

# Introduction to Convolutional Neural Networks

**Prof. Mingkui Tan**

SCUT Machine Intelligence Laboratory (SMIL)



SMIL内部资料 请勿外泄

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

4 Backpropagation

5 Convolutional Neural Networks

6 Applications

SMIL内部资料 请勿外泄

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

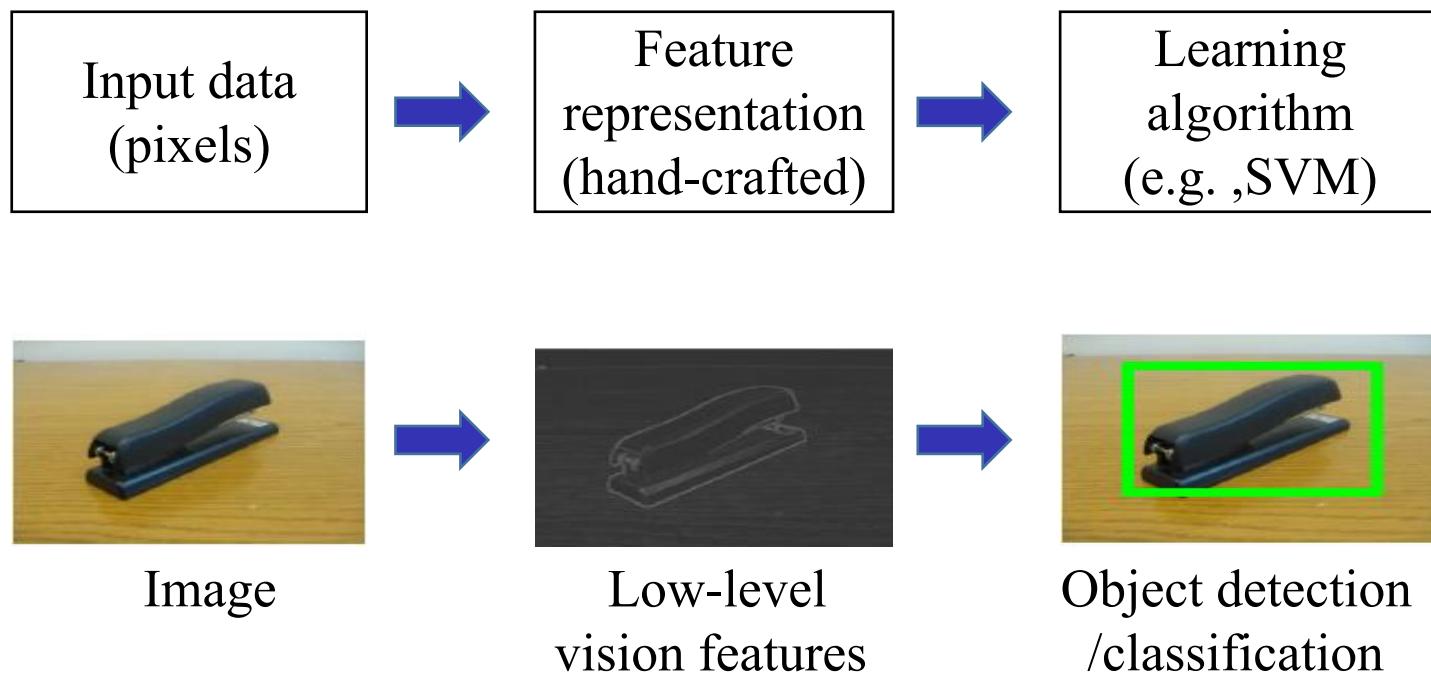
4 Backpropagation

5 Convolutional Neural Networks

6 Applications

SMIL内部资料 请勿外泄

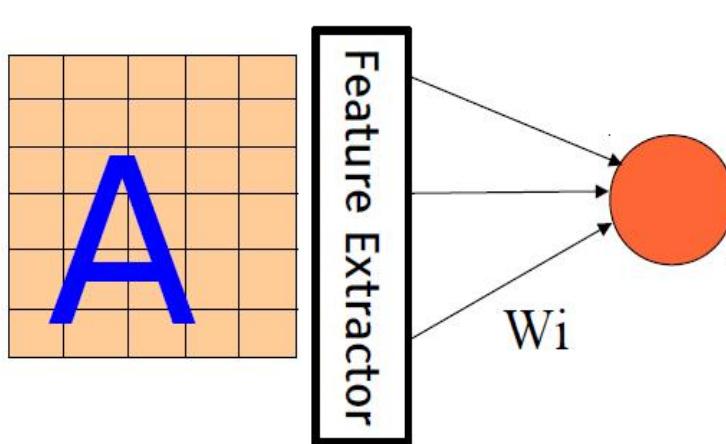
# Traditional Recognition Methods



SMIL内部资料 请勿外泄

# Traditional Recognition Methods

$$y = \text{sign} \left( \sum_{i=1}^N \mathbf{W}_i F_i(\mathbf{X}) + b \right)$$

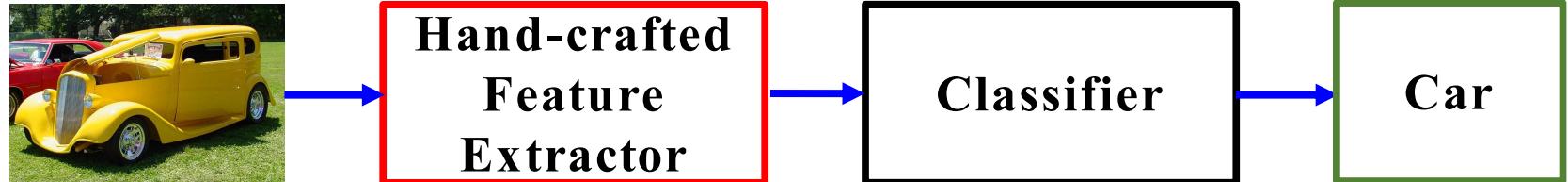


## Linear Classifier

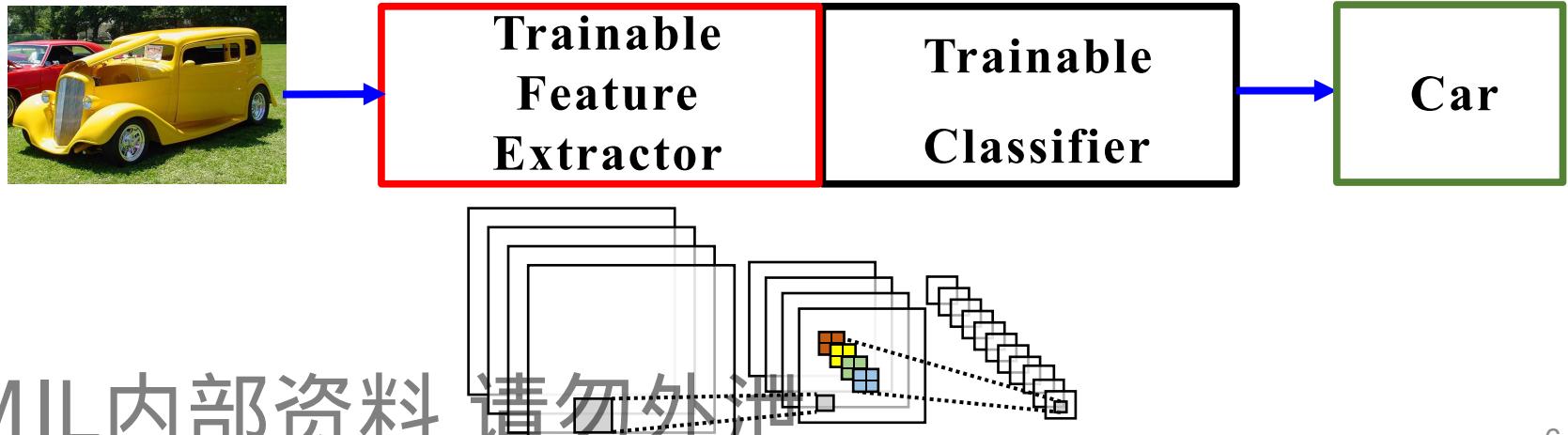
- **Linear classifier** is on top of a simple **feature extractor**
- Highly depend on the **quality of extracted feature**
- Designing a feature extractor requires **considerable efforts**

# Traditional Methods vs Deep Learning

- Traditional recognition with **hand-crafted feature**,  
e.g. SIFT and LBP feature

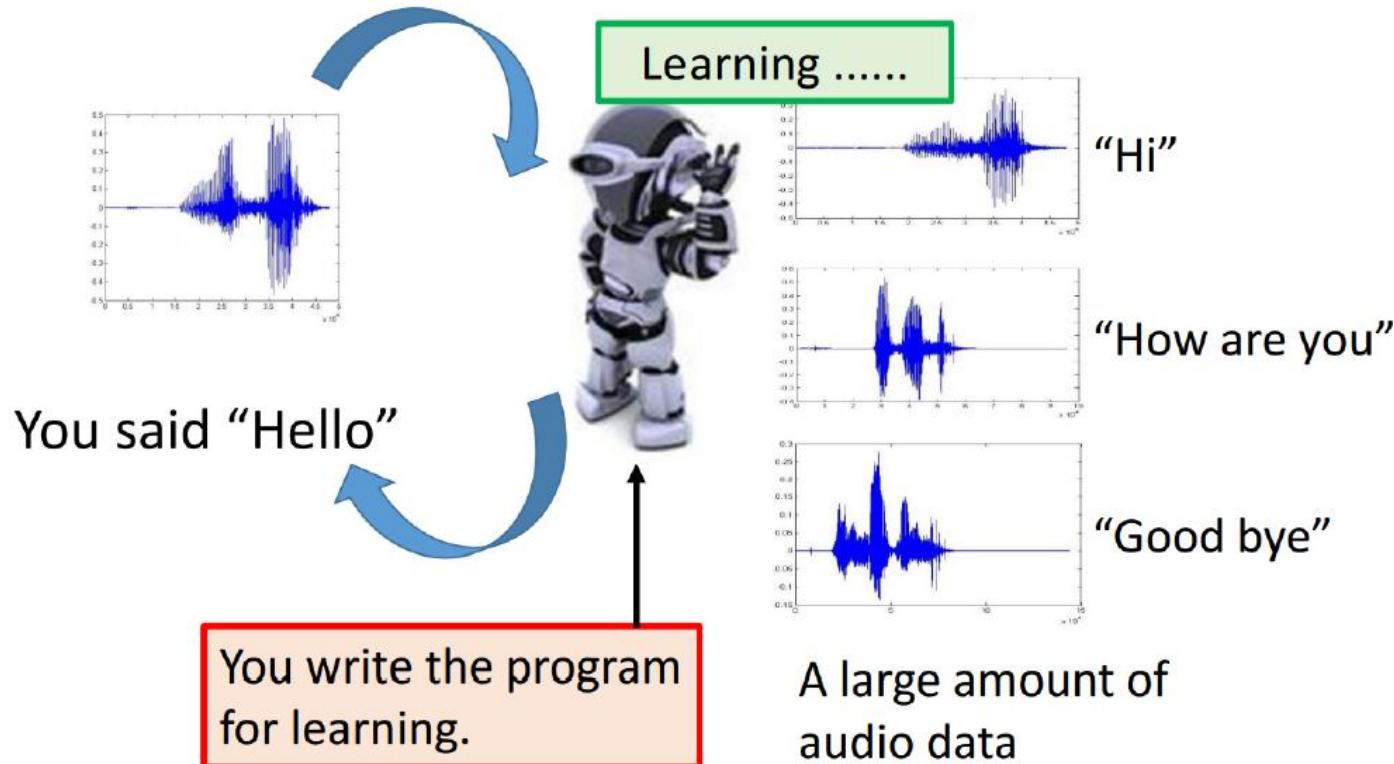


- Deep learning / Feature learning



SMIL内部资料 请勿外泄

# What is Deep Learning?



SMIL内部资料 请勿外泄

# Deep Learning

Looking for a Function  $f$ :

■ Speech Recognition

$$f\left( \begin{array}{c} \text{A spectrogram plot showing a sound waveform} \\ \text{A blue waveform on a grid from -10 to 10 ms} \end{array} \right) = \text{“How are you”}$$

■ Image Recognition

$$f\left( \begin{array}{c} \text{A photo of an orange kitten looking up} \\ \text{A small orange cat with large eyes} \end{array} \right) = \text{“Cat”}$$

■ Playing Go

$$f\left( \begin{array}{c} \text{A Go board with stones} \\ \text{A Go board with black and white stones} \end{array} \right) = \text{“5-5” (next move)}$$

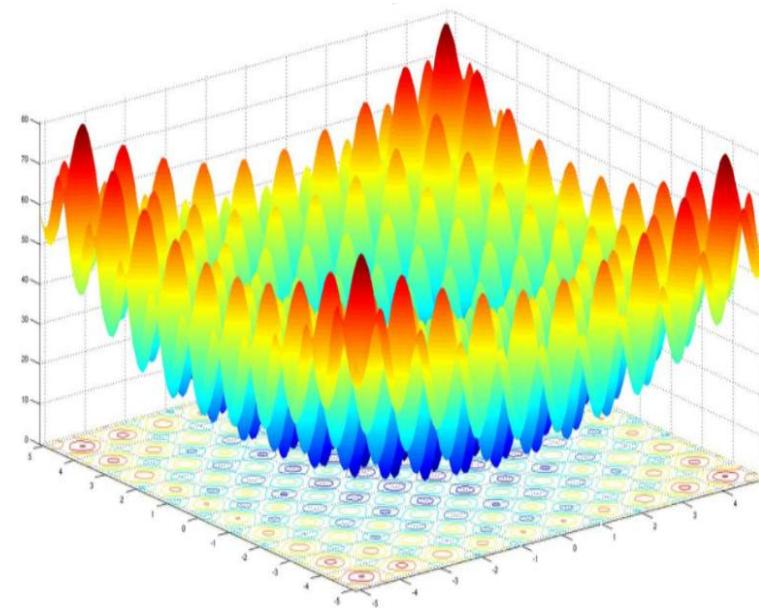
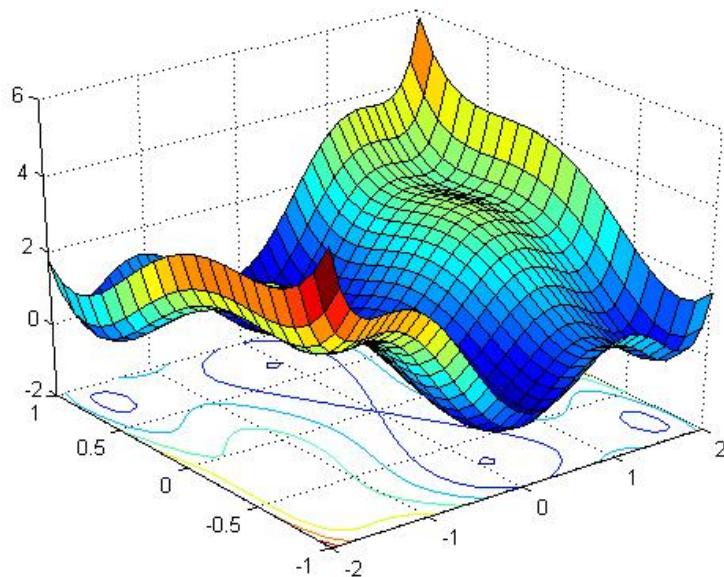
■ Dialogue System

$$f\left( \begin{array}{c} \text{“Hi”} \\ \text{(what the user said)} \end{array} \right) = \text{“Hello”} \\ \text{(system response)}$$

SMIL 内部资料 请勿外泄

# Deep Learning

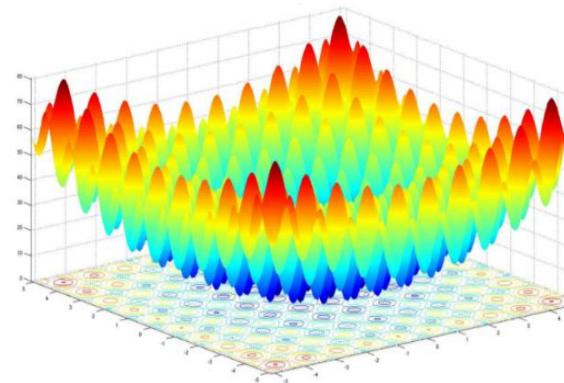
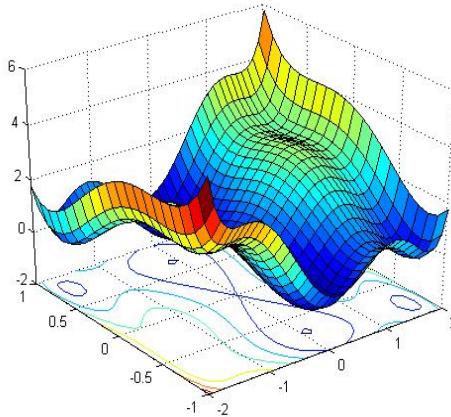
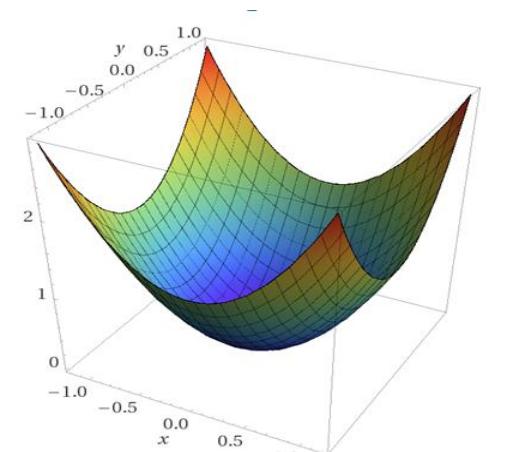
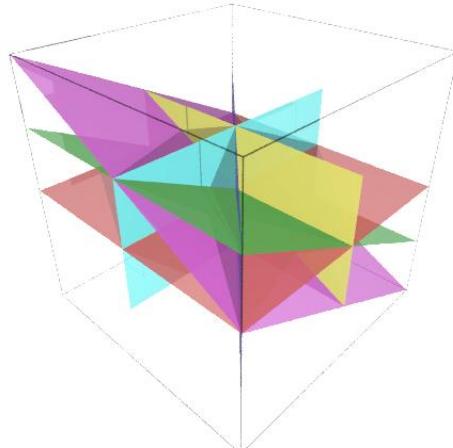
- Looking for a Function  $f$ :



SMIL内部资料 请勿外泄

# Deep Learning

- What kind of functions can we choose?



SMIL内部资料 请勿外泄

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

4 Backpropagation

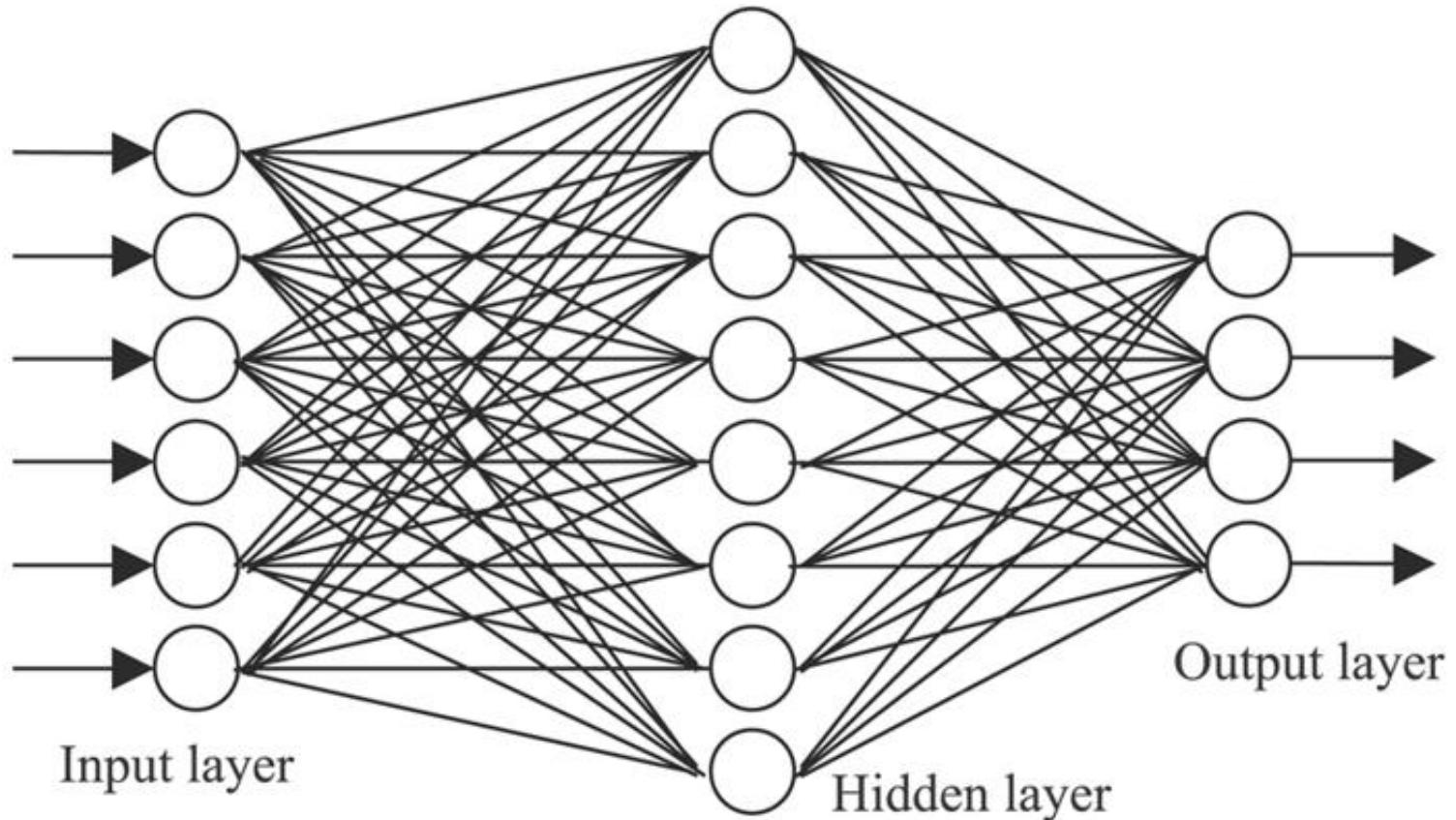
5 Convolutional Neural Networks

6 Applications

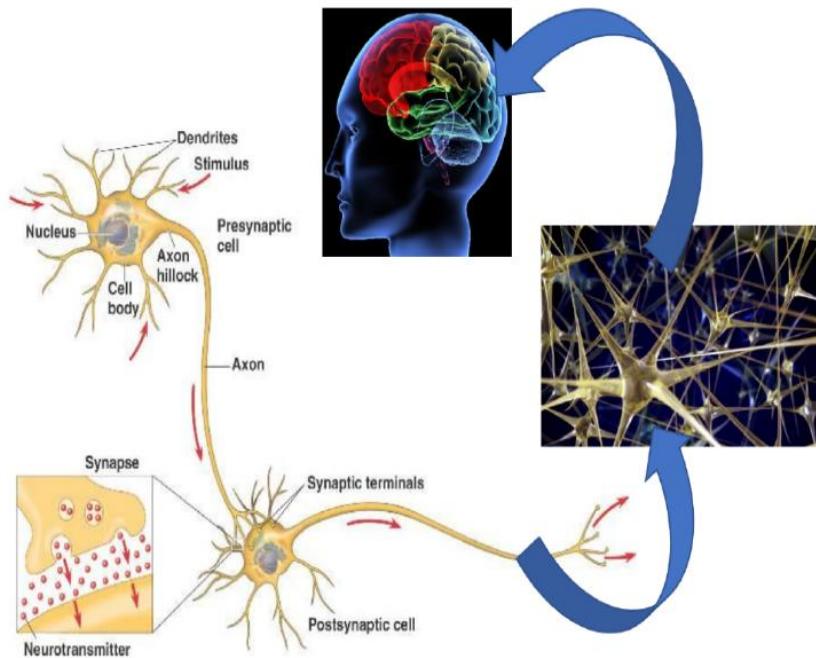
SMIL内部资料 请勿外泄

# Neural Networks

- To approximate the function, we can build **neural networks**:

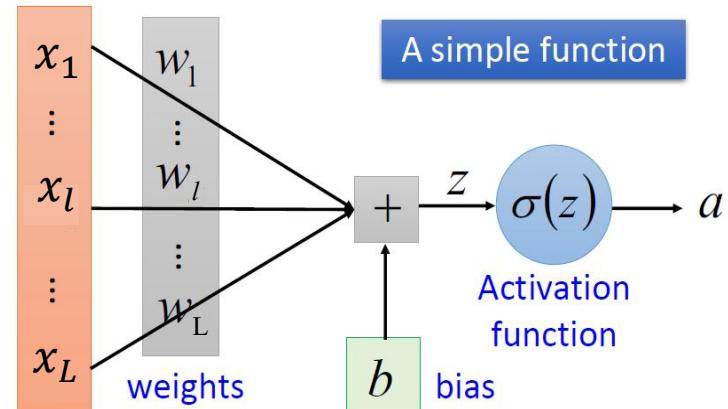


# Neurons in Human Brains



## Neuron

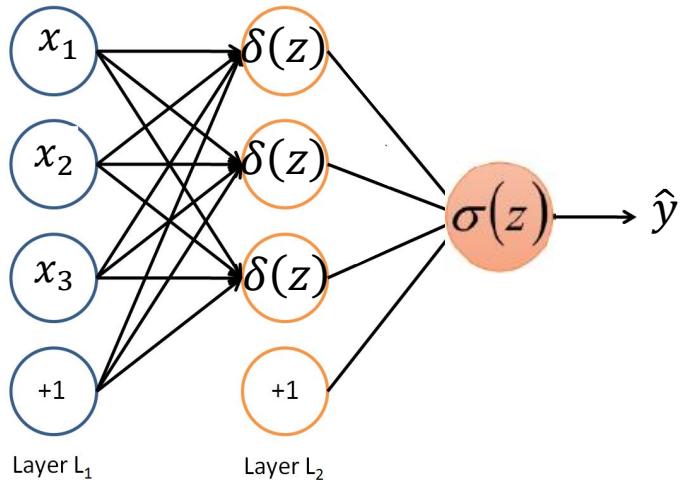
$$z = x_1 w_1 + \dots + x_l w_l + \dots + x_L w_L + b$$



