

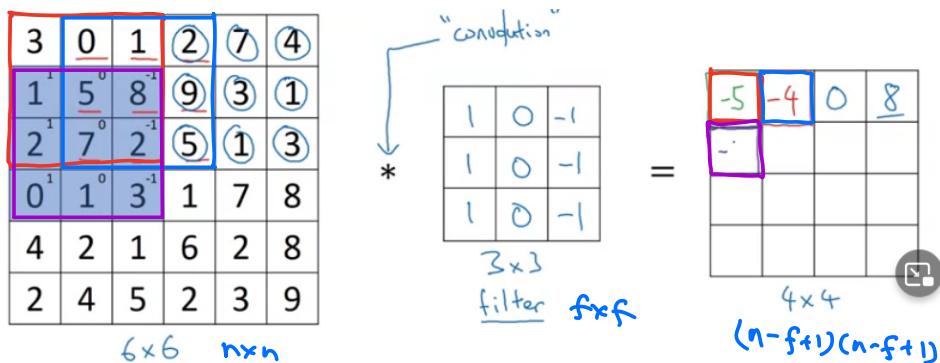
# CNN week 1

## Computer vision

- Image classification
- Object detection
- Neural Style Transfer

## Edge detection Example

Vertical edge detection



Attention: just element-wise multiple! (不是矩阵乘法!)  
not the multiplication of matrix!

Implement: python: conv\_forward

tensorflow: tf.nn.conv2d

keras: Conv2D

Different filter:

1	0	-1
2	0	-2
1	0	-1

Sobel filter

3	0	-3
0	0	-10
3	0	-3

Scharr filter

$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

treat all the number as parameter!  
And using backprop get these values

## Padding

Downsides of above method

- Picture becomes smaller and smaller after each convolution. (shrinkage output)
- The edge information can be used a few time (throw away info from edge)

Solution:

Padding

0	0	0	0	0	0	0	0
0	3	0	1	2	7	4	0
0	1	5	8	9	3	1	0
0	2	7	2	5	1	3	0
0	0	1	3	1	7	8	0
0	4	2	1	6	2	8	0
0	2	4	5	2	3	9	0
0	0	0	0	0	0	0	0

6x6

Question: how to choose the number of padding?

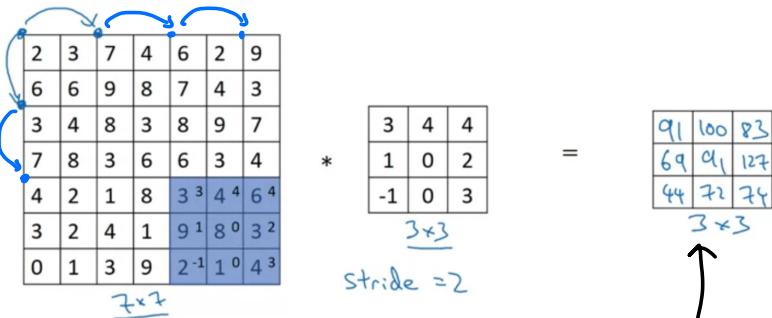
A: Valid convolution:  $n \times n * f \times f \rightarrow (n-f+1) \times (n-f+1)$

Same convolution: Pad so that output size is the same as the input size.

$$(n+2p-f+1) \times (n+2p-f+1) = n$$

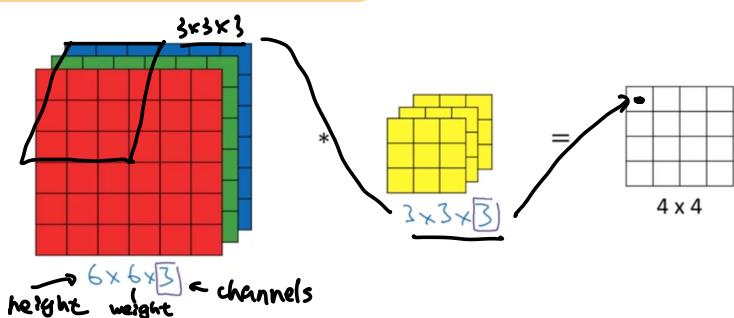
$$\Rightarrow p = f-1/2 \quad (f \text{ is usually odd})$$

### Strided Convolutions

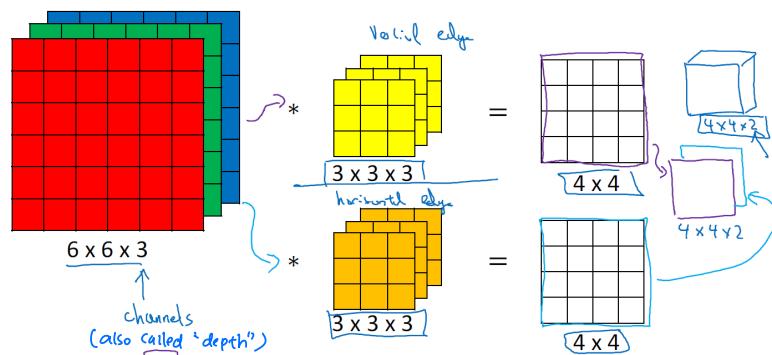


The output is  $\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$   
e.g. In above example:  $\frac{7+2 \times 0 - 3}{2} + 1 = 3$

### Convolutions Over Volume



## Multiple filters



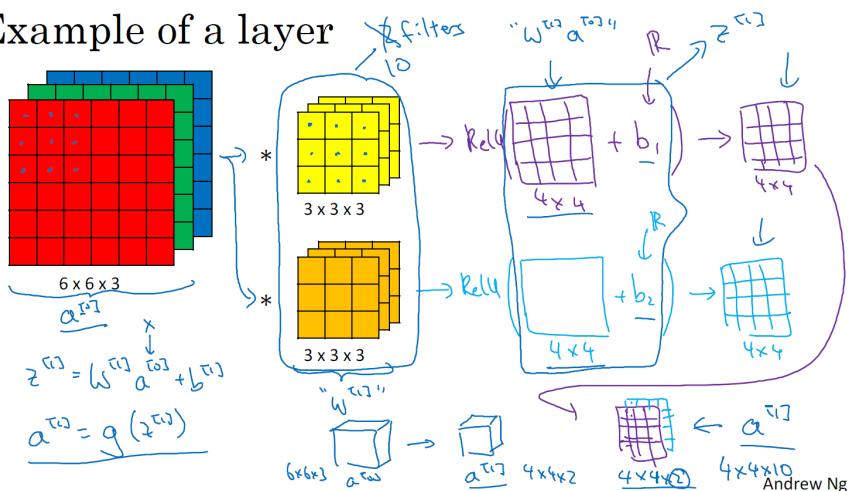
$$\text{Summary: } n \times n \times n_c \times f \times f \times n_c = (n-f+1) \times (n-f+1)$$

e.g.  $6 \times 6 \times 3 \times 4 \times 4 \times 2 = 4 \times 4$

## One layer of a convolutional Network

The comparison between neural network and CNN

Example of a layer



Question: If you have 10 filters that are  $3 \times 3 \times 3$  in one layer of a neural network, how many parameters does that layer have?

Answer:

$$3 \times 3 \times 3 + 1 = 28 \text{ parameters} \quad \times 10 = \underline{\underline{280}}$$

$\begin{matrix} \boxed{3} \\ \downarrow \\ b \end{matrix}$  - - -  $\times 10$

Summary of notation:

$$f^{(l)} = \text{filter size}$$

$$p^{(l)} = \text{padding}$$

$$s^{(l)} = \text{stride}$$

$$n_c^{(l)} = \text{number of filters}$$

$$\text{Input: } n_H^{(l-1)} \times n_W^{(l-1)} \times n_C^{(l-1)}$$

$$\text{Output: } n_H^{(l)} \times n_W^{(l)} \times n_C^{(l)}$$

$$\text{where } n_H^{(l)} = \left\lfloor \frac{n_{H^{(l-1)}} + 2p^{(l)} - f^{(l)}}{S^{(l)}} + 1 \right\rfloor$$

$$n_W^{(l)} = \left\lfloor \frac{n_{W^{(l-1)}} + 2p^{(l)} - f^{(l)}}{S^{(l)}} + 1 \right\rfloor$$

