



# 深度学习平台与应用

## 第五讲：卷积神经网络

范琦

[fanqi@nju.edu.cn](mailto:fanqi@nju.edu.cn)

2024年10月9日

# 大 纲

- 卷积神经网络-引言
- 卷积神经网络

## ■ 图像分类



假设有一个标签集合  
{人, 狗, 猫, 汽车, ...}

模型的分类预测

猫

# 课程回顾：线性分类器



输入图像

$$f(x, W) = Wx + b$$

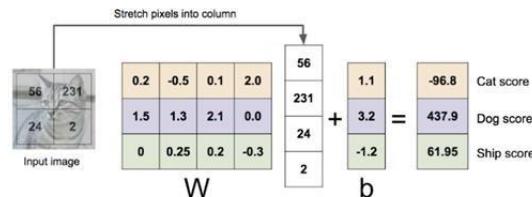
$$\longrightarrow f(x, W) \longrightarrow$$

571.3
1.25
-132.2

猫  
狗  
船

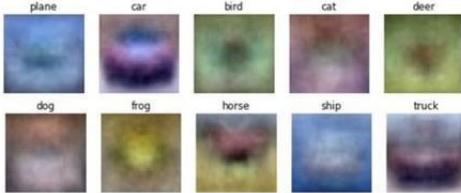
## Algebraic Viewpoint

$$f(x, W) = Wx$$



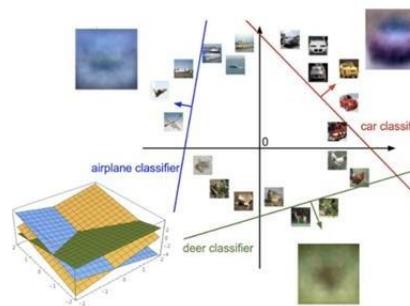
## Visual Viewpoint

One template per class



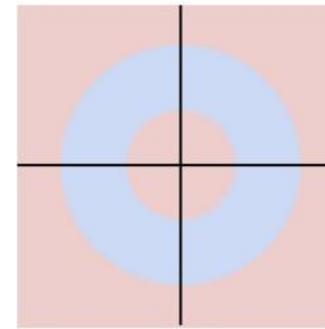
## Geometric Viewpoint

Hyperplanes cutting up space



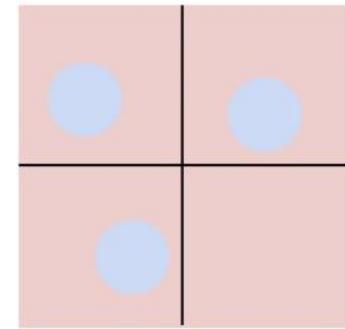
**Class 1:**  
 $1 \leq L_2 \text{ norm} \leq 2$

**Class 2:**  
Everything else



**Class 1:**  
Three modes

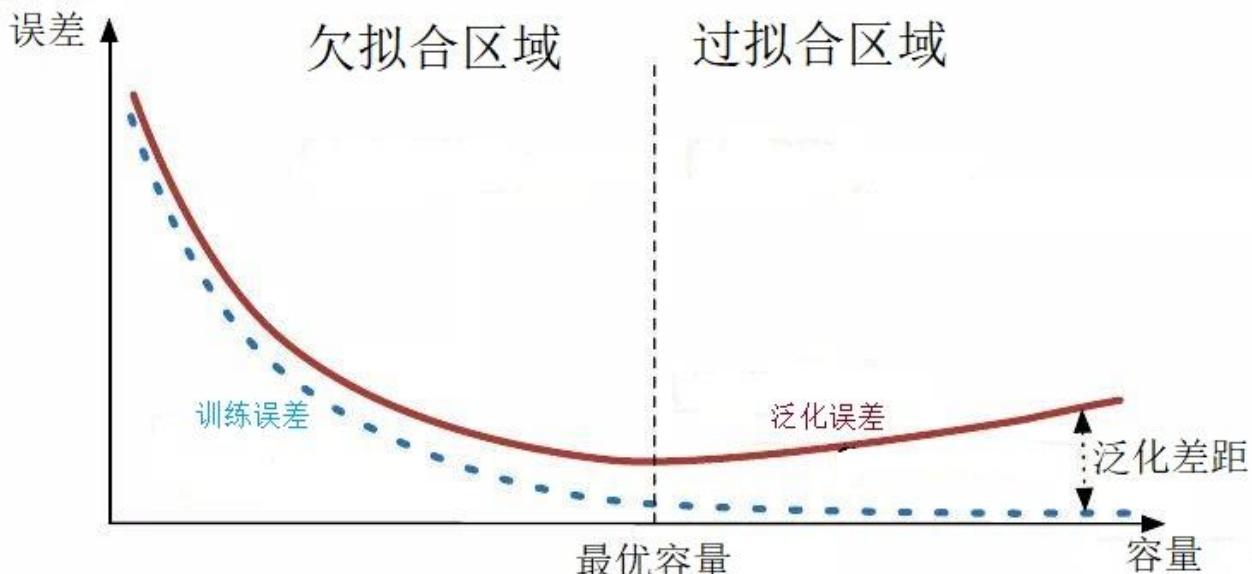
**Class 2:**  
Everything else



# 课程回顾：损失函数与优化

- 我们有训练数据:
- 我们有线性分类预测:
- 我们有损失函数:
- 优化模型来降低模型在数据上的损失: **SGD, SGD+Momentum, RMSProp, Adam**

$$\{(x_i, y_i)\}_{i=1}^N$$
$$s = f(x; W) = Wx$$
$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W)$$