SMIL内部资料 请勿外泄

# Neural Networks

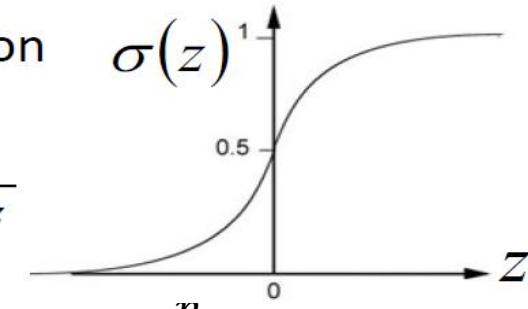
## Non – Linear Activation Functions



$$\delta(z) = \delta(W^T X + b)$$

Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



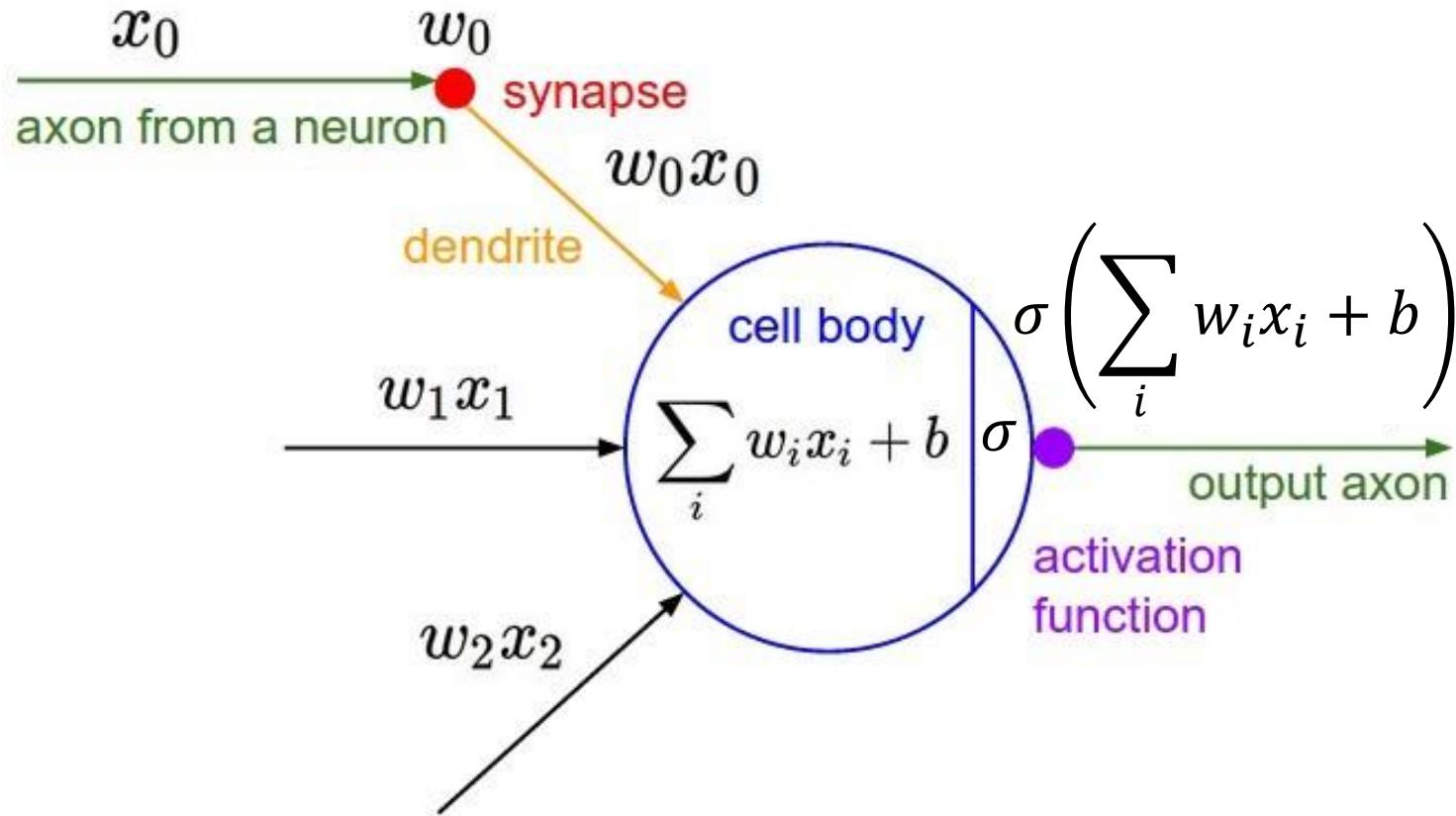
$$\text{Loss function } \mathcal{L} = \sum_{i=1}^n \|\hat{y}_i - y_i\|^2$$

## Representation Theory (Kurt Hornik, 1991)

A feed-forward neural network with **a single hidden layer** containing a **finite number of neurons** can approximate continuous function on compact subsets of  $\mathbf{R}^n$  under mild assumptions on the **activation function**

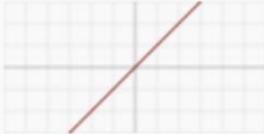
# Activation Function

- **Description:** keep value within an **acceptable** and useful range



# Activation Layer

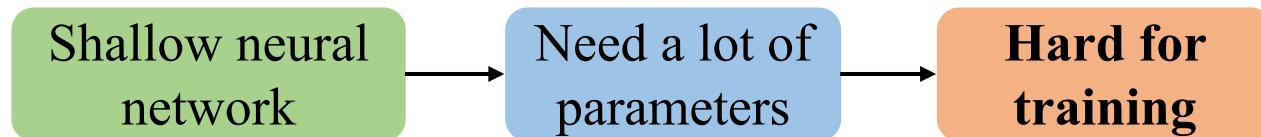
- Activation functions: **nonlinear** and **continuously differentiable**

Activation Function	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Sigmoid		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ReLU		$f(x) = \max(x, 0)$	$f'(x) = 1\{x \geq 0\}$
SoftPlus		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

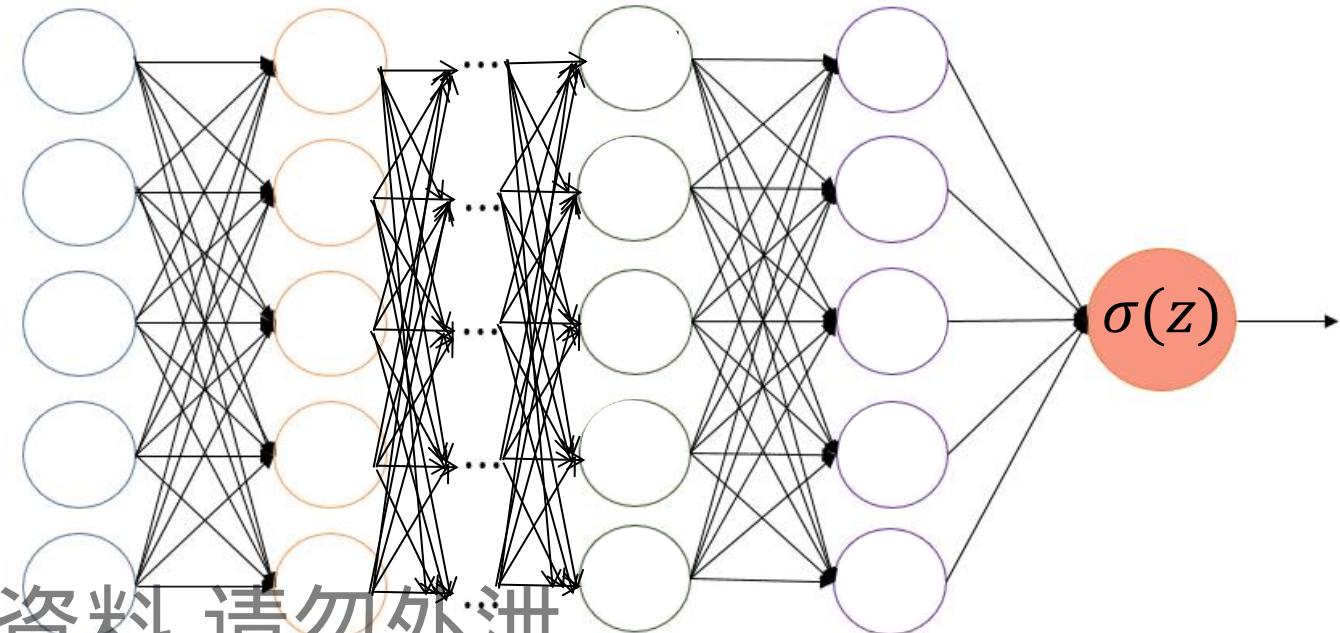
# From Shallow to Deep Neural Networks

Function Approximation Theory (Shiyu Liang, 2017)

Shallow networks require **exponentially more neurons** than a deep network to achieve the same level of accuracy for function approximation



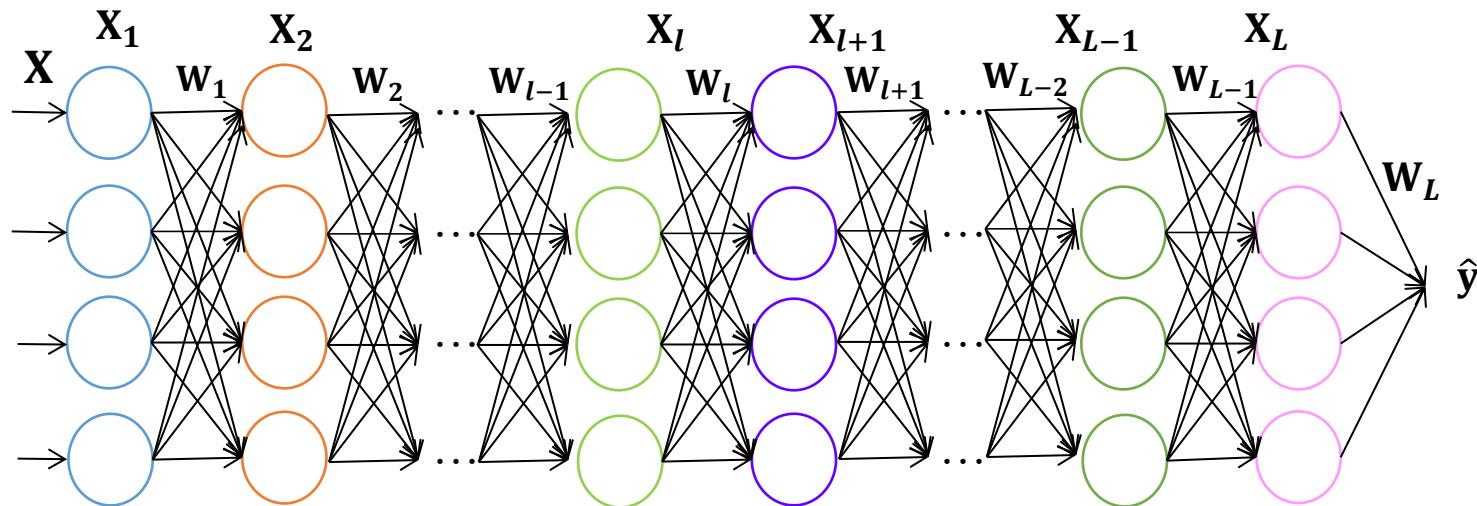
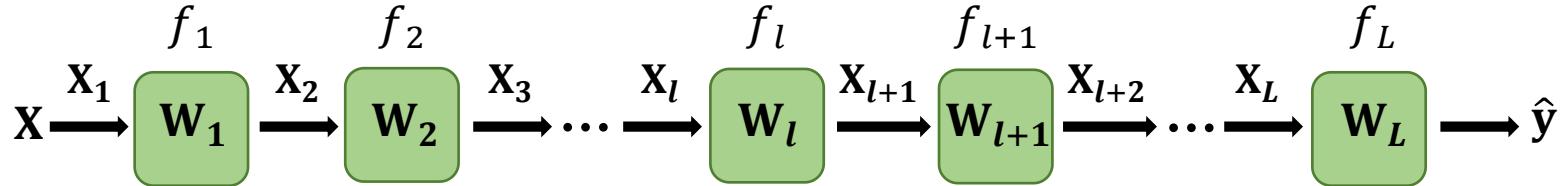
To approximate  $F(\mathbf{X})$ , we can build a **deep** neural network with less parameters:



SMIL内部资料 请勿外泄

# From Shallow to Deep Neural Networks

To approximate  $F(\mathbf{X})$ , we can build a **deep** neural network with **LESS parameters**:



$\mathbf{X}$  : input data

$\hat{\mathbf{y}}$  : prediction of the model

$\mathbf{X}_l$  : output of Layer  $l$

$f_l$  : function of Layer  $l$      $\mathbf{W}_l$ : parameters of Layer  $l$

The output of layer  $l$ , i.e.  $\mathbf{X}_{l+1}$  is the input of  $f_{l+1}$

SMIC 内部资料 请勿外泄

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

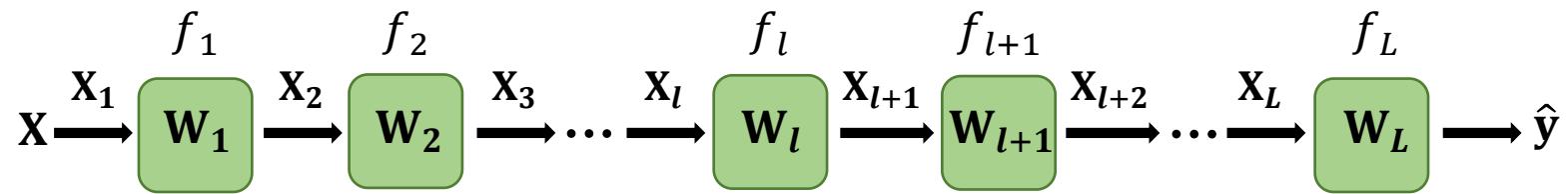
4 Backpropagation

5 Convolutional Neural Networks

6 Applications

SMIL内部资料 请勿外泄

# Forward propagation of neural networks



$$f = f_L \circ f_{L-1} \circ \cdots \circ f_2 \circ f_1$$

# Prediction for Networks $f$ : Forward Propagation

- Given input data  $\mathbf{X}$ , the forward propagation of a neural network  $f$  is to **predict** the output  $\hat{\mathbf{y}}$  of  $\mathbf{X}$
- Rewrite the network as a Composite Function

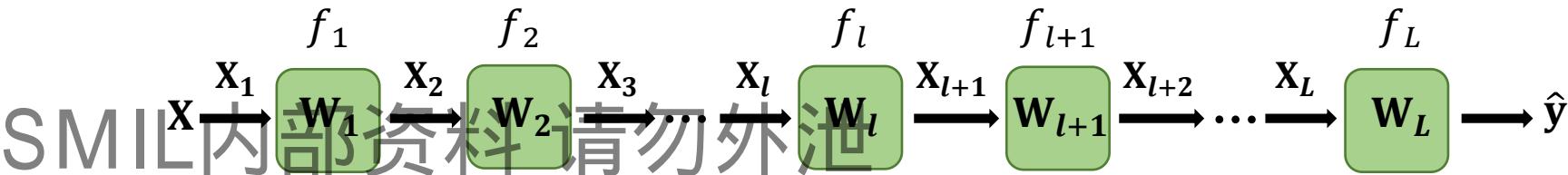
$$\hat{\mathbf{y}} = {}_W f(\mathbf{X}) = f_L \circ f_{L-1} \circ \cdots \circ f_2 \circ f_1(\mathbf{X}_1, \mathbf{W}_1)$$

- Detailed prediction process:

$$\begin{aligned}\mathbf{X}_1 &= \mathbf{X} \\ \mathbf{X}_2 &= f_1(\mathbf{X}_1, \mathbf{W}_1) \\ &\vdots \\ \mathbf{X}_{l+1} &= f_l(\mathbf{X}_l, \mathbf{W}_l) \\ &\vdots \\ \mathbf{X}_L &= f_{L-1}(\mathbf{X}_{L-1}, \mathbf{W}_{L-1}) \\ \hat{\mathbf{y}} &= f_L(\mathbf{X}_L, \mathbf{W}_L)\end{aligned}$$

$$\hat{\mathbf{y}} = f(\mathbf{X})$$

$$\begin{aligned}f_l &= \sigma(h_l(\mathbf{X}_l, \mathbf{W}_l)) \\ \text{Example:} \\ h_l &= \mathbf{W}_l * \mathbf{X}_l \\ \sigma &= \text{ReLU}(\cdot)\end{aligned}$$



# Forward Propagation: Prediction of Networks

- Detailed prediction process:

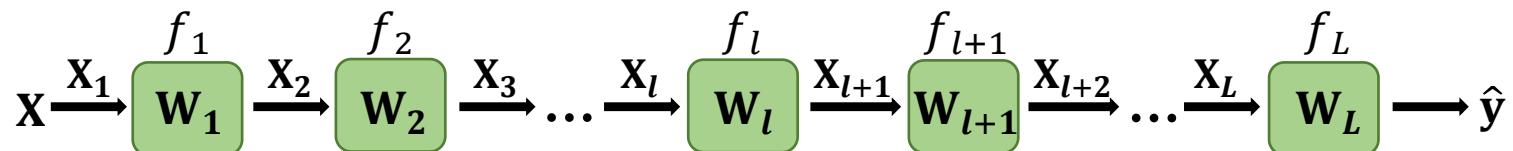
$$\begin{aligned}\mathbf{X}_1 &= \mathbf{X} \\ \mathbf{X}_2 &= f_1(\mathbf{X}_1, \mathbf{W}_1) \\ &\vdots \\ \mathbf{X}_{l+1} &= f_l(\mathbf{X}_l, \mathbf{W}_l) \\ &\vdots \\ \mathbf{X}_L &= f_{L-1}(\mathbf{X}_{L-1}, \mathbf{W}_{L-1}) \\ \hat{\mathbf{y}} &= f_L(\mathbf{X}_L, \mathbf{W}_L)\end{aligned}$$

Compact  
Form

- Forward propagation for any layer  $l$ :

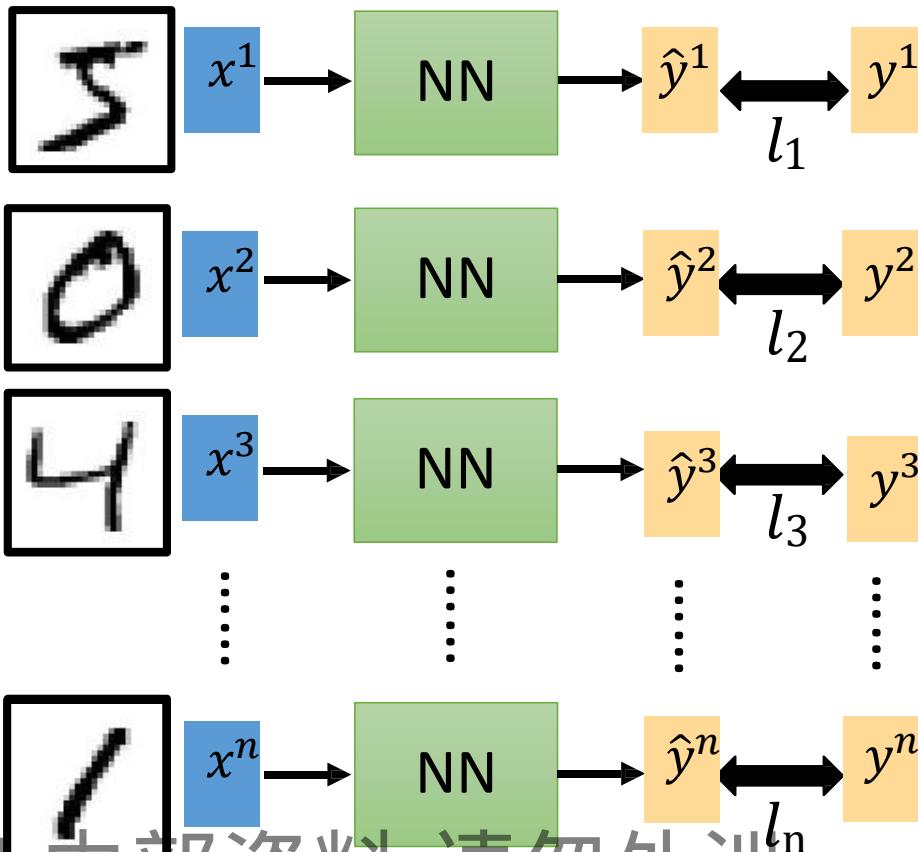
$$\begin{cases} \mathbf{Z}_l = h_l(\mathbf{X}_l, \mathbf{W}_l) \\ \mathbf{X}_{l+1} = \sigma(\mathbf{Z}_l) \end{cases}$$

$h_l$ : some linear function  
 $\sigma$ : some activation function



# Example of Forward Prediction

For all training data:



Total loss:

$$\mathcal{L} = \sum_{r=1}^n l_r$$

As small as possible

Find a function in function set that minimizes total loss  $\mathcal{L}$

Find the network parameters  $\mathbf{W}^*$  that minimize total loss  $\mathcal{L}$

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

4 Backpropagation

5 Convolutional Neural Networks

6 Applications

SMIL内部资料 请勿外泄

# **Backpropagation** of neural networks

# How to Learn Model Parameters $\{\mathbf{W}_l\}$ ?

- We minimize the following regularized loss function:

Regularized Loss function:

$$\min_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_W (\mathbf{y}_i, \hat{\mathbf{y}}_i) + \frac{1}{2} \lambda \|\mathbf{W}\|_F^2$$

$\mathbf{W}$  – parameters

$\mathcal{L}$  – loss function

$\mathbf{y}_i$  – ground-truth label

$\hat{\mathbf{y}}_i$  – predicted label

Connection to previous models?

- Linear Regression?
- Support Vector Machine?
- Logistic Regression?
- Softmax Regression?
- AdaBoost?
- GBDT?

- Gradient Descent: update the parameters  $\mathbf{W}_l$  by

$$\mathbf{W}'_l = \mathbf{W}_l - \eta \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{W}_l} \quad (1)$$

where  $\eta$  is the learning rate

**How to compute the gradient  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_l}$ ?**

# How to Compute Gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_l}$ ?

- Recall the loss function can be rewritten as the Composite Function:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \mathcal{L} \circ f_L \circ f_{L-1} \circ \cdots \circ f_l(\mathbf{X}_l, \mathbf{W}_l) \circ \cdots \circ f_2 \circ f_1$$

**Computing the gradient  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_l}$  is extremely difficult!**

- Let us resort to **Chain Rule** to compute the gradient

$$z = x \circ s \Rightarrow$$

$$\frac{\partial z}{\partial s} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial s}$$

?

# Chain Rule Revisited

- Suppose we have the following differentiable functions:

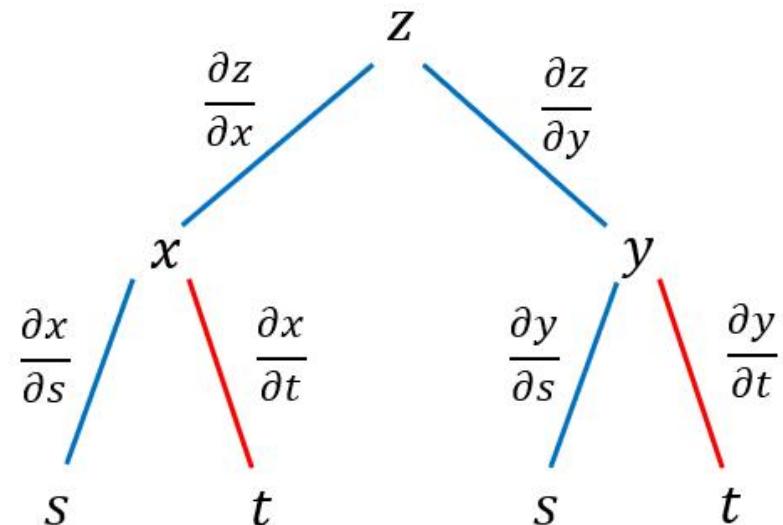
$$\begin{cases} z = f(x, y) \\ x = g(s, t) \\ y = h(s, t) \end{cases}$$

How to compute  $\frac{\partial z}{\partial s}$  and  $\frac{\partial z}{\partial t}$ ?

- We have **Chain Rule** as below:

$$\frac{\partial z}{\partial s} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial s}$$

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial t}$$



# Chain Rule Revisited

- Vector product derivative example:

functions:  $Y = X^T Q X$

$Q$  is Symmetric Matrices

How to compute  $\frac{\partial Y}{\partial X}$ ?

$$\frac{\partial Y}{\partial x} = \frac{\partial((Q^T X_c)^T X)}{\partial x} + \frac{\partial((Q X_c)^T X)}{\partial x}$$

$$= Q^T X_c + Q X_c$$

$$= 2QX$$

$$(X_c^T Q X)^T = (Q X_c)^T X$$

$$Q^T = Q$$

# Gradient Computation of Neural Network $f$

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \mathcal{L} \circ f_L \circ f_{L-1} \circ \dots \circ f_l \circ \dots \circ f_2 \circ f_1(\mathbf{X}_1, \mathbf{W}_1)$$

How to compute the gradient  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1}$ ?