Each filter is  $f^{(l)} \times f^{(l)} \times n_C^{(l-1)}$

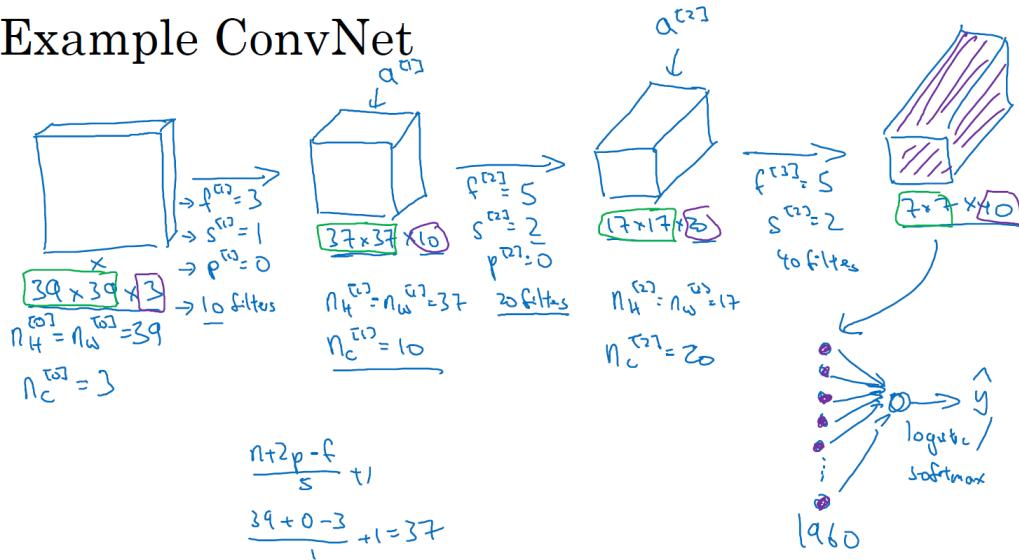
Activation:  $a^{(l)} \rightarrow n_H^{(l)} \times n_W^{(l)} \times n_C^{(l)}$

Weights:  $f^{(l)} \times f^{(l)} \times n_C^{(l-1)} \times n_C^{(l)}$

bias:  $n_C^{(l)} - (1, 1, 1, n_C^{(l)})$

### A simple convolution network example

#### Example ConvNet



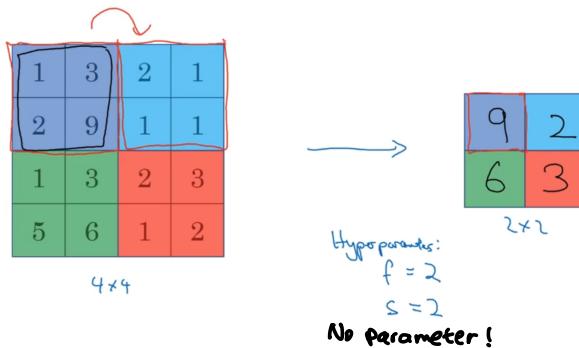
Note that  $\begin{cases} \text{the size of image } \downarrow \\ \text{the number of channels } \uparrow \end{cases}$

Type of layer in a convolution network:

- Convolution
- Pooling
- Fully connected

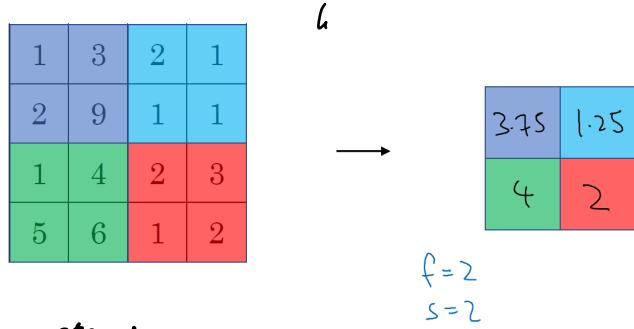
## Pooling layers

### Pooling layer: Max pooling



Note: the calculation of maximum is independent on each channel.

### Average pooling



Hyperparameters:

$f$ : filter size

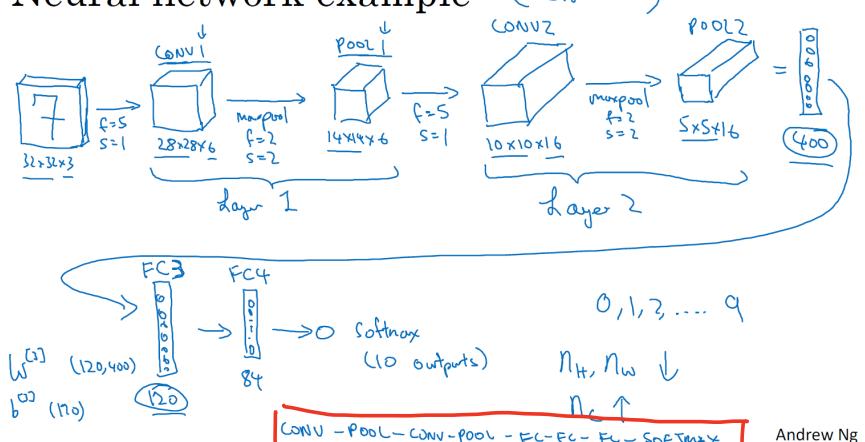
$s$ : stride

Max or average pooling

not use padding in common

## CNN Example

### Neural network example (LeNet-5)



## Why Convolutions?

- Parameter Sharing: A feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.
- Sparsity of Connection: In each layer, each output value depends only on a small number of inputs.

Some question:

1. the benefits of padding.

## CNN Week 2 Case Study

Outline:

- Classical networks:

- LeNet-5
- AlexNet
- VGG

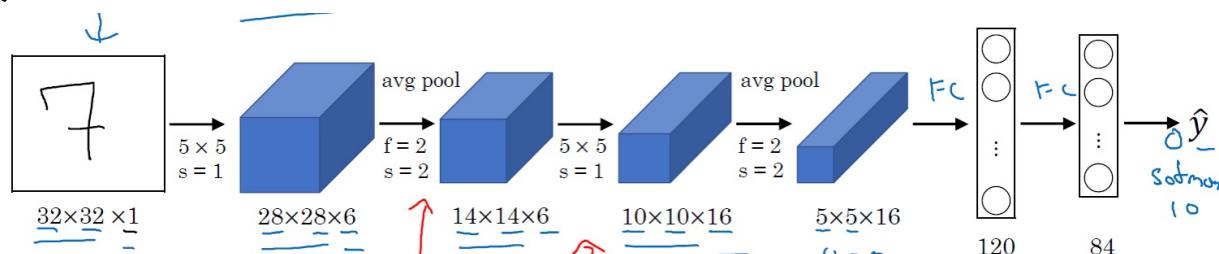
- ResNet

- Inception

- MobileNet

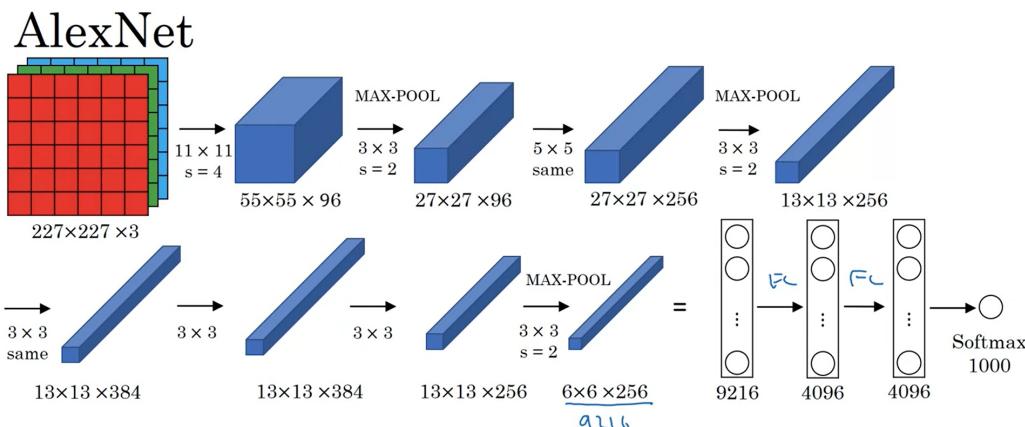
classical network

### LeNet-5



classical architecture: conv  $\rightarrow$  pool  $\rightarrow$  conv  $\rightarrow$  pool  $\rightarrow$  fc  $\rightarrow$  fc  $\rightarrow$  output