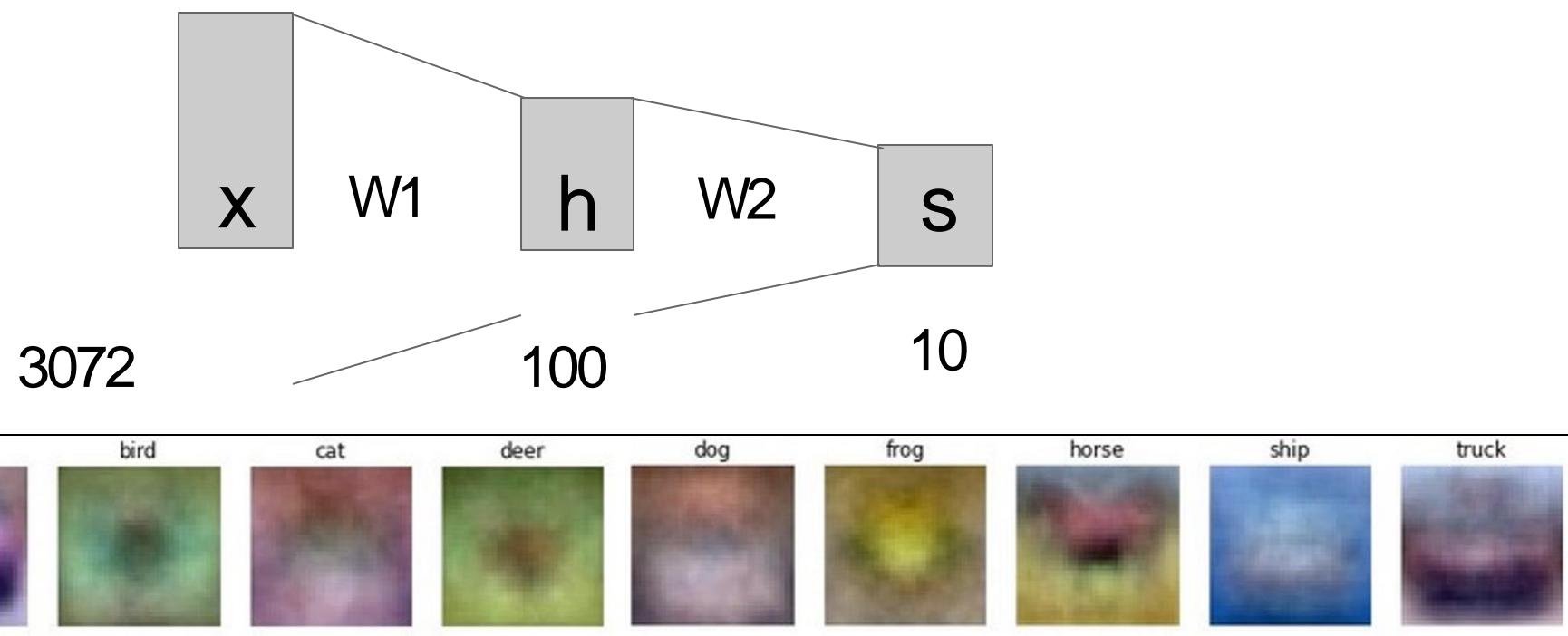


- 线性分类器

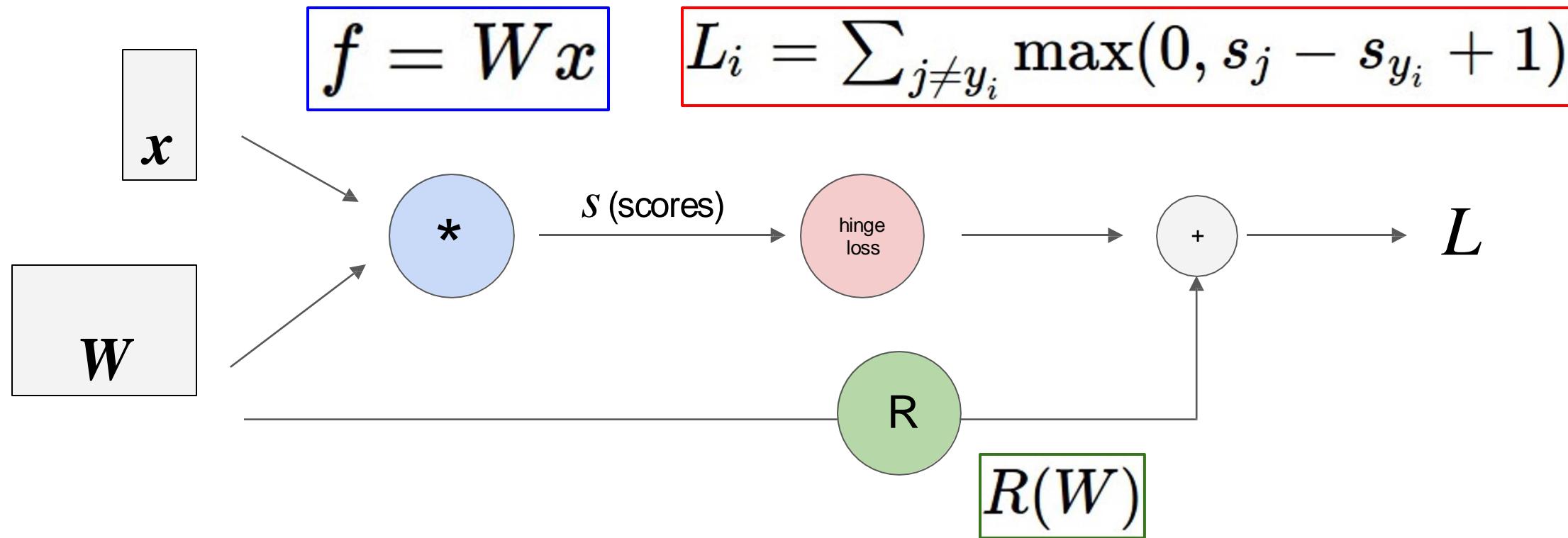
$$f = Wx$$

- 2 层神经网络

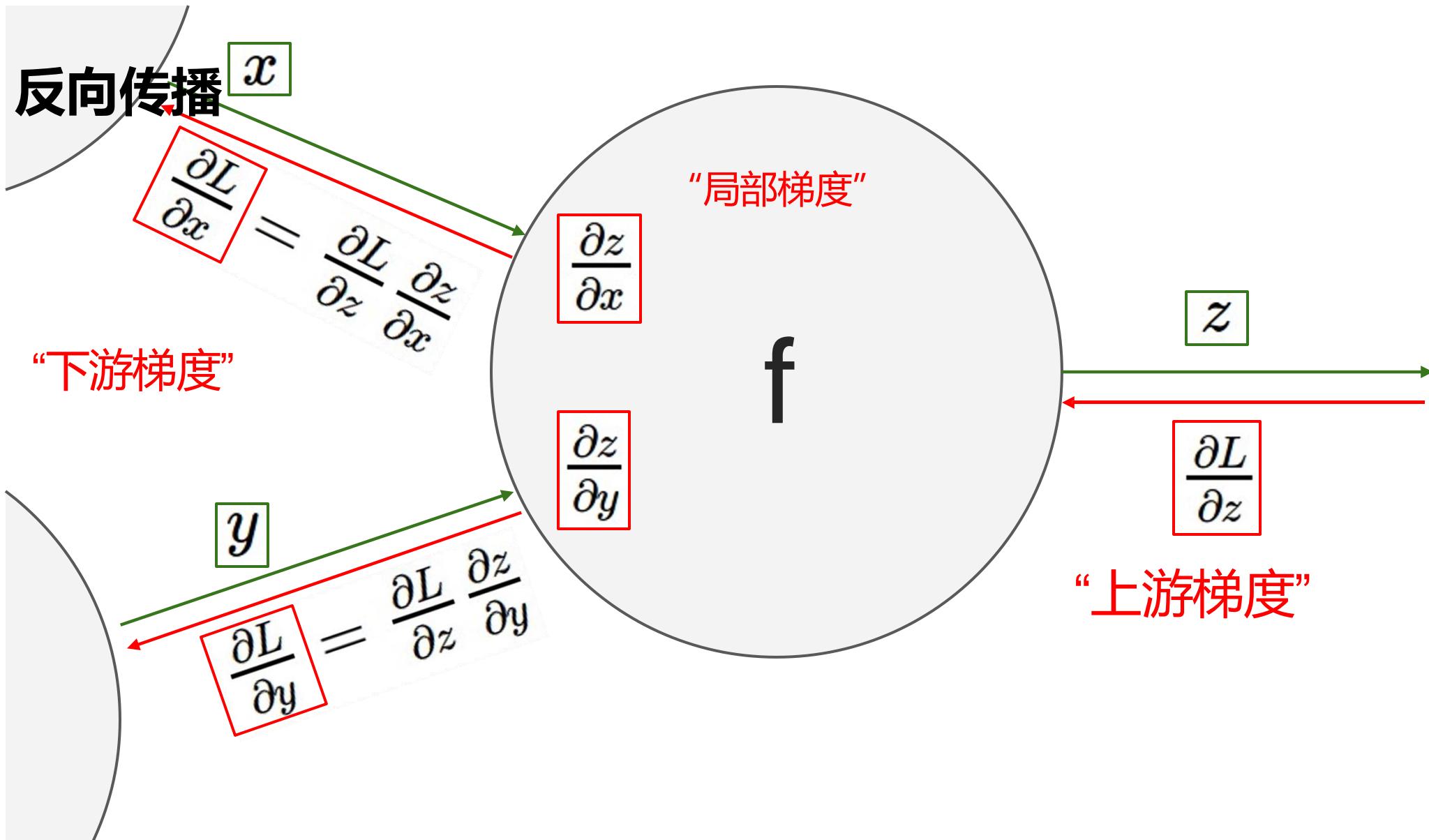
$$f = W_2 \max(0, W_1 x)$$



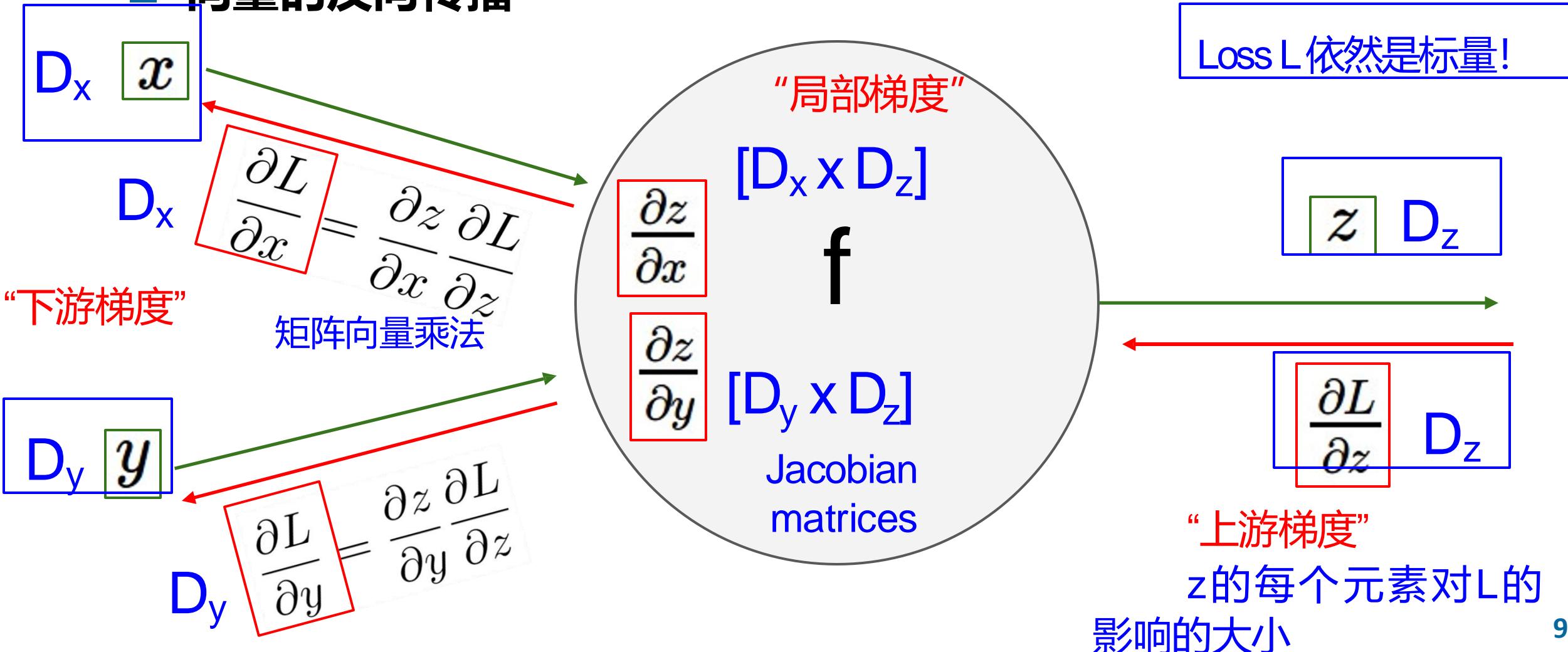
## ■ 计算图+反向传播



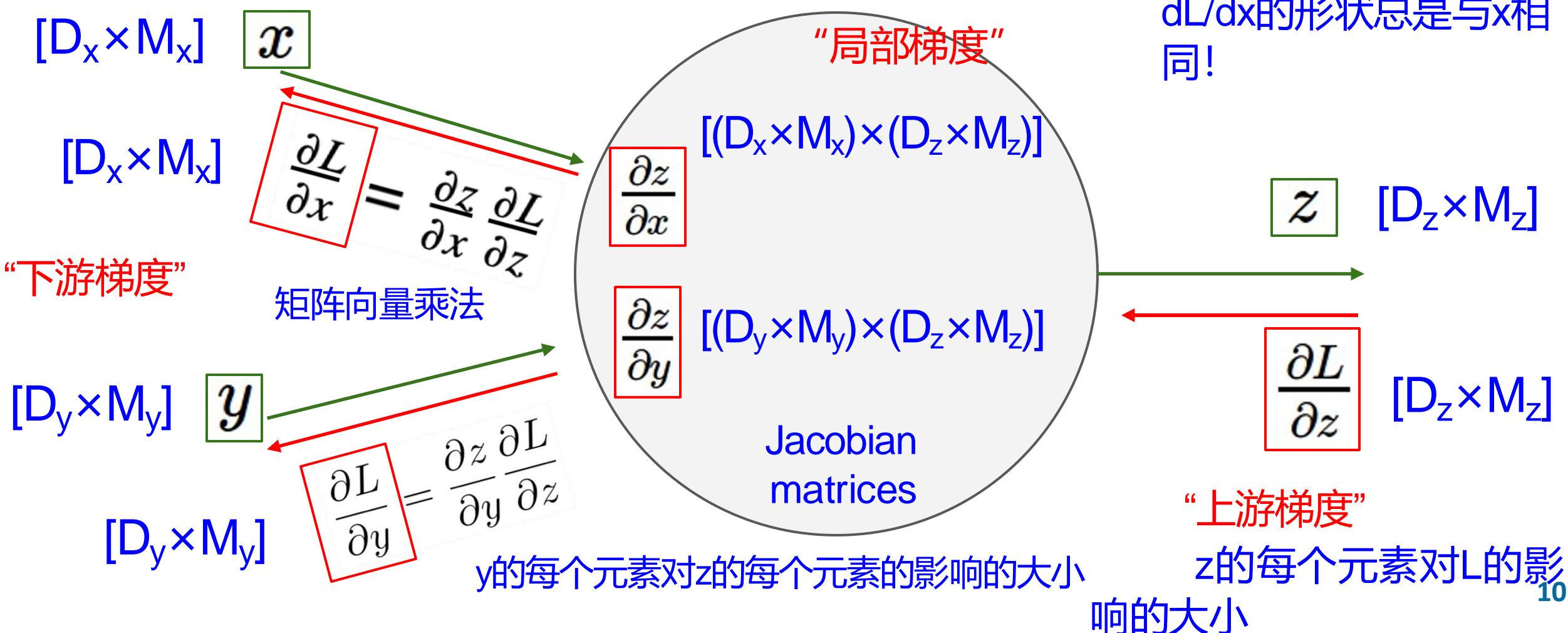
- 



## ■ 向量的反向传播



## ■ 矩阵的反向传播



## ■ 图像分类是最核心的计算机视觉任务

假设有一个标签集合  
{人, 狗, 猫, 汽车, ...}



模型的分类预测

Class  
Scores

## ■ 基于图像像素的图像分类

假设有一个标签集合  
{飞机, 狗, 猫, 汽车, ...}



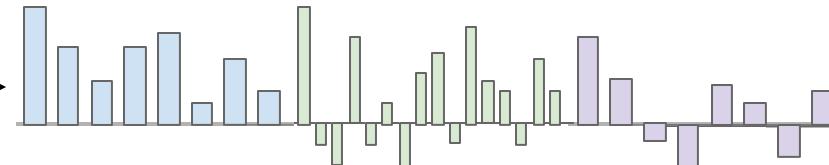
$$f(x) = Wx$$

Class  
Scores



## ■ 基于图像特征的图像分类

假设有一个标签集合  
{飞机, 狗, 猫, 汽车, ...}



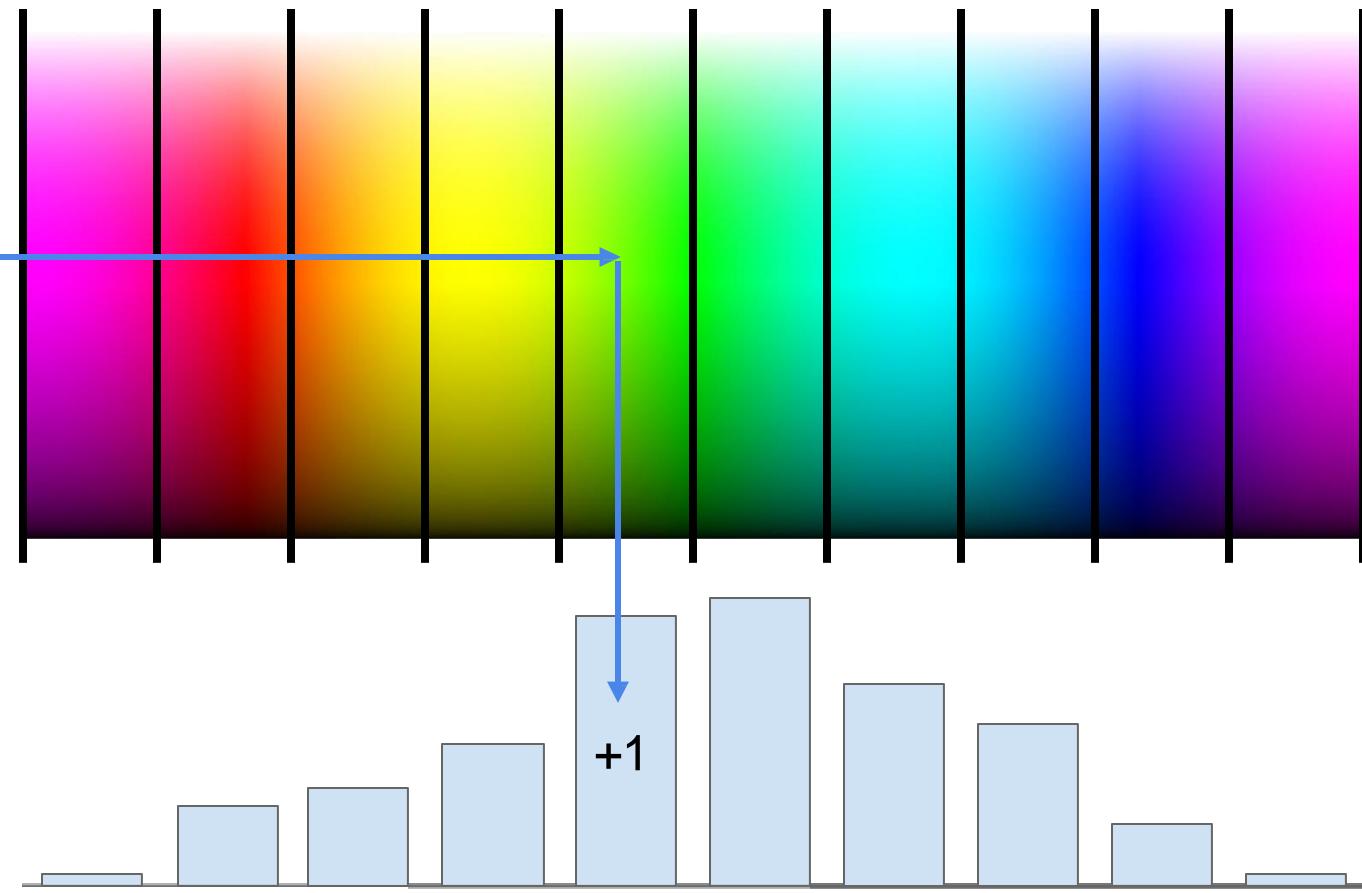
特征表示

$$f(x) = Wx$$



Class  
Scores

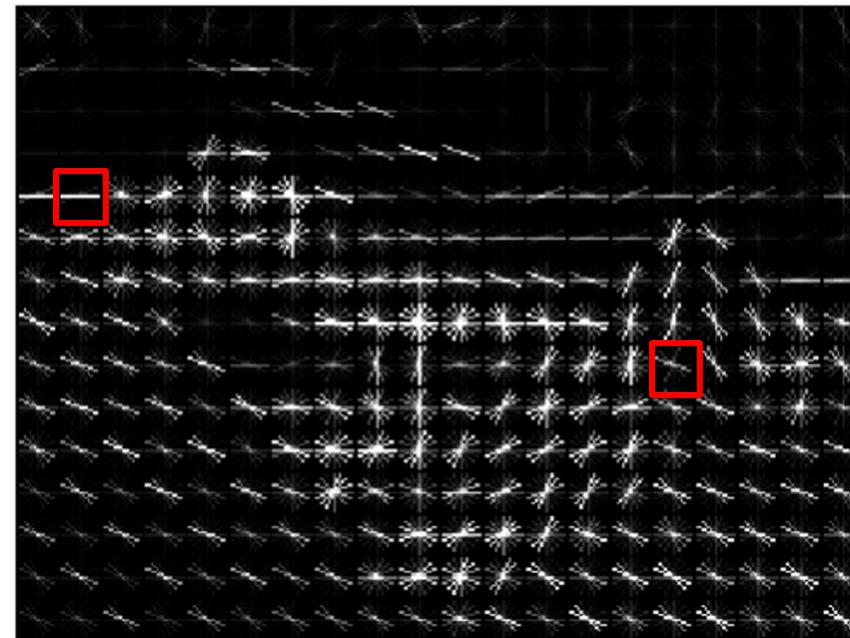
## ■ 特征提取：颜色直方图（Color Histogram）



## ■ 方向梯度直方图 (Histogram of Oriented Gradients, HoG)



**HoG:** 将图像以 $8 \times 8$ 为基础区域进行划分。在每个基础区域内，将边缘方向量化为9个数字



**例子：**320x240的图像被划分为 $40 \times 30$ 的区域；每个区域中有9个数字，因此特征向量有 $30 \times 40 \times 9 = 10800$ 个数字

## ■ 词袋 (Bag of Words)

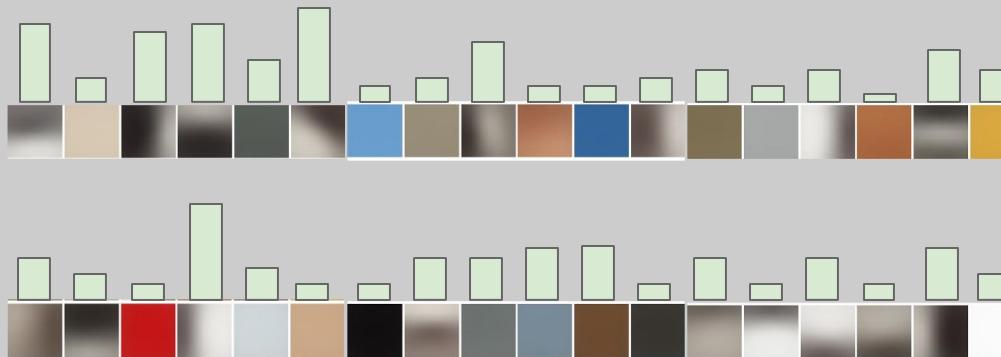
### Step 1: 构建 codebook



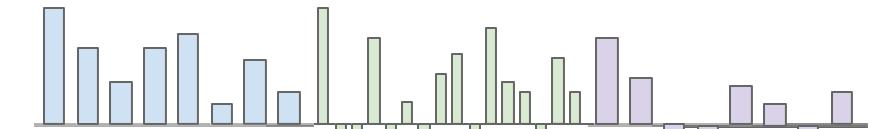
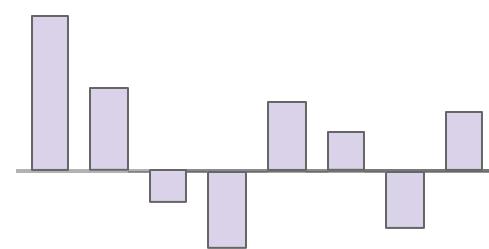
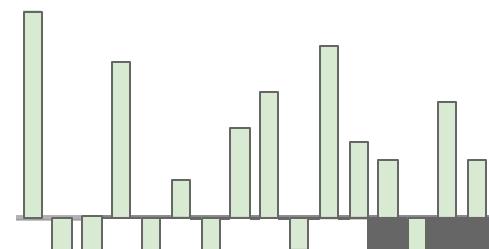
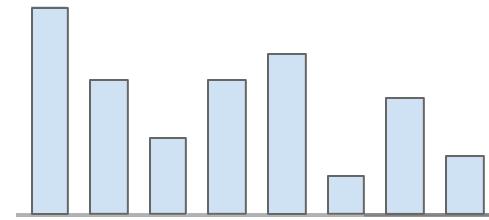
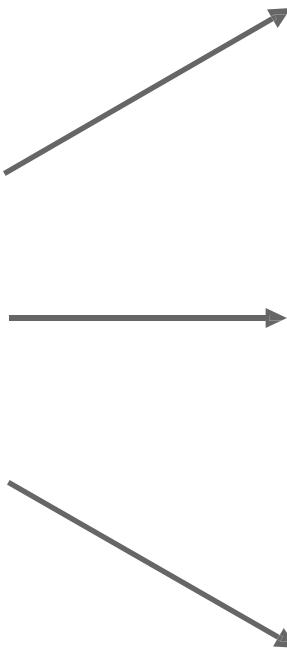
提取随机图像块



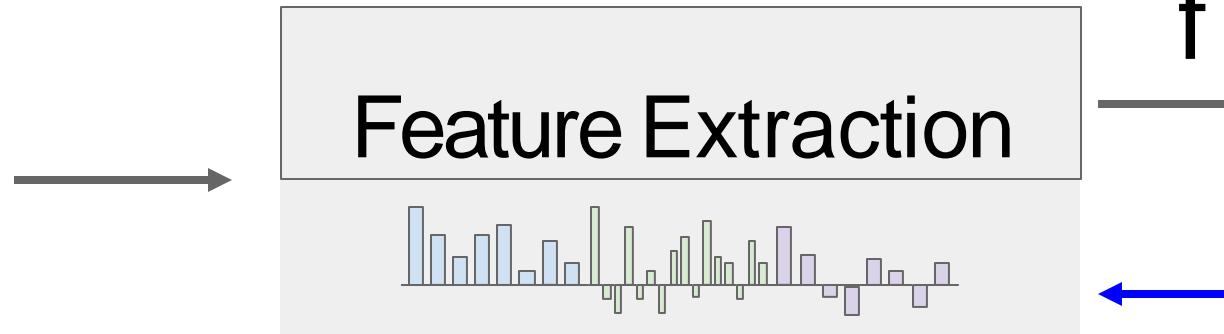
聚类图像块，形成  
“visual words”  
的codebook



## ■ 图像特征

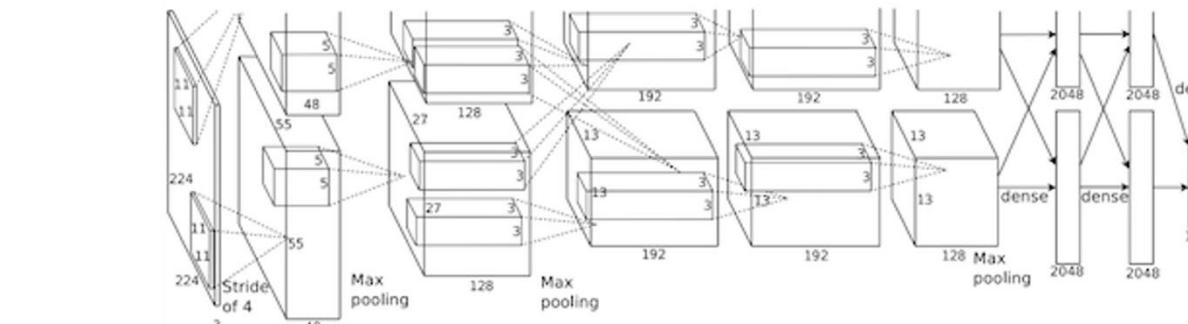


## ■ 手工图像特征与卷积神经网络



分类得分

← training



Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012.  
Figure copyright Krizhevsky, Sutskever, and Hinton, 2012.  
Reproduced with permission.

分类得分

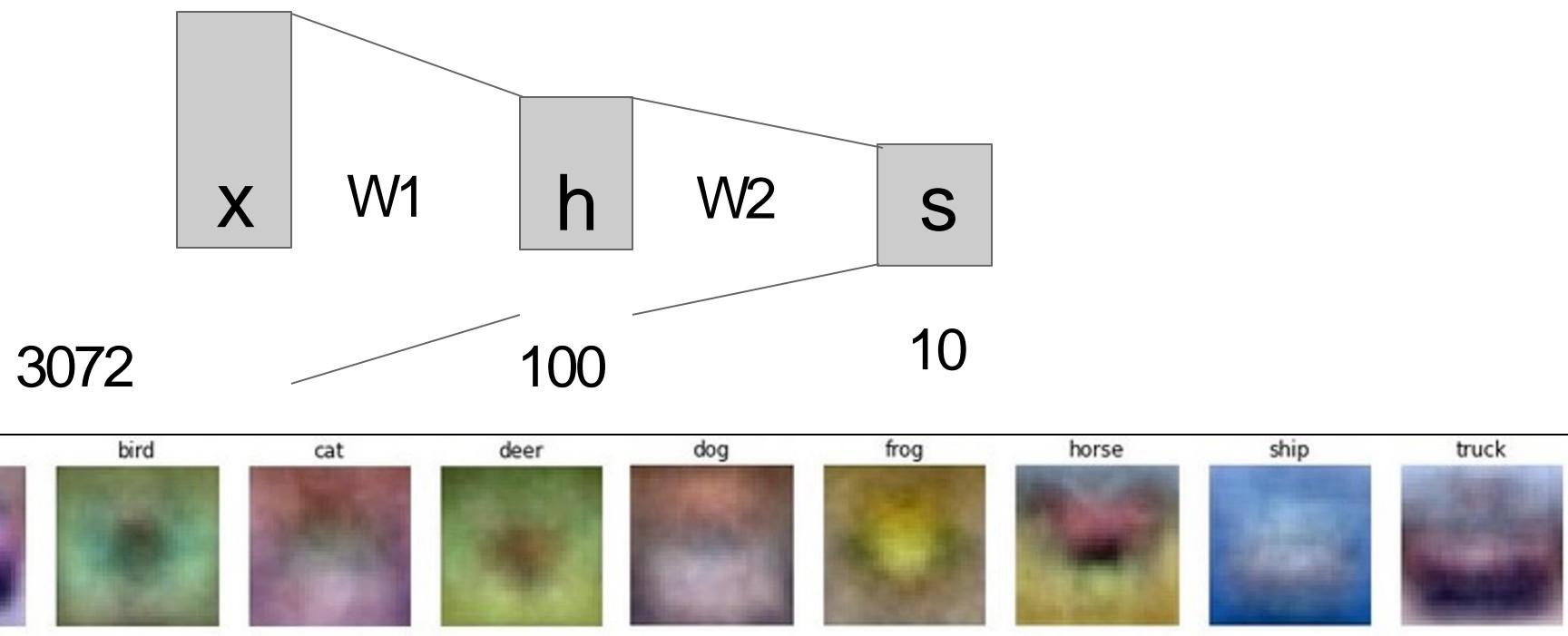
← training

- 线性分类器

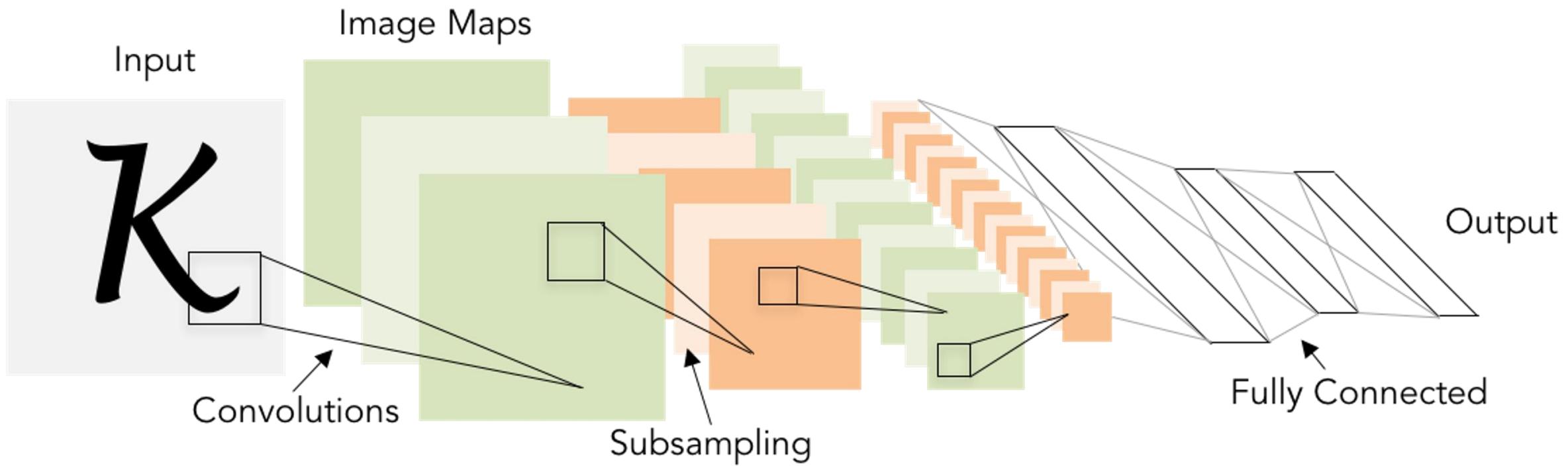
$$f = Wx$$

- 2 层神经网络

$$f = W_2 \max(0, W_1 x)$$



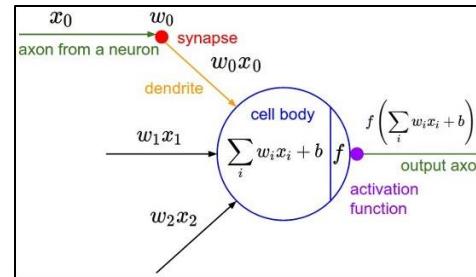
# 卷积神经网络



- 神经网络的发展历史
- 第一个感知机：Mark I 感知机

- 识别字母：

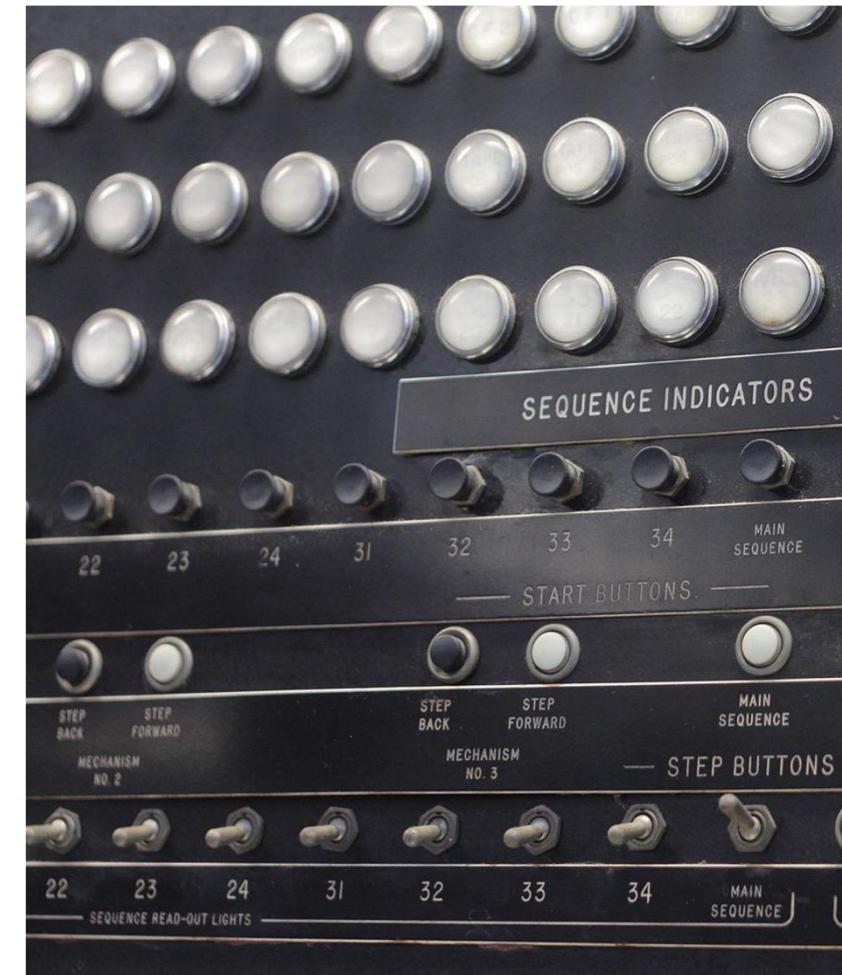
$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$