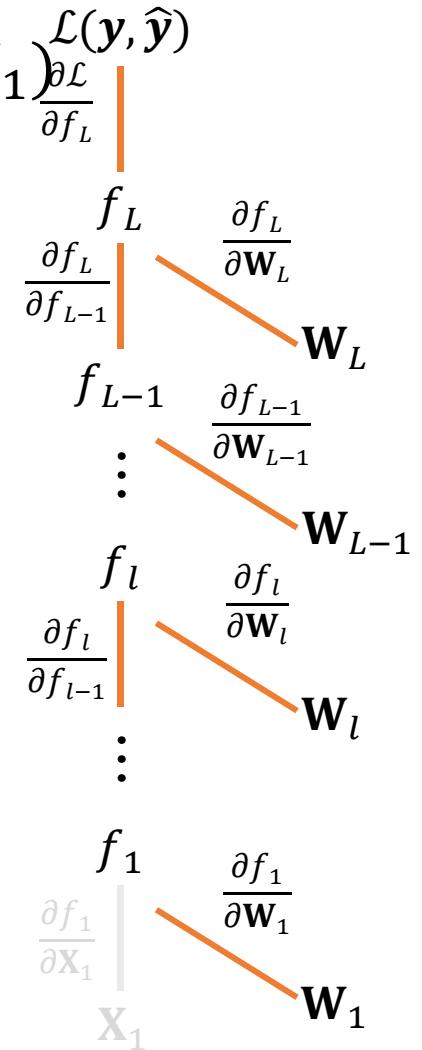
■ According to **Chain Rule**, we have

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_1} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \cdot \dots \cdot \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial \mathbf{W}_1}}$$



Gradient Chain

SMIL内部资料 请勿外泄



# Gradient for Any Layer $l$

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \mathcal{L} \circ f_L \circ f_{L-1} \circ \cdots \circ f_l(\mathbf{X}_l, \mathbf{W}_l) \circ \cdots \circ f_2 \circ f_1$$

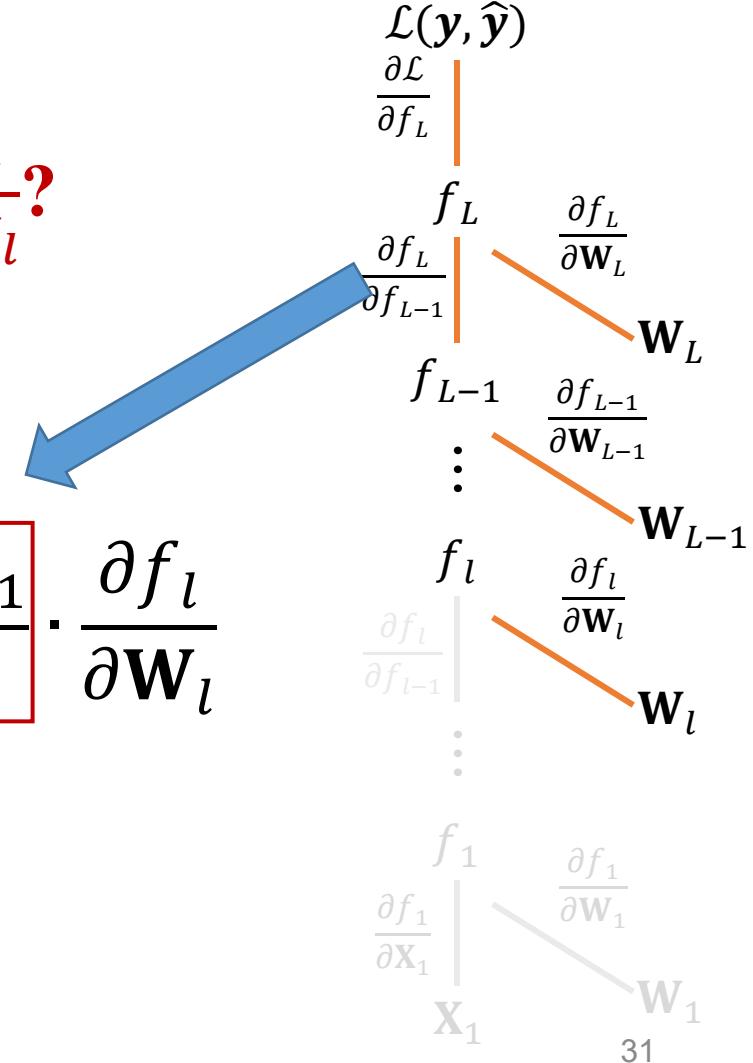
How to compute the gradient  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_l}$ ?

- According to **Chain Rule**, we have

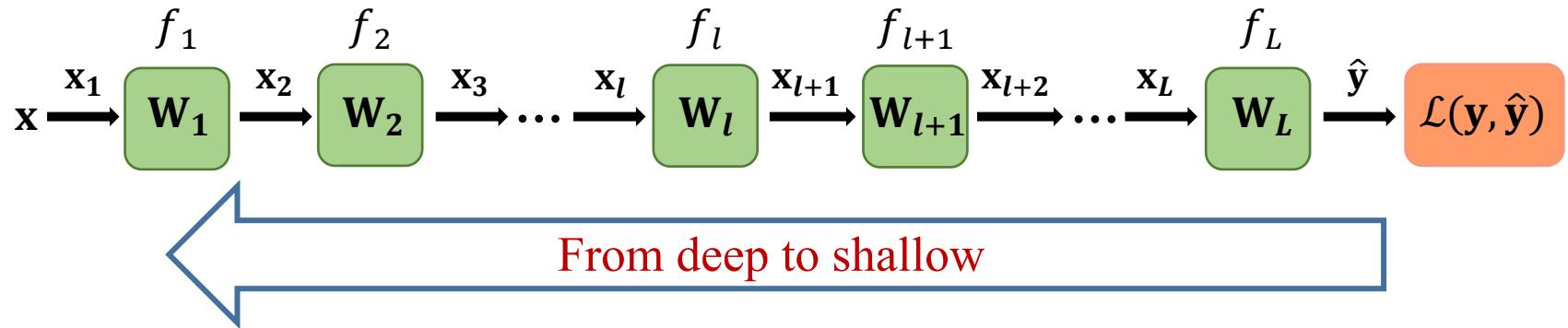
$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_l} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \cdot \cdots \cdot \frac{\partial f_{l+1}}{\partial f_l}} \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$



Gradient Chain



# How to Compute Gradient in Practice



For Layer  $L$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_L} = \frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial \mathbf{W}_L}$$

For Layer  $L - 1$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_{L-1}} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L}} \frac{\partial f_L}{\partial f_{L-1}} \cdot \frac{\partial f_{L-1}}{\partial \mathbf{W}_{L-1}}$$

⋮

For Layer  $l$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_l} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \dots \frac{\partial f_{l+2}}{\partial f_{l+1}}} \frac{\partial f_{l+1}}{\partial f_l} \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$

⋮

For Layer 1

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_1} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \dots \frac{\partial f_{l+2}}{\partial f_{l+1}} \cdot \frac{\partial f_{l+1}}{\partial f_l} \dots \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial \mathbf{W}_1}}$$

# How to Compute Gradient in Practice

Introduce intermediate variable  $\{G_l\}_{l=1}^L$

For Layer  $L$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_L} = \frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial \mathbf{W}_L}$$

For Layer  $L - 1$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_{L-1}} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L}} \cdot \frac{\partial f_L}{\partial f_{L-1}} \cdot \frac{\partial f_{L-1}}{\partial \mathbf{W}_{L-1}}$$

⋮

For Layer  $l$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_l} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \cdot \dots \cdot \frac{\partial f_{l+2}}{\partial f_{l+1}}} \cdot \frac{\partial f_{l+1}}{\partial f_l} \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$

⋮

For Layer 1

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_1} = \boxed{\frac{\partial \mathcal{L}}{\partial f_L} \cdot \frac{\partial f_L}{\partial f_{L-1}} \cdot \dots \cdot \frac{\partial f_3}{\partial f_2}} \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial \mathbf{W}_1}$$

$$G_L = \frac{\partial \mathcal{L}}{\partial f_L}$$

$$G_{L-1} = G_L \cdot \frac{\partial f_L}{\partial f_{L-1}}$$

⋮

$$G_l = G_{l+1} \cdot \frac{\partial f_{l+1}}{\partial f_l}$$

⋮

$$G_1 = G_2 \cdot \frac{\partial f_2}{\partial f_1}$$

# How to Compute Gradient in Practice

For Layer  $L$

For Layer  $L - 1$

$\vdots$

For Layer  $l$

$\vdots$

For Layer 1

$$G_L = \frac{\partial \mathcal{L}}{\partial f_L}$$

$$G_{L-1} = G_L \cdot \frac{\partial f_L}{\partial f_{L-1}}$$

$\vdots$

$$G_l = G_{l+1} \cdot \frac{\partial f_{l+1}}{\partial f_l}$$

$\vdots$

$$G_1 = G_2 \cdot \frac{\partial f_2}{\partial f_1}$$



$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_L} = G_L \cdot \frac{\partial f_L}{\partial \mathbf{W}_L}$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_{L-1}} = G_{L-1} \cdot \frac{\partial f_{L-1}}{\partial \mathbf{W}_{L-1}}$$

$\vdots$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_l} = G_l \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$

$\vdots$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_1} = G_1 \cdot \frac{\partial f_1}{\partial \mathbf{W}_1}$$

# How to Compute Gradient in Practice

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_L} = G_L \cdot \frac{\partial f_L}{\partial \mathbf{W}_L}$$

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_{L-1}} = G_{L-1} \cdot \frac{\partial f_{L-1}}{\partial \mathbf{W}_{L-1}}$$

⋮

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_l} = G_l \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$

⋮

$$\frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_1} = G_1 \cdot \frac{\partial f_1}{\partial \mathbf{W}_1}$$



$$G_L = \frac{\partial \mathcal{L}}{\partial f_L}$$

$$G_{L-1} = G_L \cdot \frac{\partial f_L}{\partial f_{L-1}}$$

⋮

$$G_l = G_{l+1} \cdot \frac{\partial f_{l+1}}{\partial f_l}$$

⋮

$$G_1 = G_2 \cdot \frac{\partial f_2}{\partial f_1}$$

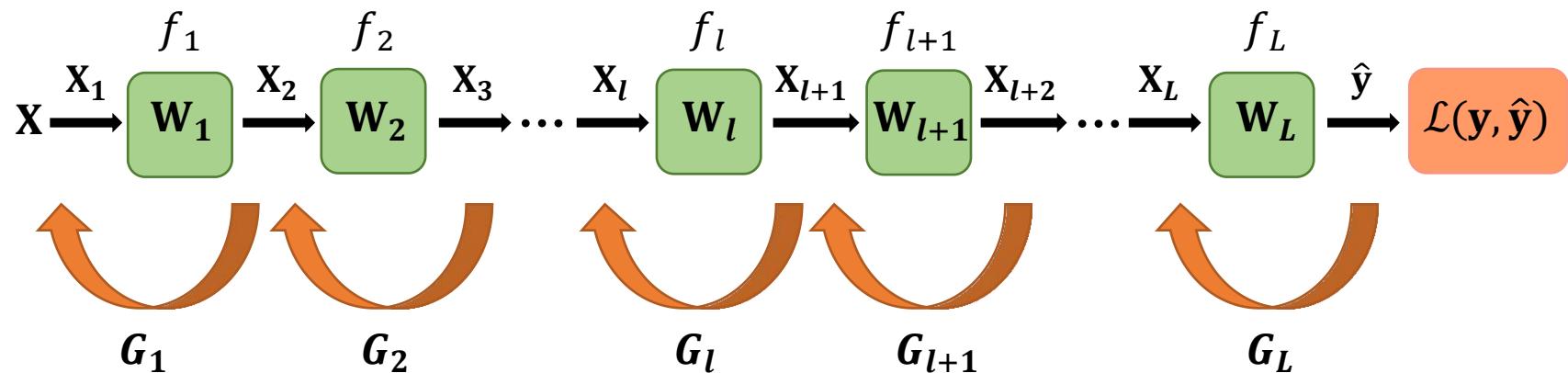
$$f_l = \sigma(h_l(\mathbf{X}_l, \mathbf{W}_l)) \implies \frac{\partial f_l}{\partial \mathbf{W}_l} = \frac{\partial \sigma}{\partial h_l} \cdot \frac{\partial h_l}{\partial \mathbf{W}_l} \quad \frac{\partial f_l}{\partial f_{l-1}} = \frac{\partial \sigma}{\partial h_l} \cdot \frac{\partial h_l}{\partial f_{l-1}}$$

SMIL 内部资料 请勿外泄

# Backpropagation for Neural Networks

## Backpropagation:

**Backpropagation** propagates the **error/loss** from the **last layer** to the **first layer** to compute the **gradient** for each layer



# Forward Propagation and Backpropagation

## Algorithm 1 Forward and Backward Propagation for an $L$ -layer Neural Network

**Require:** Input data  $\mathbf{X}$ , label  $\mathbf{y}$ , layer functions  $\{f_l\}_{l=1}^L$ , layer parameters  $\{\mathbf{W}_l\}_{l=1}^L$

### Step1: Conduct forward propagation

Compute  $\hat{\mathbf{y}} = f_W(\mathbf{X}) = f_L \circ f_{L-1} \circ \dots \circ f_l \circ \dots \circ f_2 \circ f_1(\mathbf{X}_1, \mathbf{W}_1)$

### Step2: Compute loss $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$

### Step3: Conduct backpropagation

Compute  $G = \frac{\partial \mathcal{L}}{\partial f_L}$  and  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_L} = G \cdot \frac{\partial f_L}{\partial \mathbf{W}_L}$

**For**  $l = L - 1$  to 1 **do**:

propagate  $G_{l+1}$  to current layer: 
$$G = G \cdot \frac{\partial f_{l+1}}{\partial f_l}$$

compute gradient of current layer: 
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_l} = G \cdot \frac{\partial f_l}{\partial \mathbf{W}_l}$$

**End**

Update parameter  $\mathbf{W}_l$  for any layer  $l$  with (**Stochastic**) Gradient Descent

$$\mathbf{W}_l := \mathbf{W}_l - \eta \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{W}_l}$$

# Comparison of Training Algorithms

Method	Advantages	Disadvantages
Gradient Descent (GD)	<ul style="list-style-type: none"><li>1. Accurate search direction</li><li>2. Suitable for parallelization</li></ul>	<ul style="list-style-type: none"><li>1. Slow convergence speed</li><li>2. Need huge memory for big data</li></ul>
Stochastic Gradient Descent (SGD)	<ul style="list-style-type: none"><li>1. Cheap for each step</li><li>2. Suitable for big data</li></ul>	<ul style="list-style-type: none"><li>1. Need more steps</li><li>2. Not stable</li><li>3. Hard for parallelization</li></ul>
Mini-Batch Gradient Descent (MBGD)	<ul style="list-style-type: none"><li>1. Converge faster than GD</li><li>2. More accurate update than SGD</li><li>3. Suitable for parallelization</li><li>4. Suitable for large data</li></ul>	<ul style="list-style-type: none"><li>1. Require adding learning-decay to decrease the learning rate</li><li>2. Less stable than GD</li></ul>

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

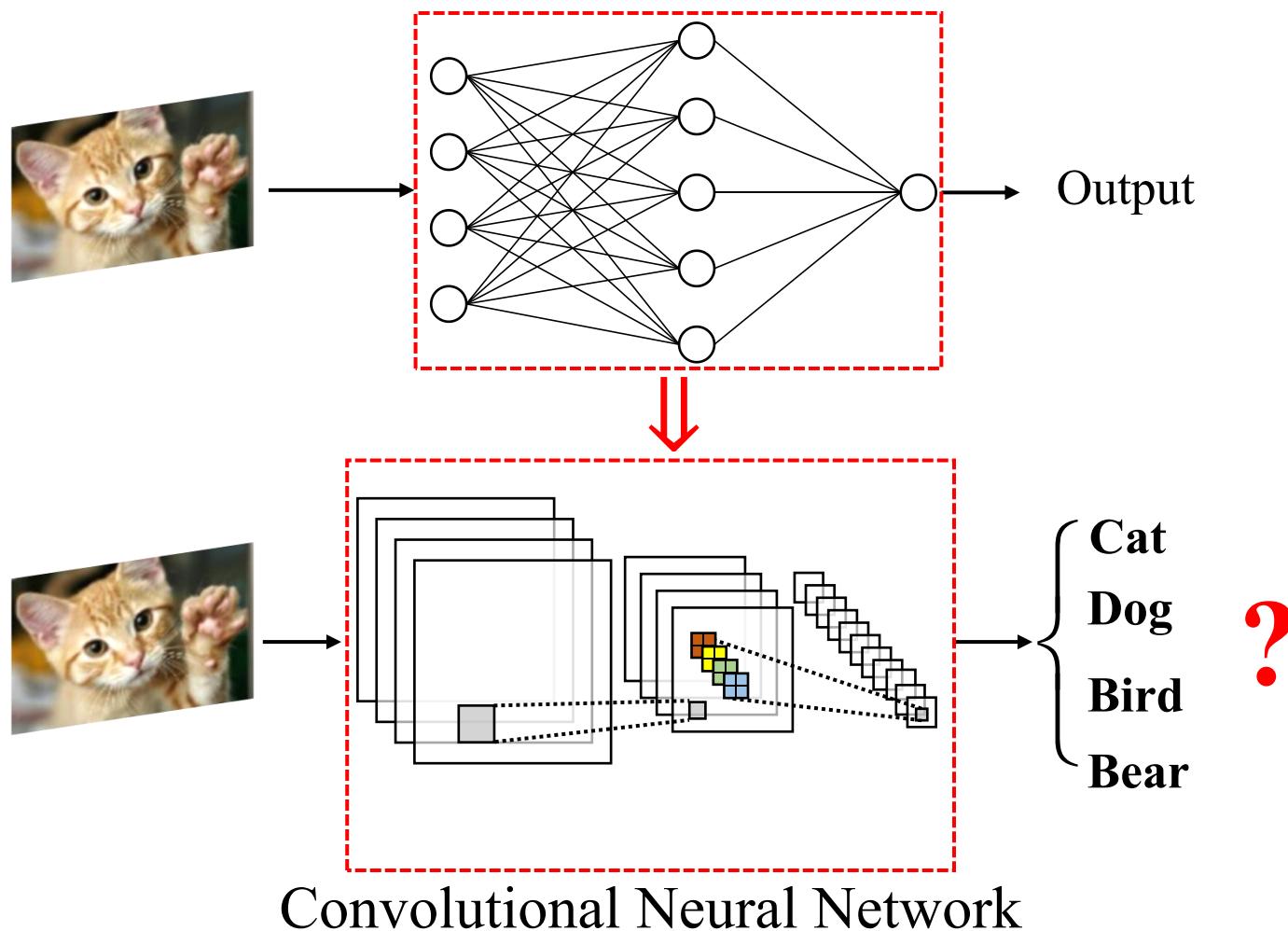
4 Backpropagation

5 Convolutional Neural Networks

6 Applications

SMIL内部资料 请勿外泄

# From Neural Networks to Convolutional Neural Networks

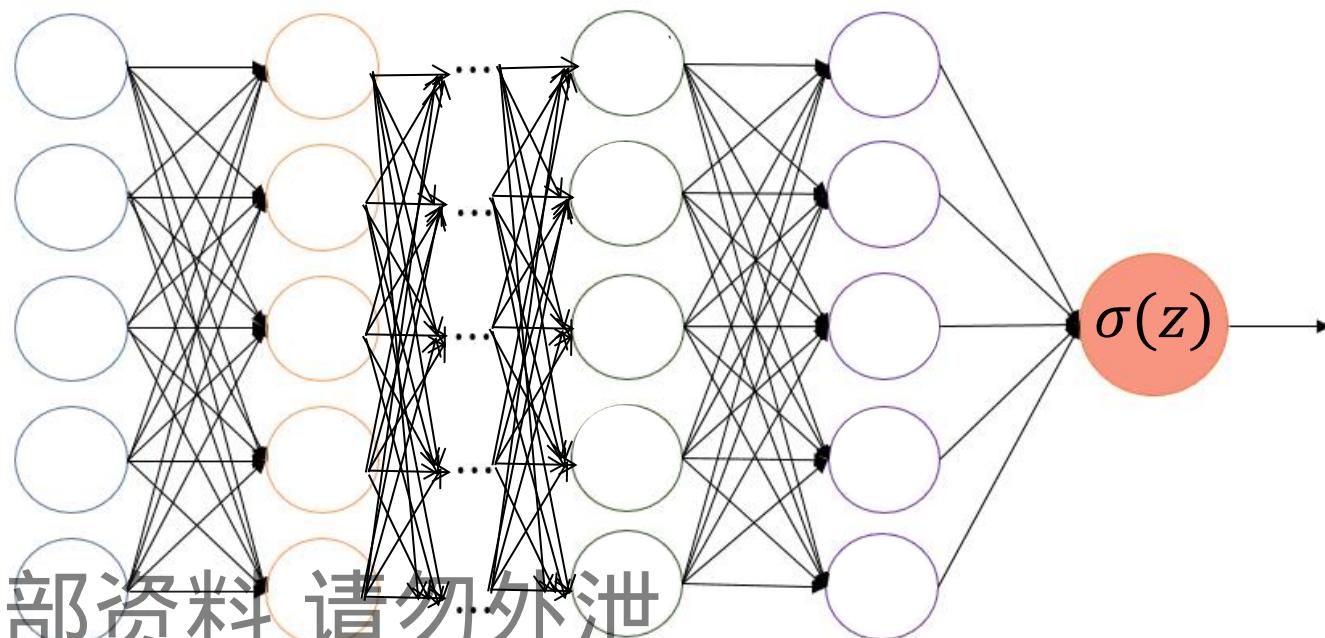
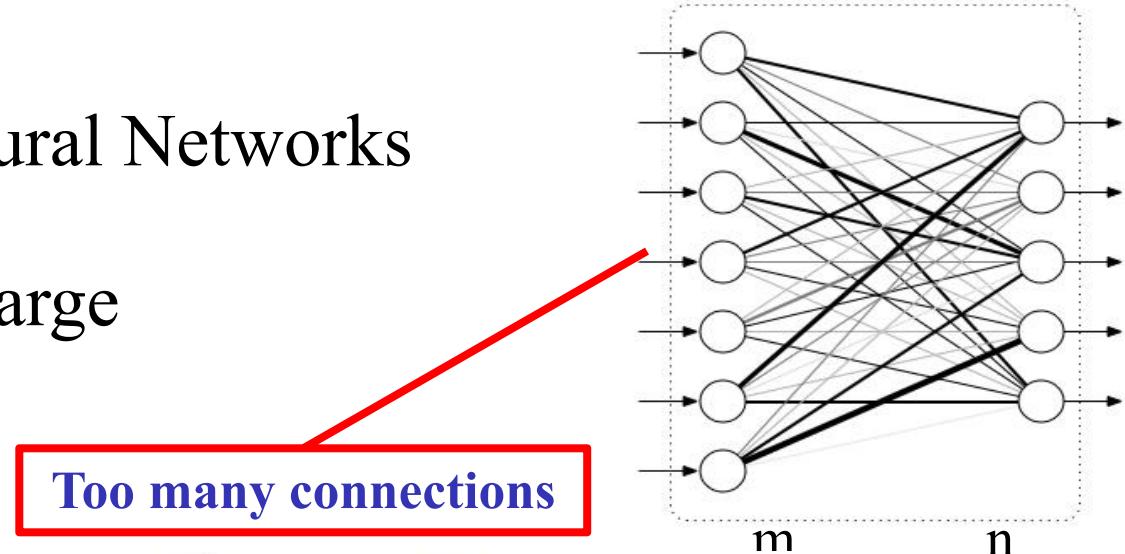


Convolutional Neural Network

SMIL内部资料 请勿外泄

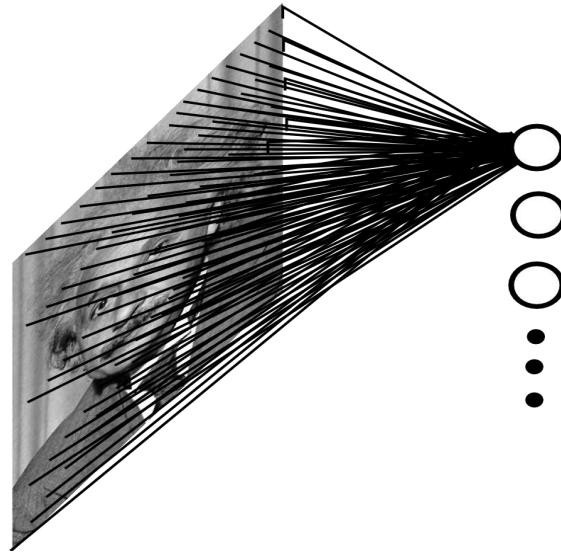
# Why Convolutional Layer

- Fully Connected Neural Networks
- $m \times n$  connections
- $m \times n$  can be very large