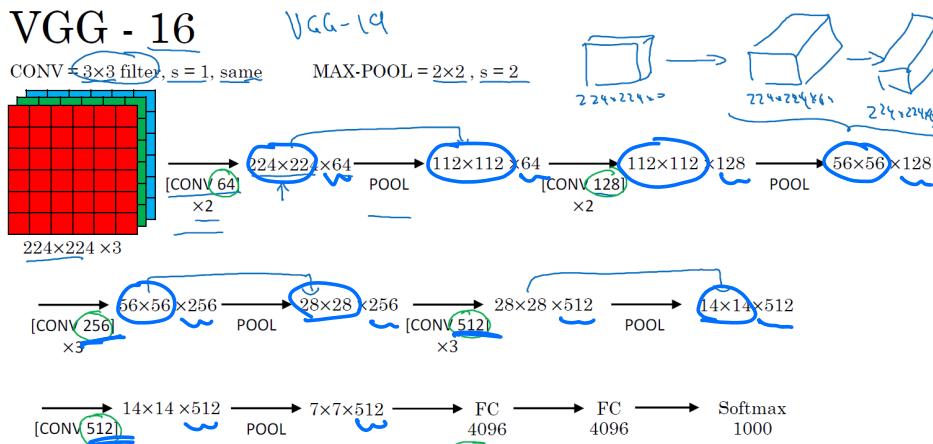
### AlexNet



- Similar to LeNet, but much bigger

- ReLU

## VGG - 16



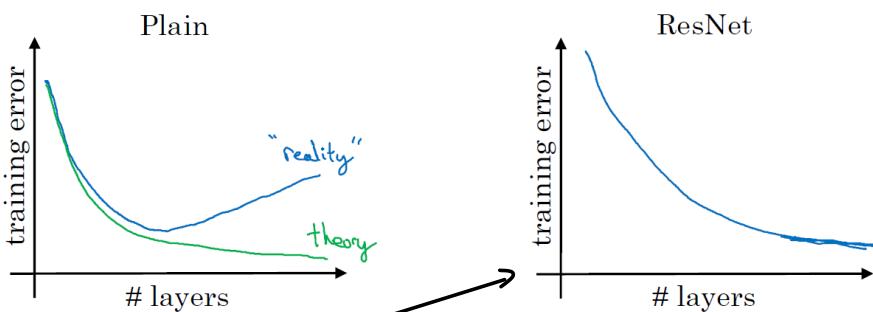
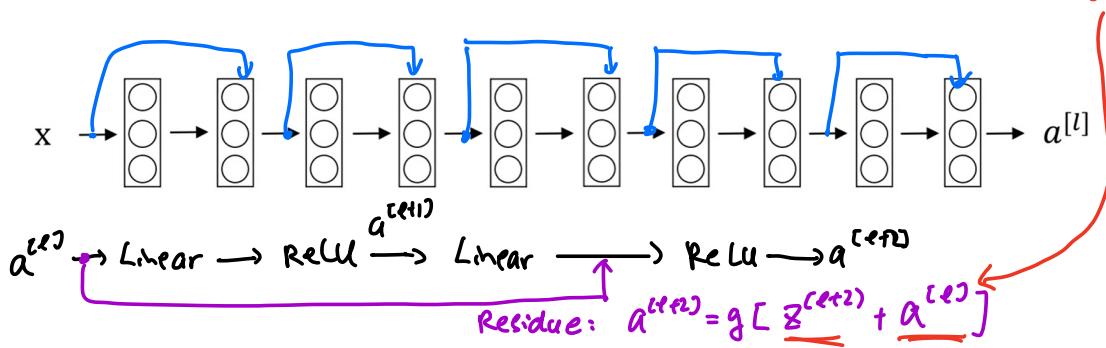
- Uniformity

- drawback: too many parameters

## Res Nets (Residual Networks)

question:

① If different dimension after mapping?



Why ResNets Work?  $\Delta$

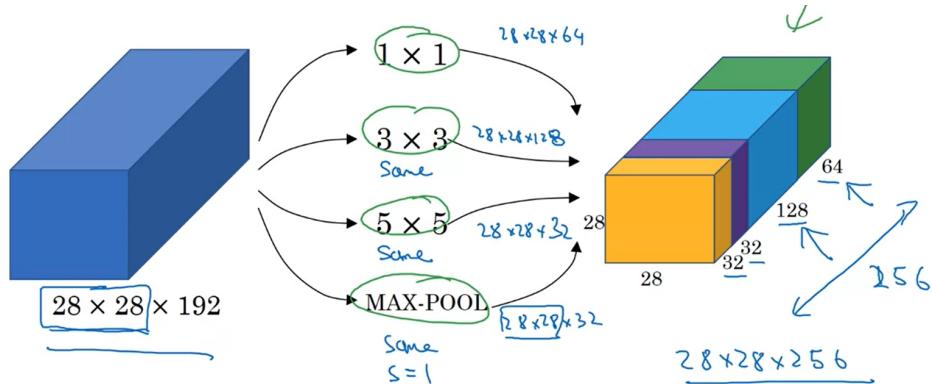
$$\begin{aligned}
 a^{[l+2]} &= g[z^{[l+2]} + a^{[l]}] \\
 &= g(W^{[l+2]}a^{[l+1]} + b^{[l+2]} + W_l a^{[l]}) = g(a^{[l]}) \\
 &\text{If } (W^{[l+2]}, b^{[l+2]}) = 0 \quad = a^{[l]}
 \end{aligned}$$

## One by One convolution (network in network)

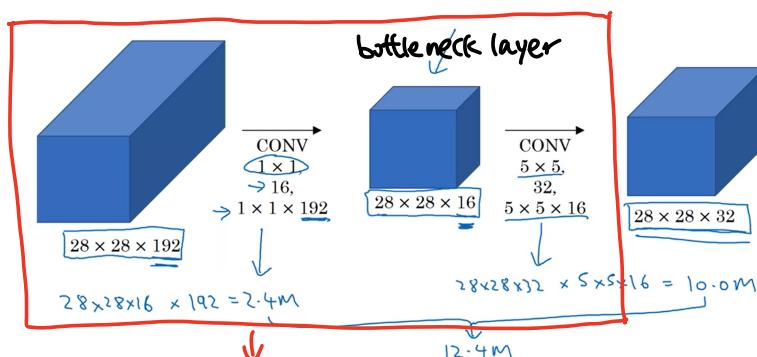
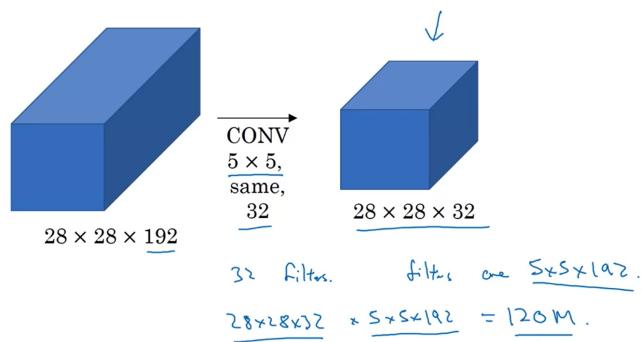
- Shrink the number of channels

## Inception Network

Intuition: Use what sizes of conv layer or pooling layer? Do them together

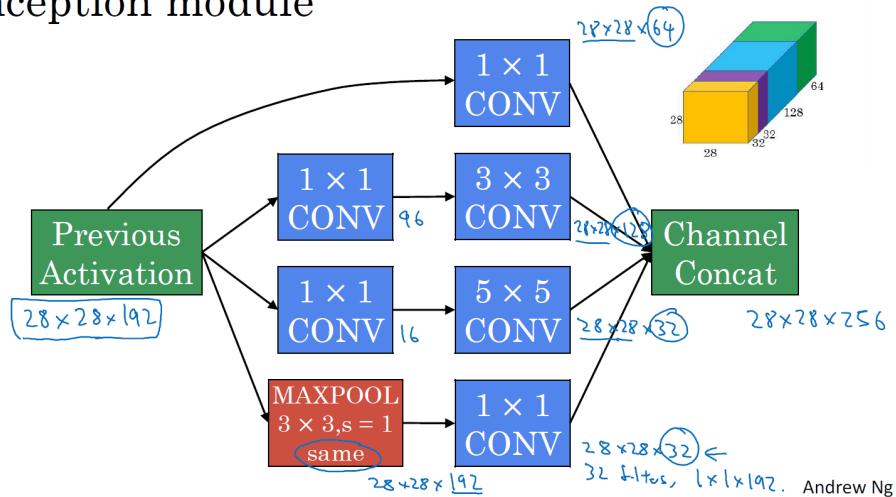


The computational cost

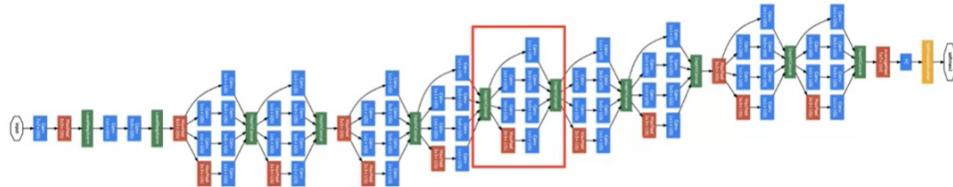


Does it influence the performance?

## Inception module



## Inception Network



## Mobile Net

Motivation:

- Low computational cost at deployment
- Useful for mobile and embedded vision application
- Key idea: Normal vs depthwise-separable convolutions

Normal Convolution.