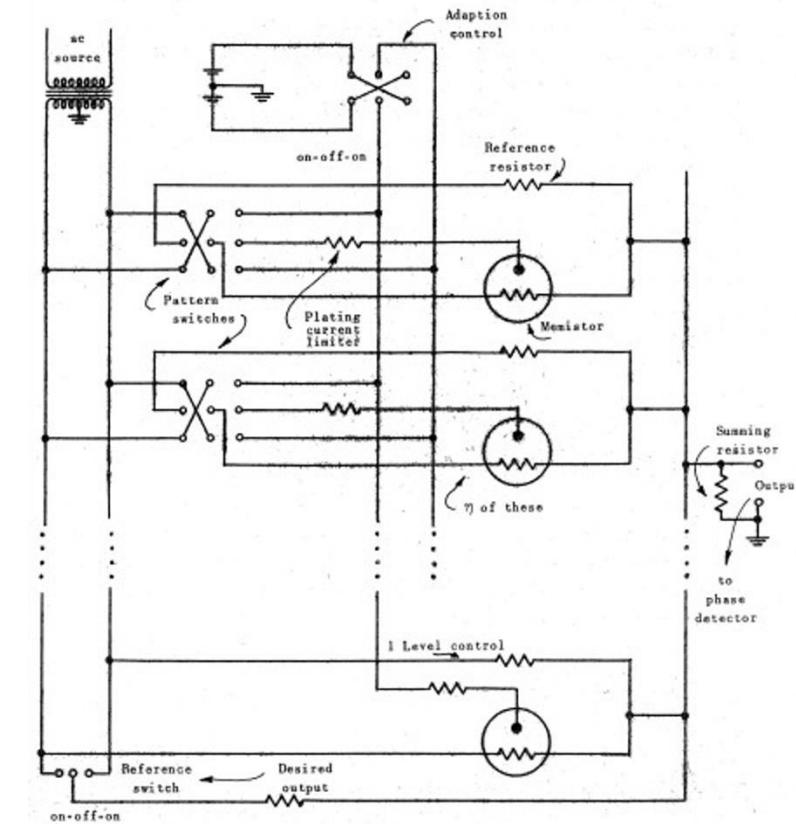
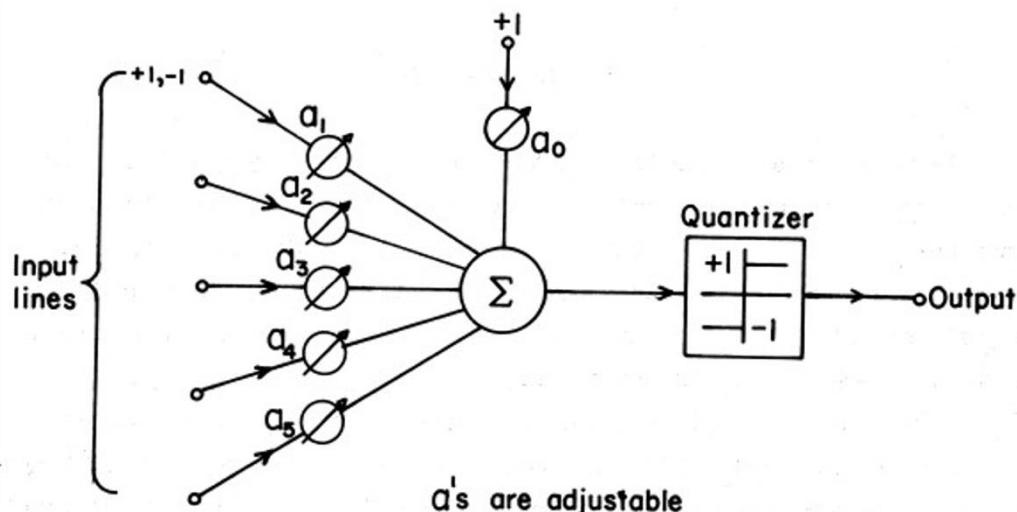
- 网络权重更新规则：

$$w_i(t+1) = w_i(t) + \alpha(d_j - y_j(t))x_{j,i},$$

- Frank Rosenblatt, ~1957: Perceptron

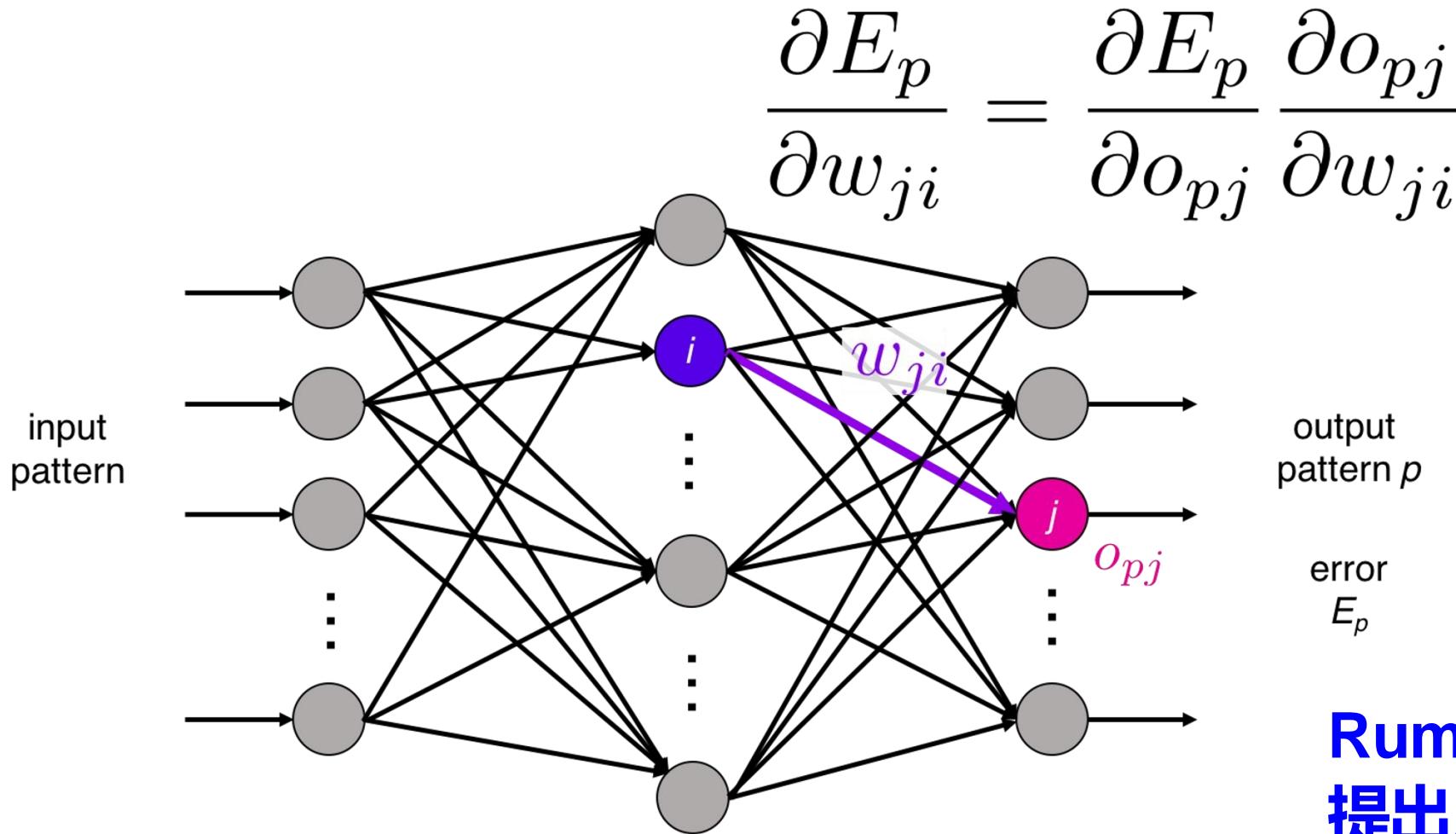


## ■ 神经网络的发展历史



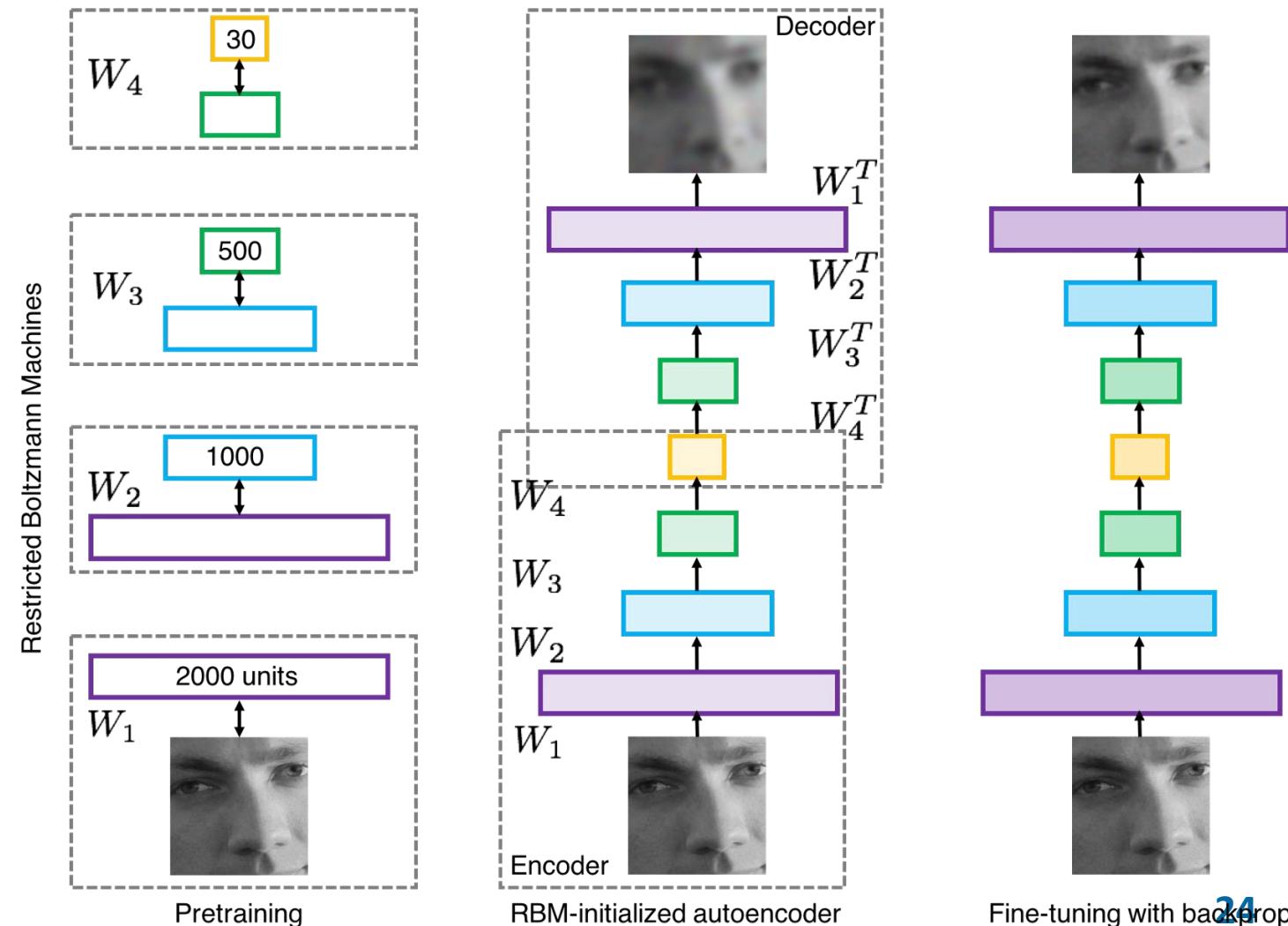
Widrow and Hoff, ~1960 : Adaline/Madaline

## ■ 神经网络的发展历史



Rumelhart et al., 1986:  
提出反向传播

- 神经网络的发展历史
- [Hinton and Salakhutdinov 2006]
- 受限玻尔兹曼机
- 重振深度学习研究



## ■ 神经网络的发展历史

Acoustic Modeling using Deep Belief Networks

Abdel-rahman Mohamed, George Dahl, Geoffrey Hinton, 2010

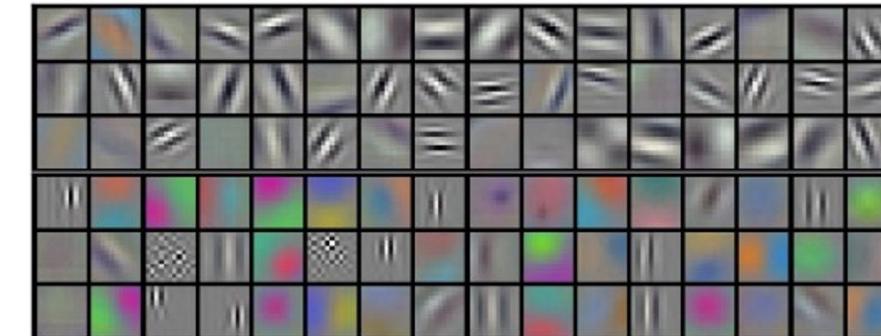
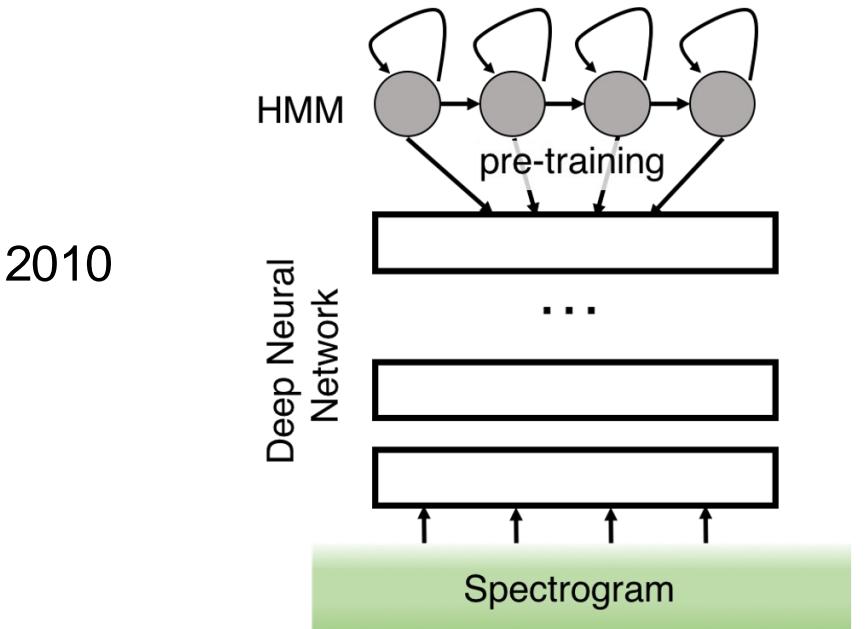
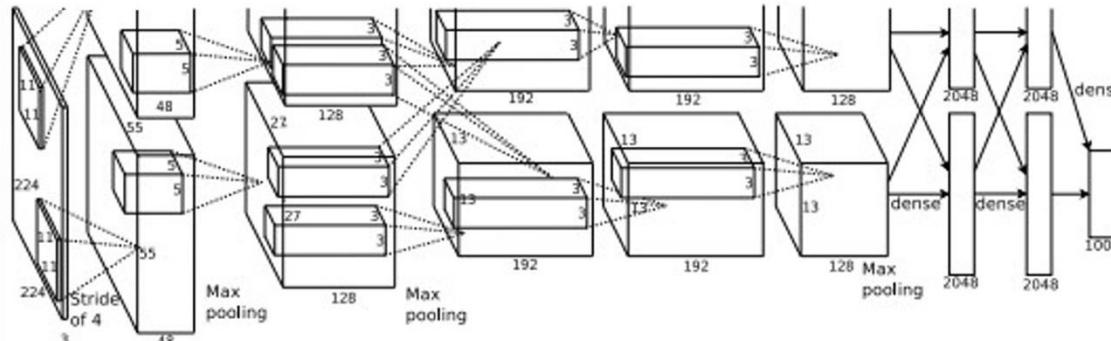
Context-Dependent Pre-trained Deep Neural Networks

for Large Vocabulary Speech Recognition

George Dahl, Dong Yu, Li Deng, Alex Acero, 2012

Imagenet classification with deep convolutional  
neural networks

Alex Krizhevsky, Ilya Sutskever, Geoffrey E Hinton, 2012



## ■ 卷积神经网络的发展历史

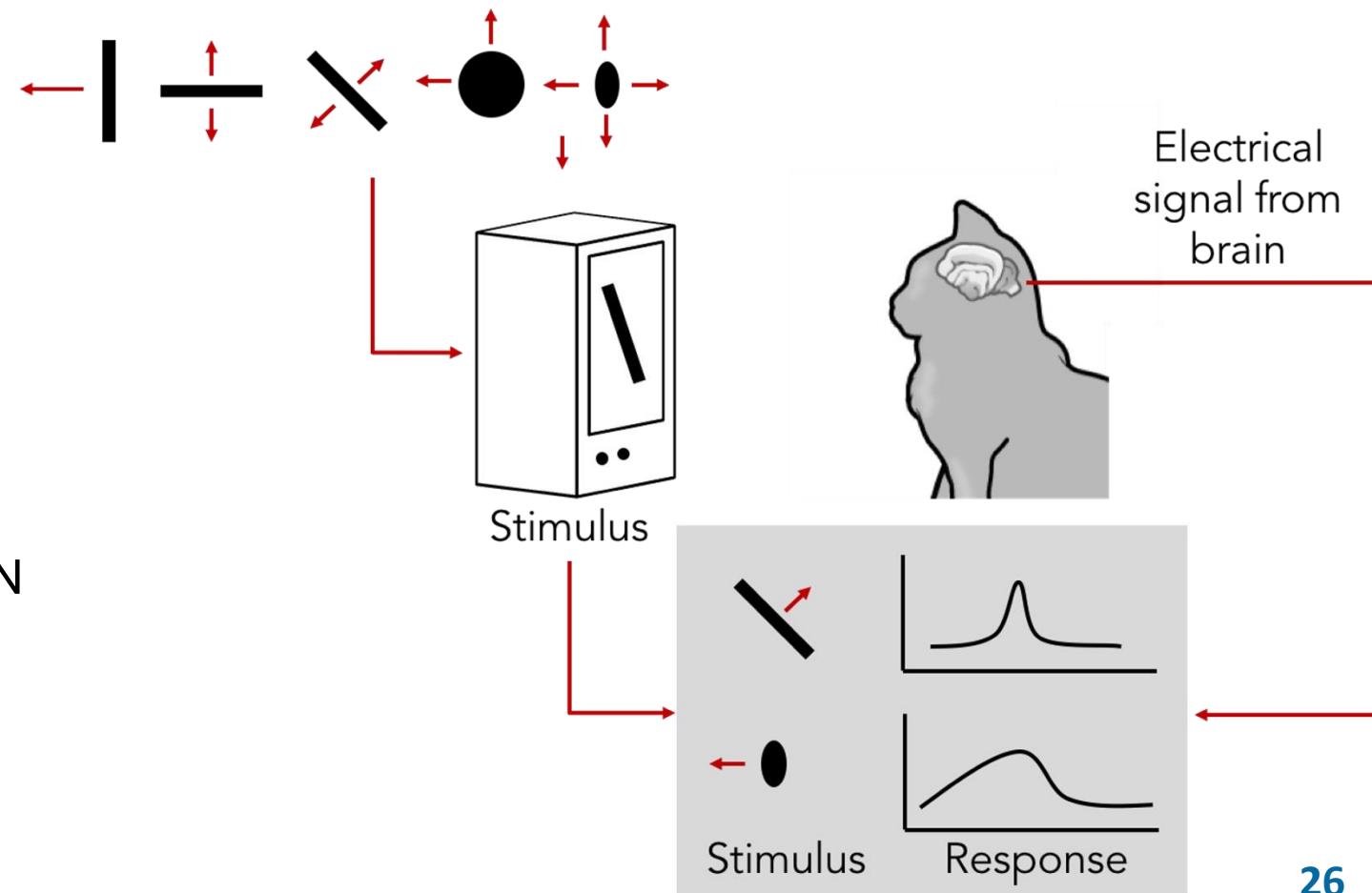
Hubel & Wiesel,  
1959

RECEPTIVE FIELDS OF SINGLE NEURONES IN  
THE CAT'S STRIATE CORTEX

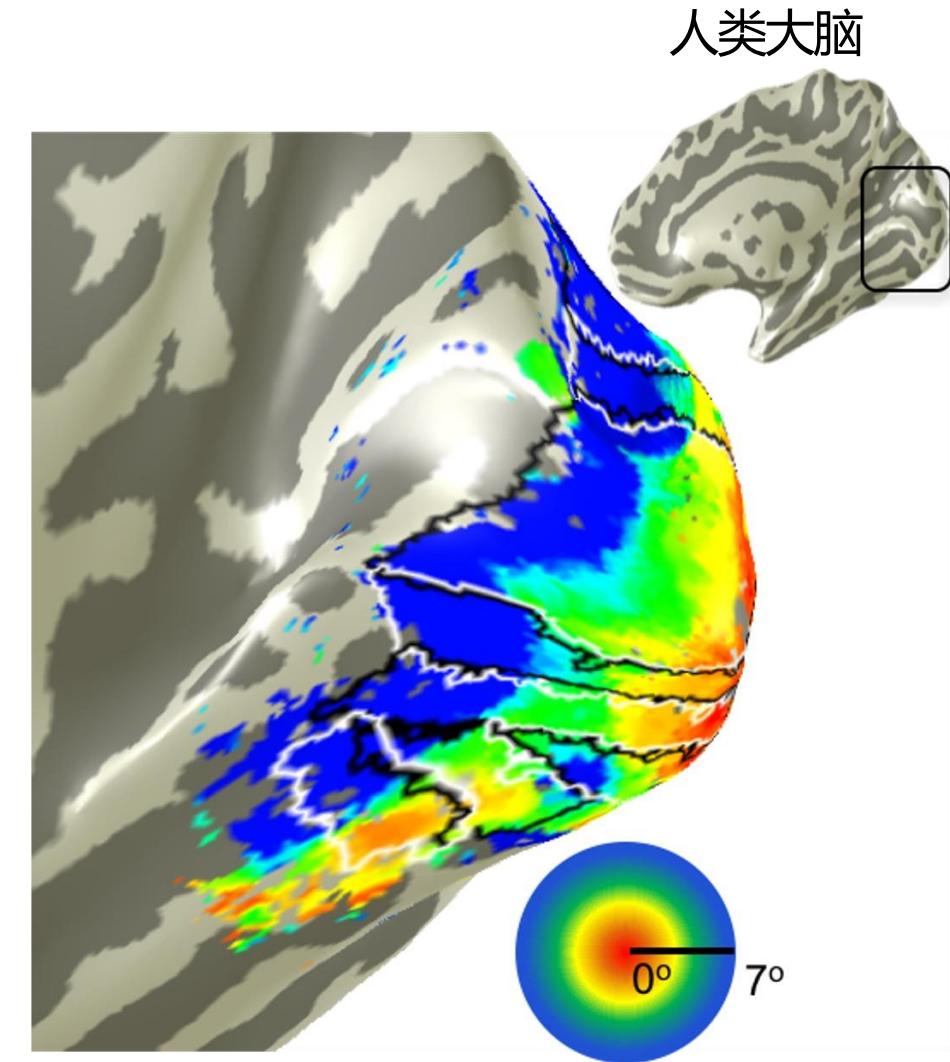
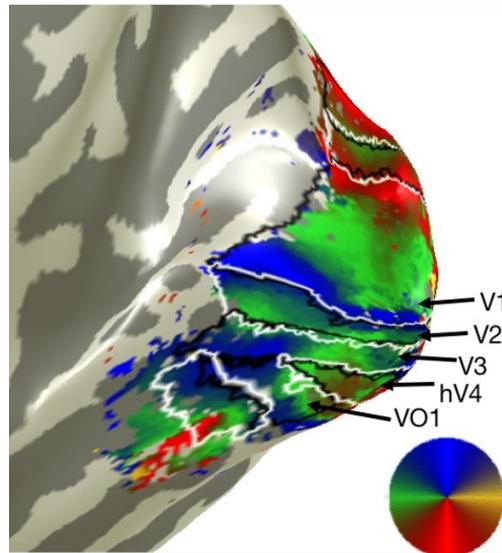
1962

RECEPTIVE FIELDS, BINOCULAR INTERACTION  
AND FUNCTIONAL ARCHITECTURE IN  
THE CAT'S VISUAL CORTEX

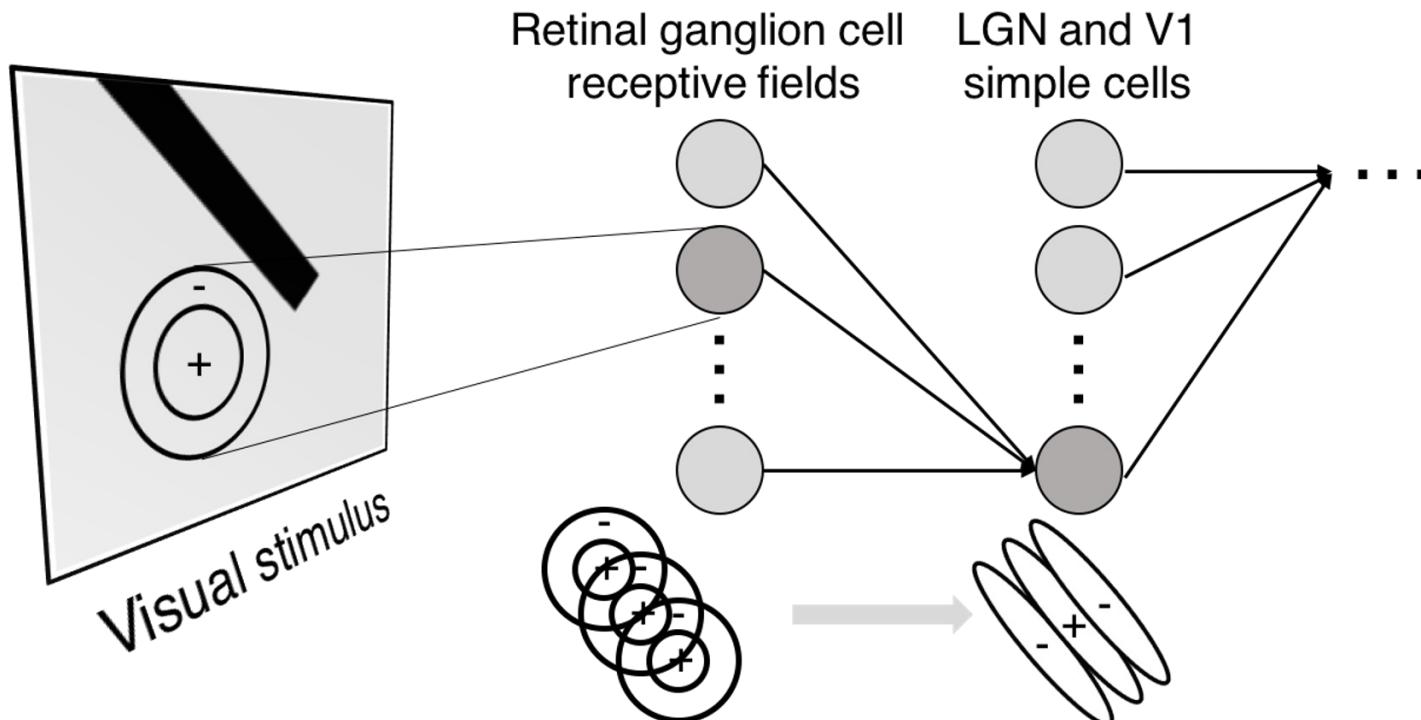
1968...