SMIL内部资料 请勿外泄

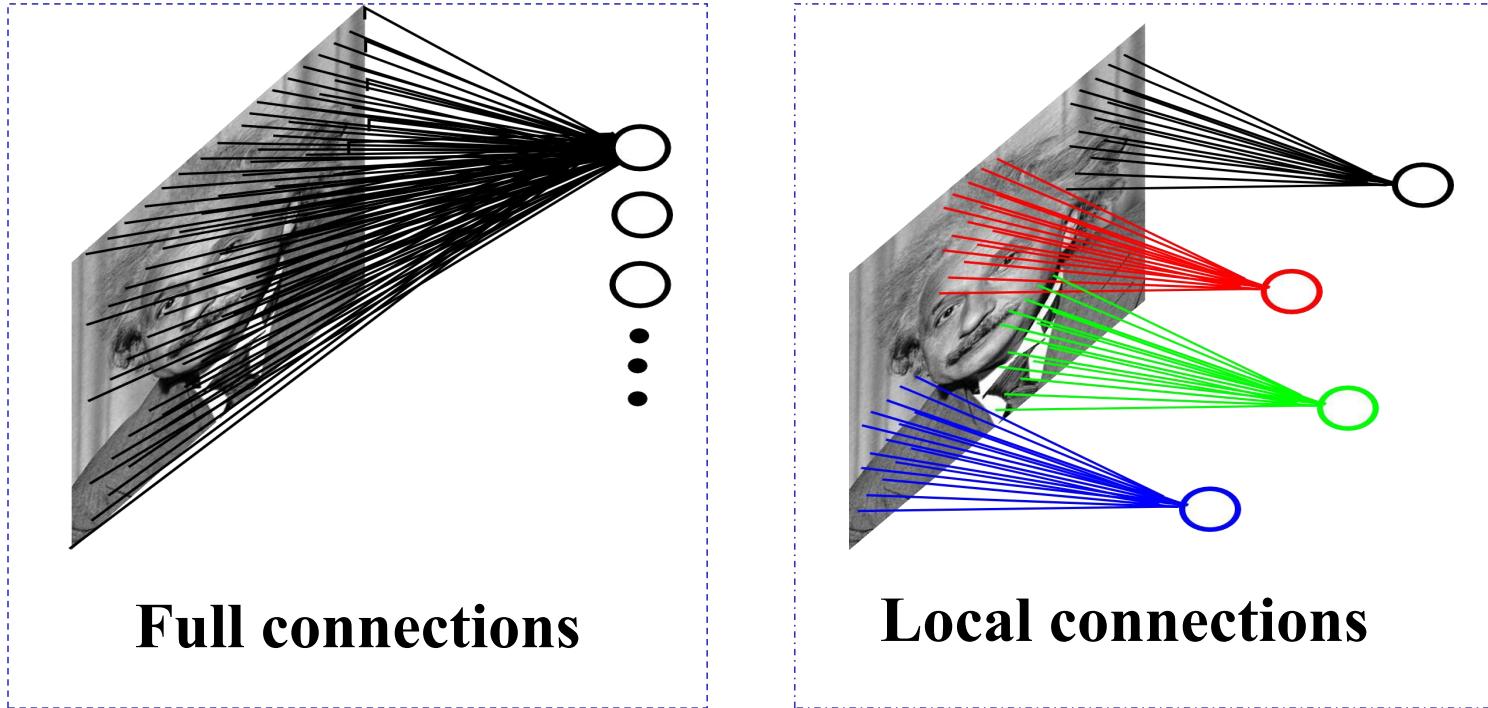
# Too Many Connections and Parameters



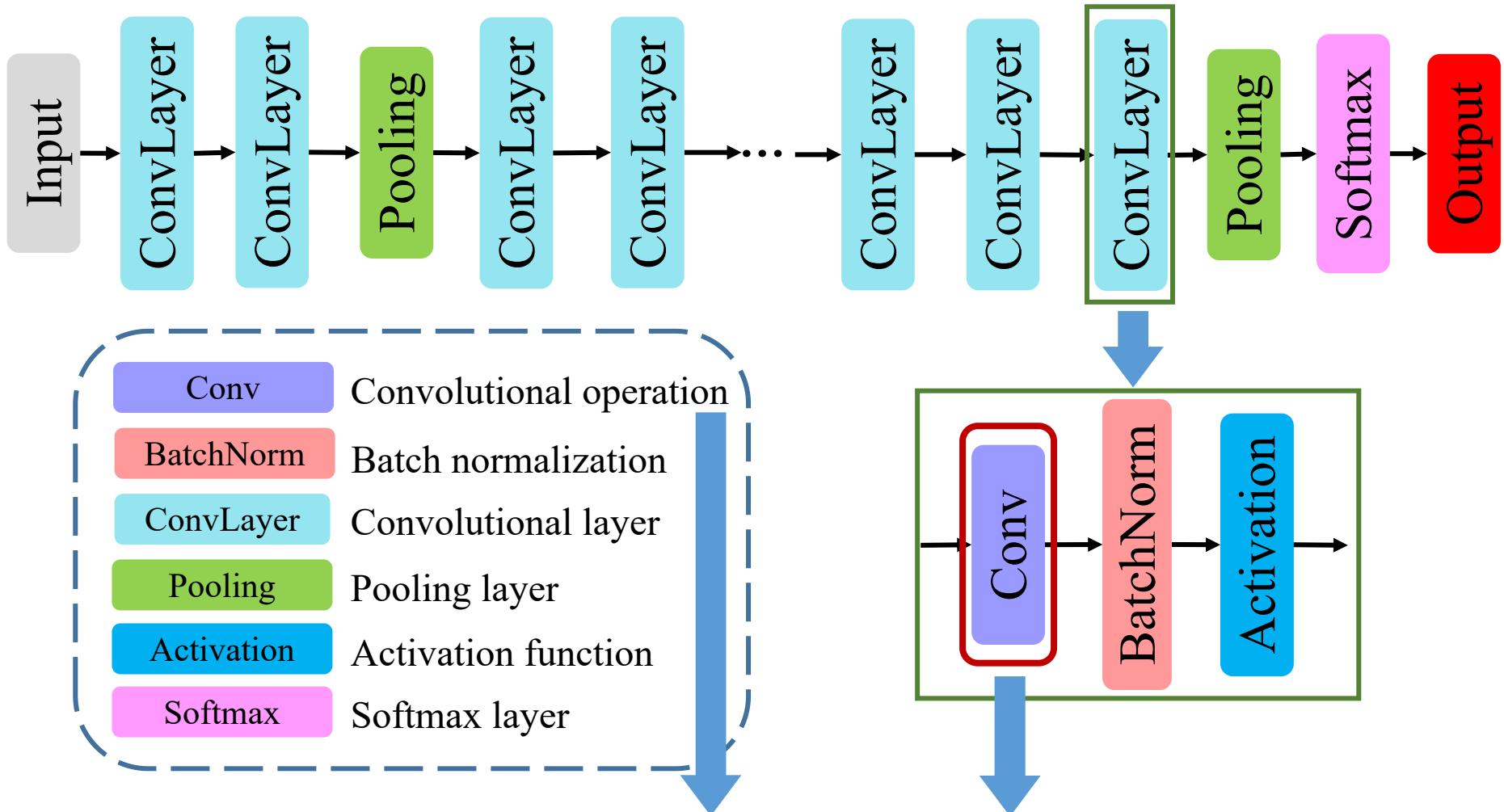
Input Size	#hidden units	#Connections
3x3	n	9n
32x32	n	1024n
224x224	n	50176n
1000x1000	n	$10^6 n$

Too many parameters

# From Full Connections to Local Connections



# Basic Concepts of Convolutional Neural Networks



Convolutional operation

SMIL内部资料 请勿外泄

# The kernel in convolutional layer

## Kernel:

A **kernel / filter** is a matrix of values, called weights, that are **trained** to detect specific features.

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

1	0	1
0	1	0
1	0	1

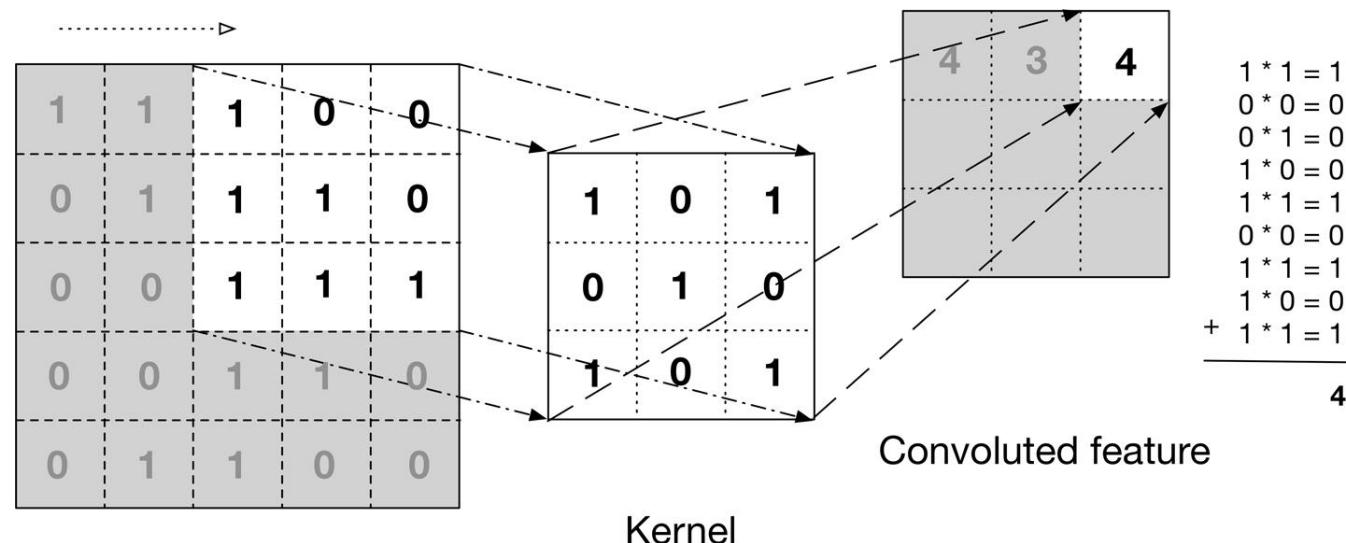
SMIL 内部资料 <sup>Input</sup> 请勿外泄

Filter / Kernel

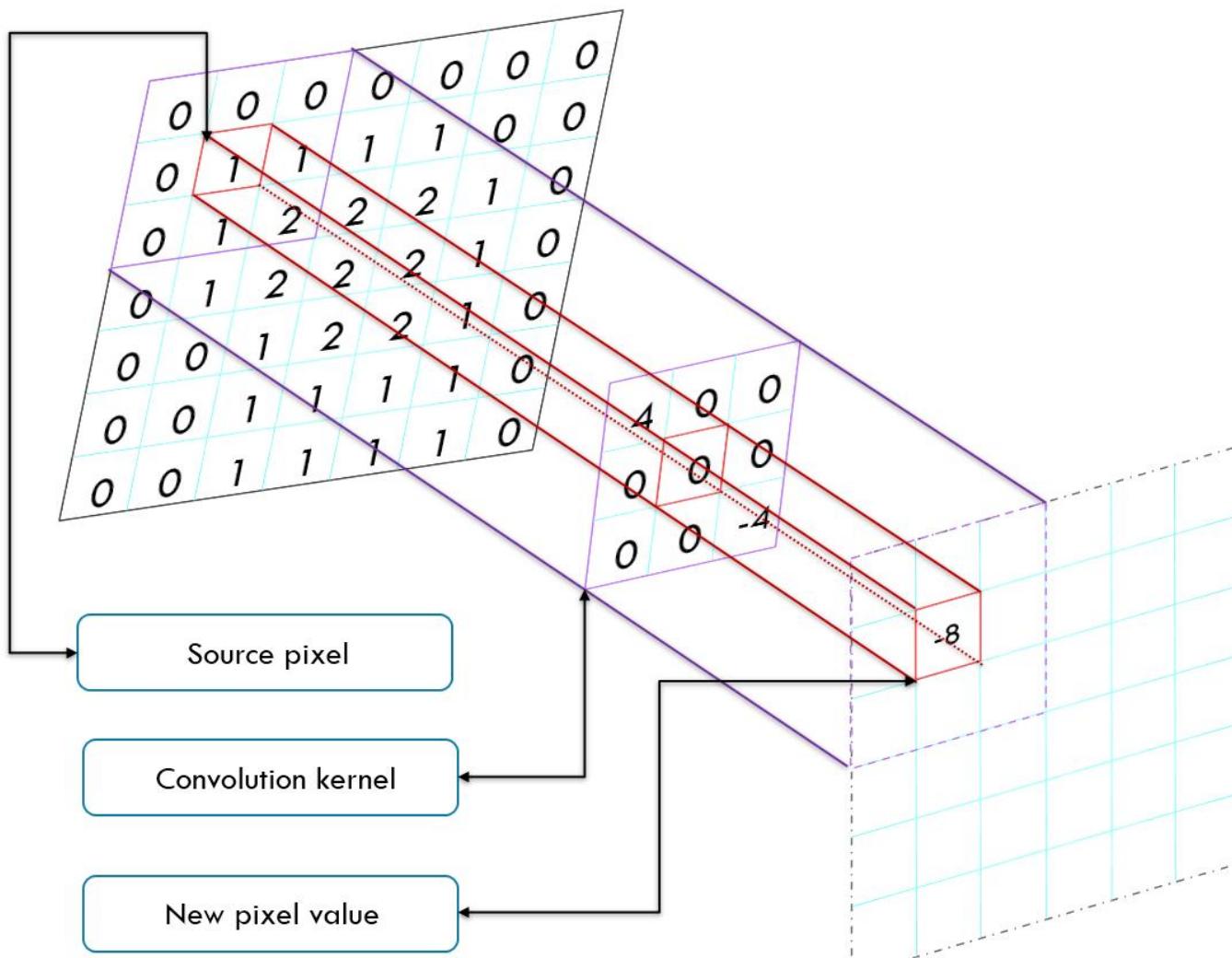
# The convolution operation in convolutional layer

## Convolution operation:

**Convolution operation** is an element-wise product and sum between two matrices.



# Convolution as Local Connections

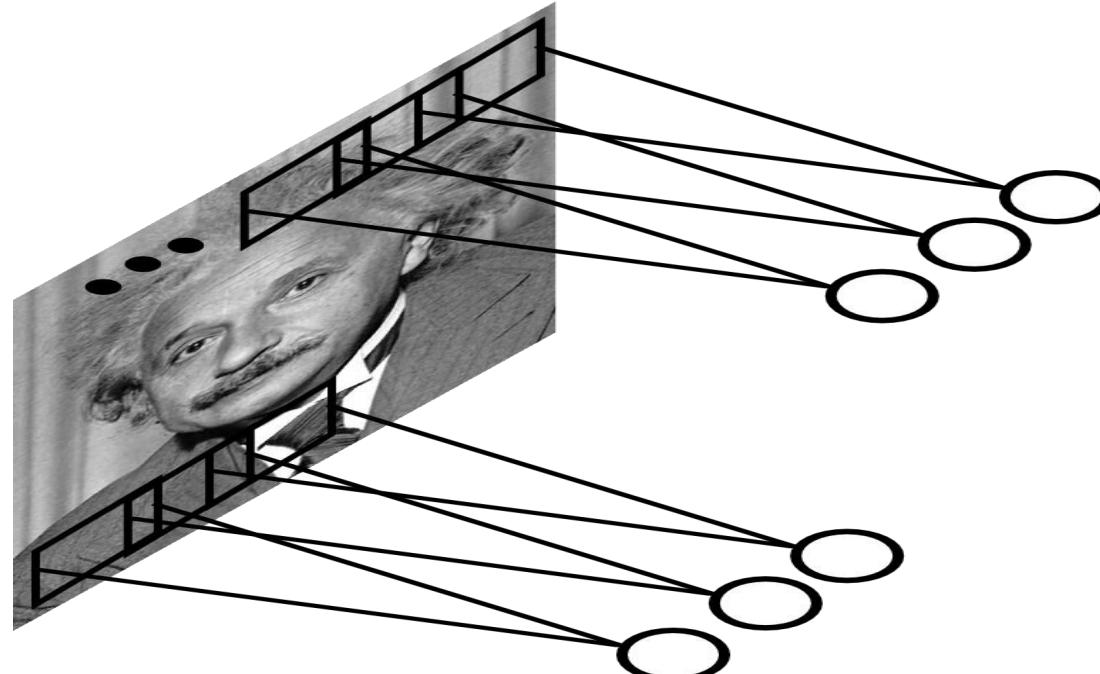


$$\begin{array}{r} 0 \times 4 = 0 \\ 0 \times 0 = 0 \\ 0 \times 0 = 0 \\ 0 \times 0 = 0 \\ 1 \times 0 = 0 \\ 1 \times 0 = 0 \\ 0 \times 0 = 0 \\ + 1 \times 0 = 0 \\ 2 \times -4 = -8 \\ \hline -8 \end{array}$$

SMIL内部资料 请勿外泄

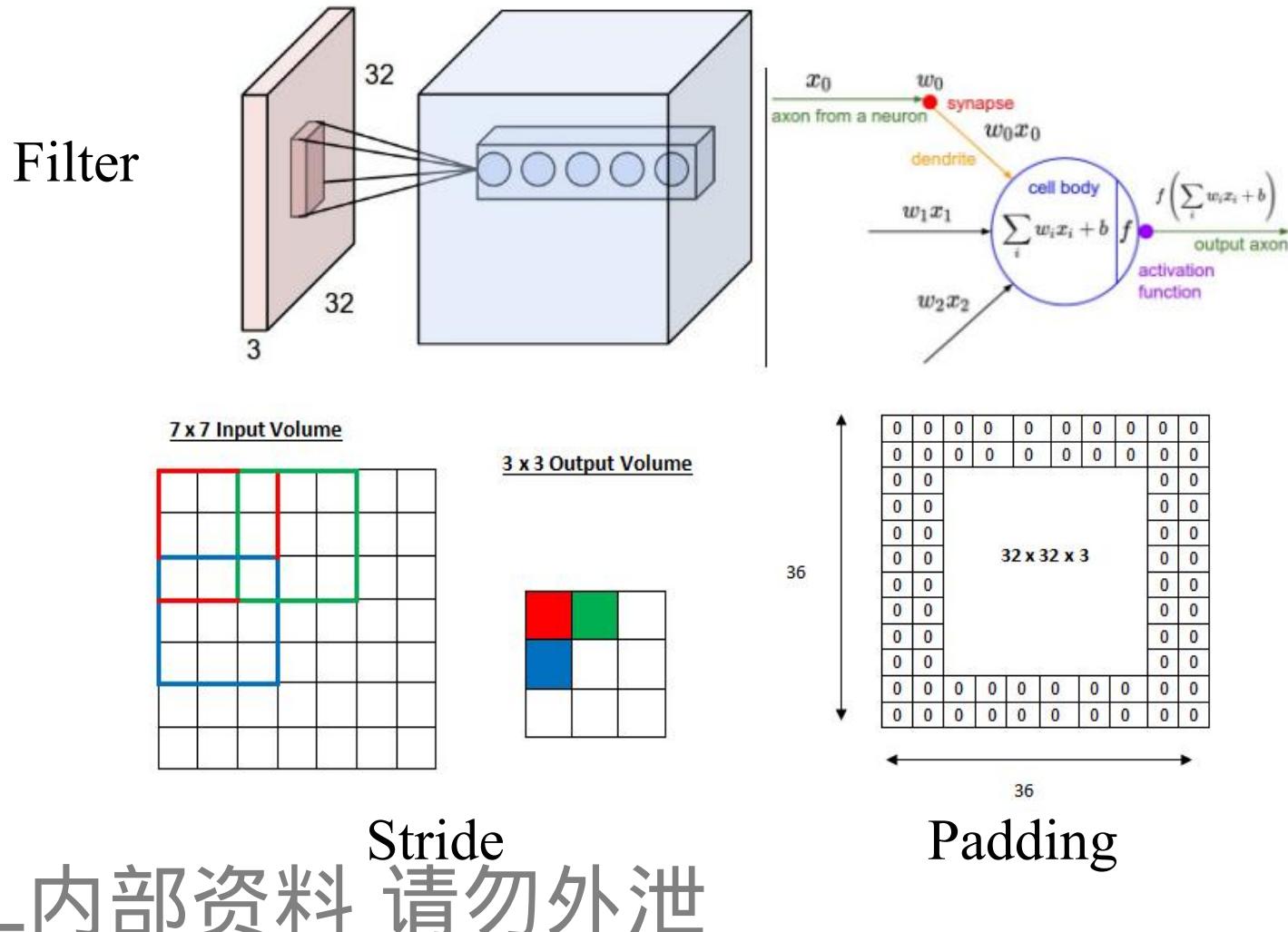
# Further Reduce Parameters: Weights Sharing

- Apply **the same filter** across the whole image
- Apply many kernels
- Each kernel is associated with one feature map



SMIL内部资料 请勿外泄

# Convolutional Layer



SMIL内部资料 请勿外泄

# Example: Convolution with Scanning

The kernel is

1	0	1
0	1	0
1	0	1

Example:

$$\sum \begin{array}{|c|c|c|} \hline 1x1=1 & 1x0=0 & 0x1=0 \\ \hline 1x0=0 & 1x1=1 & 1x0=0 \\ \hline 0x1=0 & 1x0=0 & 1x1=1 \\ \hline \end{array} = 3$$

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3
2		

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3
2	3	

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

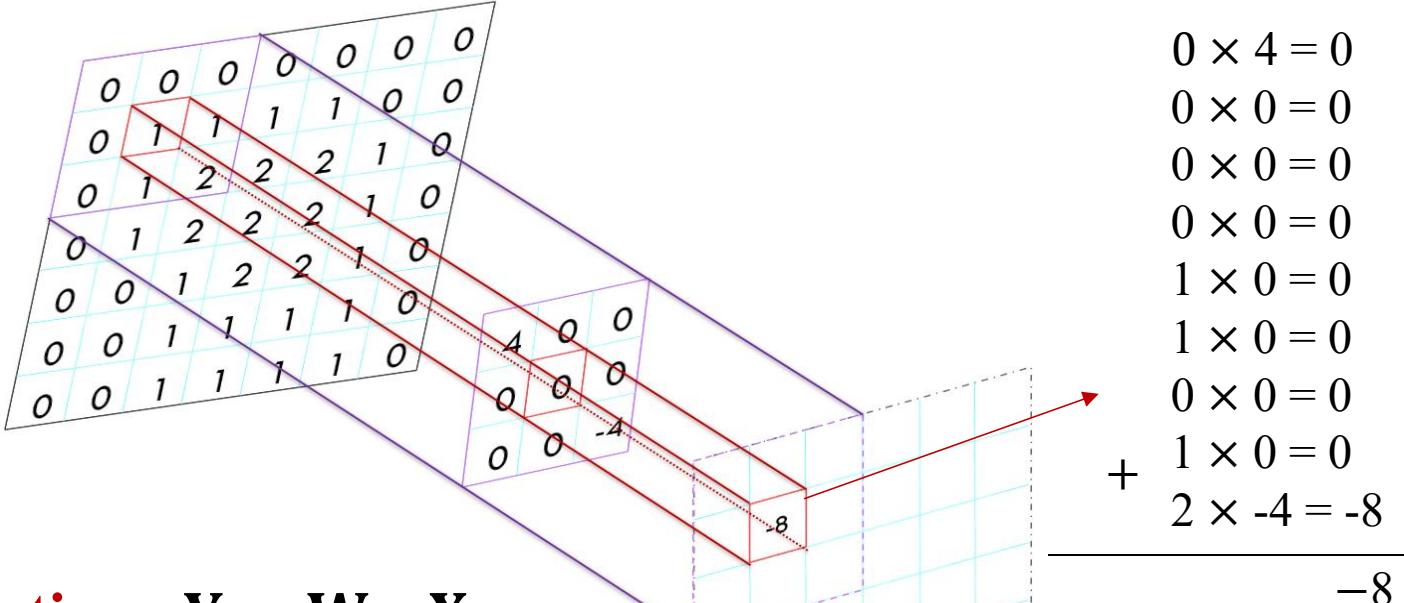
Image

4	3	4
2	4	3
2	3	4

Convolved Feature

SMIL内部资料 请勿外泄

# Properties of Convolution



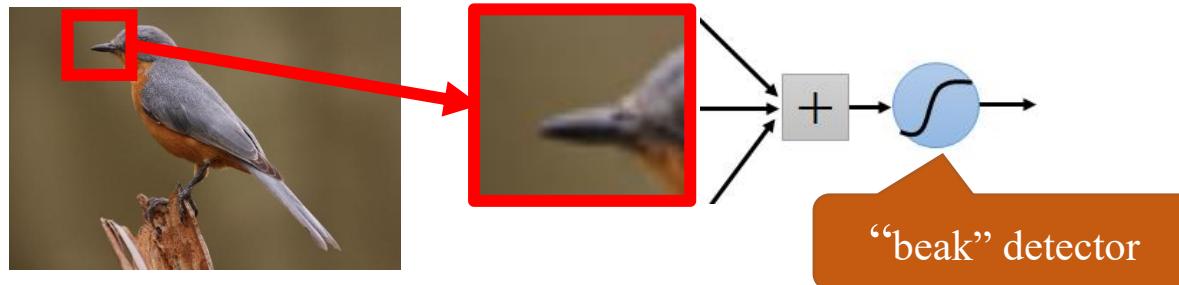
- **Linear Operation** :  $\mathbf{Y} = \mathbf{W} * \mathbf{X}$
- Shift-invariant
- Size: odd by odd ( $3 \times 3, 5 \times 5, 7 \times 7$ )

# Why Local Connections and Parameter Sharing

- Some patterns are much smaller than the whole image

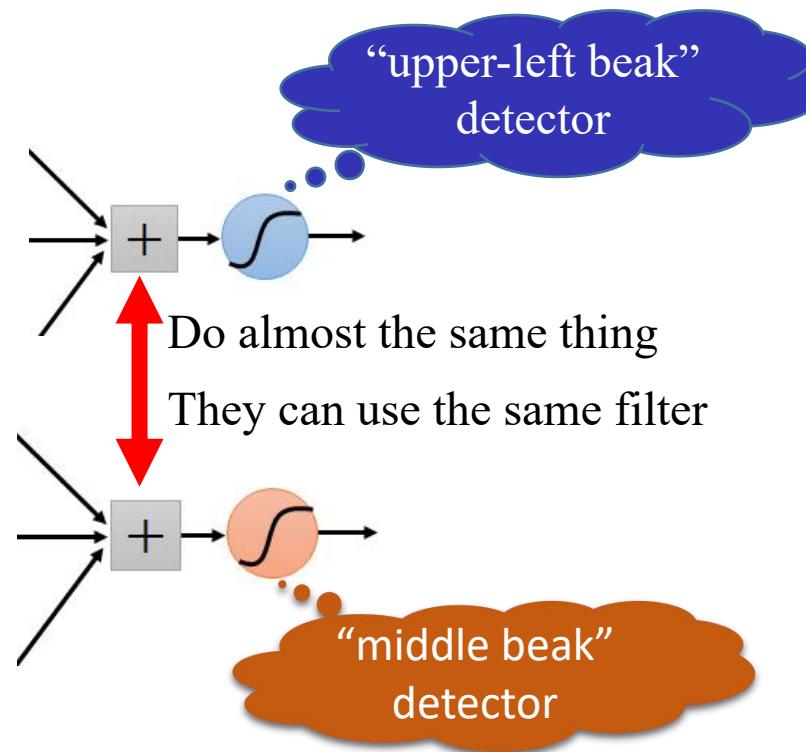
A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



# Why Local Connections

- The same patterns appear in different regions.