Computation cost = # filter params × # filter positions × # of filters

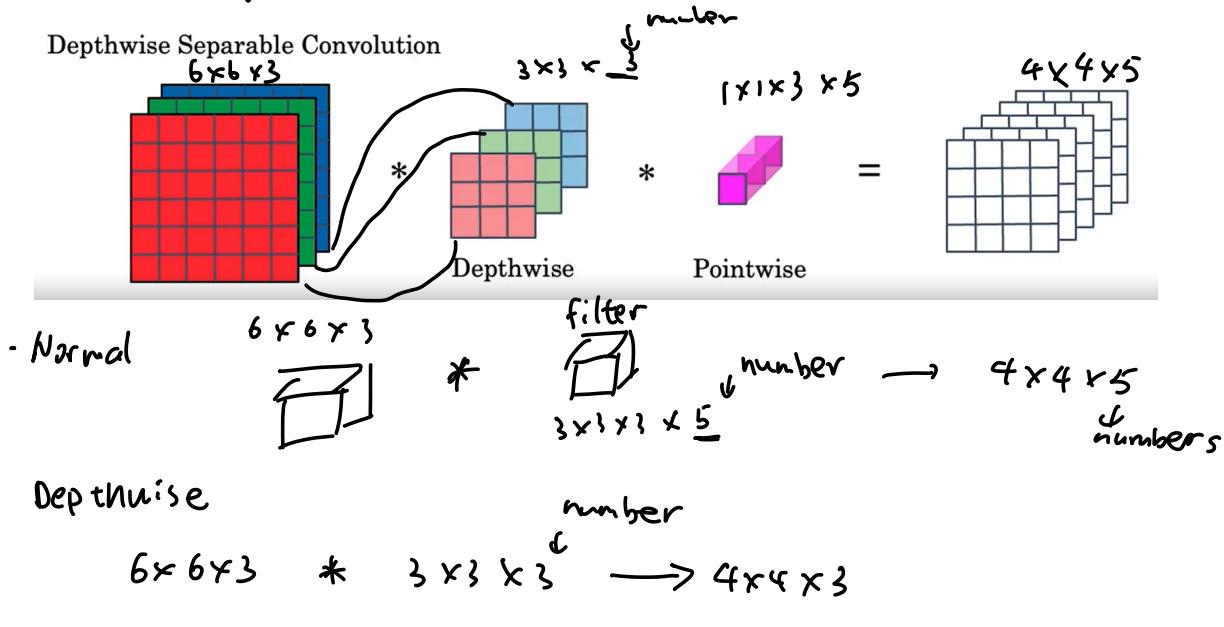
E.g.

$$6 \times 6 \times 3 \quad * \quad 3 \times 3 \times 3 \quad \stackrel{n_c' \text{ filter}}{\leftarrow} \quad \underset{\substack{\text{Stride}=1 \\ \text{no padding}}}{=} \quad 4 \times 4 \times n_c'$$

Assume  $n_c' = 5$

$$\text{Computational cost} = 3 \times 3 \times 3 \times 4 \times 4 \times 5 = 2160$$

## Depthwise Separable Convolution



$$\text{Depthwise Computation cost} = 3 \times 3 \times 3 \times 4 \times 4 = 432$$

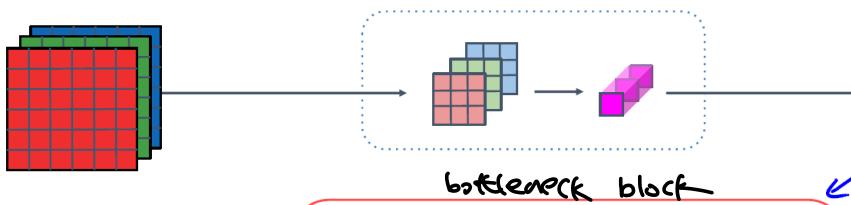
$$\text{Pointwise computation cost} = 1 \times 1 \times 3 \times 4 \times 4 \times 5 = 240$$

$$\begin{aligned} \text{Total: } & 672 \\ \text{Normal: } & 2160 \end{aligned}$$

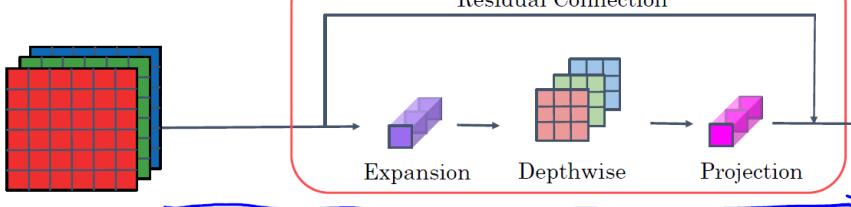
} in general case:  $\frac{1}{nc'} + \frac{1}{f^2}$

## MobileNet Architecture

MobileNet v1



MobileNet v2



\* Computationally Efficient!

## Efficient Net

- Scale up or down the neural network based on the resources of a device

## Practical Advice for Using ConvNets

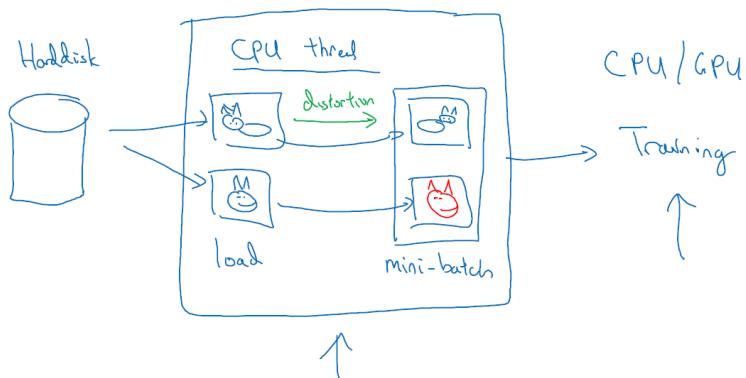
### • Transfer learning ▲

$L_2$  and  $T_{\text{dist}}$  are important hyperparameters for the initialization of your own network. →  $T_{\text{dist}} + \text{weight}$

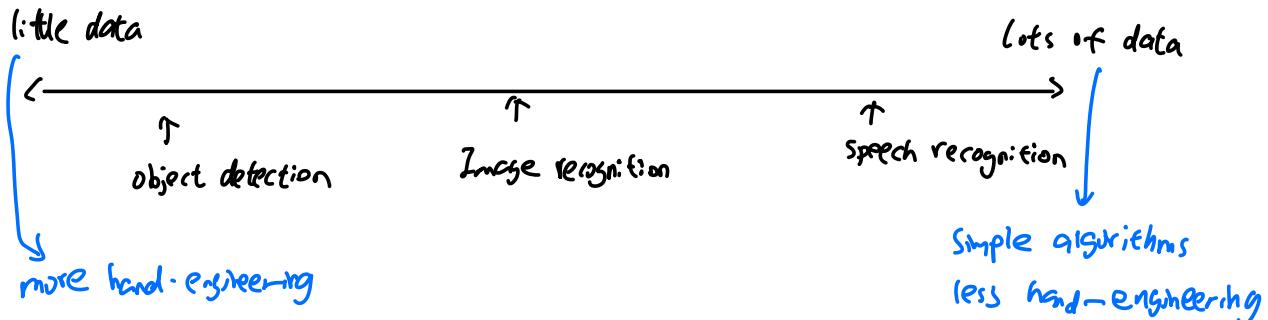
### • Data Augmentation

(CV) Mirror △  
Random Cropping △  
Rotation  
⋮

Color Shifting  
(method: PCA color augmentation)



### • State of CV



Two sources of knowledge

- Labeled data
- Hand engineered features/network architecture

Δ Tips for doing well on benchmarks / winning competitions

Ensembling keep each network available anytime, cost a lot of memory

- Train several networks independently and average their outputs

multi-crop at test time

- Run classifier on multiple versions of test image and average results

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset

## Supplement

### ① Transfer learning

Scenario 1 – Size of the Data set is small while the Data similarity is very high – In this case, since the data similarity is very high, we do not need to retrain the model. All we need to do is to customize and modify the output layers according to our problem statement. We use the pretrained model as a feature extractor. Suppose we decide to use models trained on Imagenet to identify if the new set of images have cats or dogs. Here the images we need to identify would be similar to imangenet, however we just need two categories as my output – cats or dogs. In this case all we do is just modify the dense layers and the final softmax layer to output 2 categories instead of a 1000.

Scenario 2 – Size of the data is small as well as data similarity is very low – In this case we can freeze the initial (let's say k) layers of the pretrained model and train just the remaining(n-k) layers again. The top layers would then be customized to the new data set. Since the new data set has low similarity it is significant to retrain and customize the higher layers according to the new dataset. The small size of the data set is compensated by the fact that the initial layers are kept pretrained(which have been trained on a large dataset previously) and the weights for those layers are frozen.