- 卷积神经网络的发展历史
- 皮层地形图：大脑皮层中的相邻细胞被视野中的相邻区域激活



- 卷积神经网络的发展历史
- 视觉分级感知



**Simple cells:**  
Response to light orientation

**Complex cells:**  
Response to light orientation and movement

**Hypercomplex cells:**  
response to movement with an end point



No response



Response (end point)

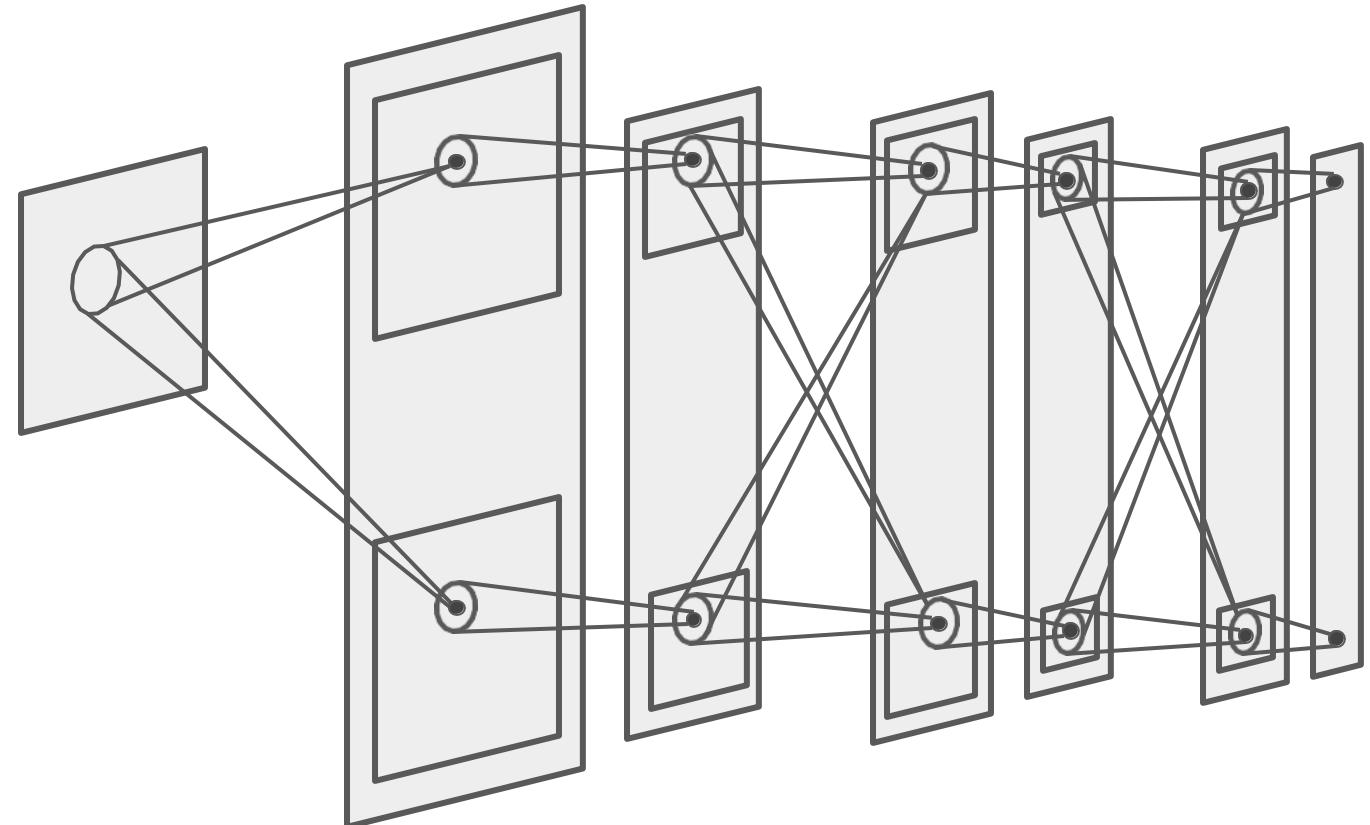
## ■ 卷积神经网络的发展历史

Neocognitron  
[Fukushima 1980]

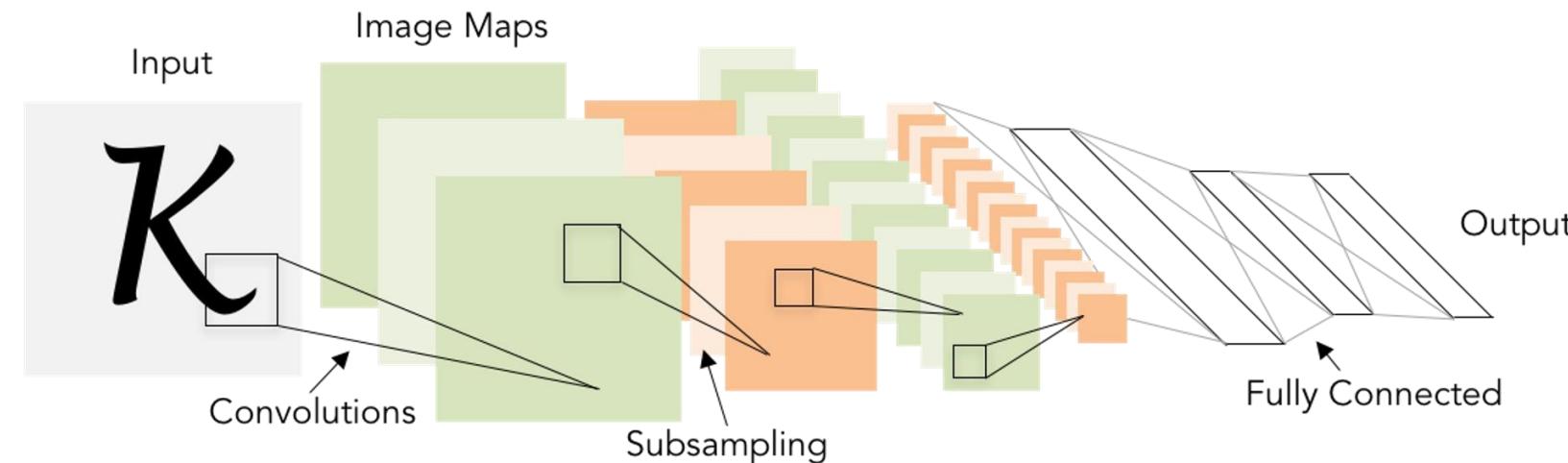
三明治结构 (SCSCSC...)

simple cells: modifiable parameters

complex cells: perform pooling



- 卷积神经网络的发展历史
- 将基于梯度优化的卷积网络应用到文档识别
- [LeCun, Bottou, Bengio, Haffner 1998]



## ■ 卷积神经网络的发展历史

**ImageNet Classification with Deep Convolutional Neural Networks**  
[Krizhevsky, Sutskever, Hinton, 2012]

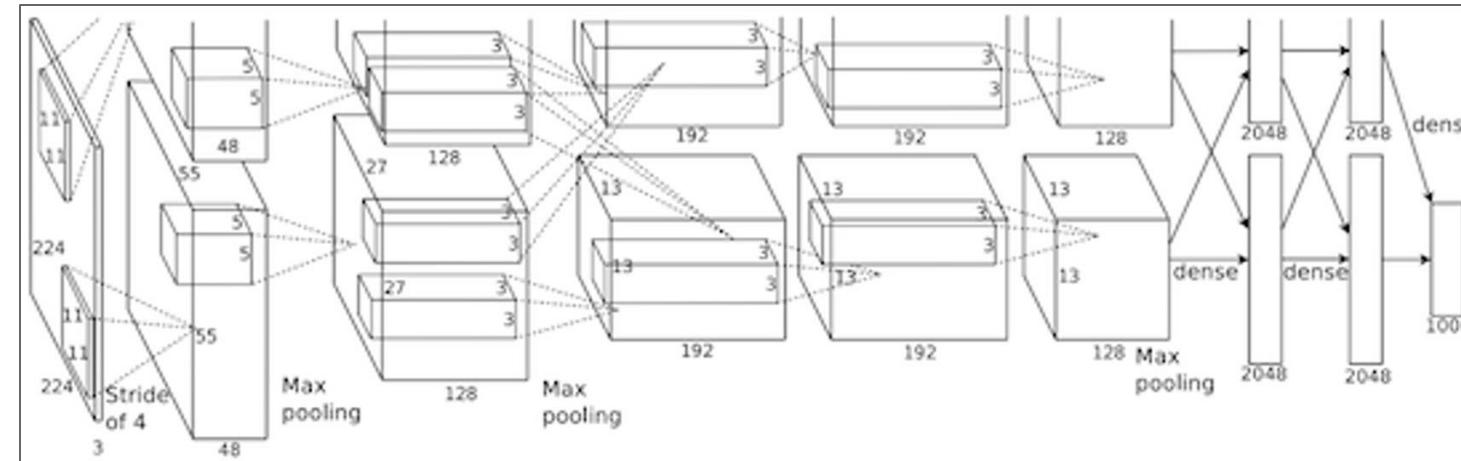


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

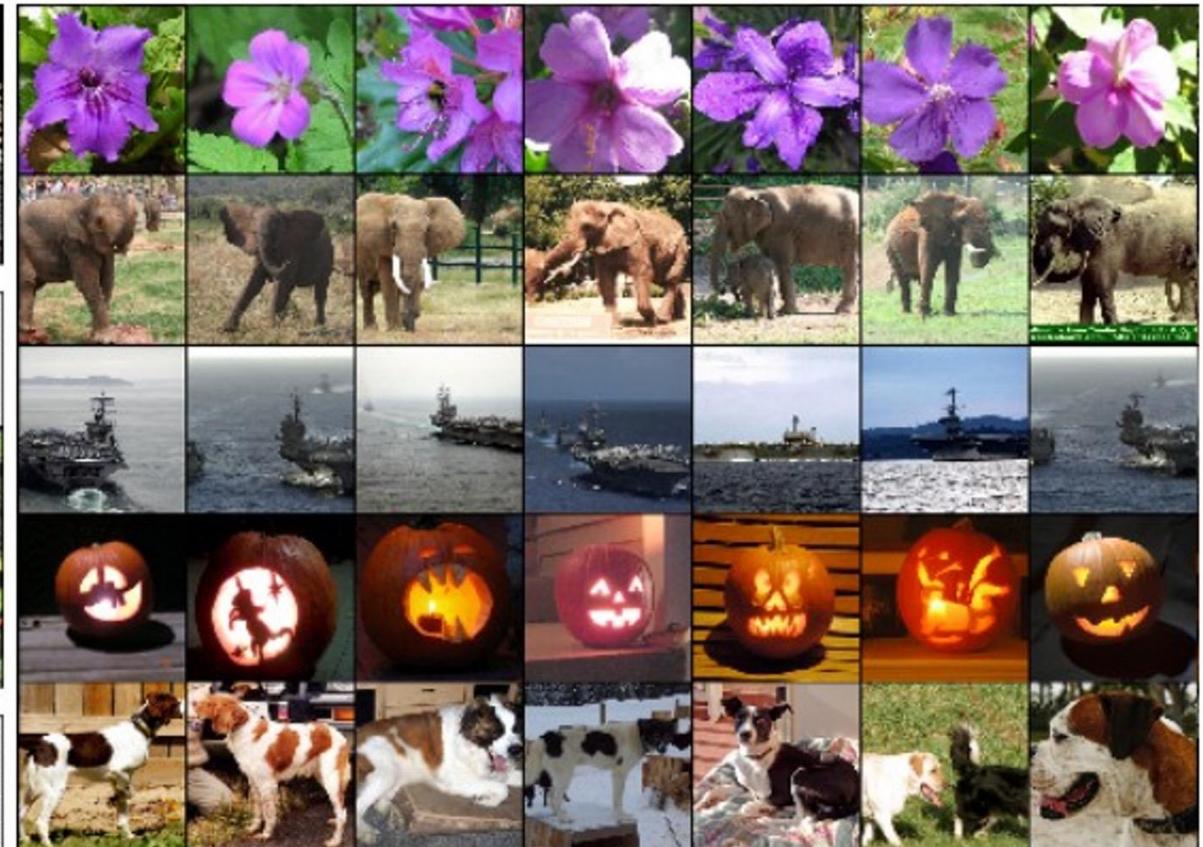
“AlexNet”

## ■ 卷积神经网络的广泛应用

Classification



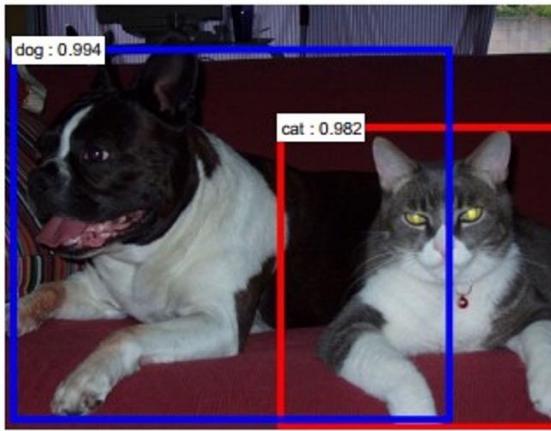
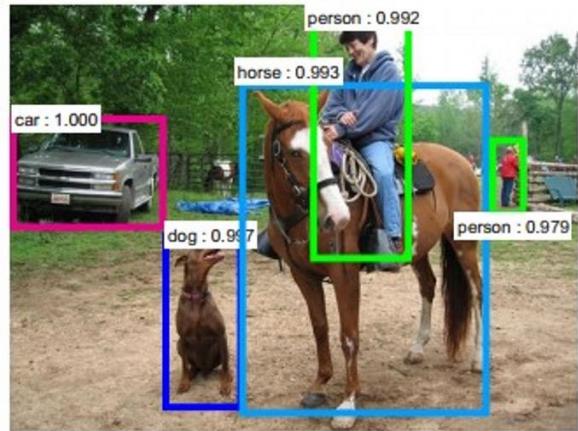
Retrieval



# 卷积神经网络

## ■ 卷积神经网络的广泛应用

# Detection

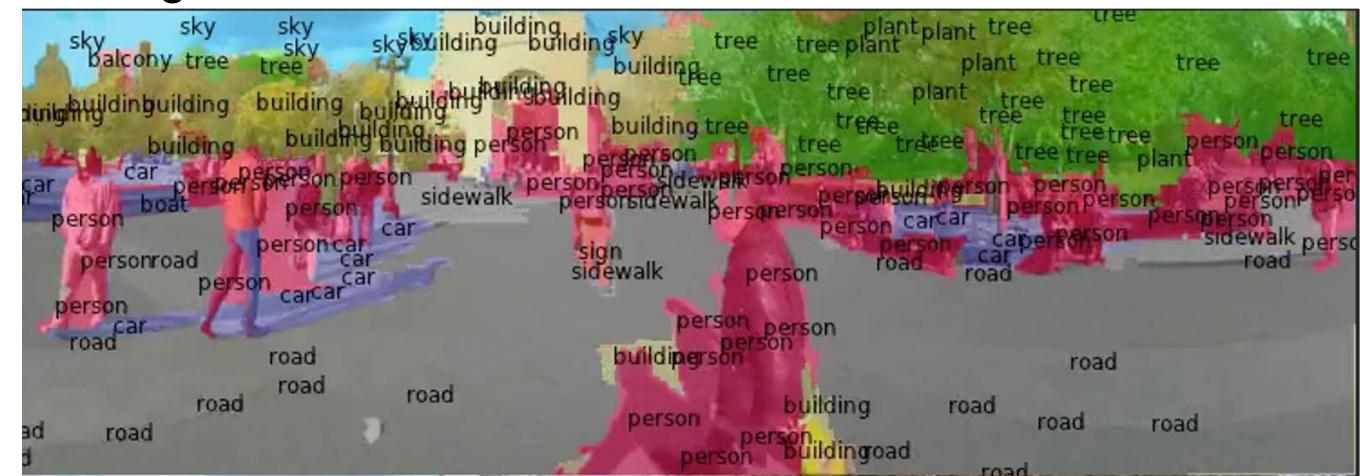


bus : 0.996

person : 0.736



## Segmentation



## ■ 卷积神经网络的广泛应用



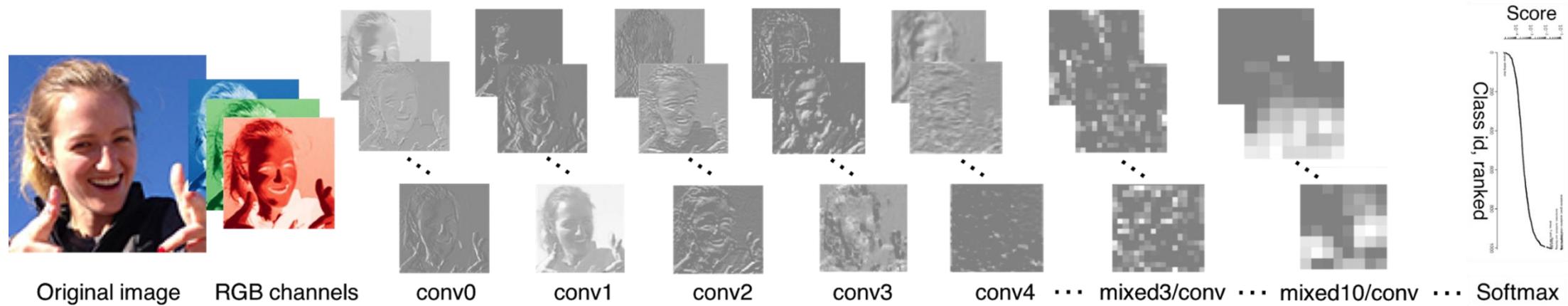
自动驾驶汽车



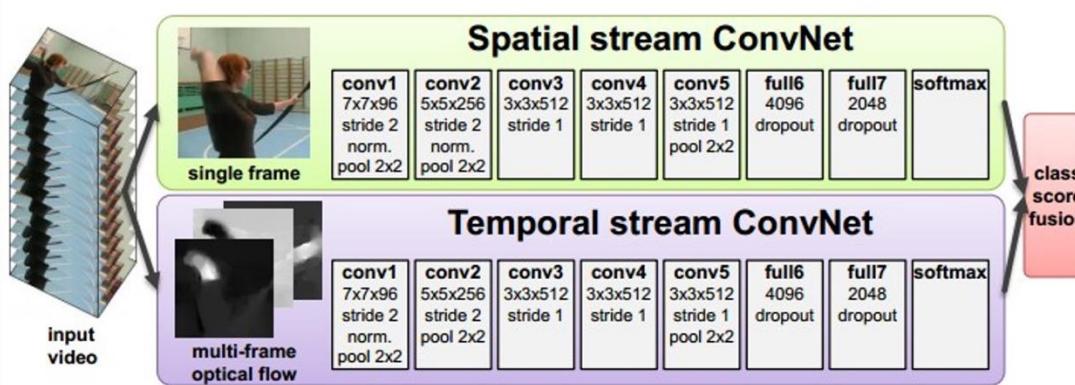
NVIDIA Tesla 系列显卡

# 卷积神经网络

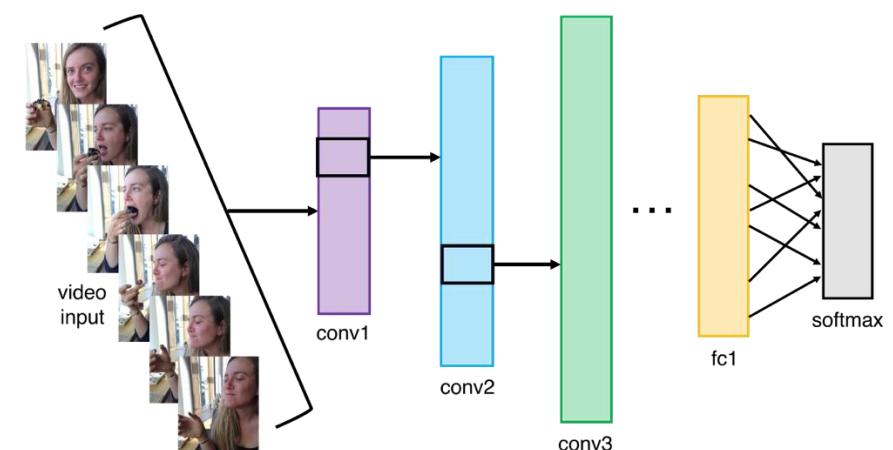
## ■ 卷积神经网络的广泛应用



[Taigman et al. 2014]