# Typical Convolution Examples

Original



$$\begin{matrix} * & \begin{matrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{matrix} & = \end{matrix}$$



# Typical Convolution Examples

## Sharpen



$$\begin{matrix} * & \begin{matrix} -1 & -1 & -1 \\ -1 & 9 & -1 \\ -1 & -1 & -1 \end{matrix} & = \end{matrix}$$



# Typical Convolution Examples

Sharpen with more details



$$\begin{array}{c} * \\ \begin{matrix} -1 & -1 & -1 & -1 & -1 \\ -1 & 2 & 2 & 2 & -1 \\ -1 & 2 & 8 & 2 & -1 \\ -1 & 2 & 2 & 2 & -1 \\ -1 & -1 & -1 & -1 & -1 \end{matrix} \\ = \end{array}$$

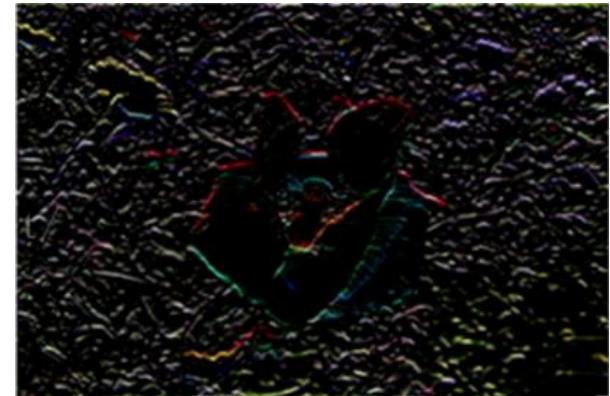


# Typical Convolution Examples

Horizontal edge detection



$$\begin{matrix} * & \begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ -1 & -1 & 2 & -1 & -1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix} \\ = & \end{matrix}$$



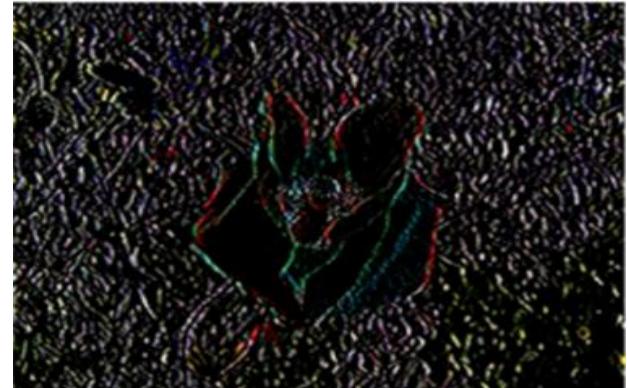
# Typical Convolution Examples

Vertical edge detection



$$\begin{matrix} * & \begin{matrix} 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \end{matrix} \end{matrix}$$

=

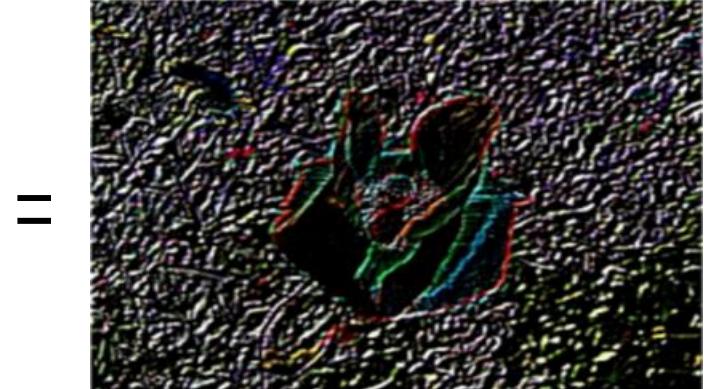


# Typical Convolution Examples

45 degree edge detection



$$\begin{matrix} * & \begin{matrix} -1 & 0 & 0 & 0 & 0 \\ 0 & -2 & 0 & 0 & 0 \\ 0 & 0 & 6 & 0 & 0 \\ 0 & 0 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & -1 \end{matrix} \end{matrix}$$

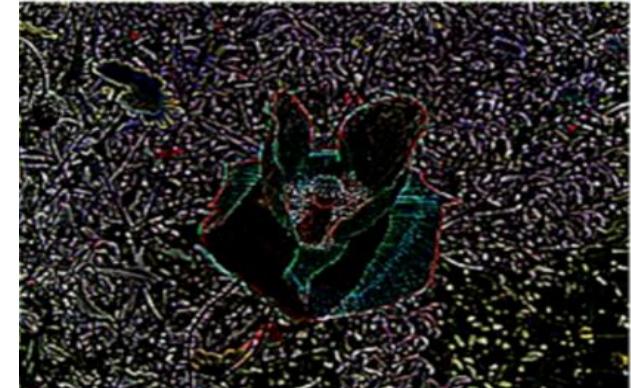


# Typical Convolution Examples

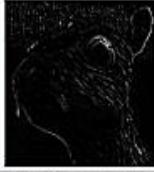
All direction edge detection



$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline -1 & 8 & -1 \\ \hline -1 & -1 & -1 \\ \hline \end{array} & = \end{matrix}$$

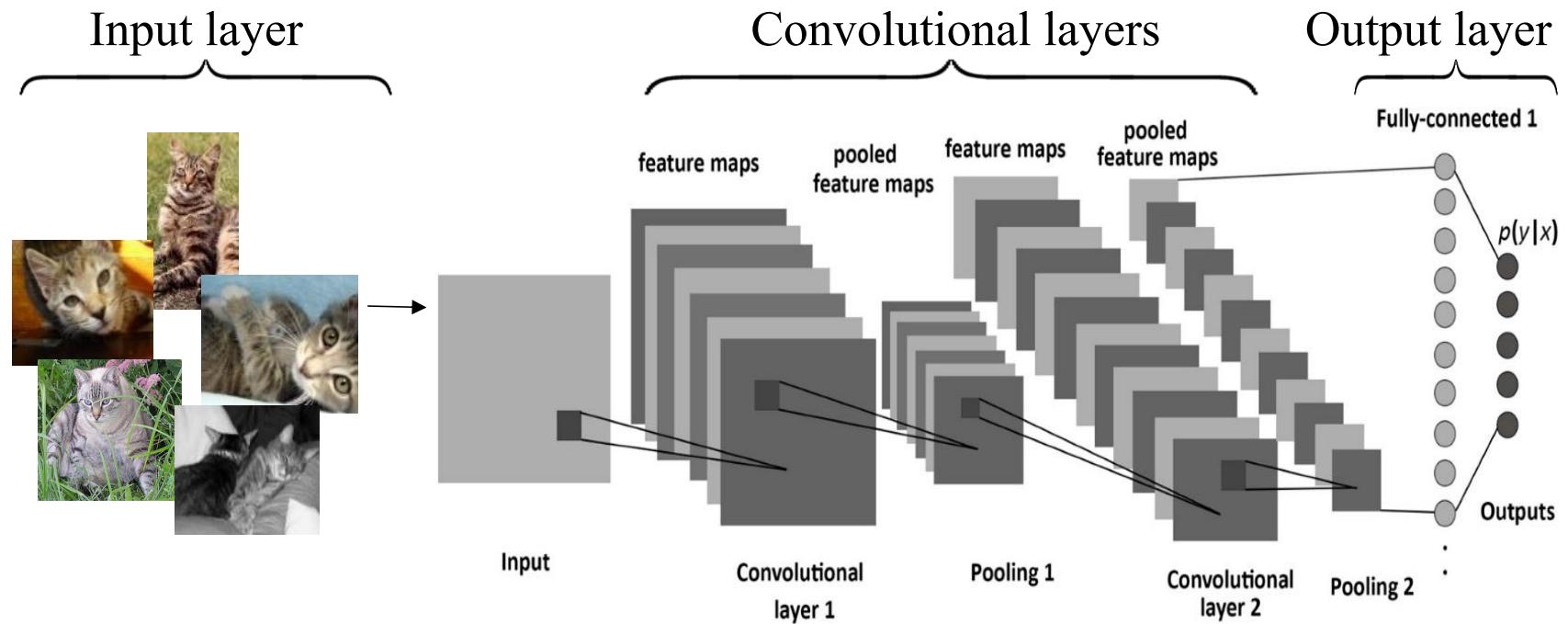


# Some Typical Kernels

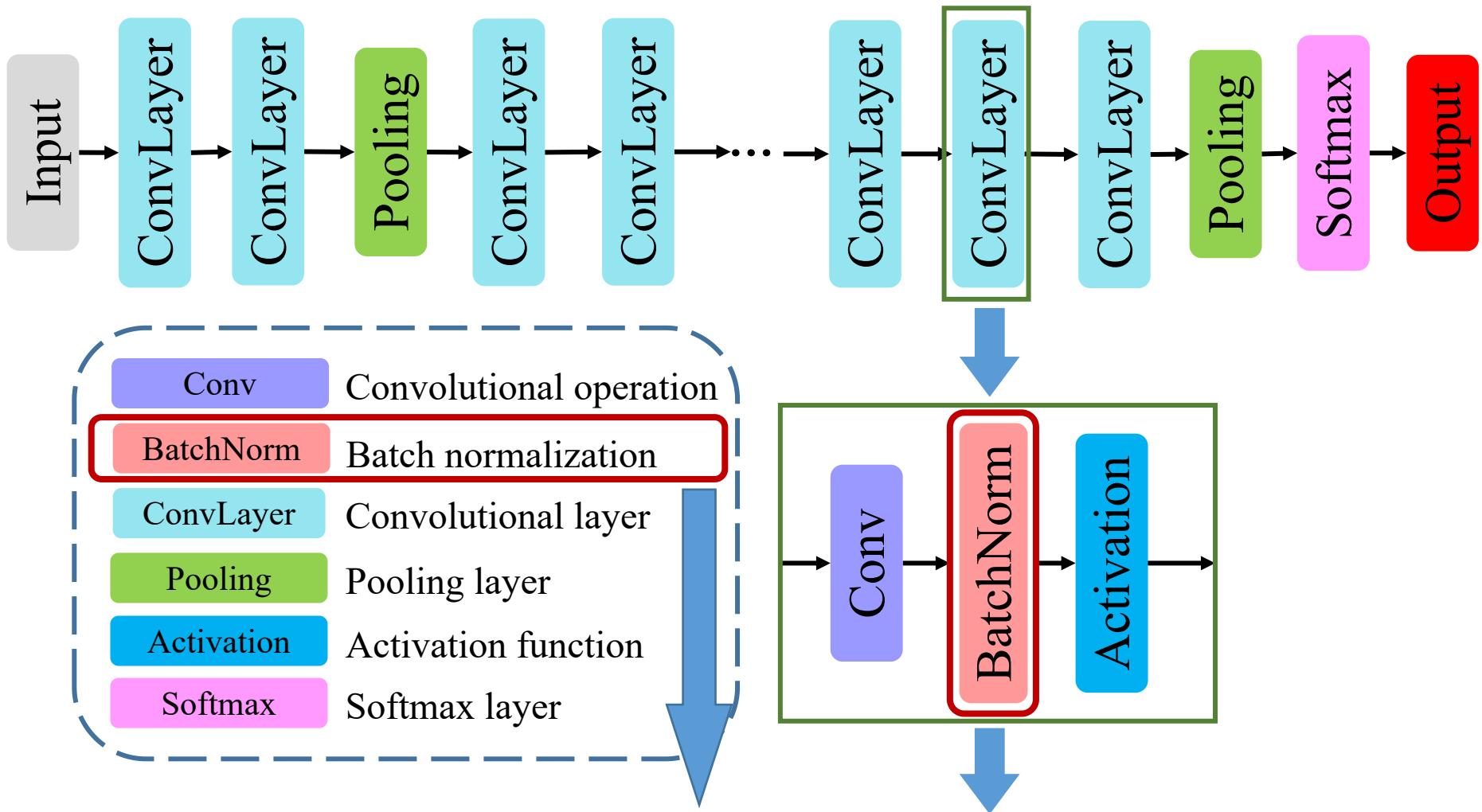
Operation	Kernel	Image result
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	

Operation	Kernel	Image result
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	

# Convolutional Neural Networks



# Batch Normalization



Batch Normalization

SMIL内部资料 请勿外泄

# Batch Normalization Layer

- **Internal Covariate Shift:**

features are normalized to have zero-mean and unit variance

- **Feature Normalization:**

distribution changes during training

- **Batch Normalization:**

**Input:** Values of  $\mathbf{x}$  :  $\mathbf{x} = \{x_1, \dots, x_m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(\mathbf{x})\}$

# Normalization in Each Channel

- Compute data **mean** per channel

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

- Compute data **variance** per channel

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

- Normalization: **Linear Transformation**

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$

- Rescale the normalized value to **restore the representation power**

$$y_i = \gamma \hat{x}_i + \beta$$

# Backpropagation for Batch Normalization Layer

$$\frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \gamma$$

$$\frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot (x_i - \mu_B) \cdot \frac{-1}{2} (\sigma_B^2 + \varepsilon)^{-3/2}$$

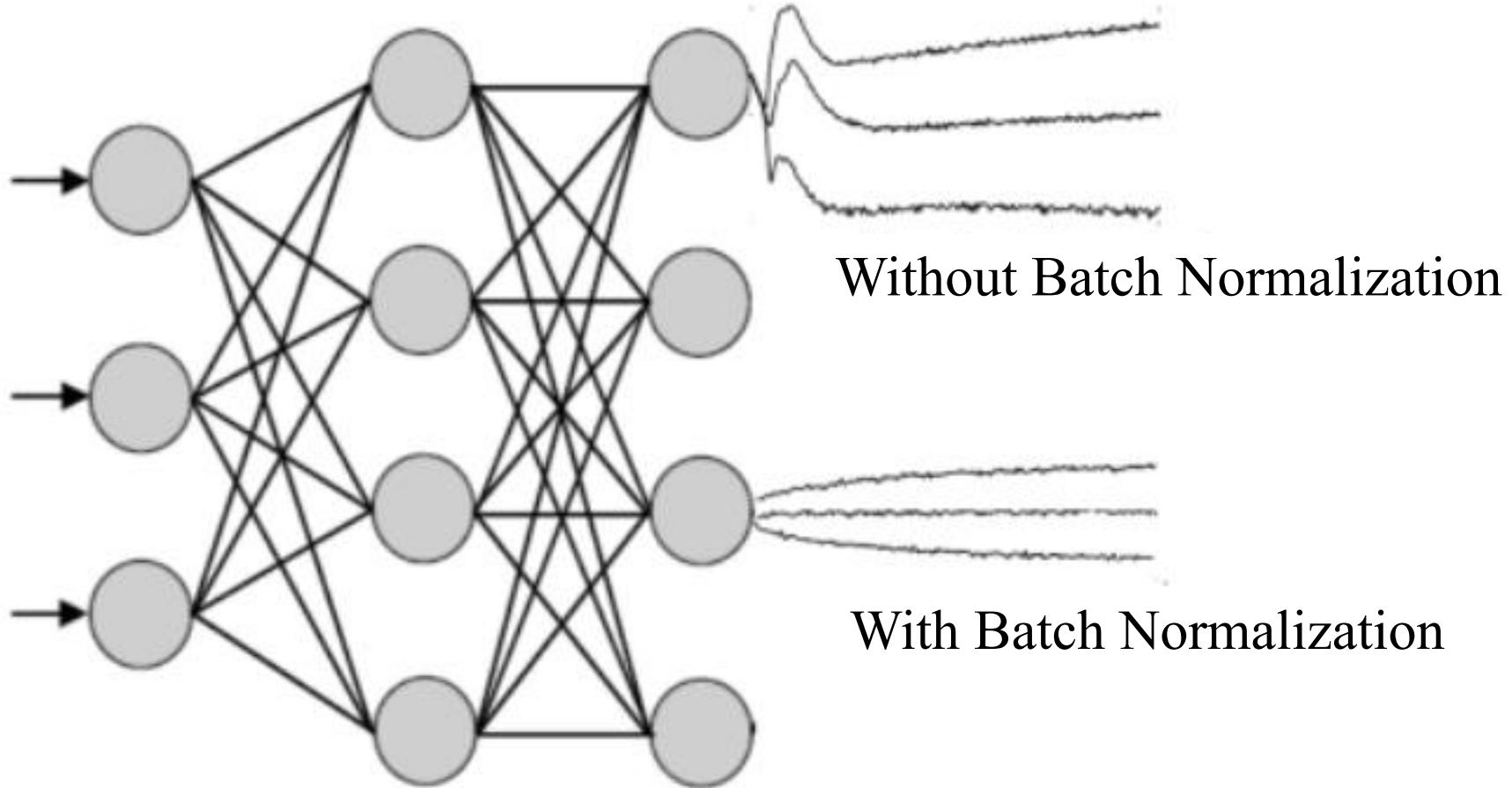
$$\frac{\partial l}{\partial \mu_B} = \left( \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \varepsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_B)}{m}$$

$$\frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \varepsilon}} + \frac{\partial l}{\partial \sigma_B^2} \cdot \frac{2(x_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \cdot \frac{1}{m}$$

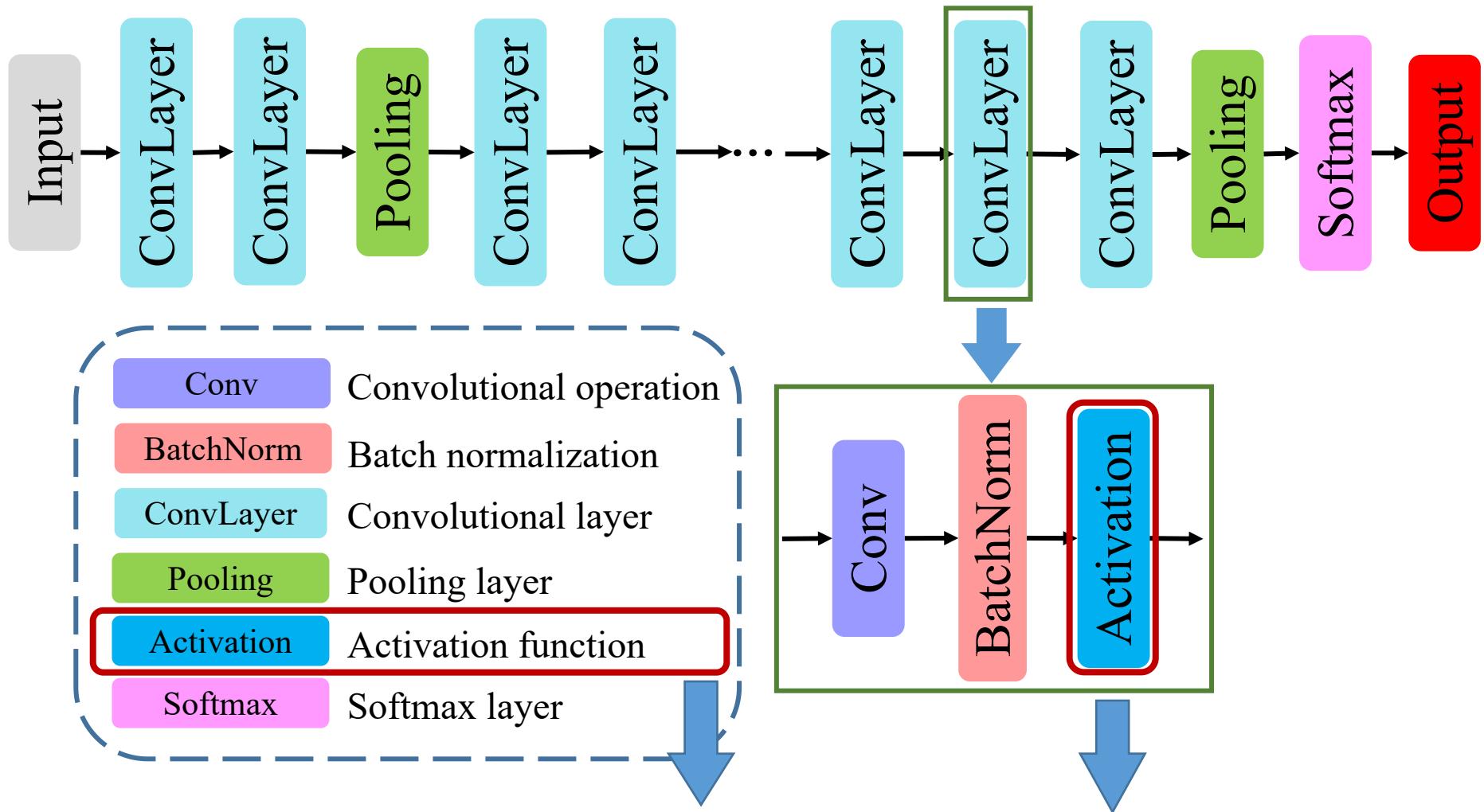
$$\frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \partial \hat{x}_i$$

$$\frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

# Effect of Batch Normalization



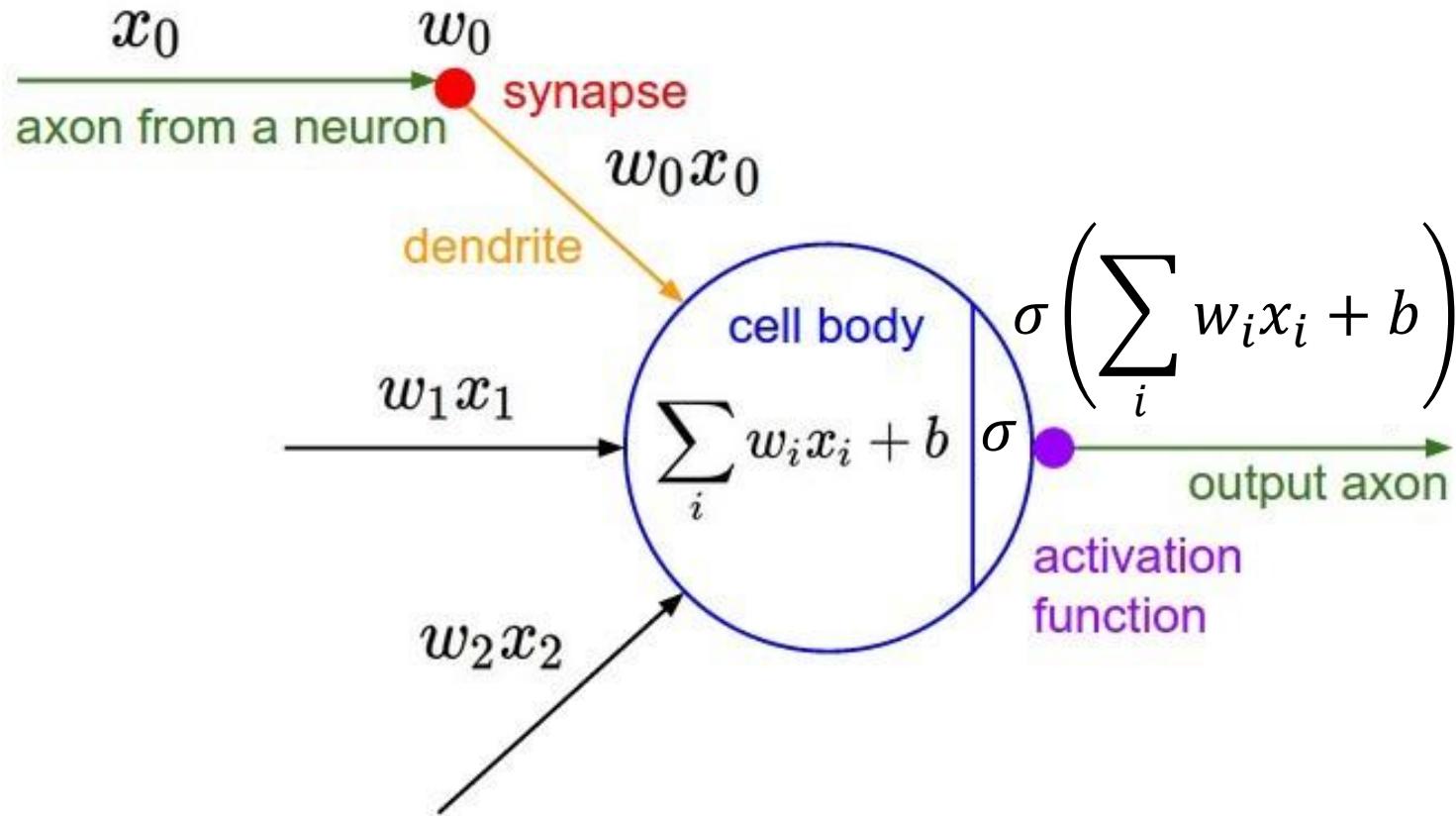
# Activation Layer



SMIL内部资料 请勿外泄 **Activation function**

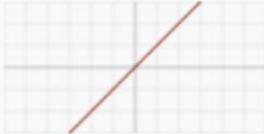
# Activation Function

- **Description:** keep value within an **acceptable** and useful range



# Activation Function

- Activation functions: **nonlinear** and **continuously differentiable**

Activation Function	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Sigmoid		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ReLU		$f(x) = \max(x, 0)$	$f'(x) = 1\{x \geq 0\}$
SoftPlus		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