Scenario 3 – Size of the data set is large however the Data similarity is very low – In this case, since we have a large dataset, our neural network training would be effective. However, since the data we have is very different as compared to the data used for training our pretrained models. The predictions made using pretrained models would not be effective. Hence, its best to train the neural network from scratch according to your data.

Scenario 4 – Size of the data is large as well as there is high data similarity – This is the ideal situation. In this case the pretrained model should be most effective. The best way to use the model is to retain the architecture of the model and the initial weights of the model. Then we can retrain this

## ② transfer learning 实际过程

### ① 划分训练集和验证集

```
image_dataset_from_directory( )
```

### ② Use prefetch() to prevent memory bottlenecks when reading from disk

### ③ data augment

```
data_augmentation = tf.keras.Sequential(  
    [layers.RandomFlip('horizontal_and_vertical'),  
     layers.RandomRotation(0.2),  
     ...  
    ])
```

④

```
def data_augmenter():  
    ...  
    Create a Sequential model composed of 2 layers  
    Returns:  
        tf.keras.Sequential  
    ...  
    ### START CODE HERE  
    data_augmentation = tf.keras.Sequential()  
    data_augmentation.add(RandomFlip('horizontal'))  
    data_augmentation.add(RandomRotation(0.2))  
    ### END CODE HERE  
  
    return data_augmentation
```

## ④ Use a model for transfer learning

## ⑤ Layer freezing

### 1. Delete the top layer (the classification layer)

- Set `include_top` in `base_model` as False

### 2. Add a new classifier layer

- Train only one layer by freezing the rest of the network
- As mentioned before, a single neuron is enough to solve a binary classification problem.

### 3. Freeze the base model and train the newly-created classifier layer

- Set `base_model.trainable=False` to avoid changing the weights and train *only* the new layer
- Set `training` in `base_model` to False to avoid keeping track of statistics in the batch norm layer

去掉 base model 的分类器，冻结 base model，输入自己的分类层（drop过的）。

比如 Global avg pooling, Dropout, predict\_layer, output

new model = tf.keras.Model(input, outputs)

## ⑥ backprop in my model.

```
base_learning_rate = 0.001
model2.compile(optimizer=tf.keras.optimizers.Adam(lr=base_learning_rate),
                loss=tf.keras.losses.BinaryCrossentropy(from_logits=True),
                metrics=['accuracy'])

initial_epochs = 5
history = model2.fit(train_dataset, validation_data=validation_dataset, epochs=initial_epochs)
```

## ⑦ 看 training and validation accuracy

## ⑧ Fine-tuning model

You could try fine-tuning the model by re-running the optimizer in the last layers to improve accuracy. When you use a smaller learning rate, you take smaller steps to adapt it a little more closely to the new data. In transfer learning, the way you achieve this is by unfreezing the layers at the end of the network, and then re-training your model on the final layers with a very low learning rate. Adapting your learning rate to go over these layers in smaller steps can yield more fine details - and higher accuracy.

The intuition for what's happening: when the network is in its earlier stages, it trains on low-level features, like edges. In the later layers, more complex, high-level features like wispy hair or pointy ears begin to emerge. For transfer learning, the low-level features can be kept the same, as they have common features for most images. When you add new data, you generally want the high-level features to adapt to it, which is rather like letting the network learn to detect features more related to your data, such as soft fur or big teeth.

To achieve this, just unfreeze the final layers and re-run the optimizer with a smaller learning rate, while keeping all the other layers frozen.

Where the final layers actually begin is a bit arbitrary, so feel free to play around with this number a bit. The important takeaway is that the later layers are the part of your network that contain the fine details (pointy ears, hairy tails) that are more specific to your problem.

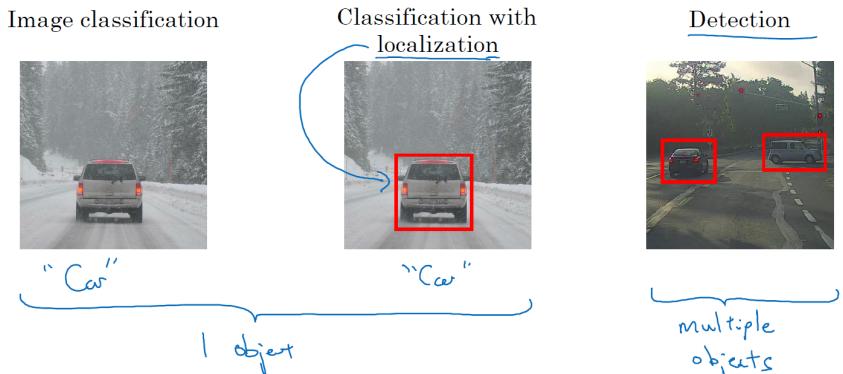
First, unfreeze the base model by setting `base_model.trainable=True`, set a layer to fine-tune from, then re-freeze all the layers before it. Run it again for another few epochs, and see if your accuracy improved!

## CNN Week 3 Detection Algorithms

### Object localization

what objects are in this image and where in the image are those objects located?

Some definitions



Defining the target label  $y$

1 - pedestrian  
2 - car ←  
3 - motorcycle  
4 - background ←

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_i - y_i)^2 + (\hat{y}_j - y_j)^2 & \text{if } y_i = 1 \\ \dots + (\hat{y}_k - y_k)^2 & \text{if } y_i = 1 \\ (\hat{y}_l - y_l)^2 & \text{if } y_i = 0 \end{cases}$$

Need to output  $b_x, b_y, b_h, b_w$ , class label (1-4)

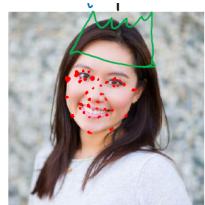
$x =$   $y =$   $(x, y)$

is there an object?

Andrew Ng

### Landmark detection

The label of landmark should be same among the all training set.



$l_{1x}, l_{1y},$   
 $l_{2x}, l_{2y},$   
 $l_{3x}, l_{3y},$   
 $l_{4x}, l_{4y},$   
 $\vdots$   
 $l_{64x}, l_{64y}$

$X, Y$

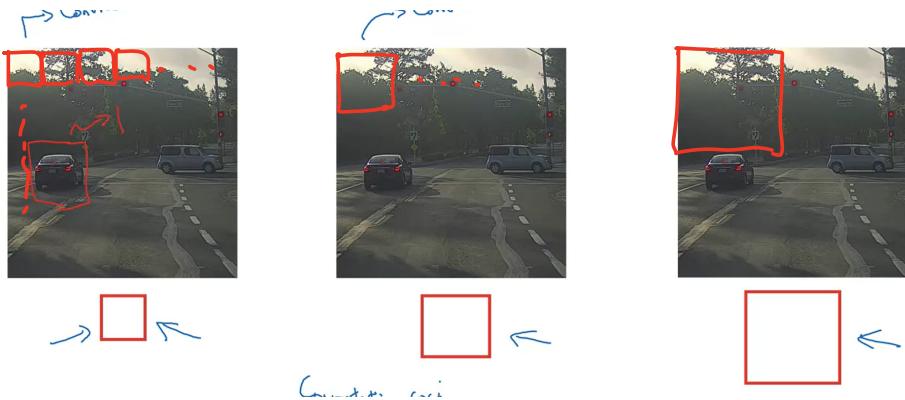
→ Conv Net

face ?

$l_{1x}, l_{1y}, l_{2x}, l_{2y}, l_{3x}, l_{3y}, l_{4x}, l_{4y}, \dots, l_{64x}, l_{64y}$

## object detection

### Sliding windows detection



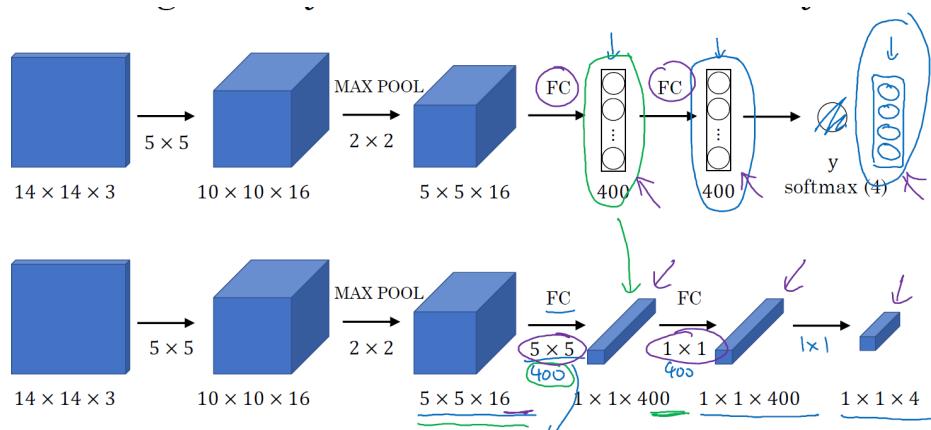
Training set:

X	y
	1
	1
	1
	0
	0

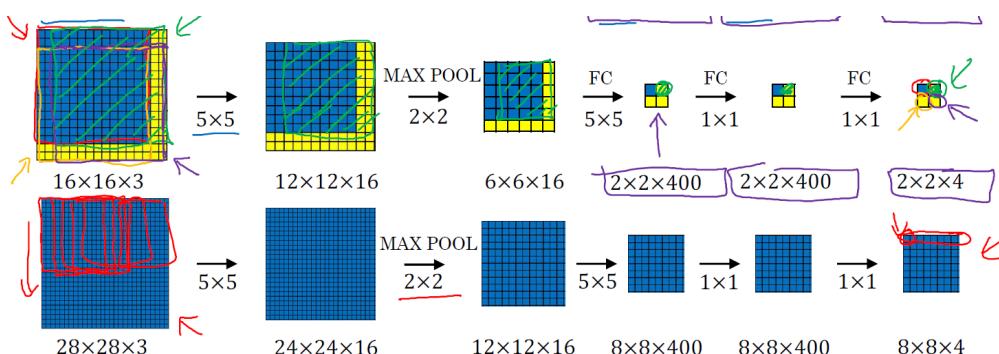
Drawback: high computational cost, unacceptable for CNNs.

## Convolutional implementation of sliding windows

### Turning FC layer into convolutional layers

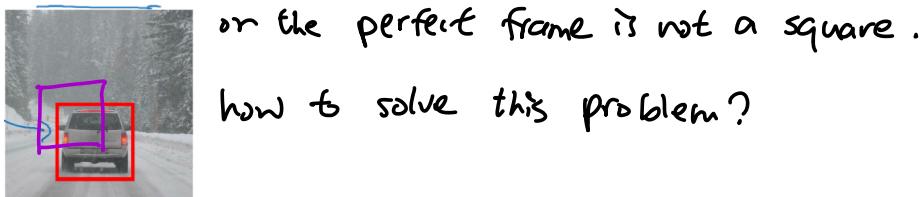


### Convolution implementation of sliding windows



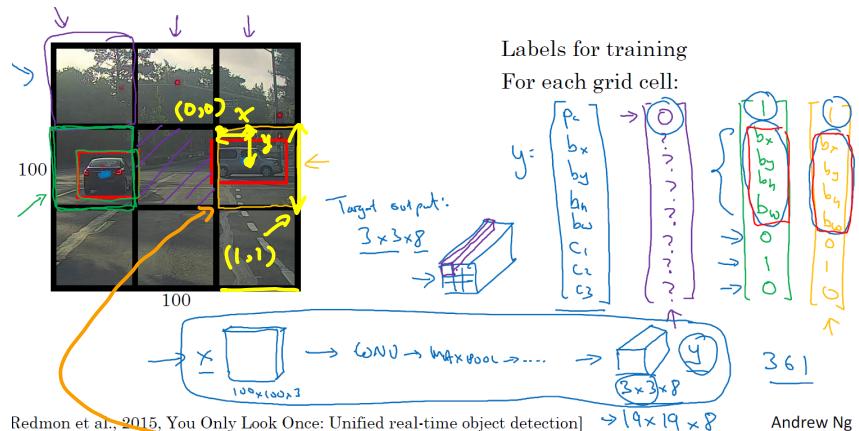
Implement all the sliding windows at the same time.

Problem: 无法精确输出边界框。(the purple frame is more common)



## Building box Predictions

### YOLO Algorithm



Advantage:

- ① 不会受到滑动窗口的干扰
- ② very fast

Question:

how to decide  $b_x, b_y, b_h, b_w$

$$\begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ \dots \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0.4 \\ 0.3 \\ 0.9 \\ 0.5 \\ \dots \\ 0 \end{bmatrix}$$

$\leftarrow$  less than 1  
 $\leftarrow$  fraction of red line  
 $\leftarrow$  could be larger than 1

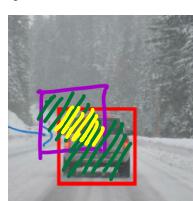
## Intersection over union

- { Evaluate the detection algorithm
- { Add another components in detection algorithm

$$IoU = \frac{\text{size of yellow}}{\text{size of green}}$$

"correct" if  $IoU \geq 0.5$

$\frac{T}{T}$   
just a convention



## Non-max suppression (解决一个物体出现多个框的问题)

- Just output the max probability one.

• Step:

Each output prediction is  $\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$

Discard all boxes with  $p_c \leq 0.6$

while there are any remaining boxes:

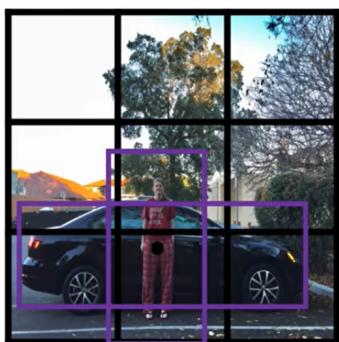
- Pick the box with the largest  $p_c$ . Output that as a prediction
- Discard any remaining box with  $I_{IoU} \geq 0.5$  with the box input in the previous step.

## Anchor box (固定尺寸特征图重叠和每个物体的框)

Previously: Each object in training image is assigned to grid cell that contains that object's midpoint

with two anchor boxes: Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest  $I_{IoU}$ .

Example:



Anchor box 1: Anchor box 2:



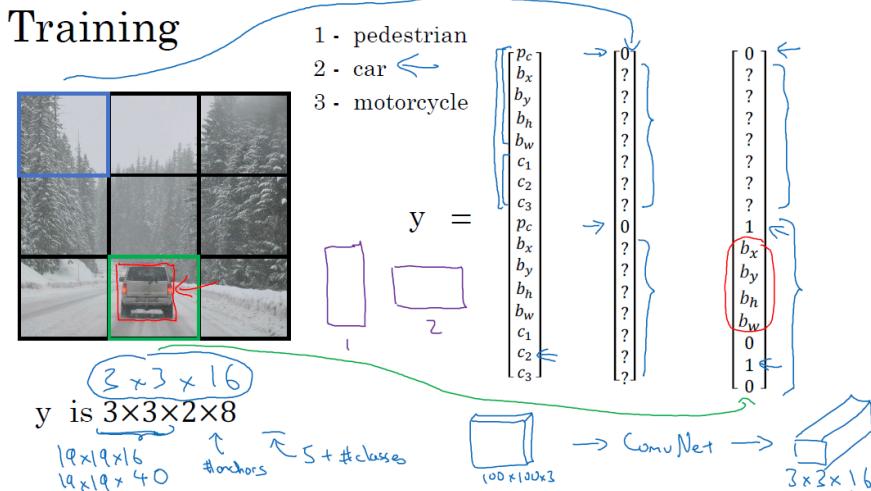
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ - \\ 0 \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ - \\ 0 \\ . \\ . \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

only car?

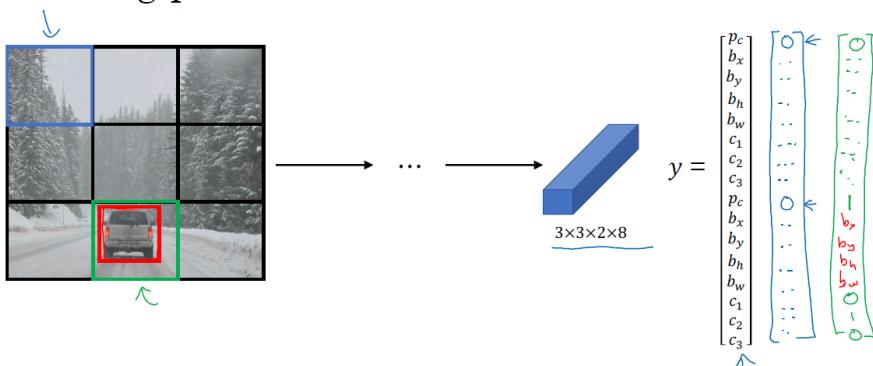
why anchor box?  
 {- 处理同一个物体出现在同一格子的逻辑  
 {- 避免过长过宽的情况 (too tall/wide ...)

how to choose anchor box? — k-means

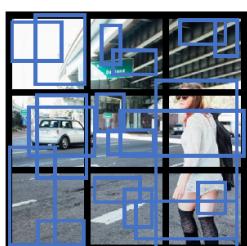
### YOLO Algorithm (update)



### Making predictions



### Outputting the non-max suppressed outputs



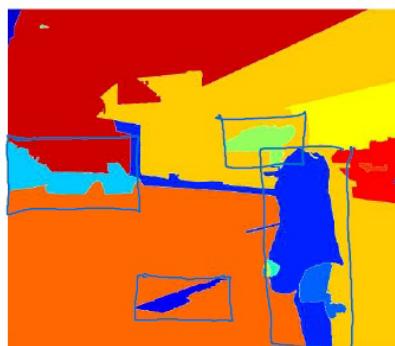
- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

## Region proposals (R-CNN)

R-CNN: Propose regions. Classify proposed regions one at a time. Output label & bounding box.