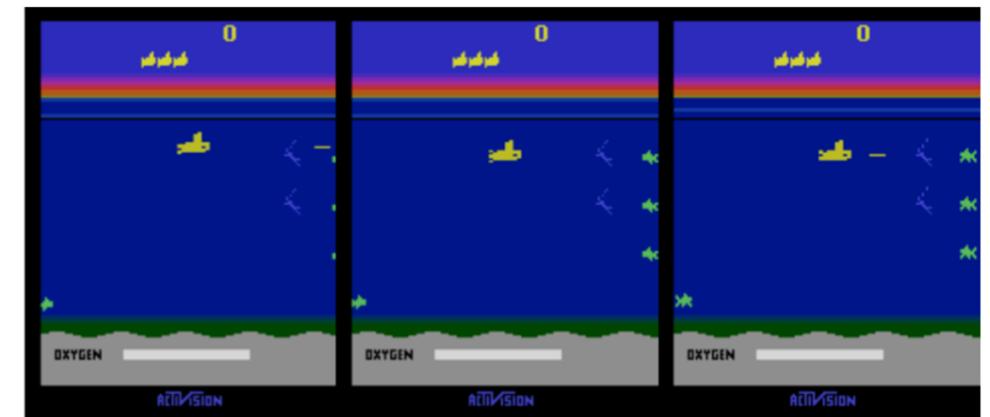
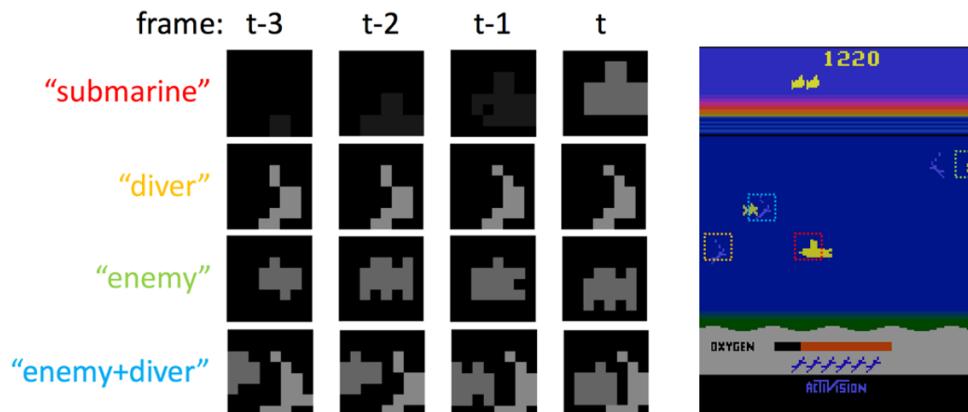
[Simonyan et al. 2014]



## ■ 卷积神经网络的广泛应用

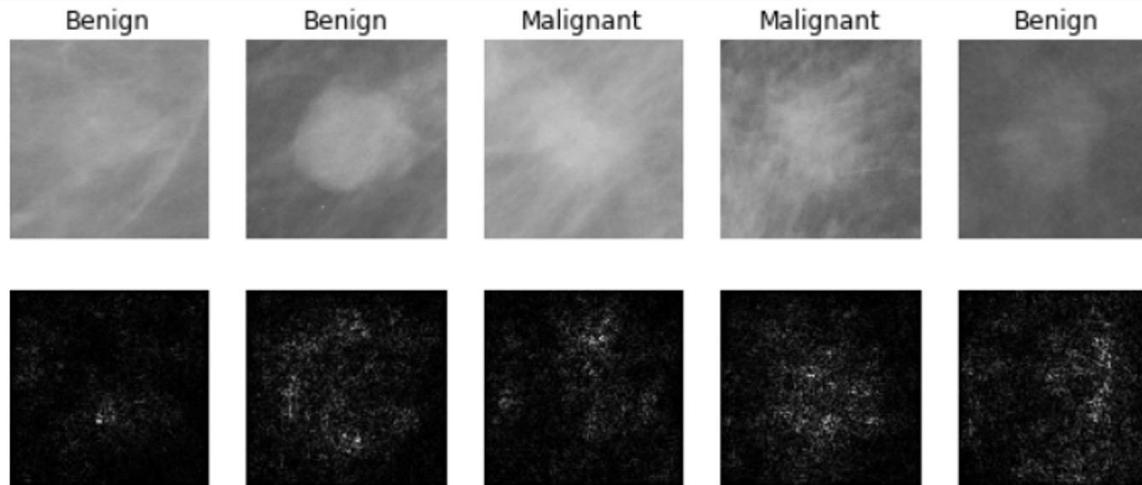


[Toshev, Szegedy 2014]



[Guo et al. 2014]

## ■ 卷积神经网络的广泛应用



[Levy et al. 2016]



[Dieleman et al. 2014]



[Sermanet et al. 2011]  
[Ciresan et al.]

## ■ 卷积神经网络的广泛应用



Kaggle Challenge



Mnih and Hinton, 2010

## ■ 卷积神经网络的广泛应用

No errors



A white teddy bear  
sitting in the grass



A man riding a wave on  
top of a surfboard

Minor errors



A man in a baseball  
uniform throwing a ball



A cat sitting on a  
suitcase on the floor

Somewhat related



A woman is holding a  
cat in her hand

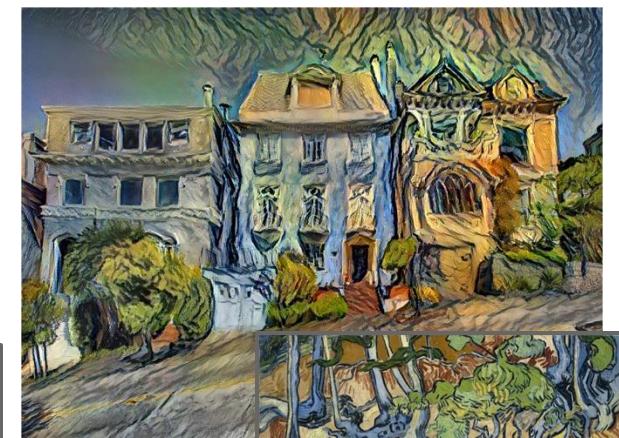
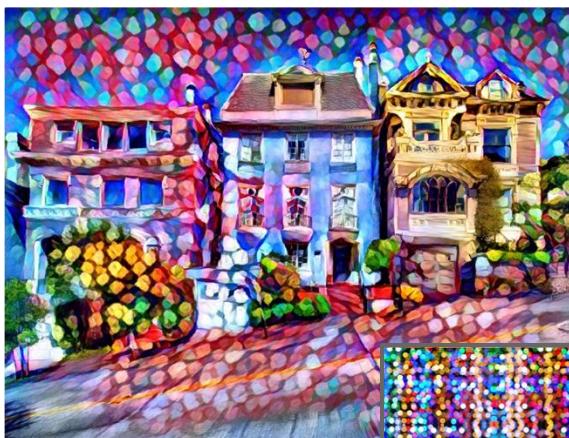
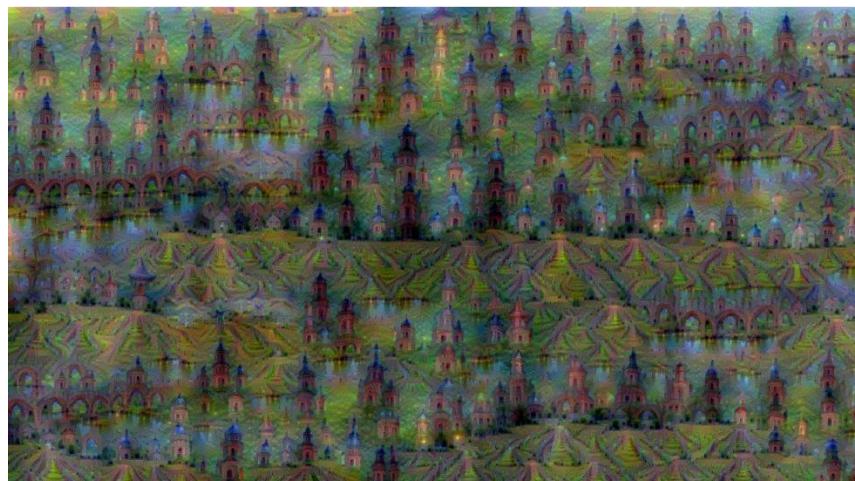
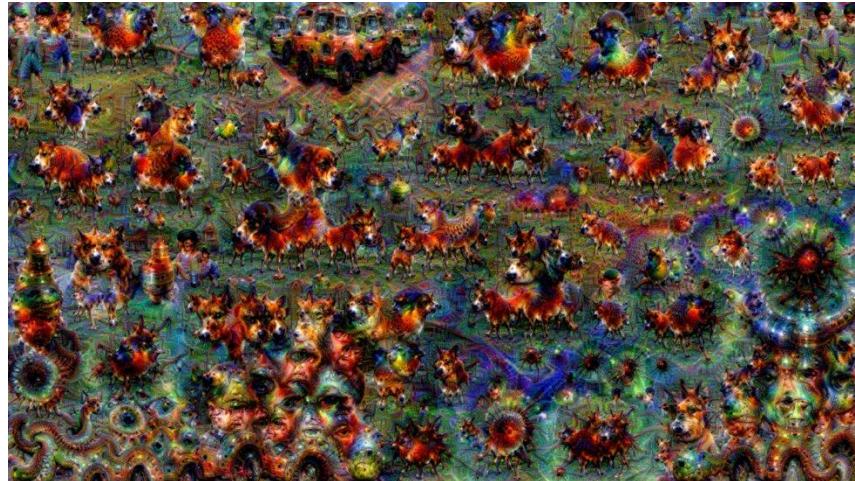


A woman standing on a  
beach holding a surfboard

## 图像描述

[Vinyals et al., 2015]  
[Karpathy and Fei-Fei,  
2015]