SMIL内部资料 请勿外泄

# Rectified Linear Unit (ReLU)

## ■ ReLU:

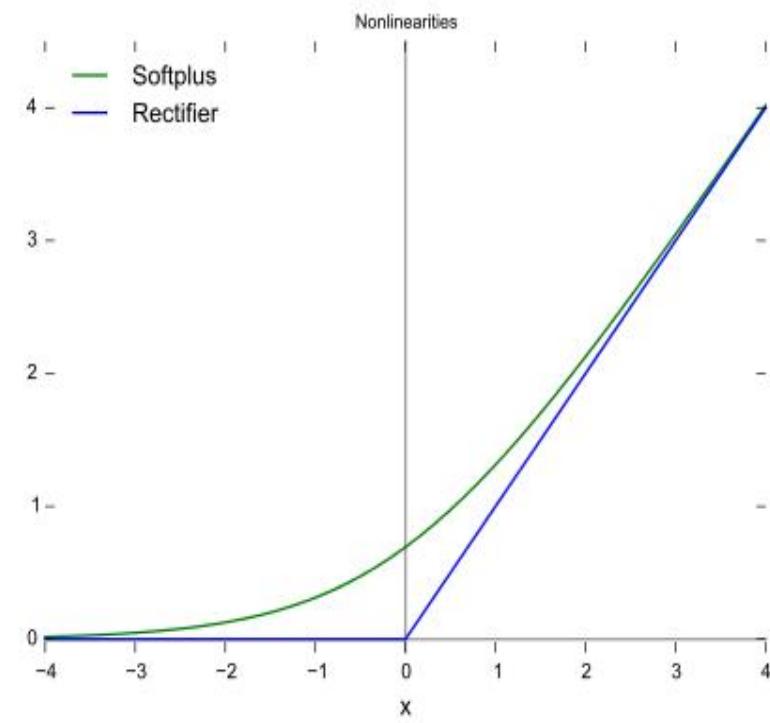
**Forward**  $y = \max(x, 0)$

**Backward**  $\frac{\partial y}{\partial x} = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$

## ■ Softplus:

**Forward**  $y = \ln[1 + e^x]$

**Backward**  $\frac{\partial y}{\partial x} = \frac{1}{1 + e^{-x}}$



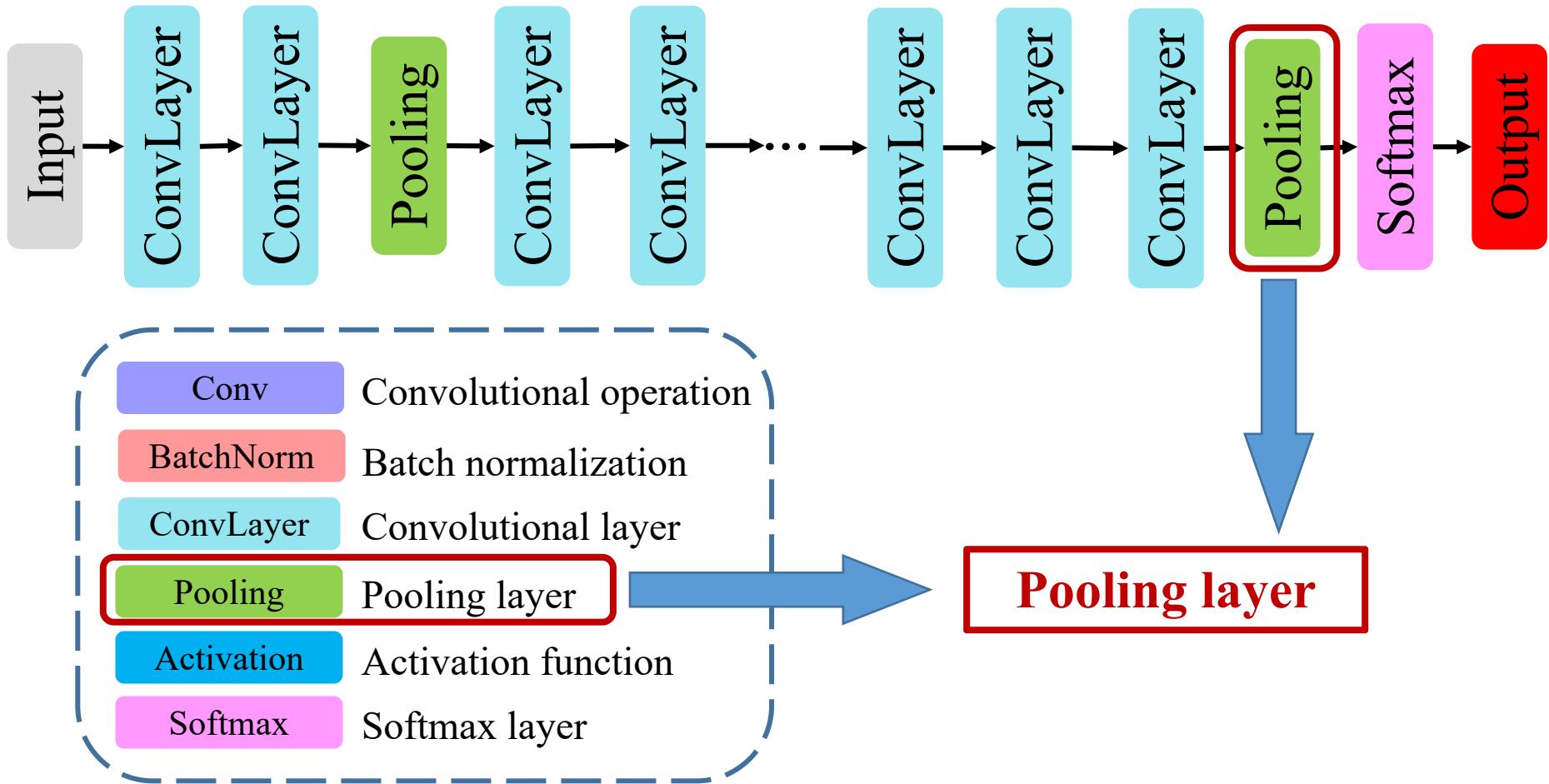
# Leaky ReLU

- Guarantee a **small** and **non-zero** gradient when the unit is not active.

**Forward**  $y = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases} \Leftrightarrow y = \max(x, ax)$   
where  $a < 1$

**Backward**  $y = \begin{cases} 1 & \text{if } x > 0 \\ a & \text{otherwise} \end{cases}$

# Pooling Layer

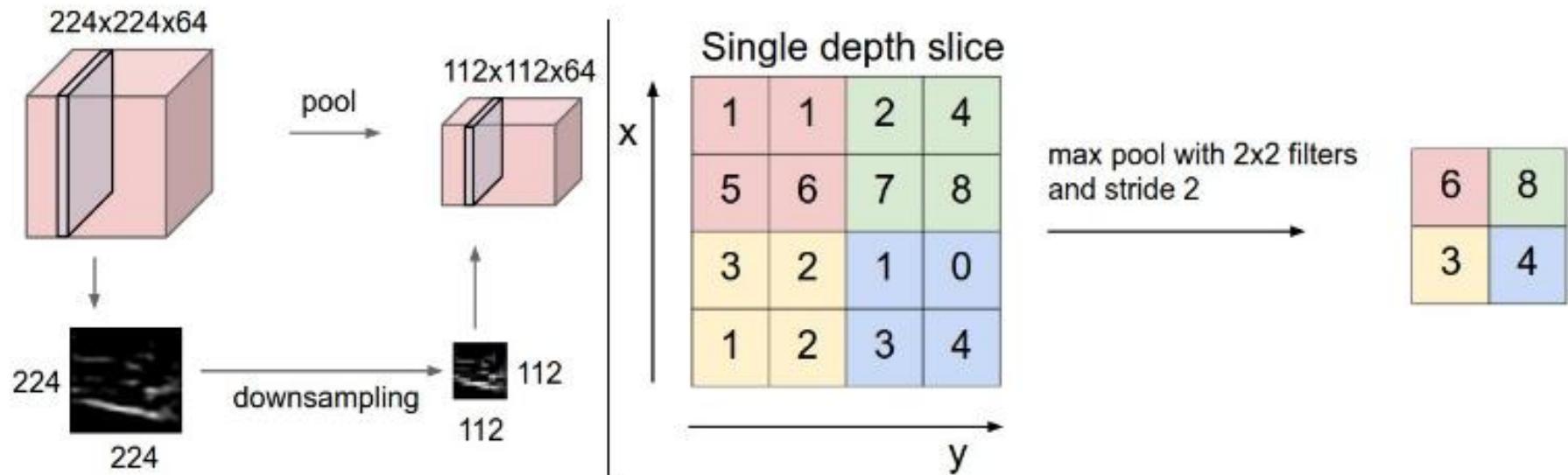


SMIL内部资料 请勿外泄

# Pooling Layer

## Max pooling:

- Max value of a cluster of neurons
- Nonlinear down-sampling



SMIL内部资料 请勿外泄

# Pooling Layer

Reduce **spatial size** of the representation and the **number of parameters**

## Forward

- Input data size:  $W_1 \times H_1 \times C_1$
- Output data size:  $W_2 \times H_2 \times C_2$

$$W_2 = [(W_1 - F)/S + 1]$$

$$H_2 = [(H_1 - F)/S + 1]$$

$$C_2 = C_1$$

where  $S$  is Spatial extent,  $F$  is Stride, and  $C$  is the number of channel

# **Forward propagation and backpropagation of Convolutional Neural Networks**

# Prediction for CNN $f$ : Forward Propagation

- Given input data  $\mathbf{X}$  and a convolutional neural network  $f$ , forward propagation is to compute the **prediction**  $\hat{\mathbf{y}}$  of  $\mathbf{X}$ , namely

$$\hat{\mathbf{y}} = f(\mathbf{X})$$

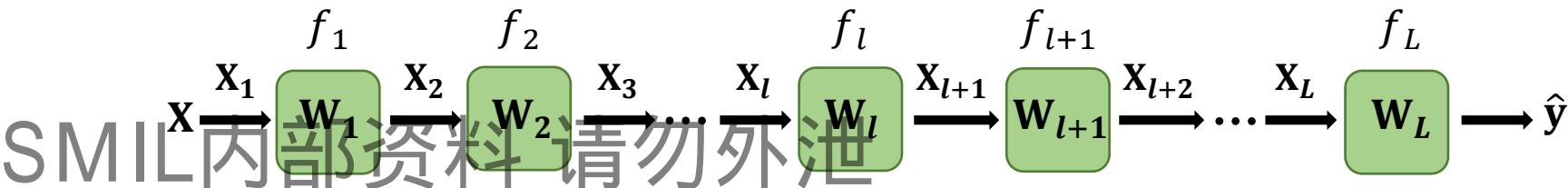
- Detailed prediction process:

$$\begin{aligned}\mathbf{X}_1 &= \mathbf{X} \\ \mathbf{X}_2 &= f_1(\mathbf{X}_1, \mathbf{W}_1) \\ &\vdots \\ \mathbf{X}_{l+1} &= f_l(\mathbf{X}_l, \mathbf{W}_l) \\ &\vdots \\ \mathbf{X}_L &= f_{L-1}(\mathbf{X}_{L-1}, \mathbf{W}_{L-1}) \\ \hat{\mathbf{y}} &= f_L(\mathbf{X}_L, \mathbf{W}_L)\end{aligned}$$

$$f_l = \sigma(h_l(\mathbf{X}_l, \mathbf{W}_l))$$

$\sigma$  is some activation function

$h_l = \mathbf{W}_l * \mathbf{X}_l$  is the **Convolutional Operation**—a linear operation



# How to Learn Parameters $\{\mathbf{W}_l\}$ in CNNs?

- We minimize the following regularized loss function:

Regularized Loss Function:

$$\min_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_W(\mathbf{y}_i, \hat{\mathbf{y}}_i) + \frac{1}{2} \lambda \|\mathbf{W}\|_F^2$$

$\mathbf{W}$  – parameters       $N$  – number of samples       $\mathcal{L}$  – loss function  
 $\mathbf{y}_i$  –  $i^{th}$  ground-truth label       $\hat{\mathbf{y}}_i$  –  $i^{th}$  predicted label

- To minimize the loss, we update the parameters  $\mathbf{W}_l$  by using **(Stochastic) Gradient Descent**:

$$\mathbf{W}'_l = \mathbf{W}_l - \eta \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{w}_l} \quad (2)$$

where  $\eta$  is the learning rate and gradient  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}_l}$  is computed

by **backpropagation**

# Contents

1 Introduction

2 Neural Networks

3 Forward Propagation

4 Backpropagation

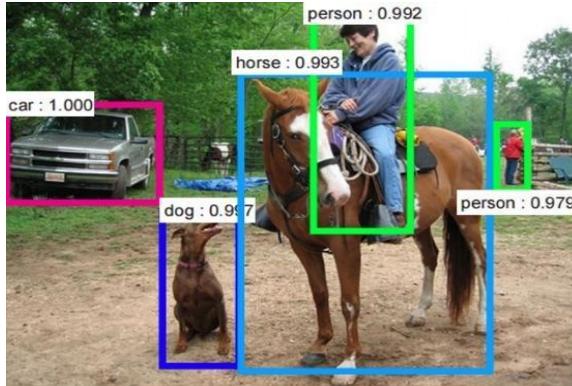
5 Convolutional Neural Networks

6 Applications

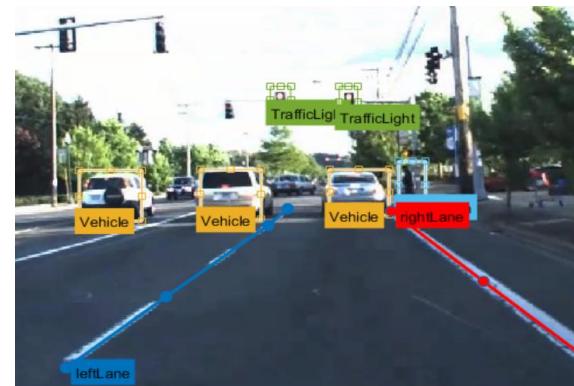
SMIL内部资料 请勿外泄

# Deep learning Applications

## ■ Image Classification; Object Detection; Image Segmentation

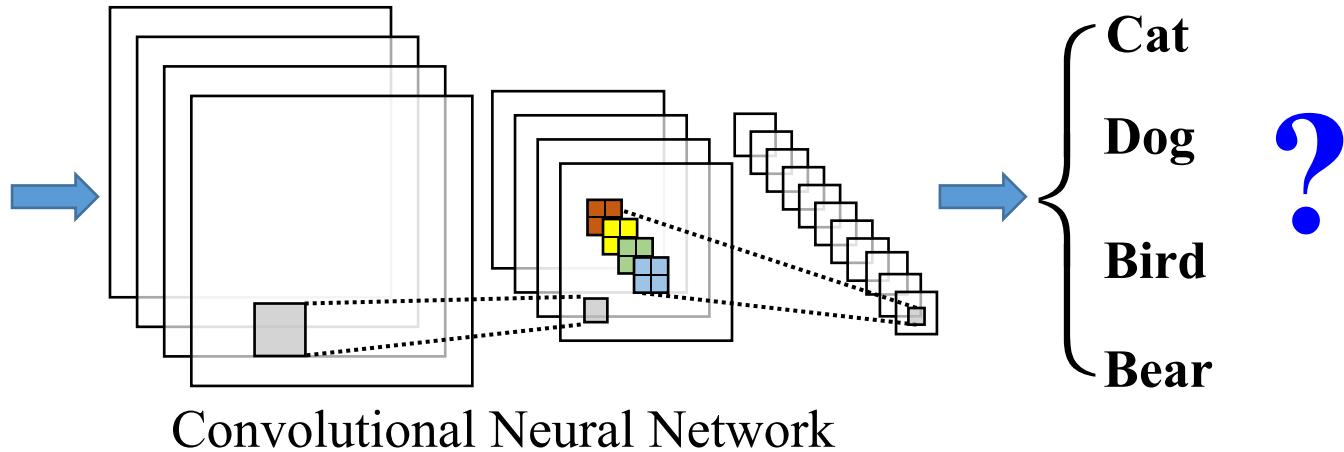


## ■ Pedestrian Detection; Lane Detection; Traffic Sign Recognition



SMIL内部资料 请勿外泄

# Image Recognition



Image

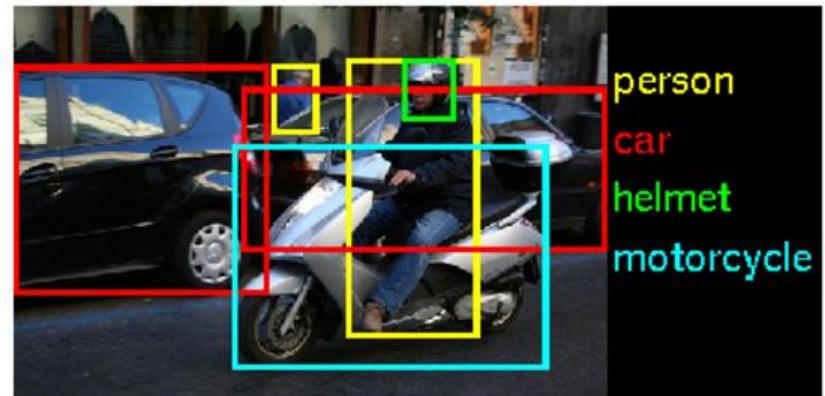
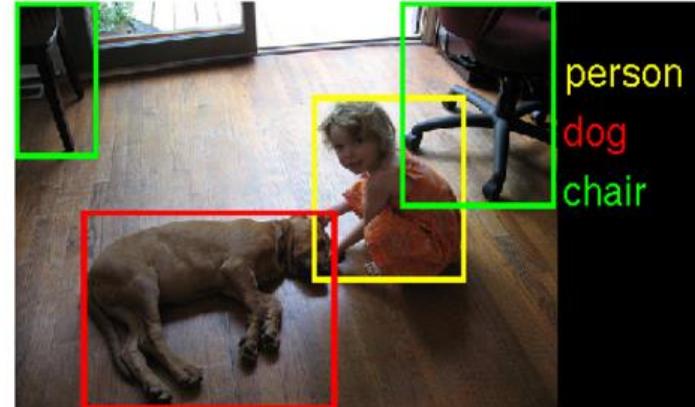
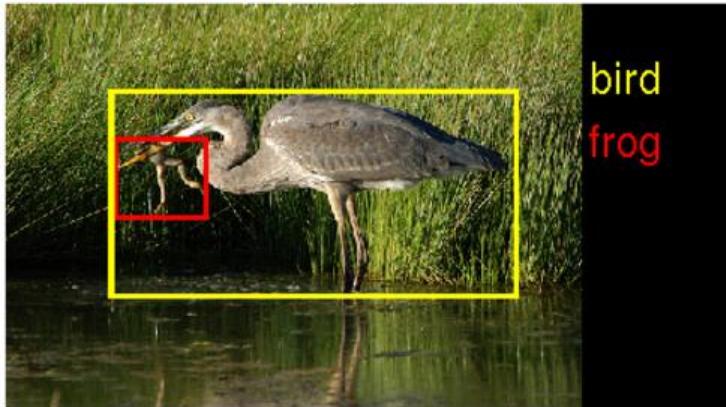


Prediction

mite	container ship	motor scooter	leopard
mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