Faster R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions

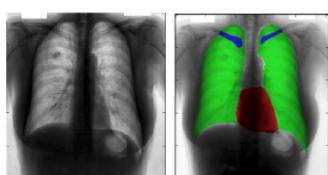
Faster R-CNN: Use convolutional network to propose regions



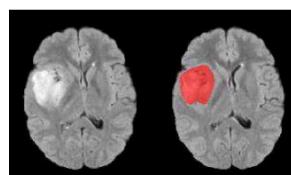
## Segmentation algorithm

## Semantic Segmentation with U-net (语义分割) ☆☆

## Motivation for U-Net

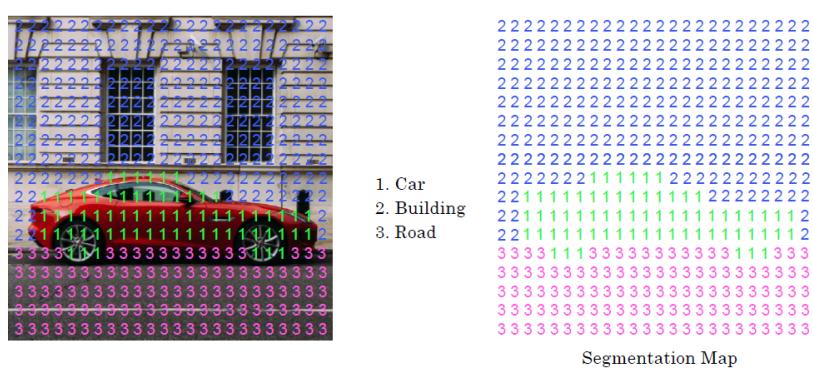


Chest X-Ray

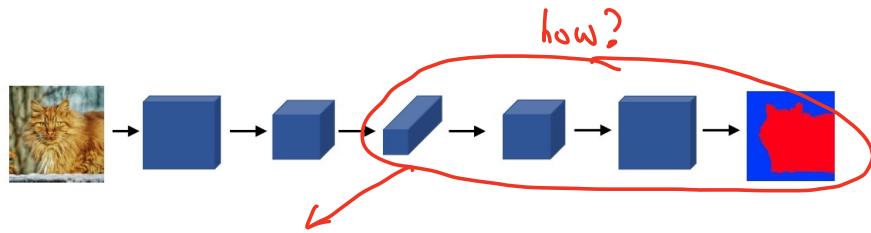


Brain MRI

### Per-pixel class labels



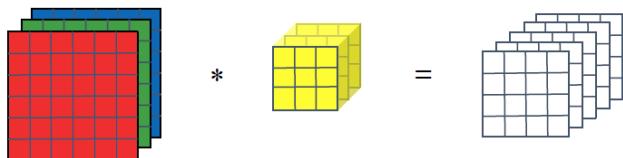
# Deep Learning for Semantic Segmentation



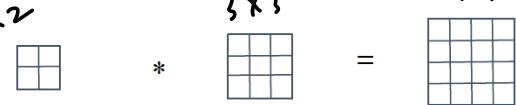
## Transpose Convolutions

### Transpose Convolution

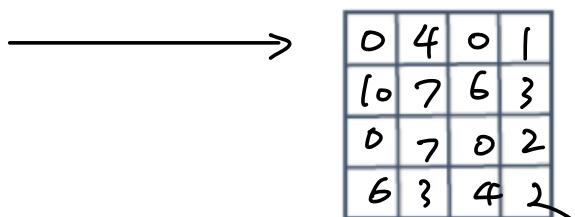
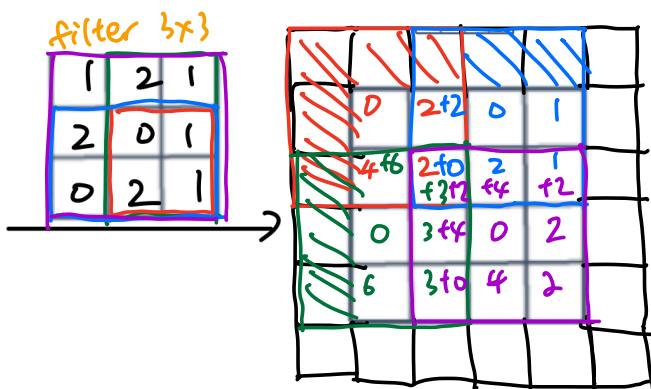
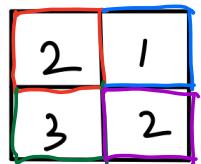
Normal Convolution



Transpose Convolution

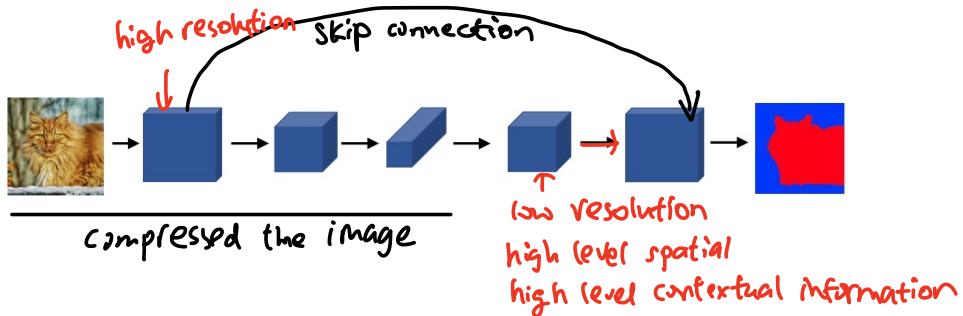


In detail.



## U-Net Architecture Intuition

### Deep Learning for Semantic Segmentation



## U-Net Architecture

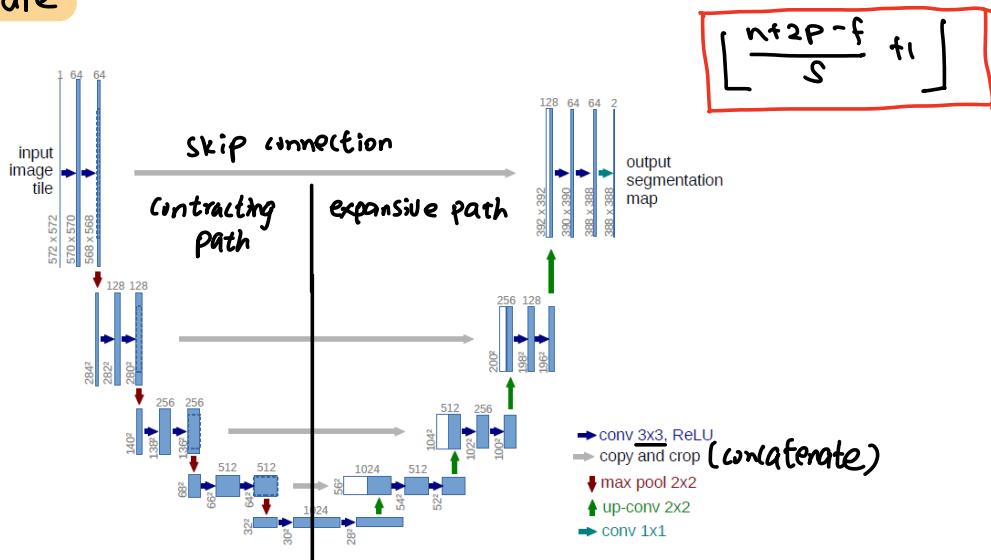


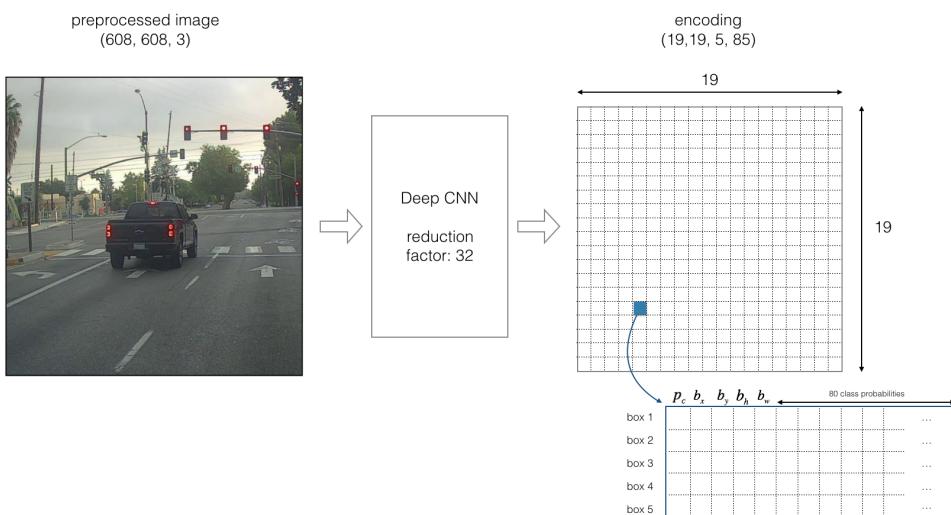
Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