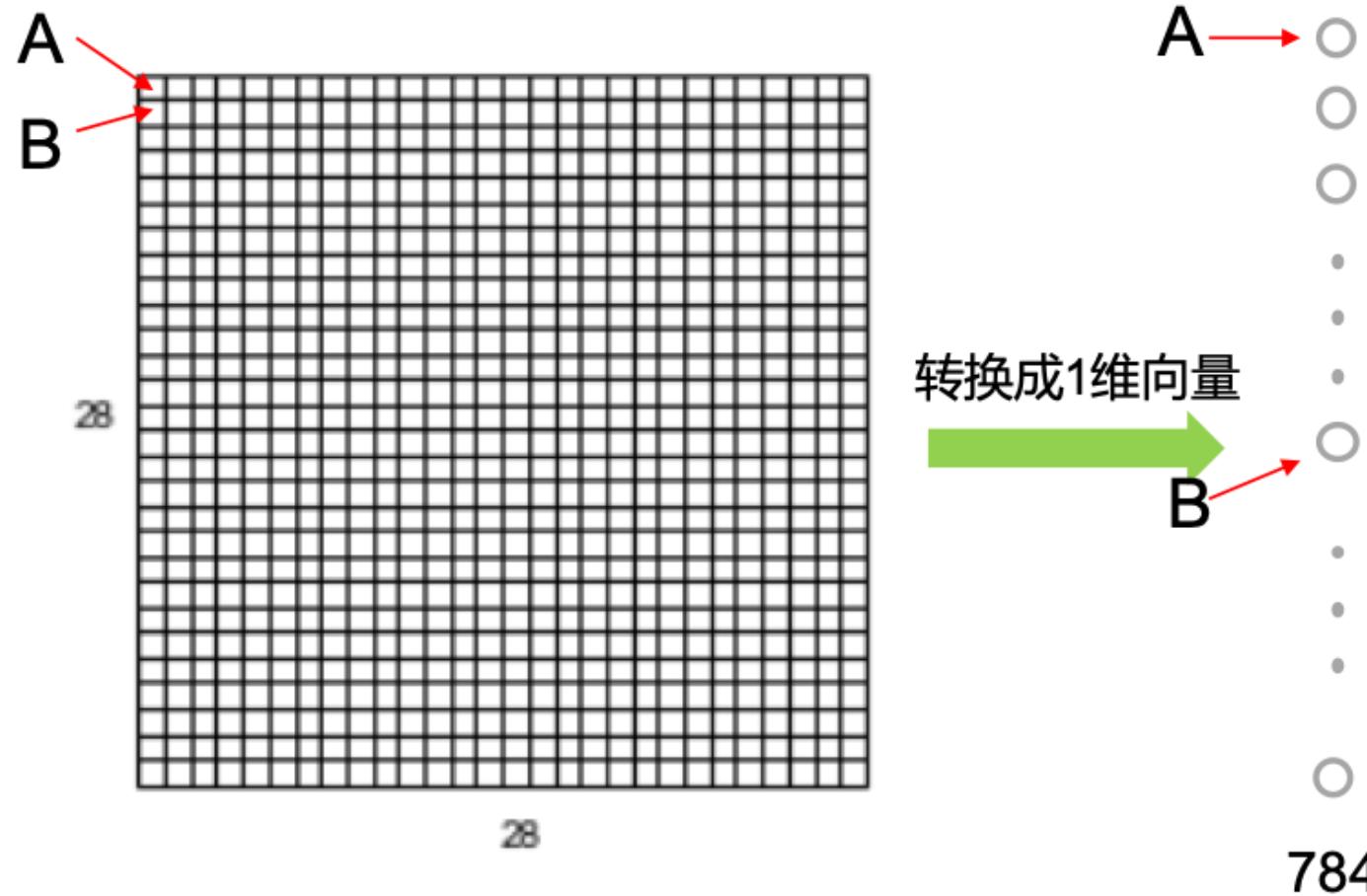
## ■ 卷积神经网络的广泛应用



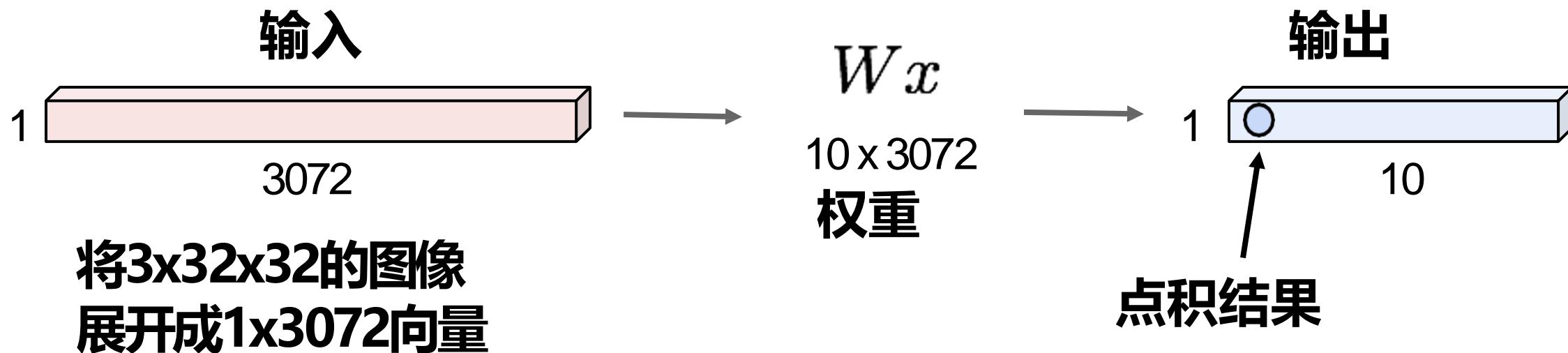
# 大 纲

- 卷积神经网络-引言
- 卷积神经网络

## ■ 全连接层（参数量，优化，空间信息）

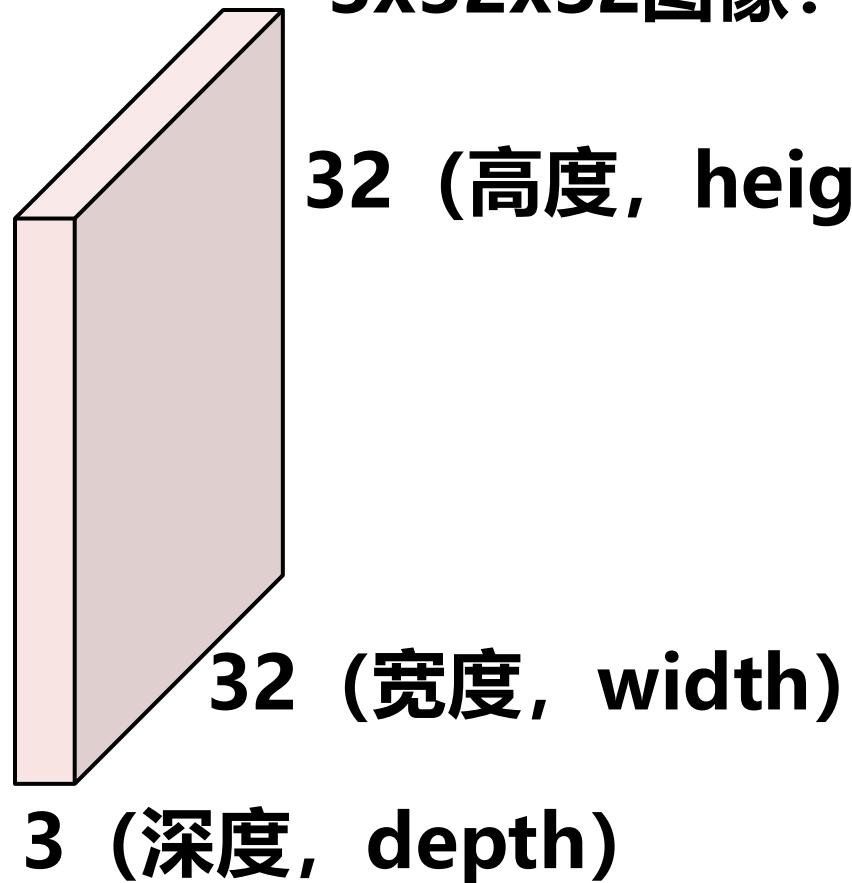


## ■ 全连接层（参数量，优化，空间信息）

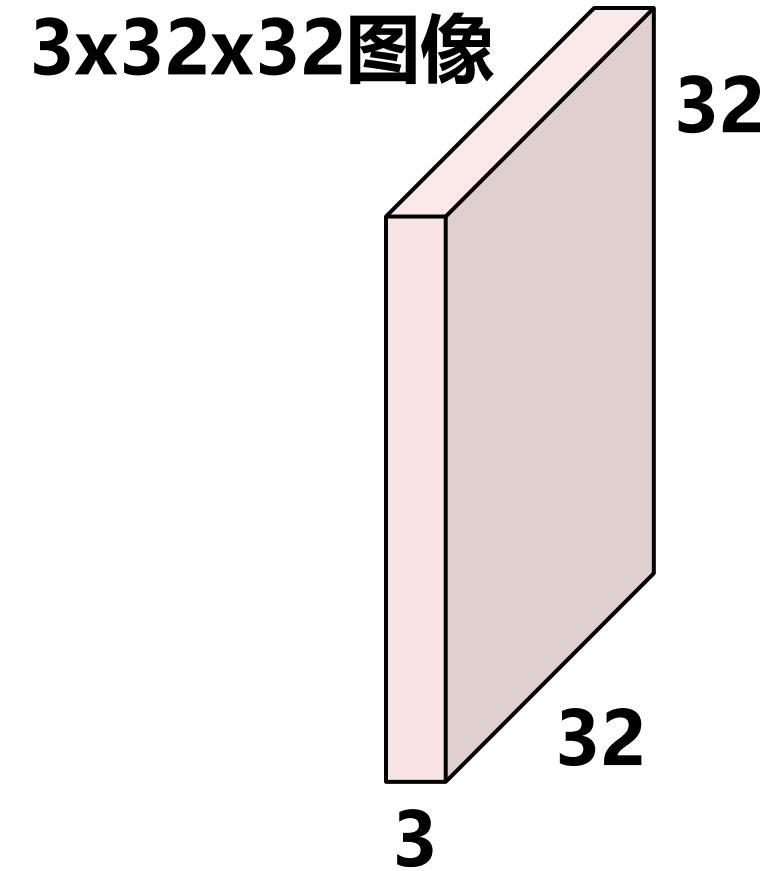


## ■ 卷积层

**3x32x32图像：保持图像的空间结构**



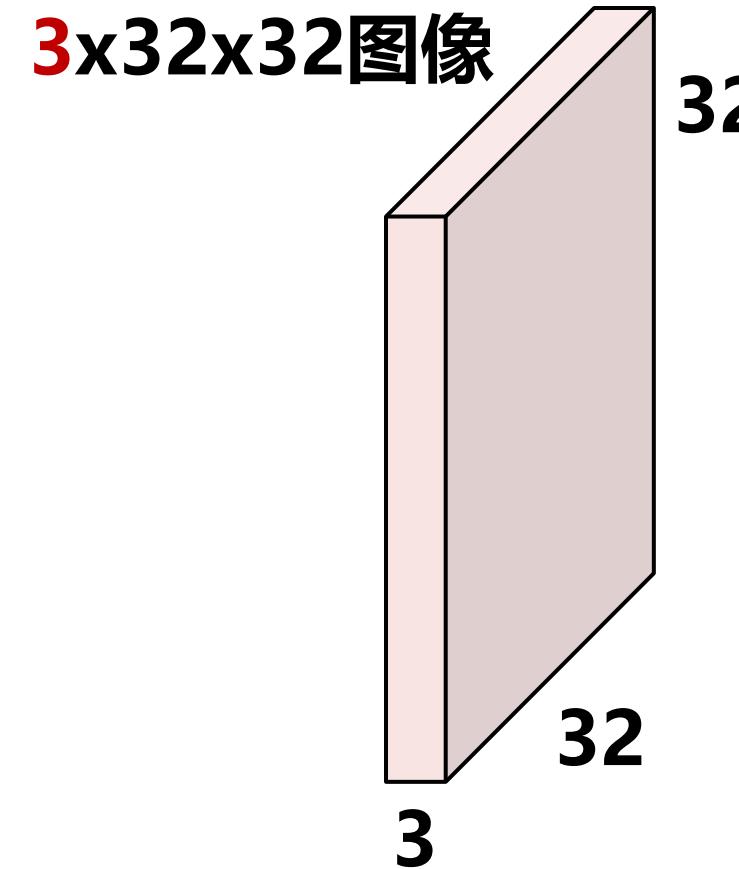
## ■ 卷积层



3x5x5滤波器

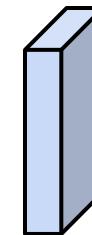


## ■ 卷积层



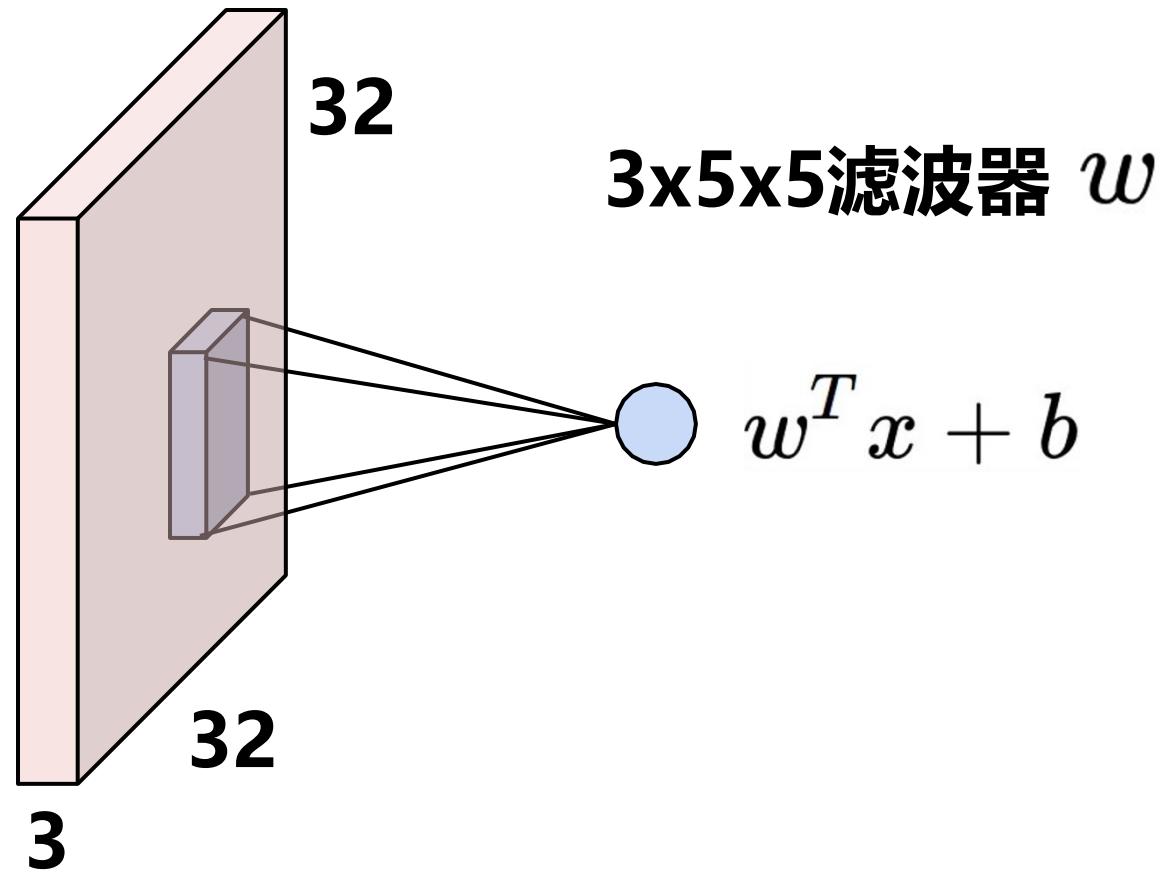
滤波器的深度始终与输入的  
深度相同

3x5x5滤波器

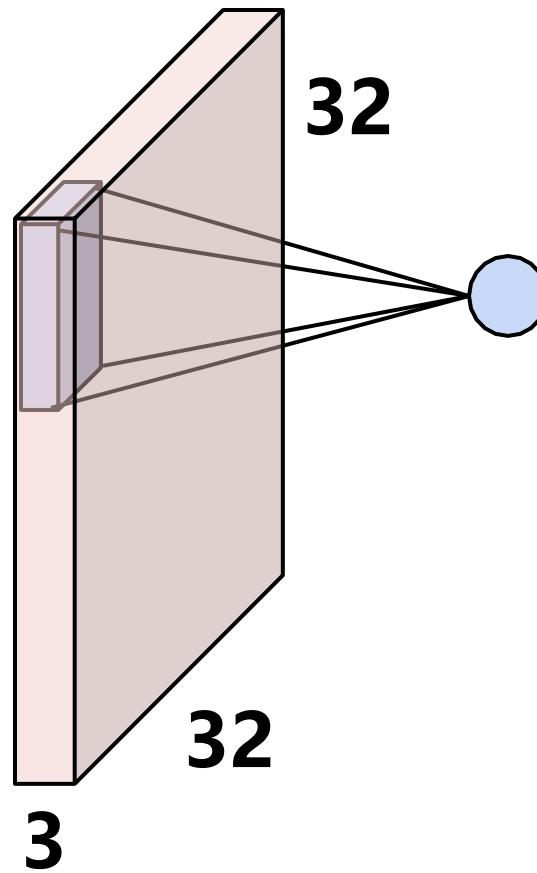


将滤波器与图像卷积，  
即“在图像上空间滑动，计算点积”

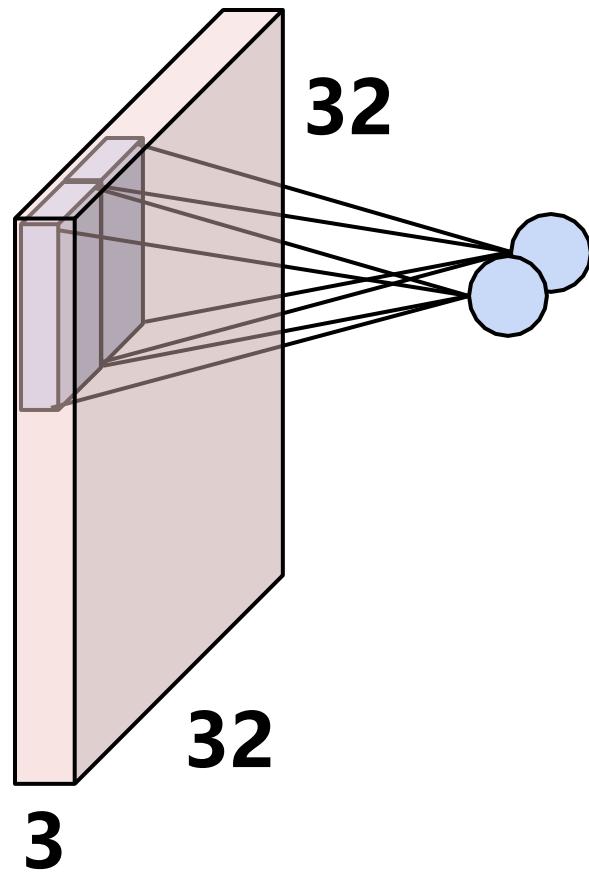
## ■ 卷积层



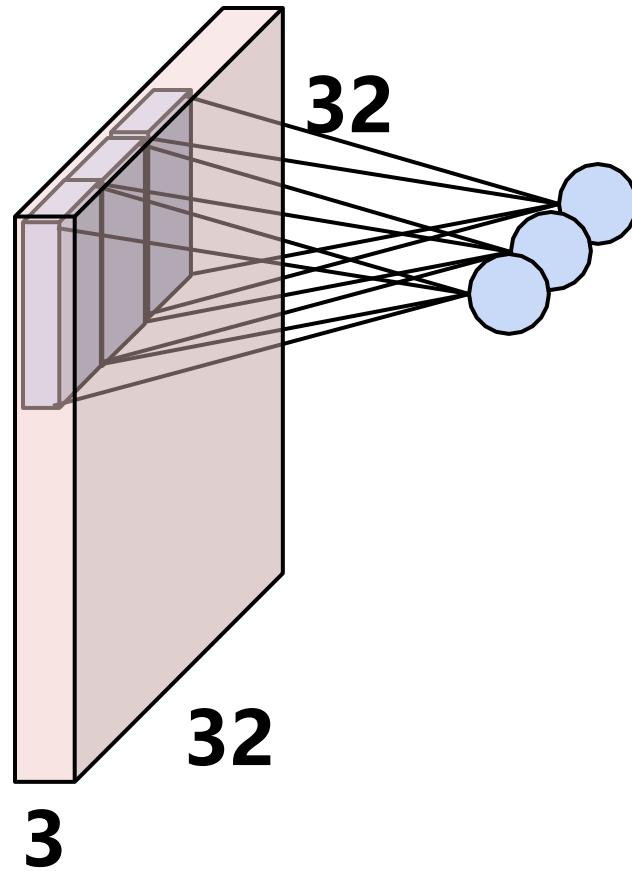
## ■ 卷积层



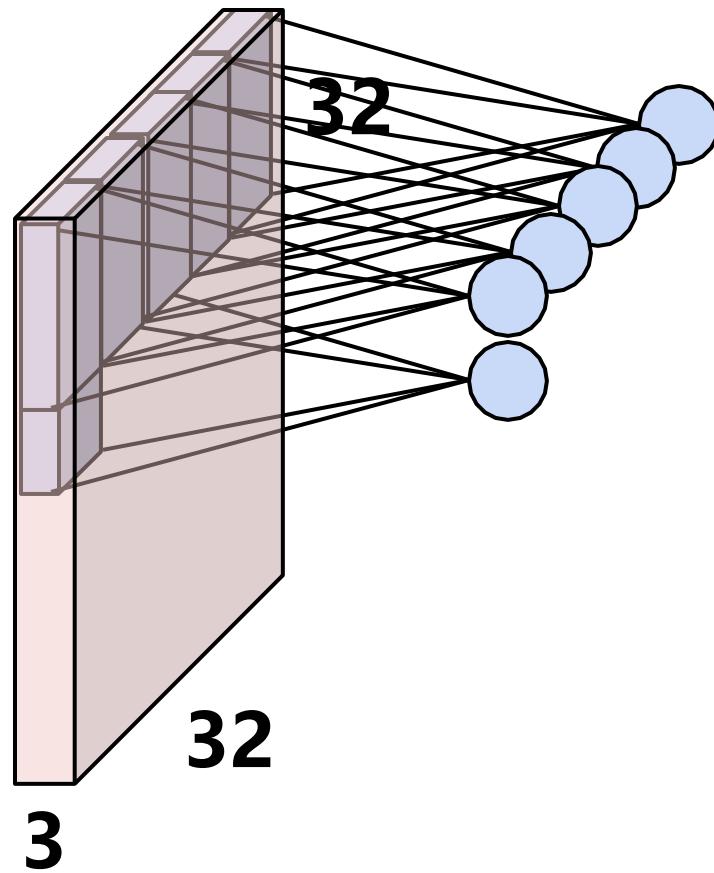
## ■ 卷积层



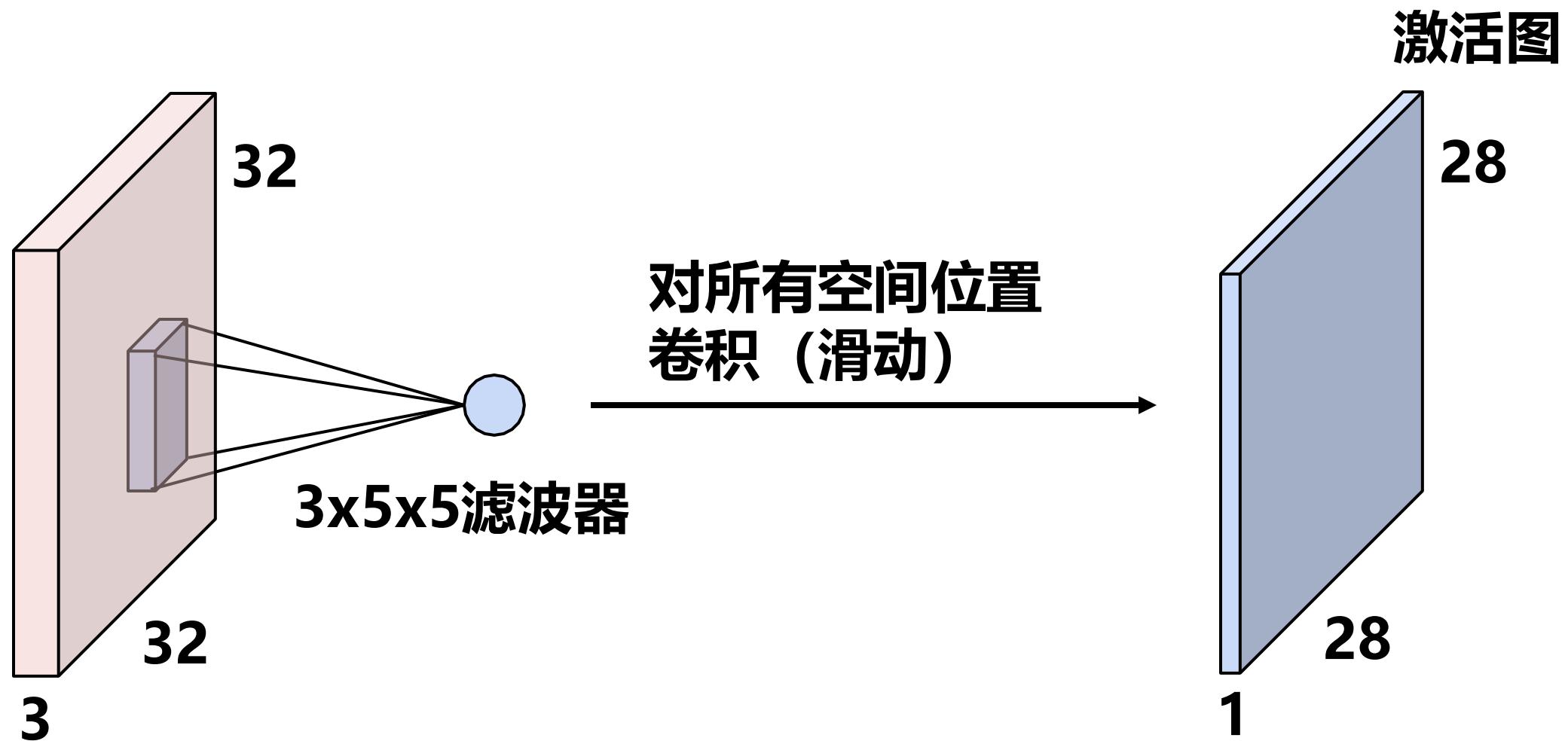
## ■ 卷积层



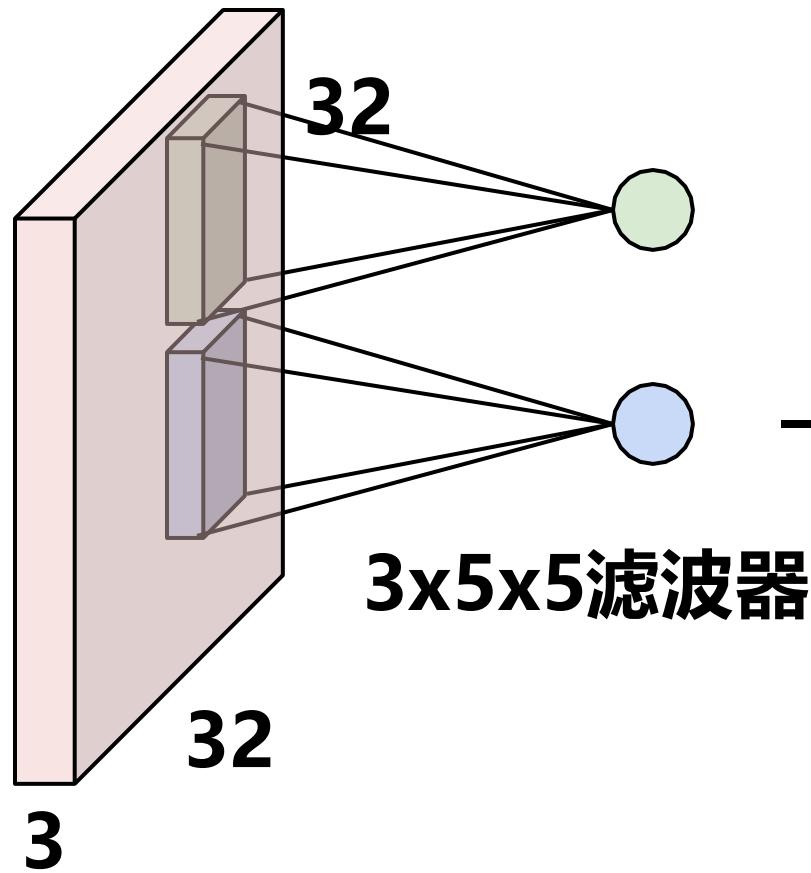
## ■ 卷积层



## ■ 卷积层



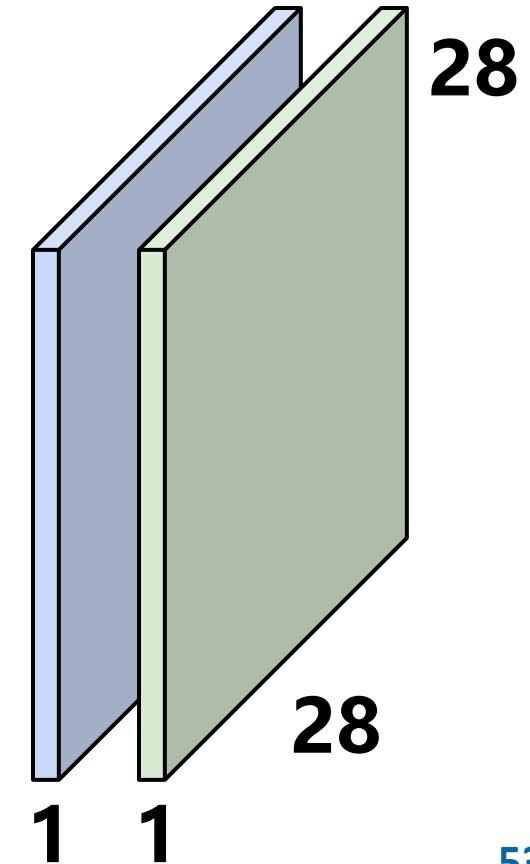
## ■ 卷积层



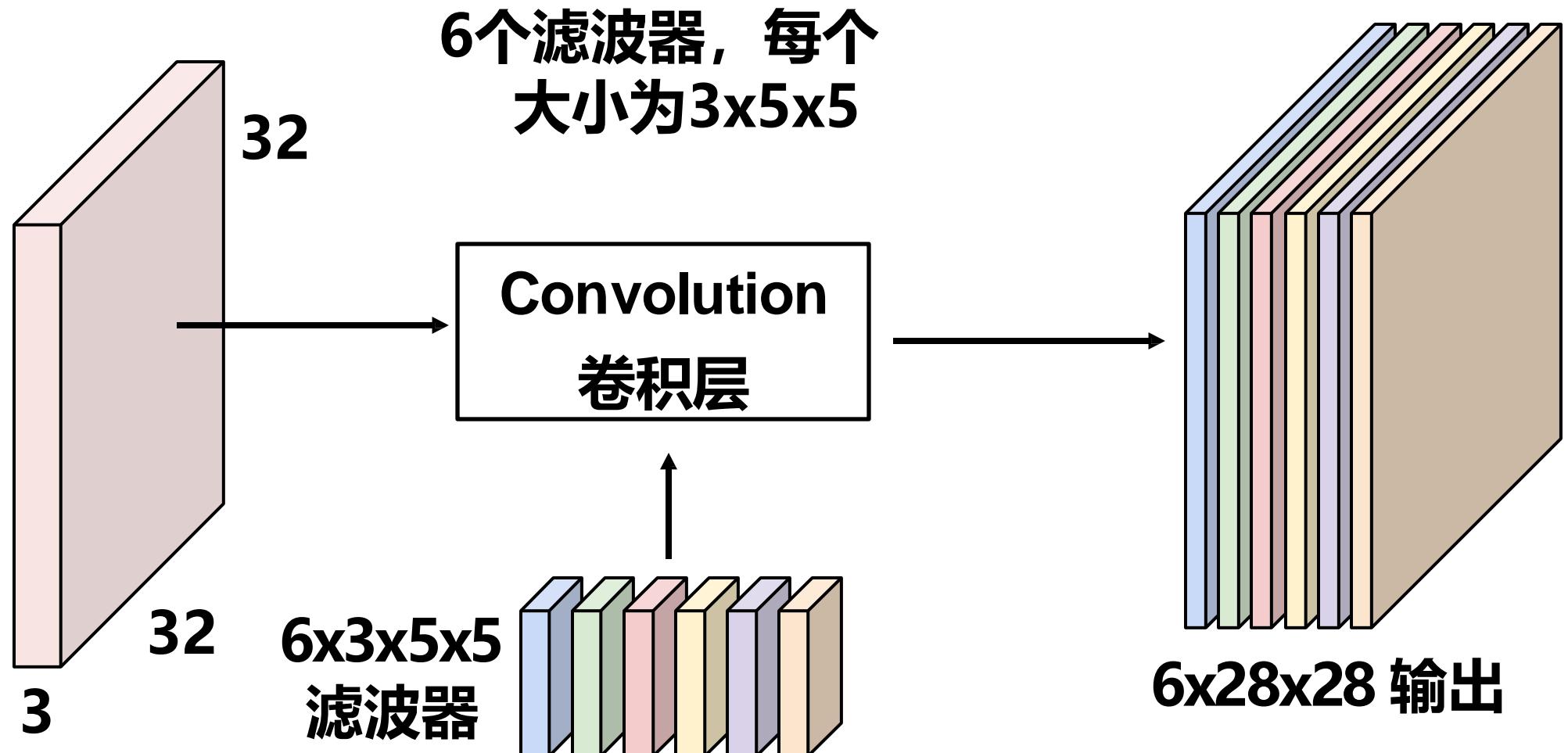
增加一个卷积 (绿色)

对所有空间位置  
卷积 (滑动)

