img[i].data.numpy().transpose(1,2,0))

```

```
34.     figs[i, 1].imshow(pred)
```



3-L

3.4 Experimental experience

Semantic segmentation, both textual and image semantic segmentation, has been a remarkable achievement in the field of artificial intelligence in the last decade, not only in many areas of industry, but also in our understanding of the extraction process and the intrinsic composition of semantic information. In real life we recognise an 'object' as if we are naturally able to extract information about its boundaries and, in the process of further understanding, to clearly distinguish the different details of the subject to which it belongs. The success of machine learning has led us to reflect on this phenomenon and to realise that it is possible to attach this phenomenon (behaviour) to a machine through the language of mathematics, and that semantic segmentation of images at pixel level is undoubtedly an important part of what can be achieved in the future with general or strong artificial intelligence. FCN, the pioneer of pixel-level semantic segmentation of images, has pioneered the introduction of deconvolution into the traditional convolutional neural network-full-connected neural network model of image segmentation processing, making it possible not only to achieve unlimited scale of the input image, but also to give new ideas for semantic edge processing of images. We have learnt about the fundamentals of image semantic segmentation, the effect of perceptual field on network structure, residual networks and other related knowledge, and become proficient in the use of frameworks such as pytorch, which provides a basis for further study. The pixel accuracy predicted by the model in this project was maintained at around 89.4%, and we collected some photos in our daily life, and the test can see that the model has better semantic segmentation results in a variety of environments.

However, the FCN-based semantic segmentation is still unable to achieve fine detail in dark areas and edge details. Future improvements are expected in terms of more detailed tuning of the network and further work on semantic segmentation algorithms such as u-net and segnet.

References

- [1] Jonathan Long, Evan Shelhamer, Trevor Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR2015 paper. 2015 3431-3440.
- [2] 李梦怡, 朱定局. 基于全卷积网络的图像语义分割方法综述. 华南师范大学计算机学院. 2019
- [3] 张建, 基于深度学习的图像语义分割方法. 电子科技大学. 2018.
- [4] 章琳, 袁非牛, 张文睿, 曾夏玲. 全卷积神经网络研究综述. 计算机工程与应用. 2020(01)
- [5] 刘健, 袁谦, 吴广, 喻晓. 卷积神经网络综述. 计算机时代. 2018(11)

Group division of labour

张梓琦: Theoretical studies, Paper preparation(1/3)

关乐平: Data set collection and testing,PPT Production,(1/3)

刘鸿嘉: Course defence, Code implementation(1/3)