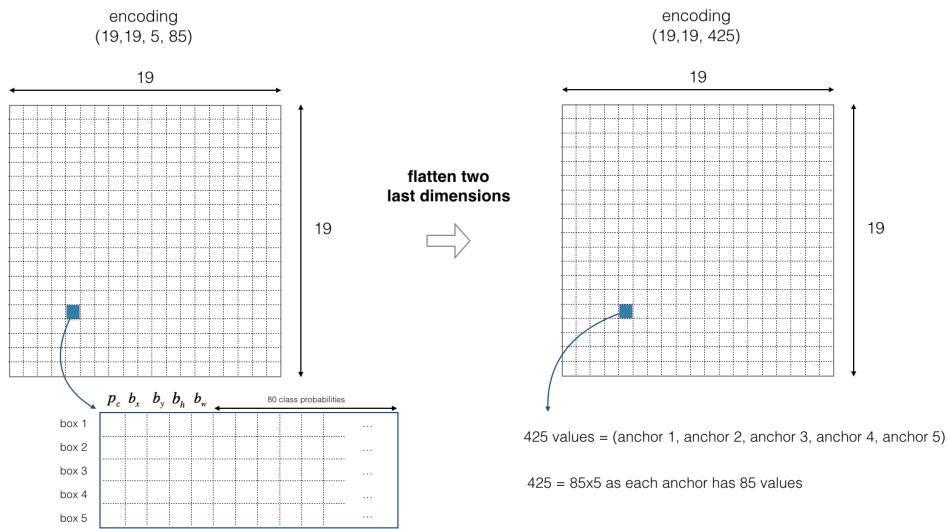
### Implementation:

① `Conv1 = conv2D(filters = 64,  
kernel_size = 3,  
padding = 'same', (Valid?)  
activation = 'relu',  
kernel_initializer = 'he_normal')`

## Summary for YOLO

- Input image (608, 608, 3)
- The input image goes through a CNN, resulting in a (19, 19, 5, 85) dimensional output.
- After flattening the last two dimensions, the output is a volume of shape (19, 19, 425)
  - Each cell in a 19x19 grid over the input image gives 425 numbers.
  - $425 = 5 \times 85$  because each cell contains predictions for 5 boxes, corresponding to 5 anchor boxes.
  - $85 = 80 + 5$  where 5 is because  $(p_c, b_x, b_y, b_h, b_w)$  has 5 numbers. and 80 is the number of classes we'd like to detect
- You then select only few boxes based on:
  - Score-thresholding: throw away boxes that have defeated a class with a score less than the threshold
  - Non-max suppression: Compute the Intersection over Union and avoid selecting overlapping boxes.
- This gives you YOLO's final output.



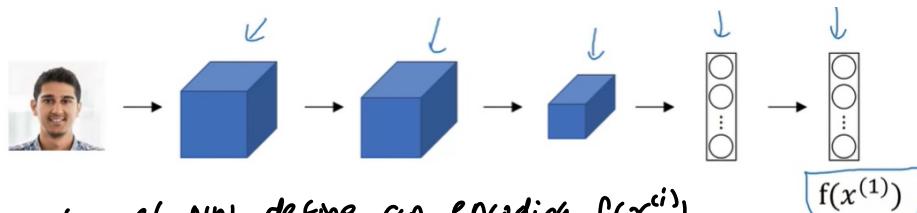


## CNN Week 4 Face Recognition

### One-Shot learning

- Learning from one example to recognize the person again
- Learning a 'similarity' function  
 $d(\text{img 1}, \text{img 2}) = \text{degree of difference between image}$   
 If  $d(\text{img 1}, \text{img 2}) \leq \tau$   
 $\quad > \tau$
- Question: How to learn 'similarity' function?

### Siamese Network



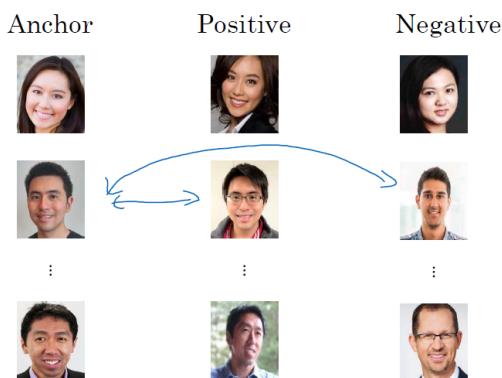
- Parameter of NN define an encoding  $f(x^{(i)})$
- Learn parameters so that:  
 If  $x^{(i)}, x^{(j)}$  are same person,  $\|f(x^{(i)}) - f(x^{(j)})\|$  is small  
 different  $\|f(x^{(i)}) - f(x^{(j)})\|$  is large

### Triplet loss

Given 3 images A, P, N:

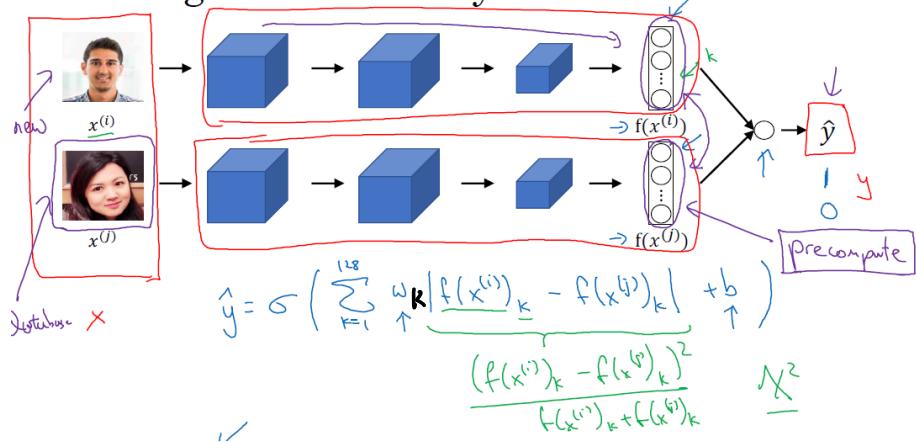
$$\ell(A, P, N) = \max ( \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \delta, 0)$$

$$J = \sum_{i=1}^m \ell(A^{(i)}, P^{(i)}, N^{(i)})$$

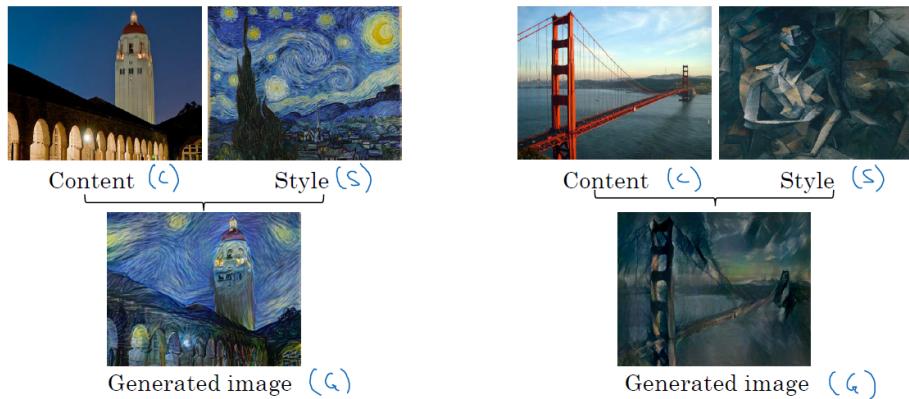


## Face verification and binary classification

### Learning the similarity function



## Neural Style Transfer



## Cost Function

$$J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

1. Initiate  $G$  randomly

$$G: 100 \times 100 \times 3$$

2. Use gradient descent to minimize  $J(G)$

### Content Cost Function

Content cost function

$$J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

- Say you use hidden layer  $l$  to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let  $a^{[l](C)}$  and  $a^{[l](G)}$  be the activation of layer  $l$  on the images
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content

$$J_{content}(C, G) = \frac{1}{2} \|a^{[l](C)} - a^{[l](G)}\|^2$$

2

- Style cost Function