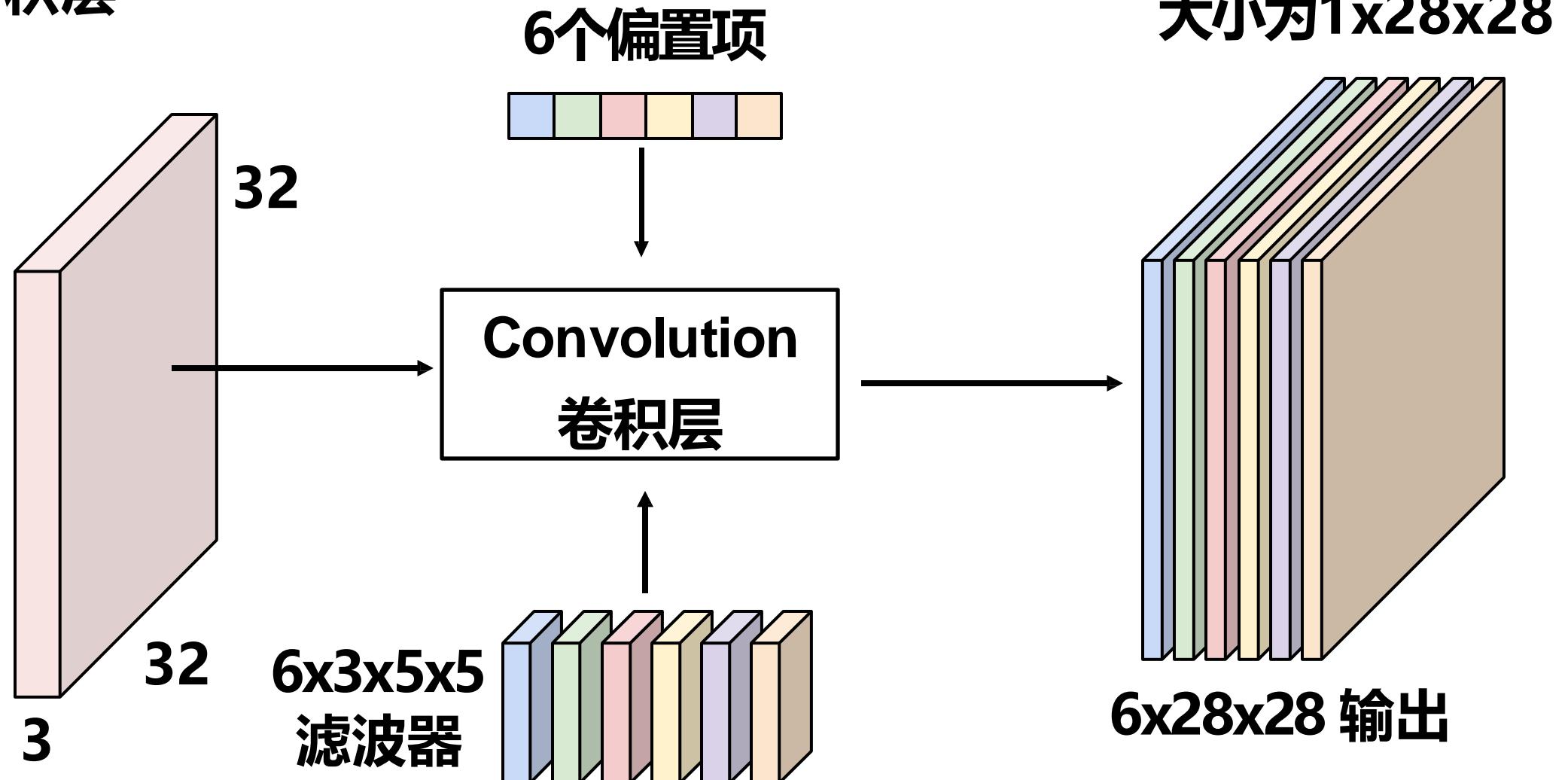
激活图



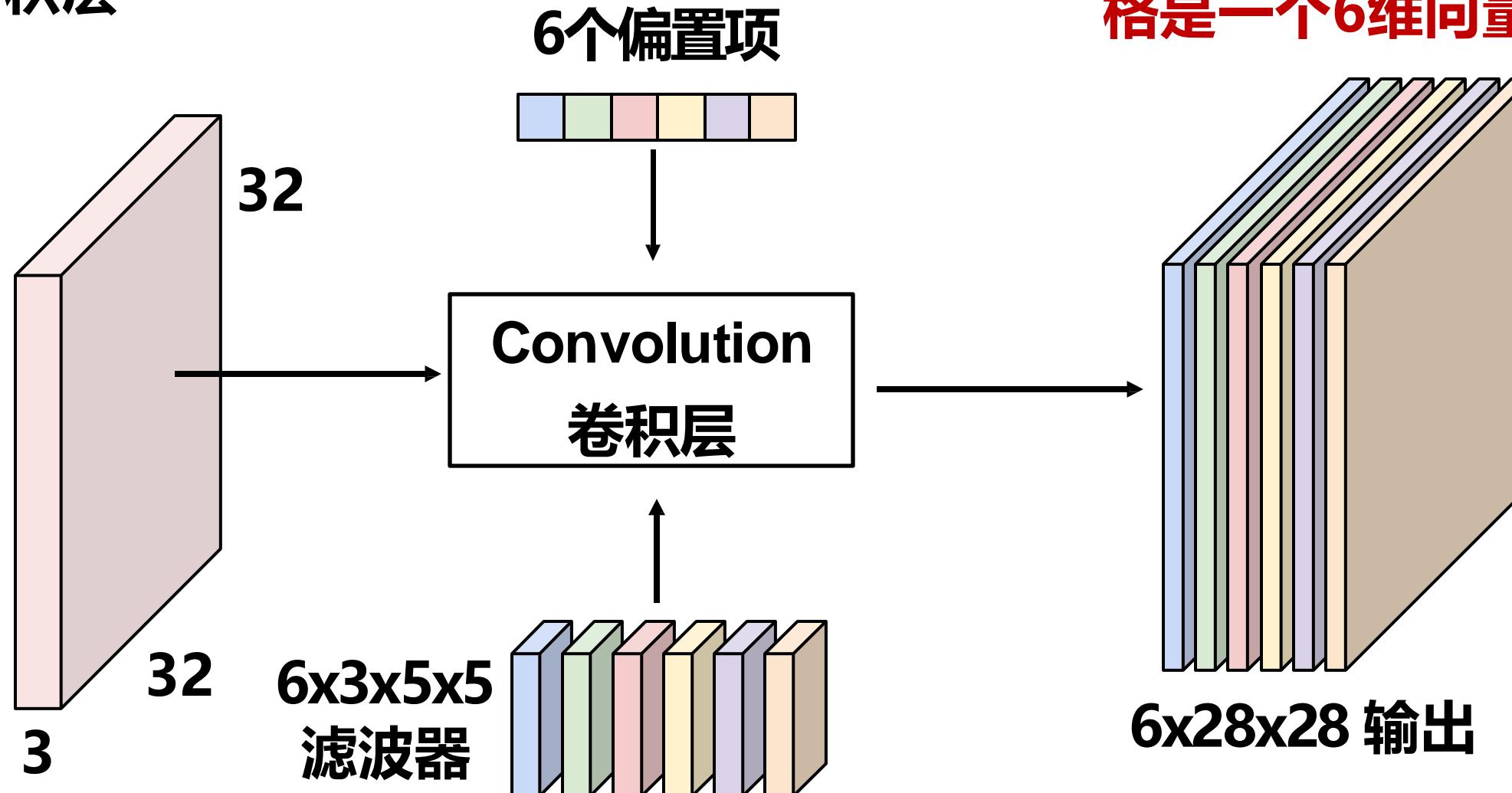
## ■ 卷积层



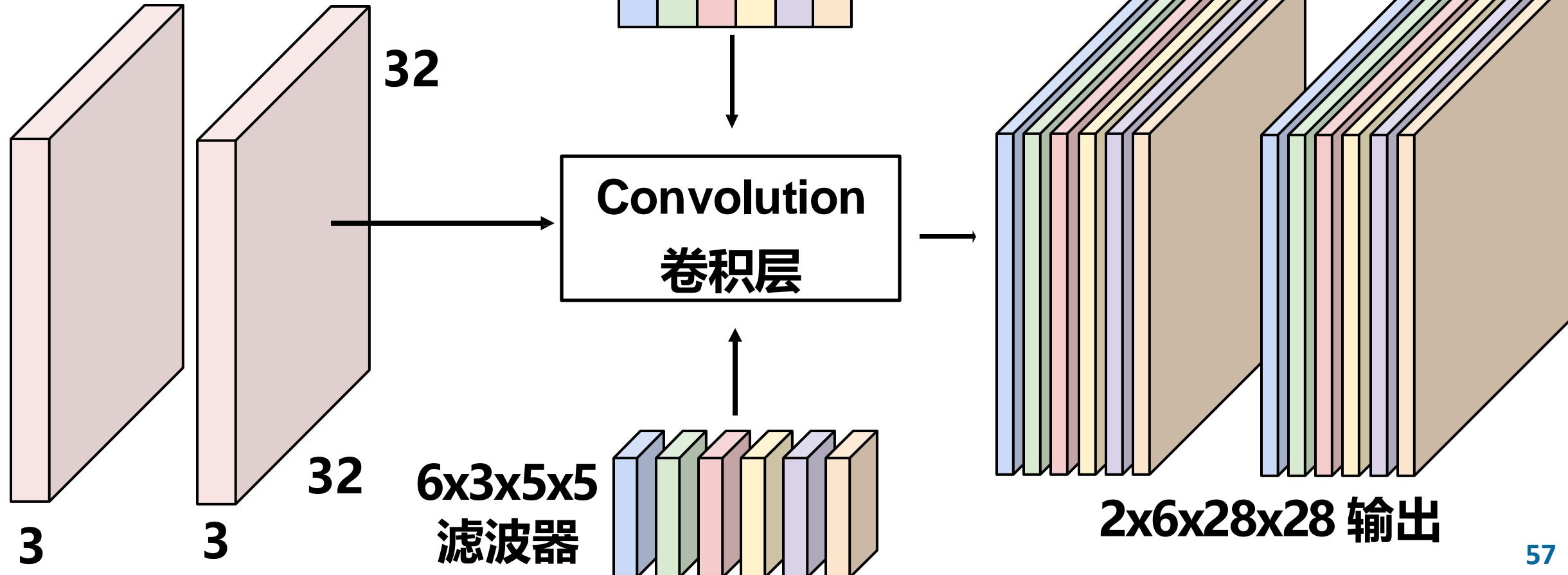
## ■ 卷积层



## ■ 卷积层

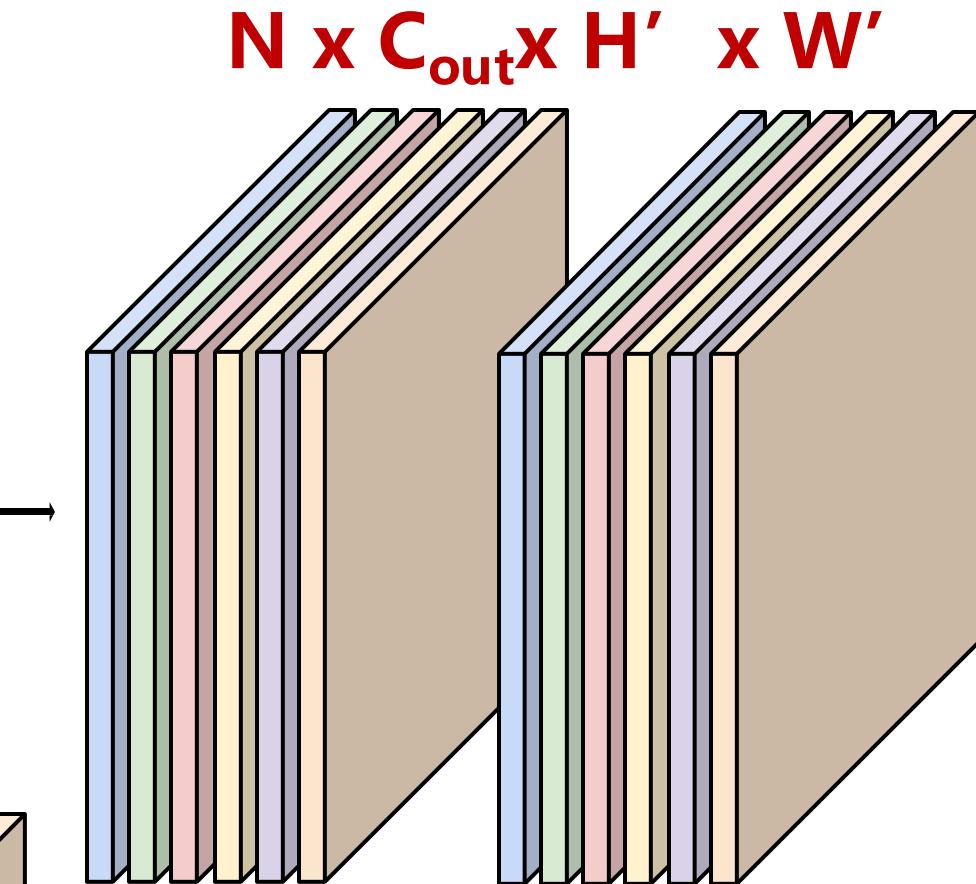
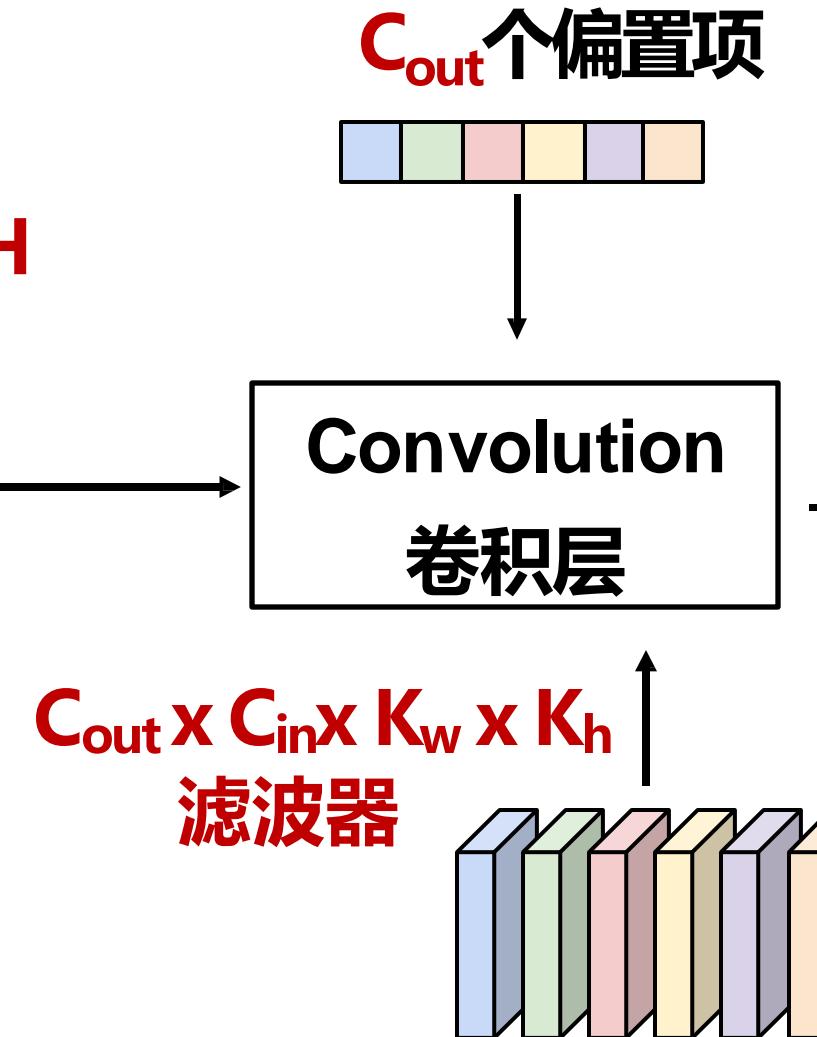
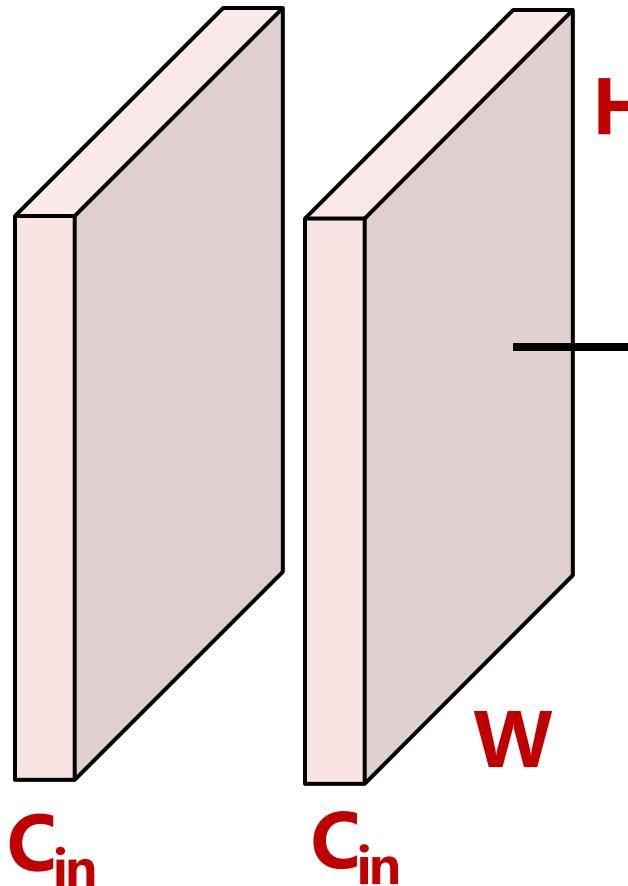