SMIL内部资料 请勿外泄

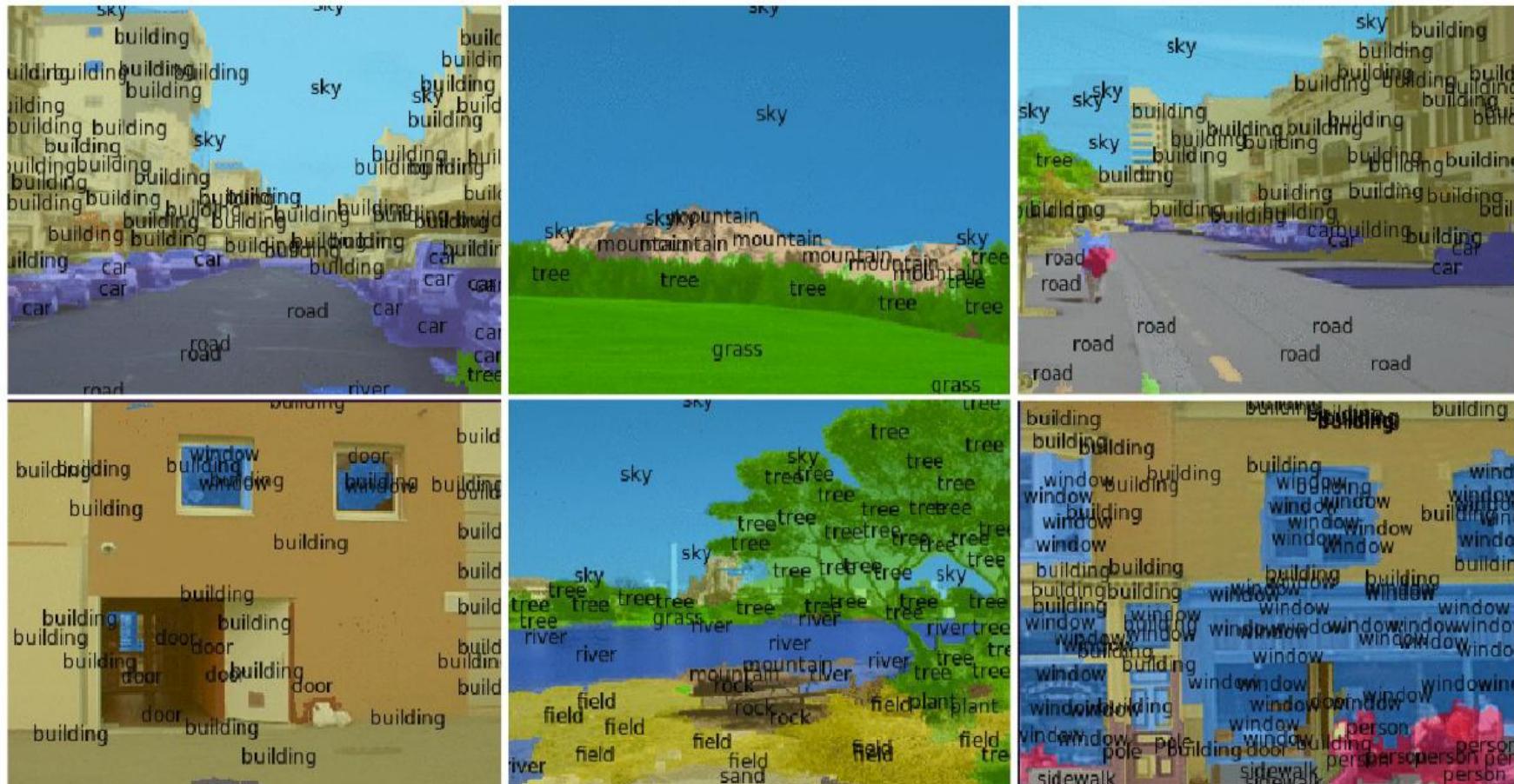
# Object Detection



SMIL内部资料 请勿外泄

# Semantic Labeling

- Labeling every pixel with the object category it belongs to



SMIL内部资料 请勿外泄

# Scene Parsing

- Each output sees a large input context



SMIL内部资料 请勿外泄

# Other Recognition Applications

## ■ Traffic Sign Recognition (GTSRB)

- ▶ German Traffic Sign Reco Bench
- ▶ 99.2% accuracy



## ■ House Number Recognition (Google)

- ▶ Street View House Numbers
- ▶ 94.3 % accuracy

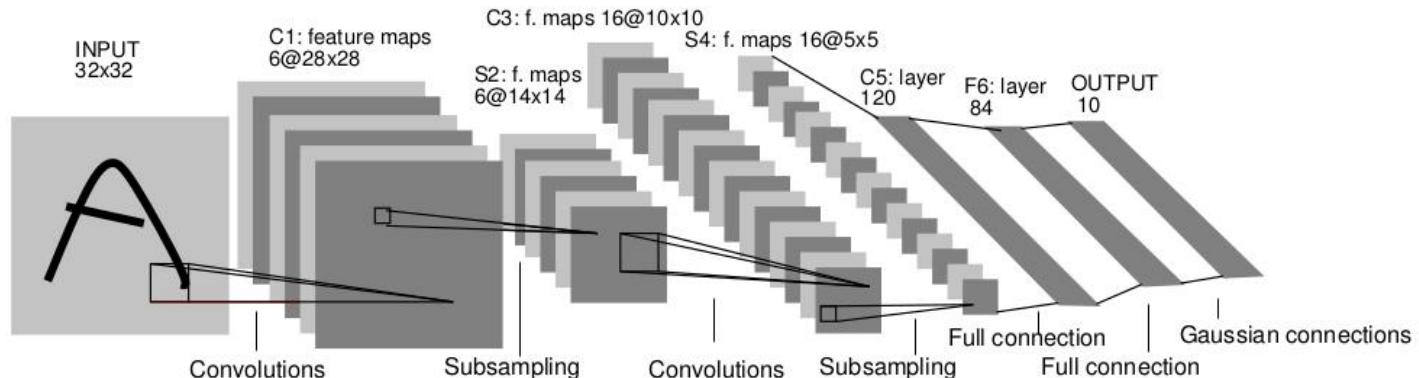


SMIL内部资料 请勿外泄

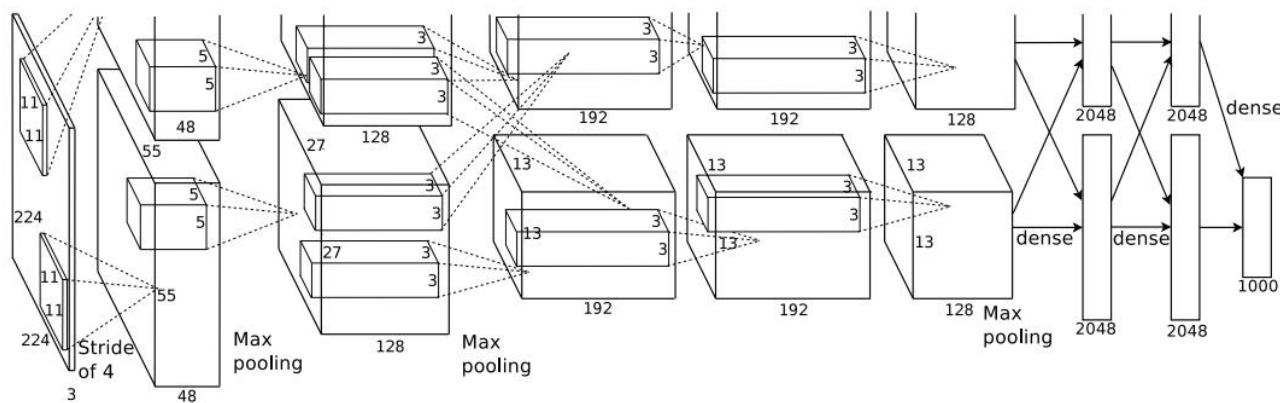
# Classical CNNs

# Classical CNNs

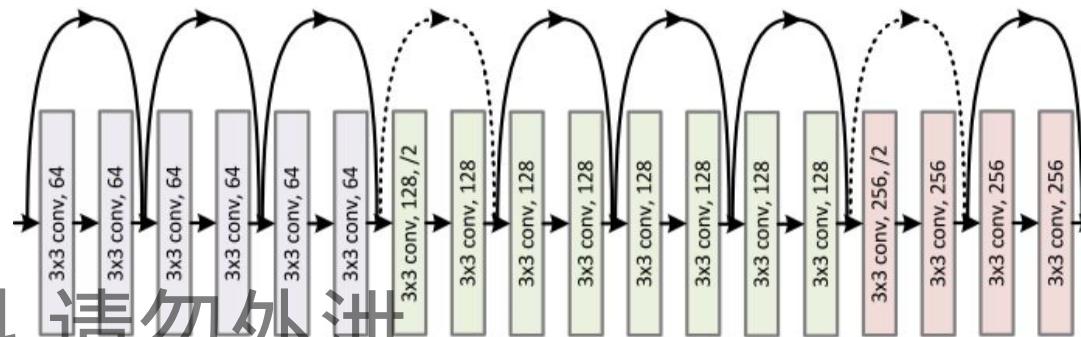
LeNet



AlexNet

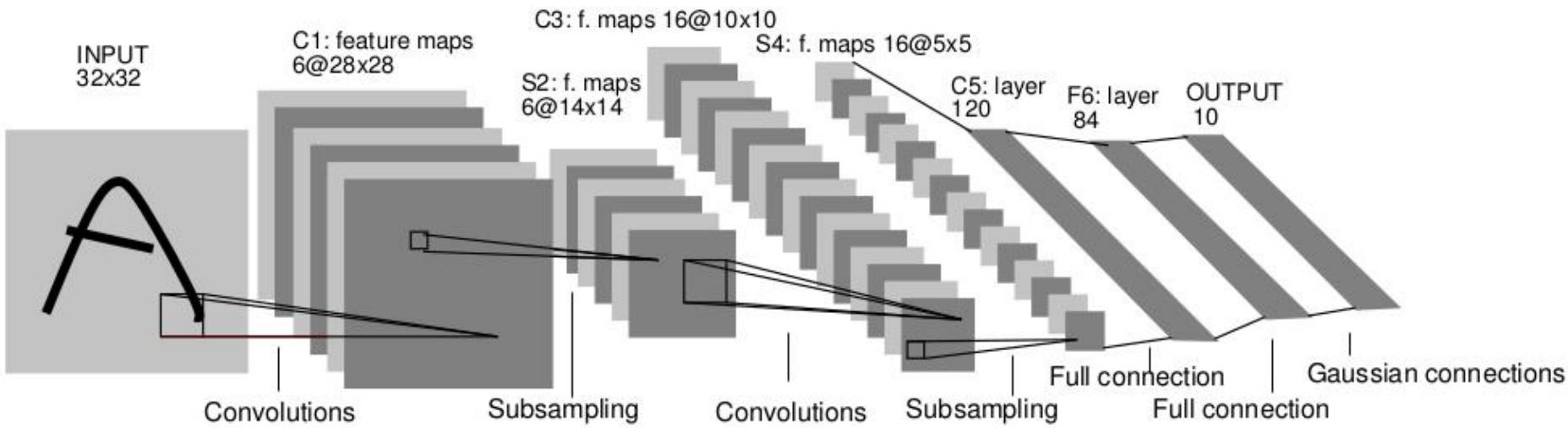


ResNet

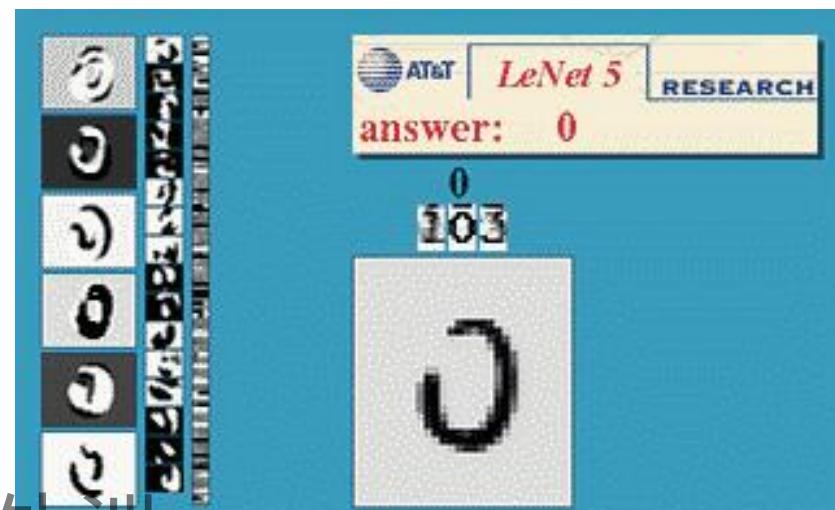


SMIL 内部资料 请勿外泄

# LeNet



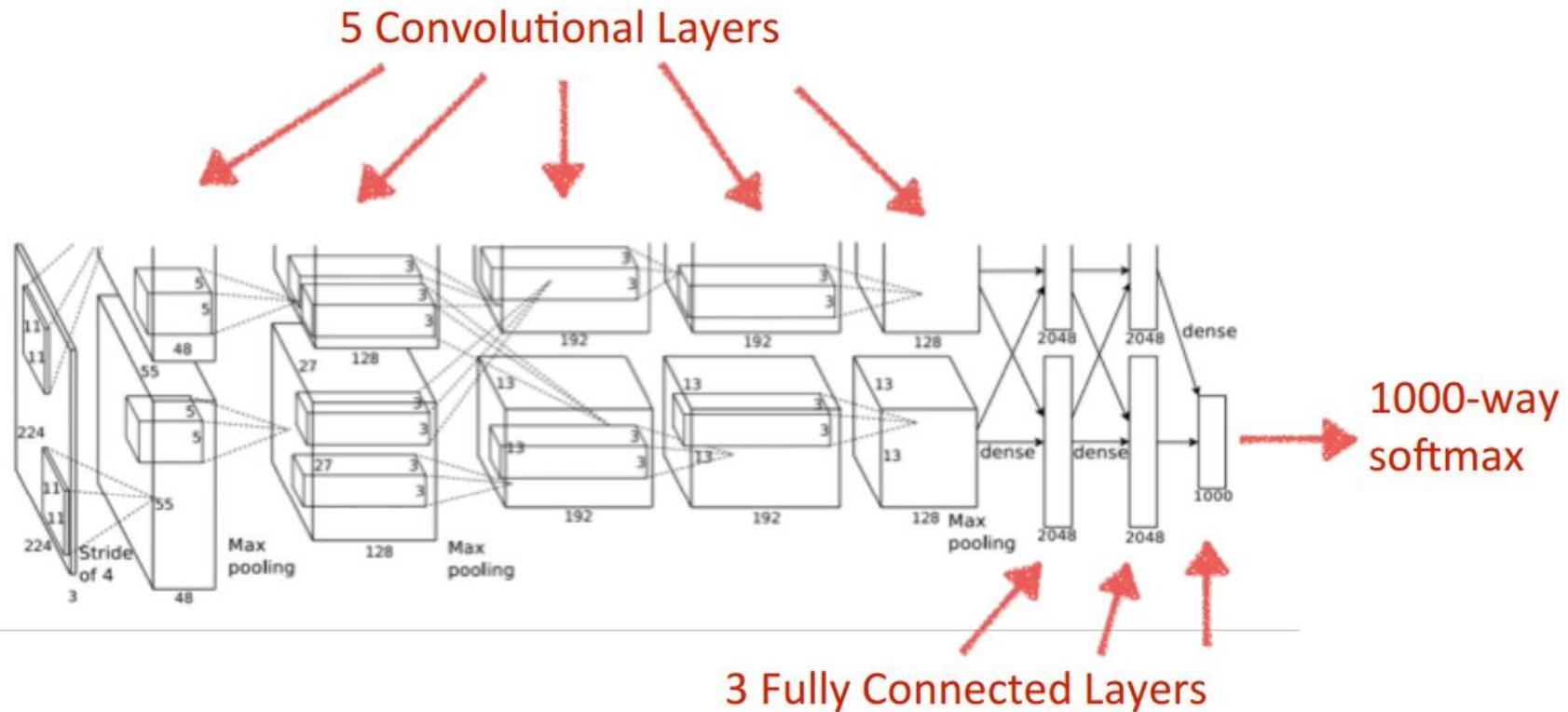
## ■ Handwriting number recognition



SMIL内部资料 请勿外泄

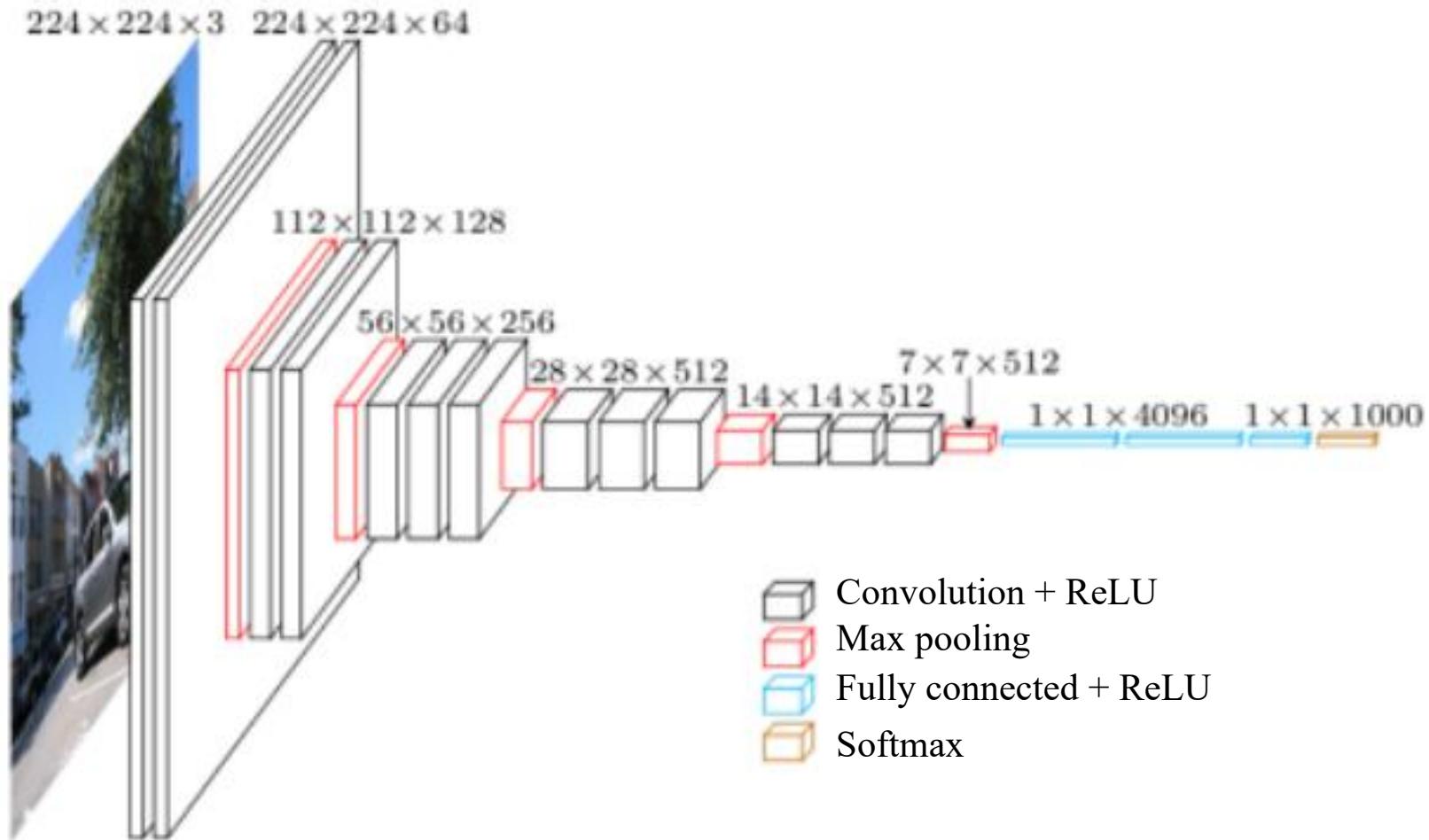
# AlexNet

- Trained the network on ImageNet.
- Used ReLU for the nonlinearity functions
- Implemented dropout layers to combat overfitting issue



SMIL内部资料 请勿外泄

# VGG



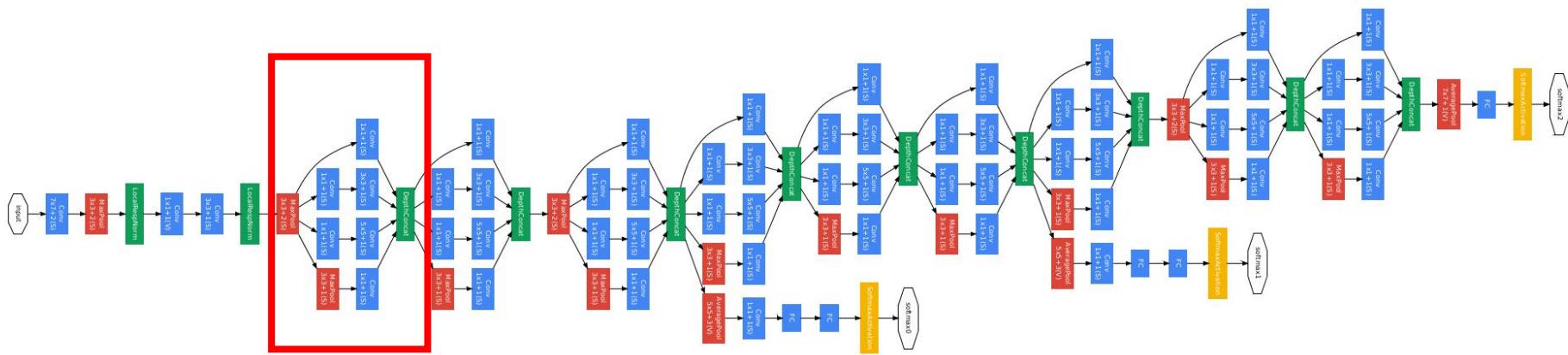
SMIL内部资料 请勿外泄

# Six Different Architectures of VGG Net

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 <b>conv3-256</b> <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

SMIL内部资料 请勿外泄

# GoogLeNet

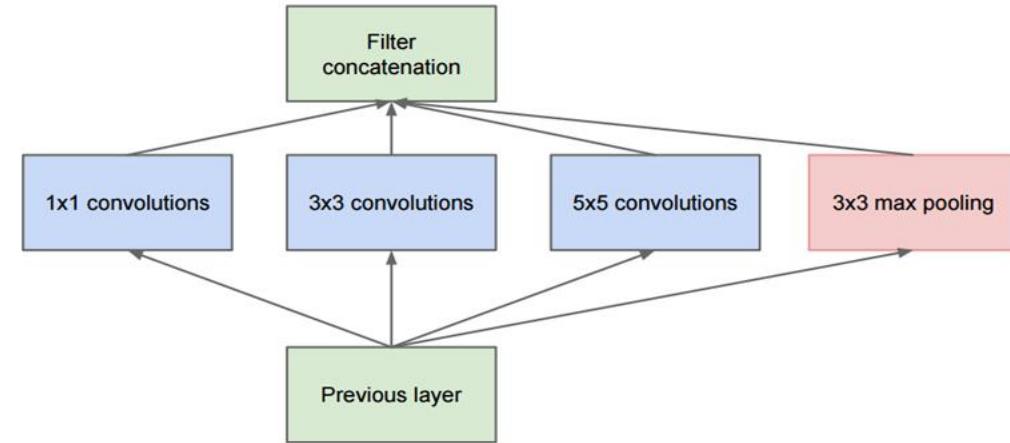


- Used 9 Inception modules, over 100 layers (very deep)
- Uses 12x fewer parameters than AlexNet
- Trained on multiple GPUs within a week

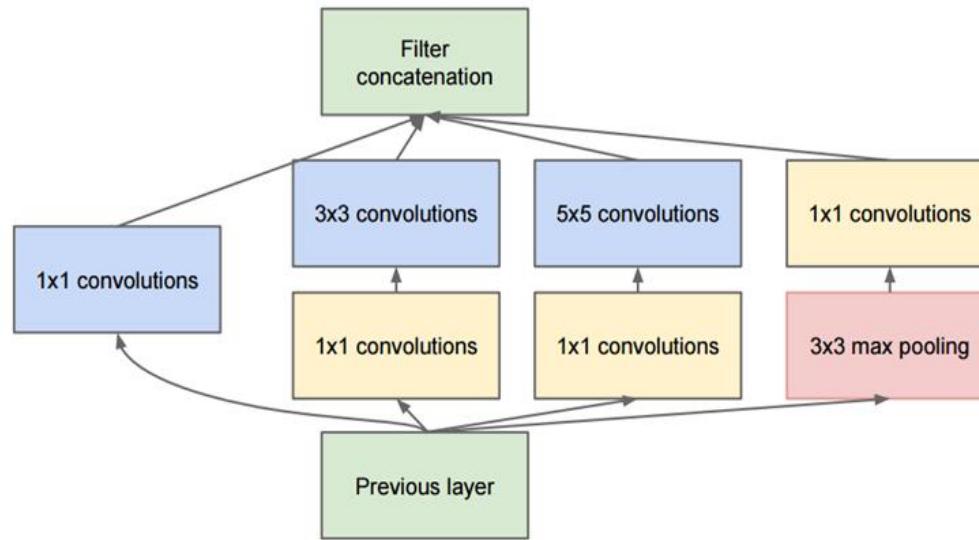
SMIL内部资料 请勿外泄

# Inception Module

## ■ Basic Inception Module

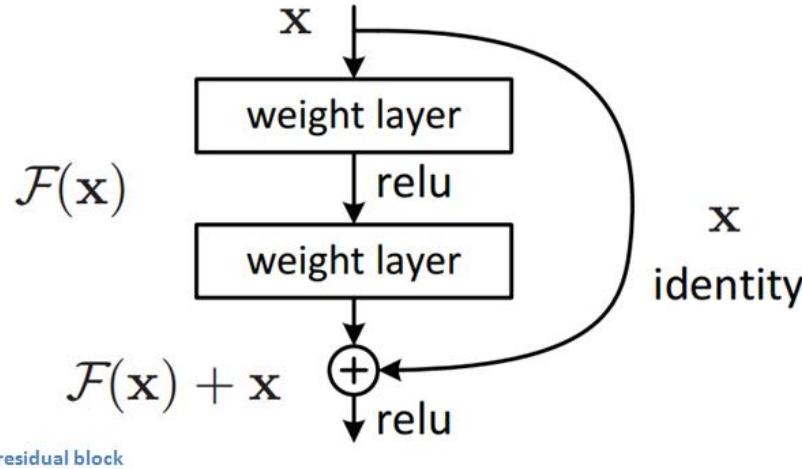


## ■ Full Inception Module



SMIL内部资料 请勿外泄

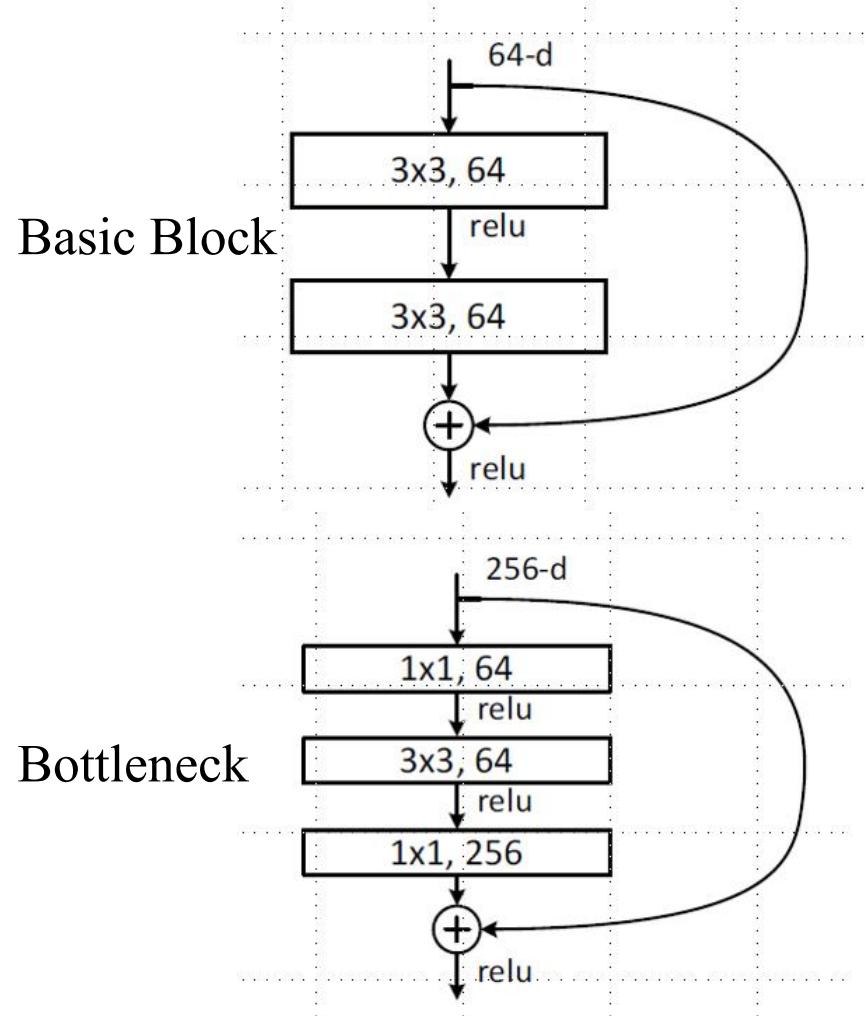
# Residual Network



A residual block

$$y = \boxed{\mathcal{F}(x)} + x$$

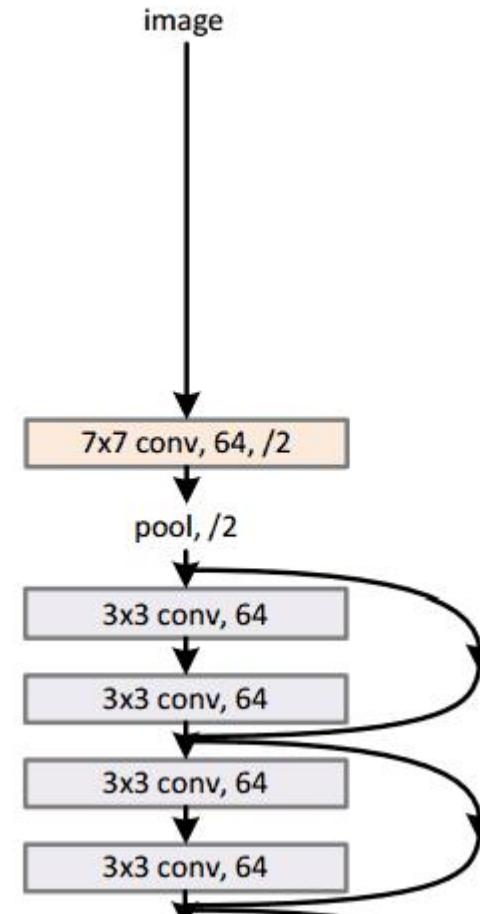
Residual Term



# Residual Network

- Residual learning
- Ultra-deep 152 and 200 layers
- Trained a **1202-layer** network
- Trained on an **8 GPUs** machine

34-layer residual



SMIL内部资料 请勿外泄

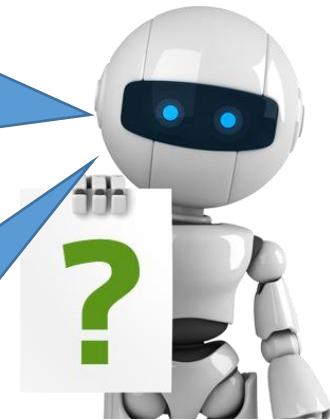
# **What does machine learn in CNNs?**

# What does Machine Learn?



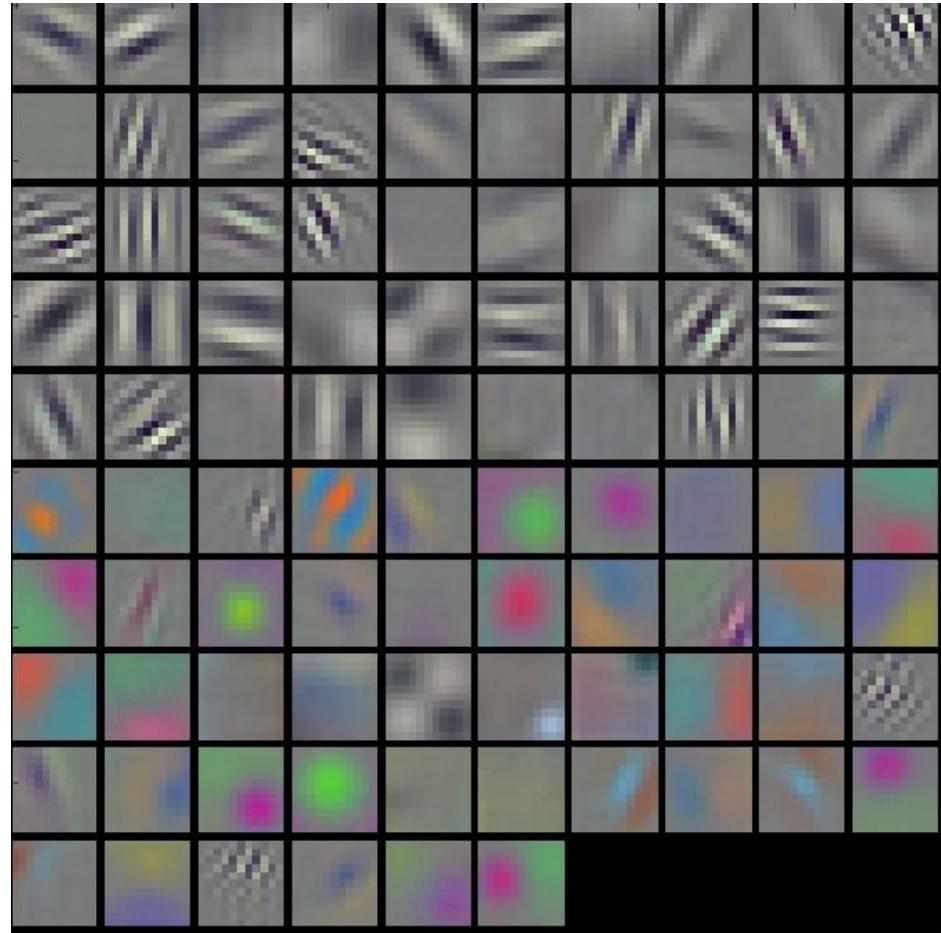
shoes

puma



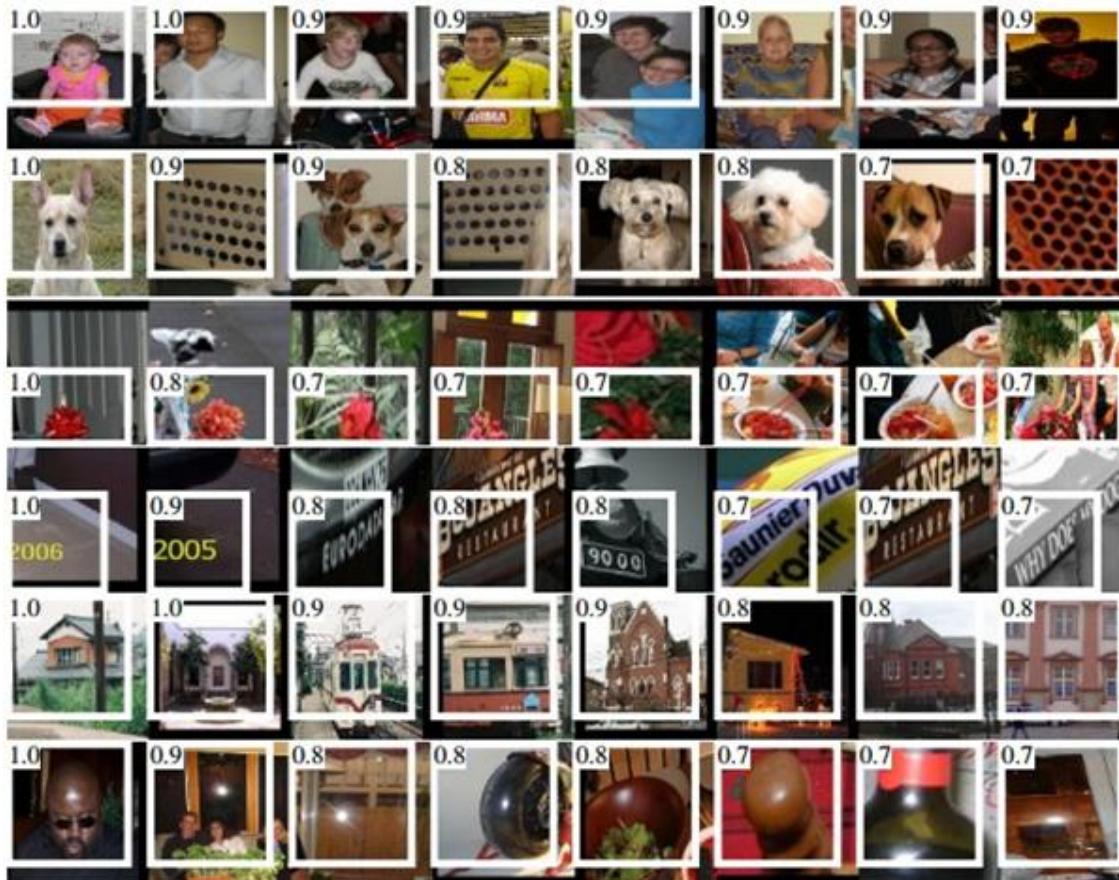
# First Convolutional Layer

- Typical-looking filters on the trained first layer



# How about Higher Layers?

- Which images make a specific neuron activate?

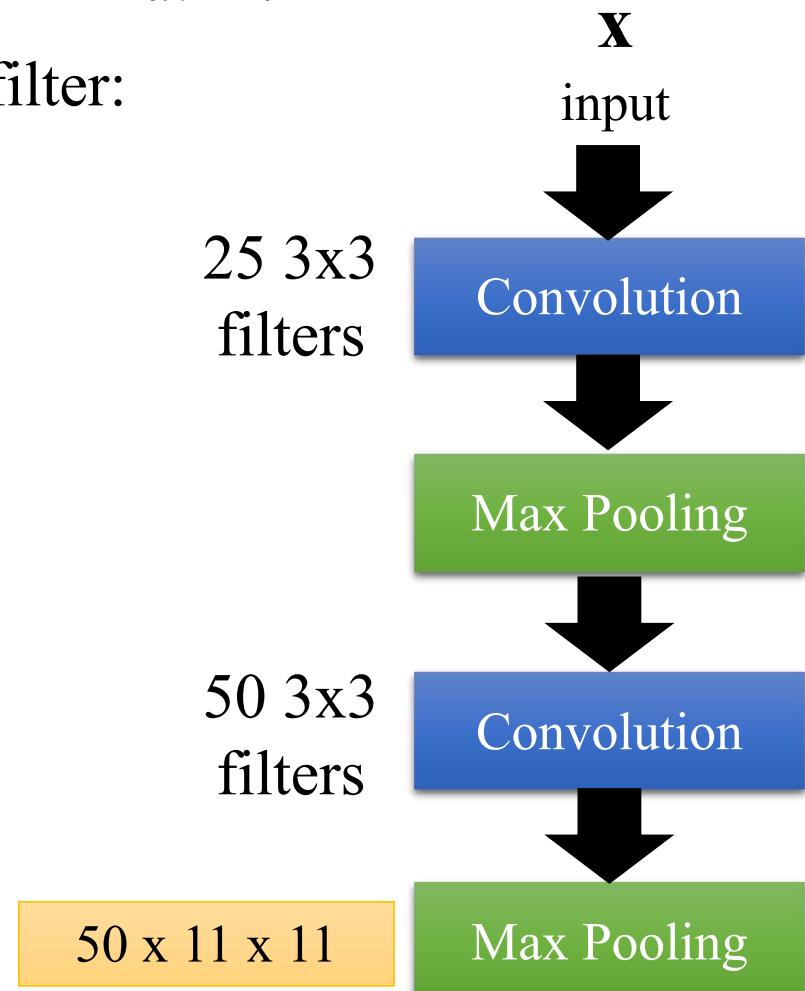
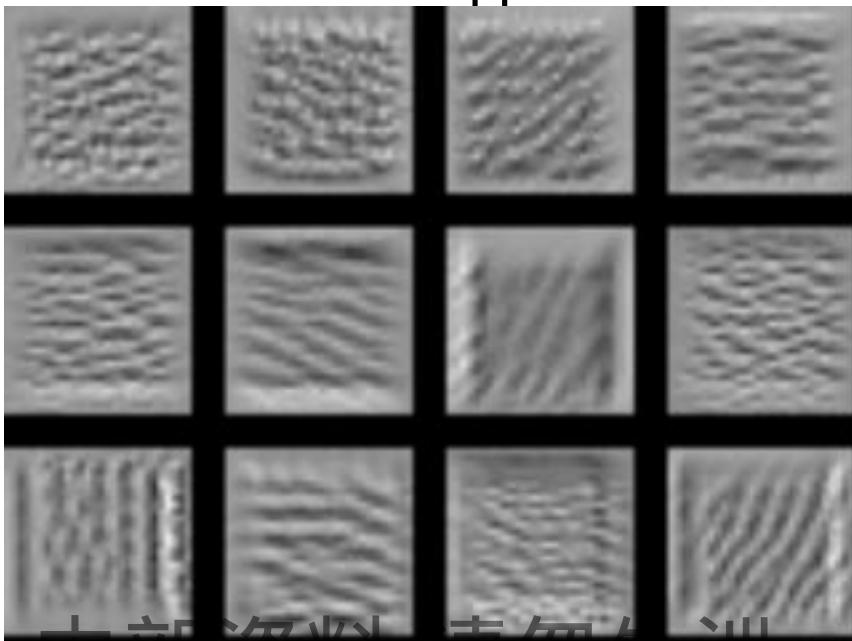


Ross Girshick, Jeff  
Donahue, Trevor  
Darrell, Jitendra Malik, “Rich  
feature hierarchies for accurate  
object detection and semantic  
segmentation”, CVPR, 2014

# What does Machine Learn?

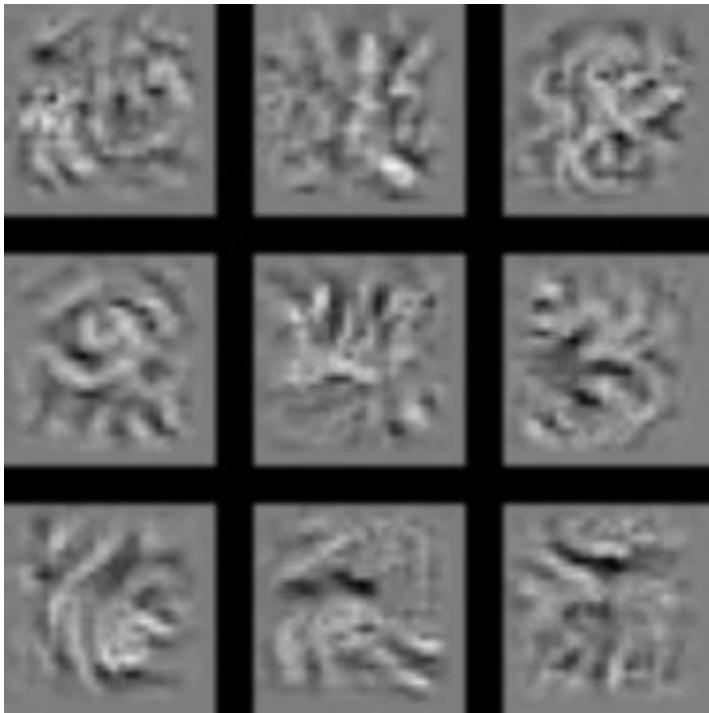
- The output of the k-th filter is a  $11 \times 11$  matrix.
- Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$



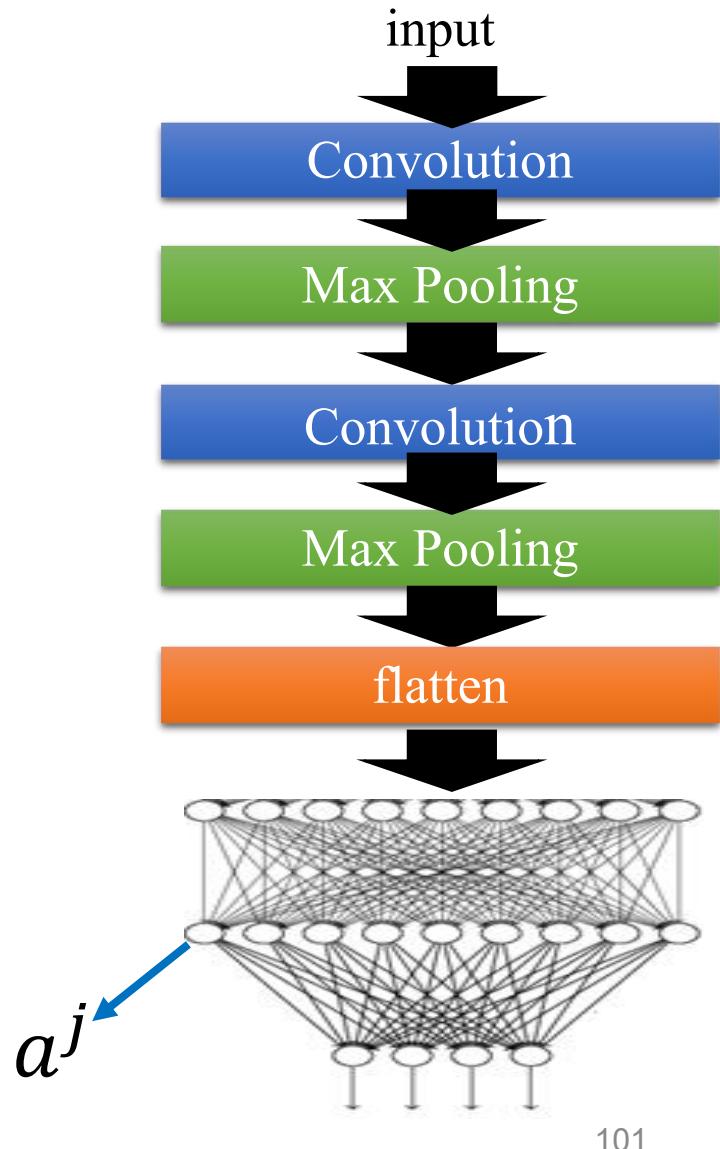
# What does Machine Learn?

- Find an image maximizing the output of neuron:  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} a^j$



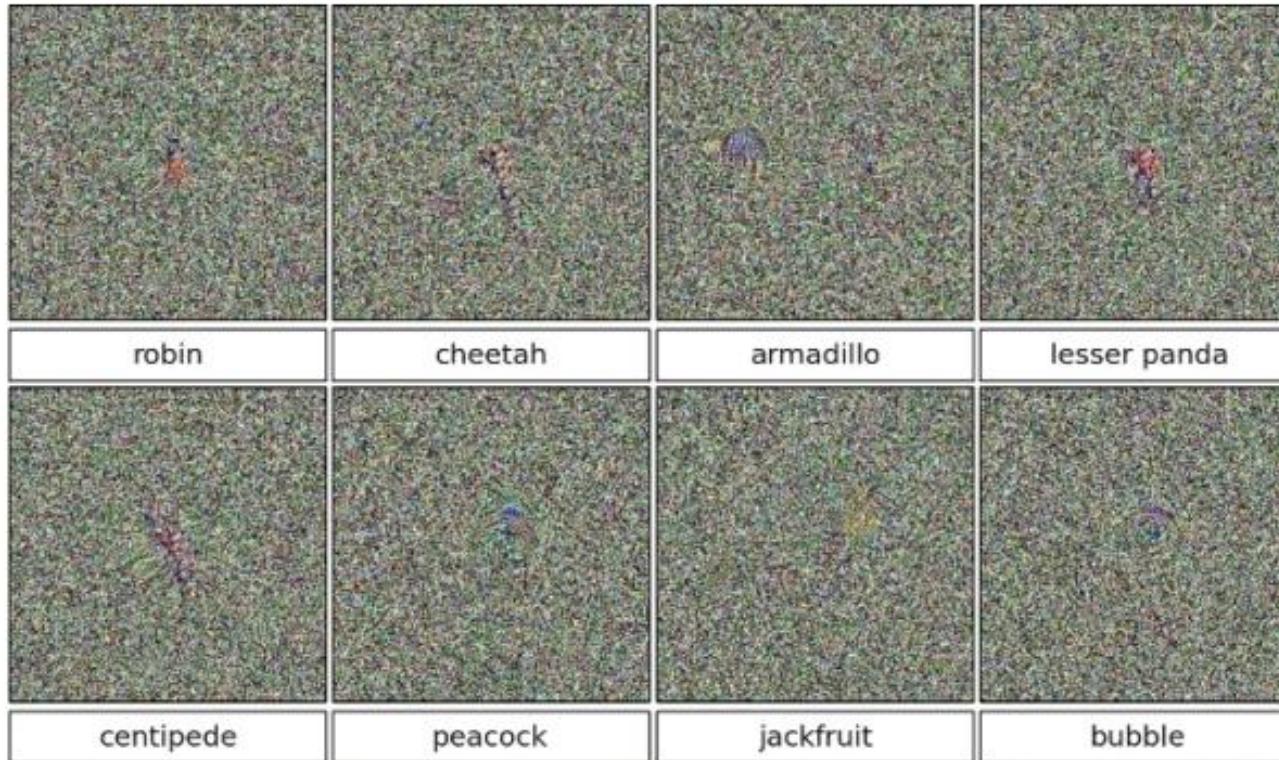
- Each figure corresponds to a neuron

SMIL 内部资料 请勿外泄



# Will Deep Neural Networks be Fooled?

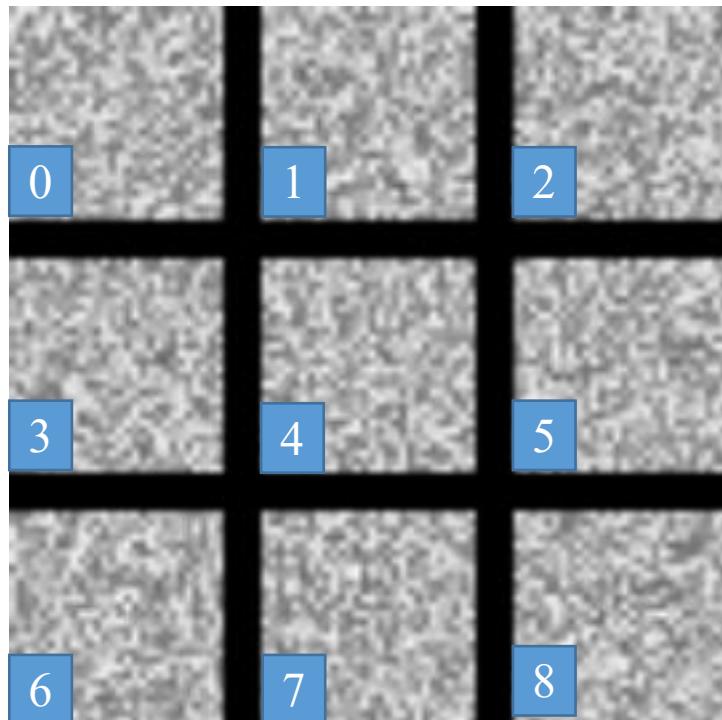
- Evolved images that are unrecognizable to humans but DNNs can recognize



(Deep Neural Networks are Easily Fooled  
<https://www.youtube.com/watch?v=M2IebCN9Ht4>)

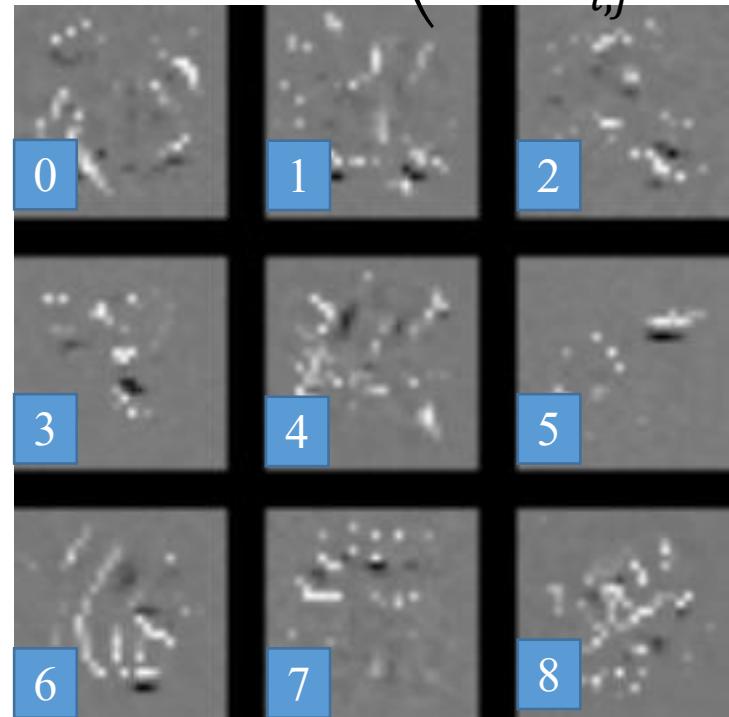
# What does Machine Learn?

$$\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} y^i$$



Over all pixel values

$$\mathbf{x}^* = \operatorname{argmax}_X \left( y^i - \sum_{l,J} |\mathbf{x}_{lj}| \right)$$



SMIL内部资料 请勿外泄

# What does Machine Learn?

- Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, “Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps”, ICLR, 2014



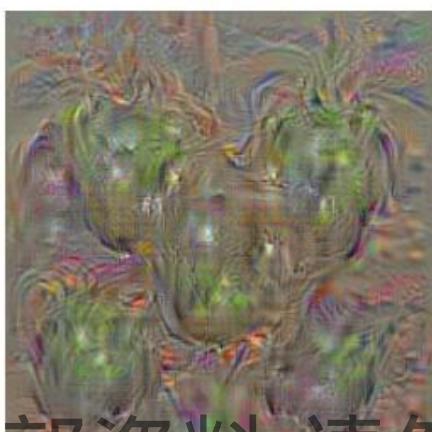
dumbbell



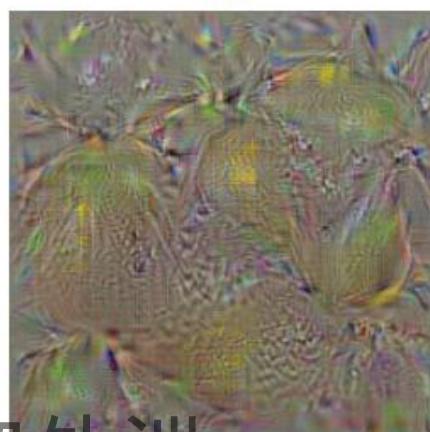
cup



dalmatian



bell pepper



lemon



husky

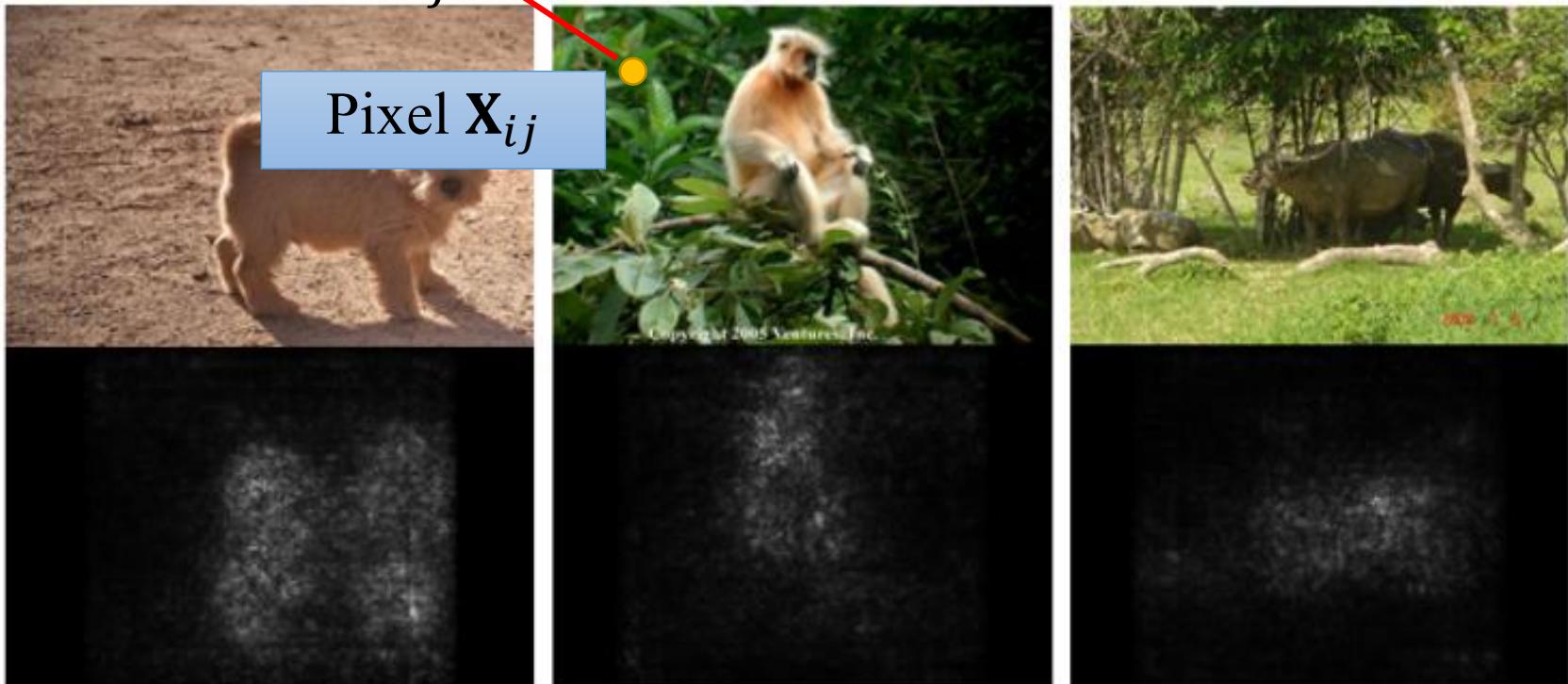
SMIL内部资料 请勿外泄

# What does Machine Learn?

- Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, “Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps”, ICLR, 2014

$$\left| \frac{\partial y_k}{\partial X_{ij}} \right|$$

$y_k$ : the predicted class  
of the model



# What does Machine Learn?

- Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014* (pp. 818-833)



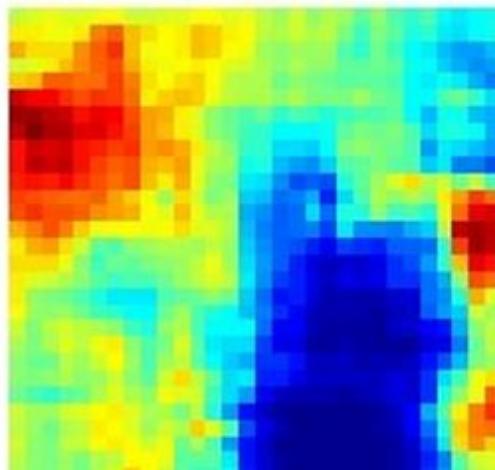
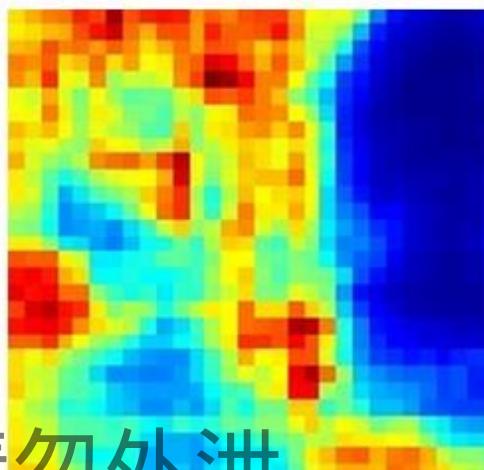
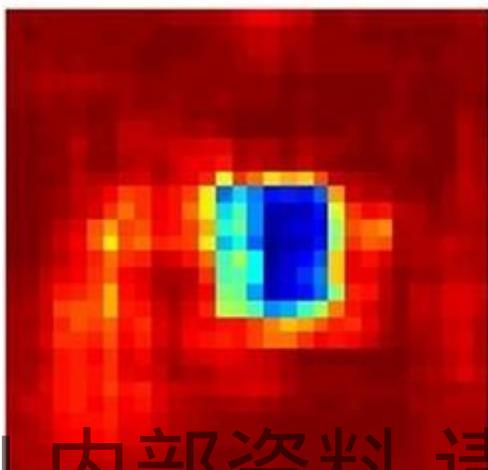
True Label: Pomeranian



True Label: Car Wheel



True Label: Afghan Hound



Thank you