■ 卷积层  
**一个batch里2张图**

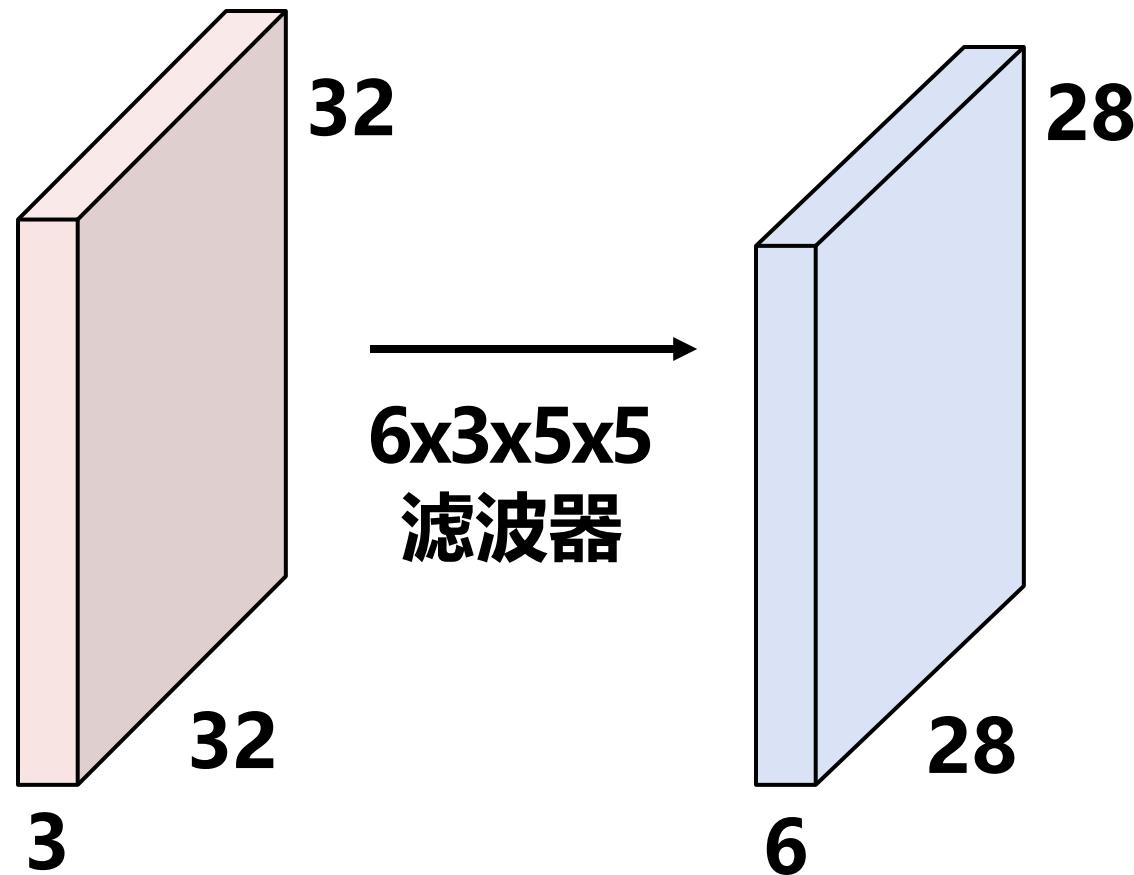


## ■ 卷积层

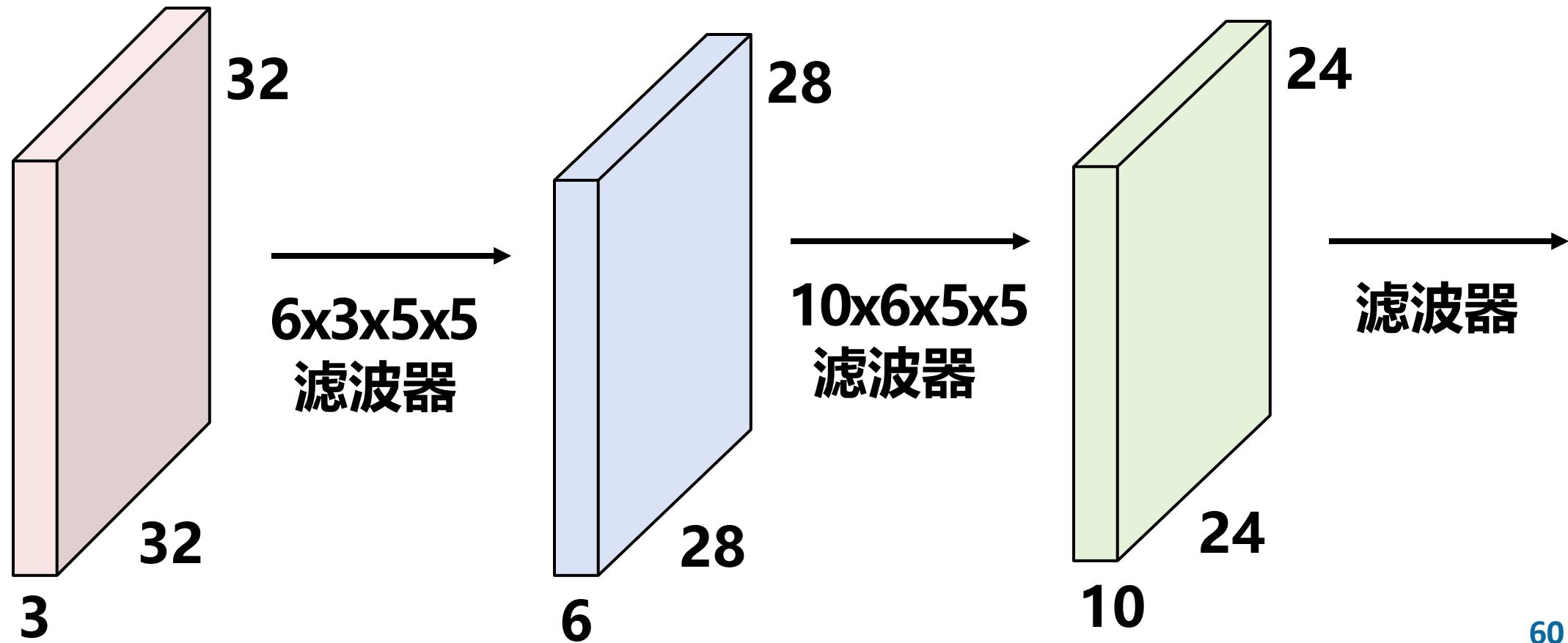
$N \times C_{in} \times H \times W$



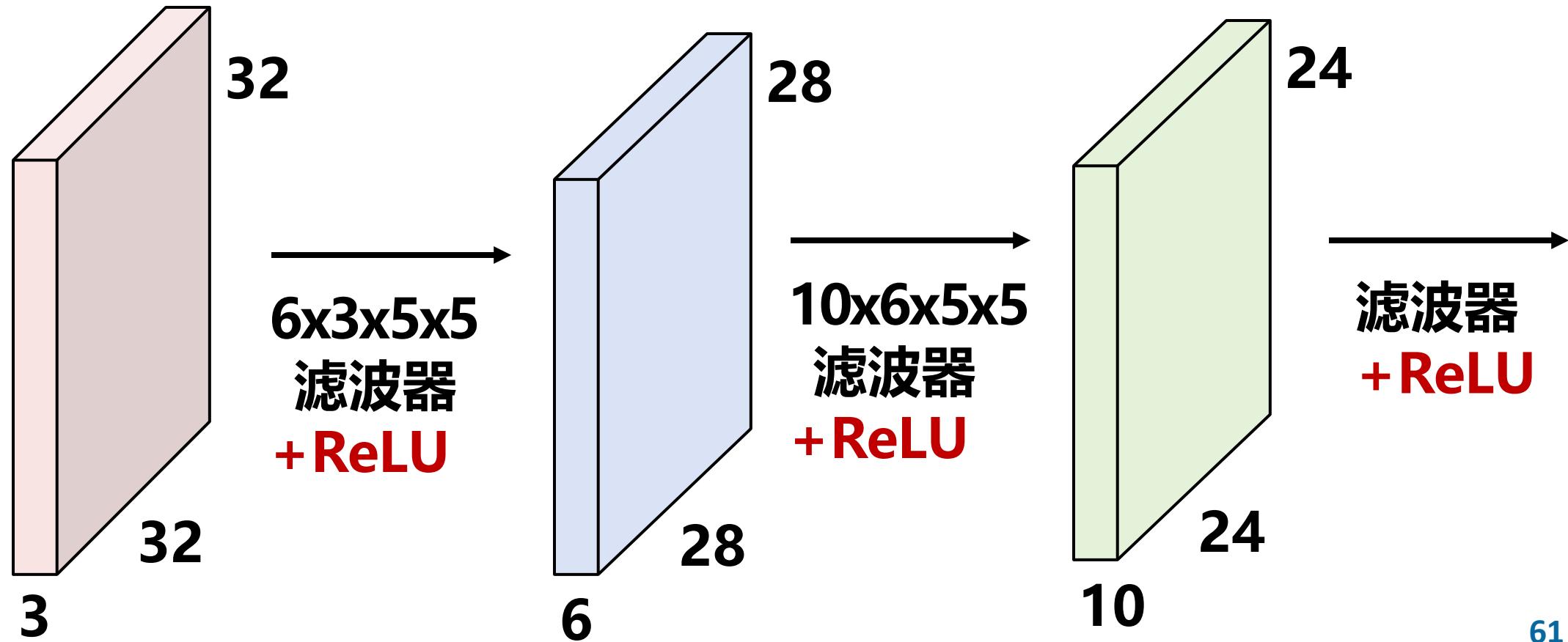
- 卷积神经网络由一系列卷积层组成



- 卷积神经网络由一系列卷积层组成



## ■ 卷积神经网络由一系列卷积层组成



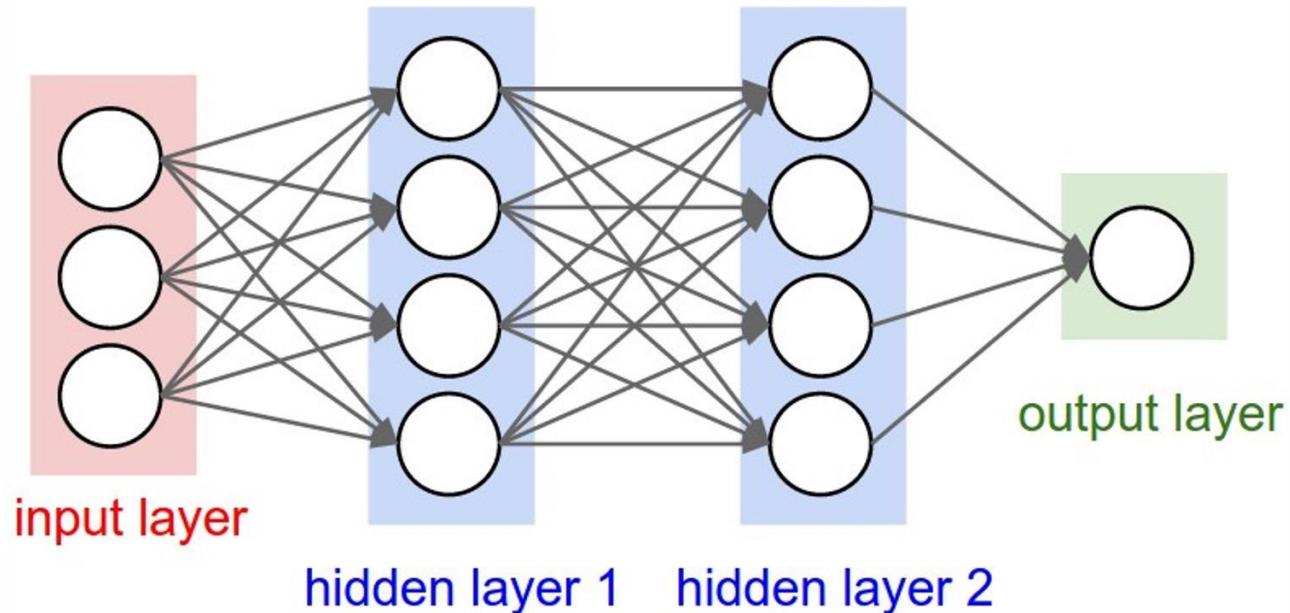
## ■ 线性分类器学到了什么?

$$f(x, W) = Wx + b$$

每个类别的模板



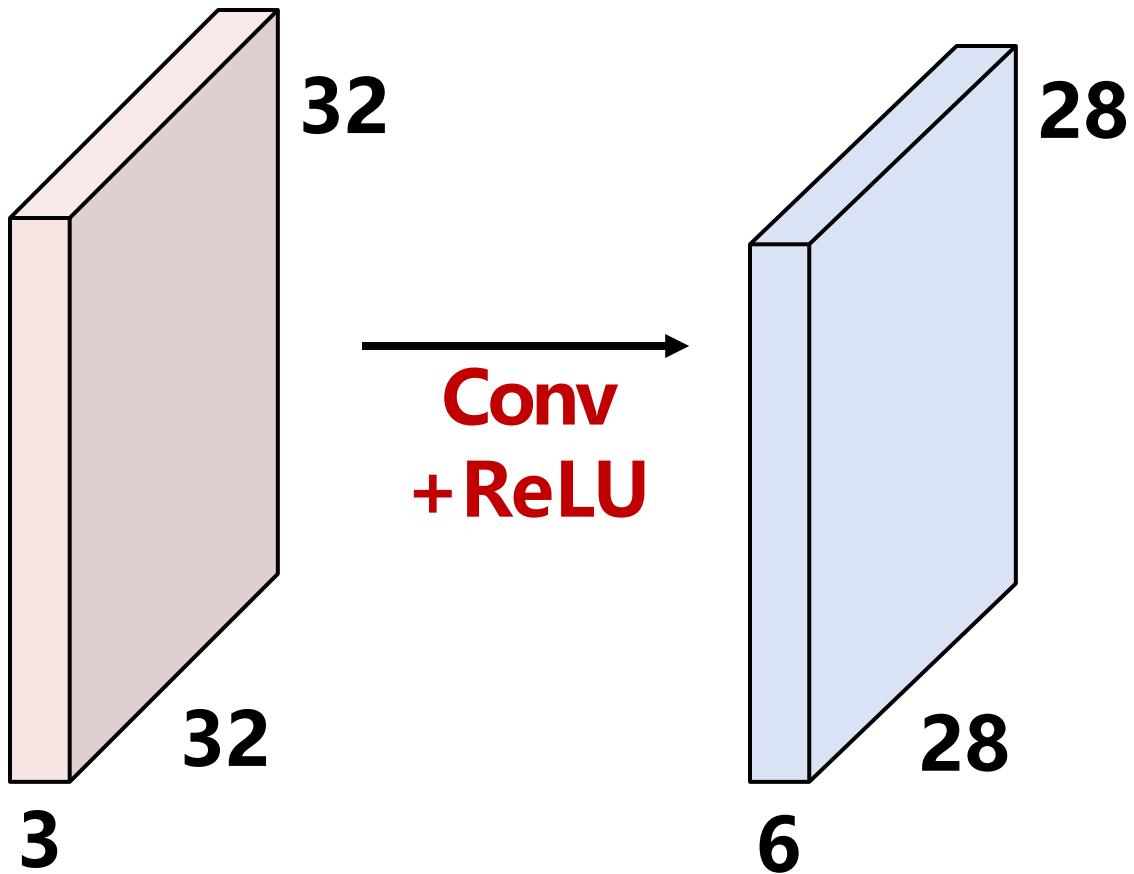
## ■ 全连接神经网络学到了什么？



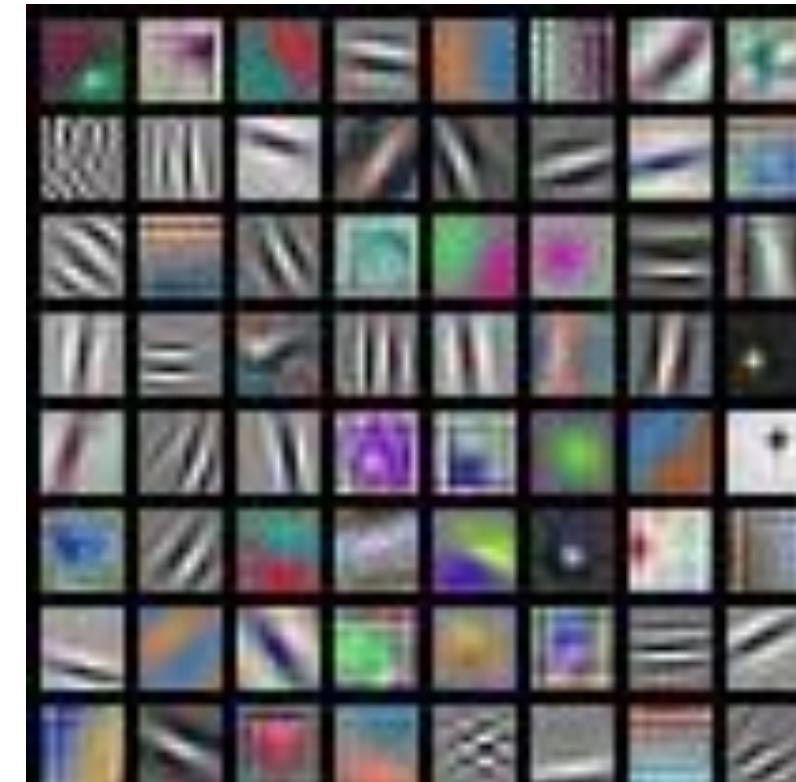
## 整张图像的的模板库



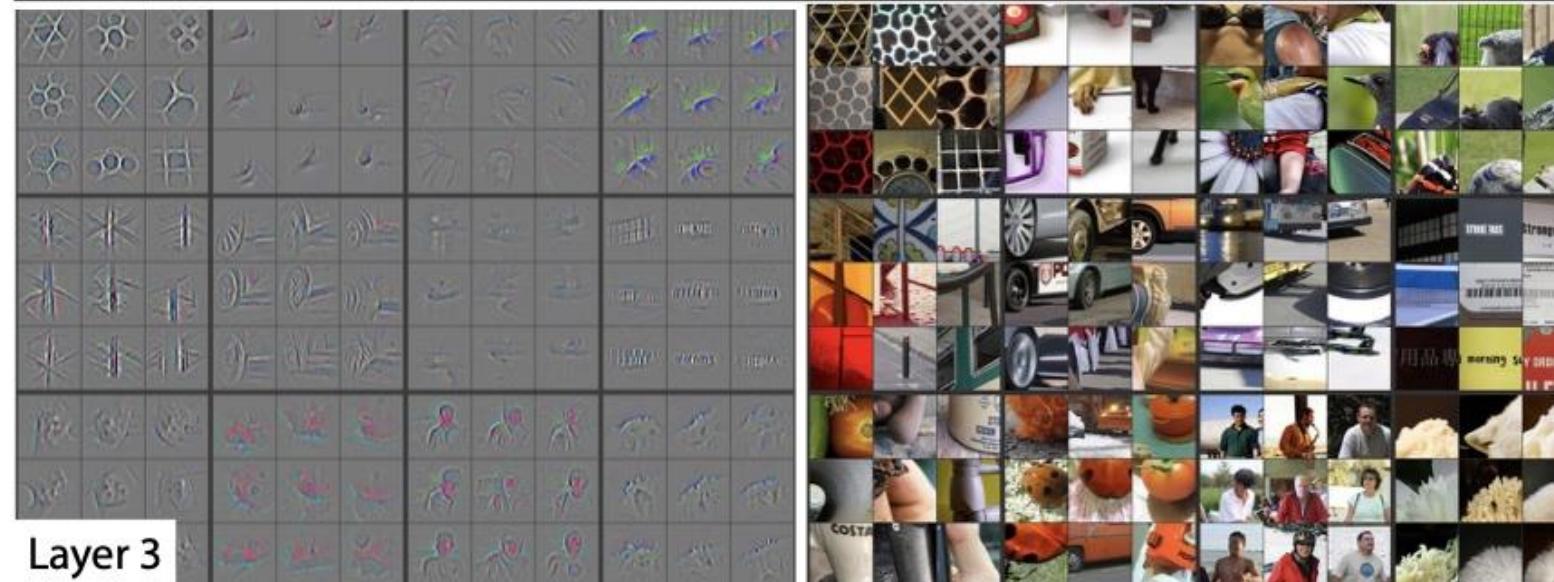
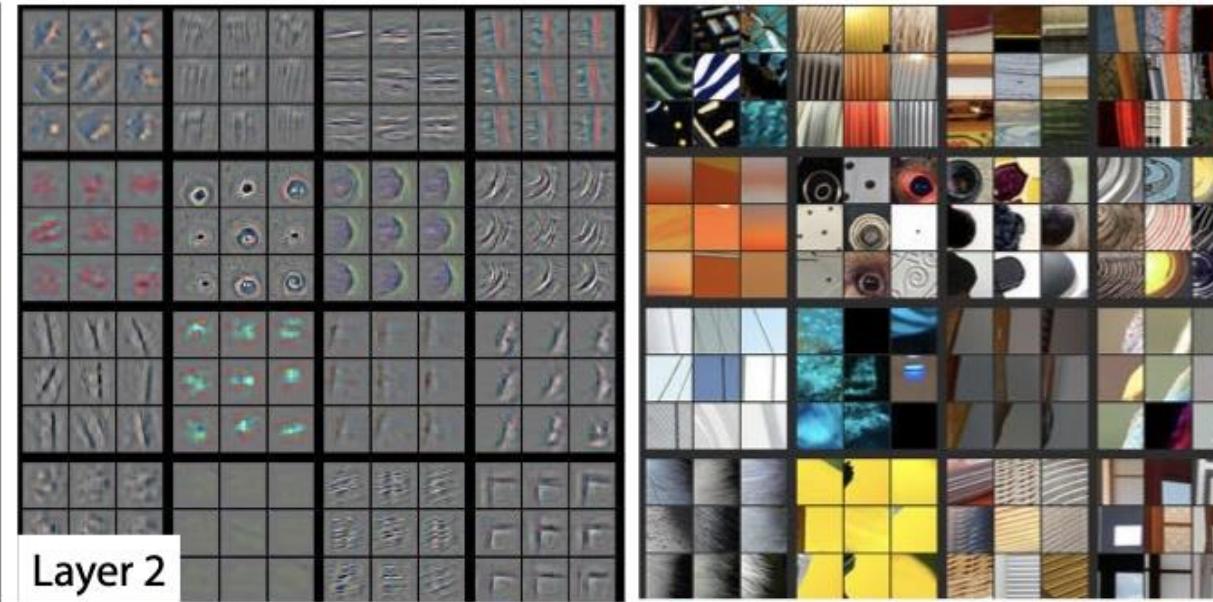
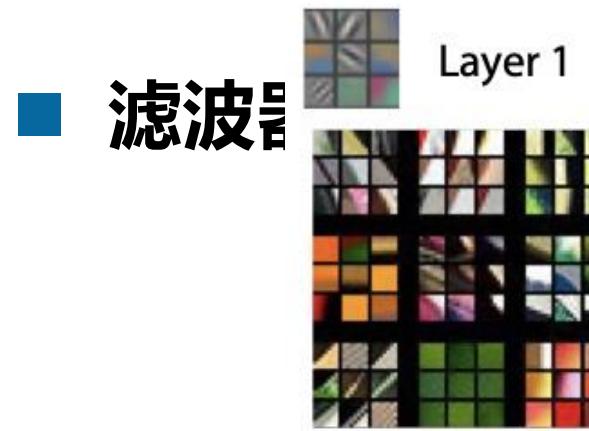
## ■ 滤波器学到了什么？



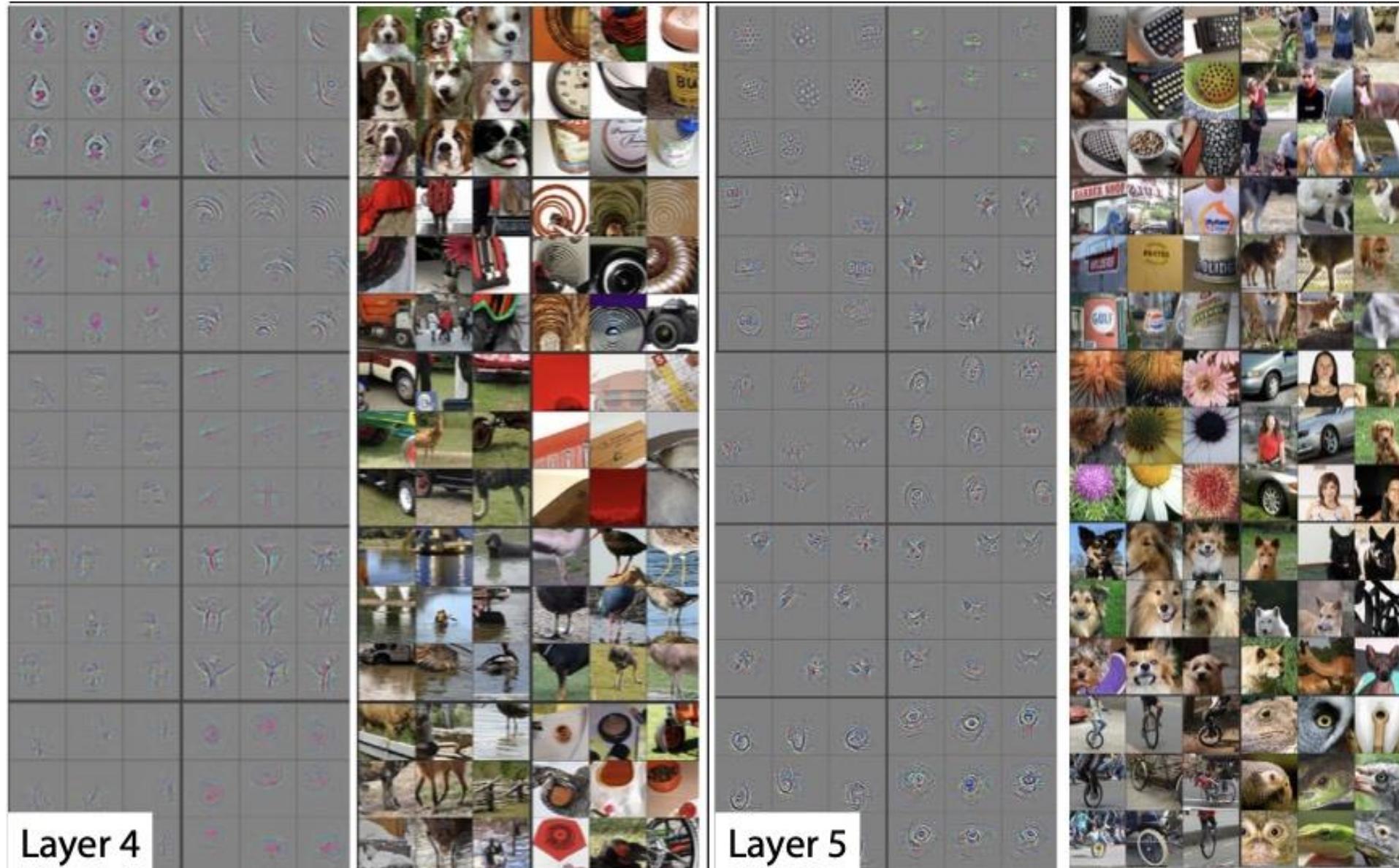
**第一层卷积滤波器：局部图像模板  
(通常学习边缘、颜色等纹理)**



# 卷积神经网络

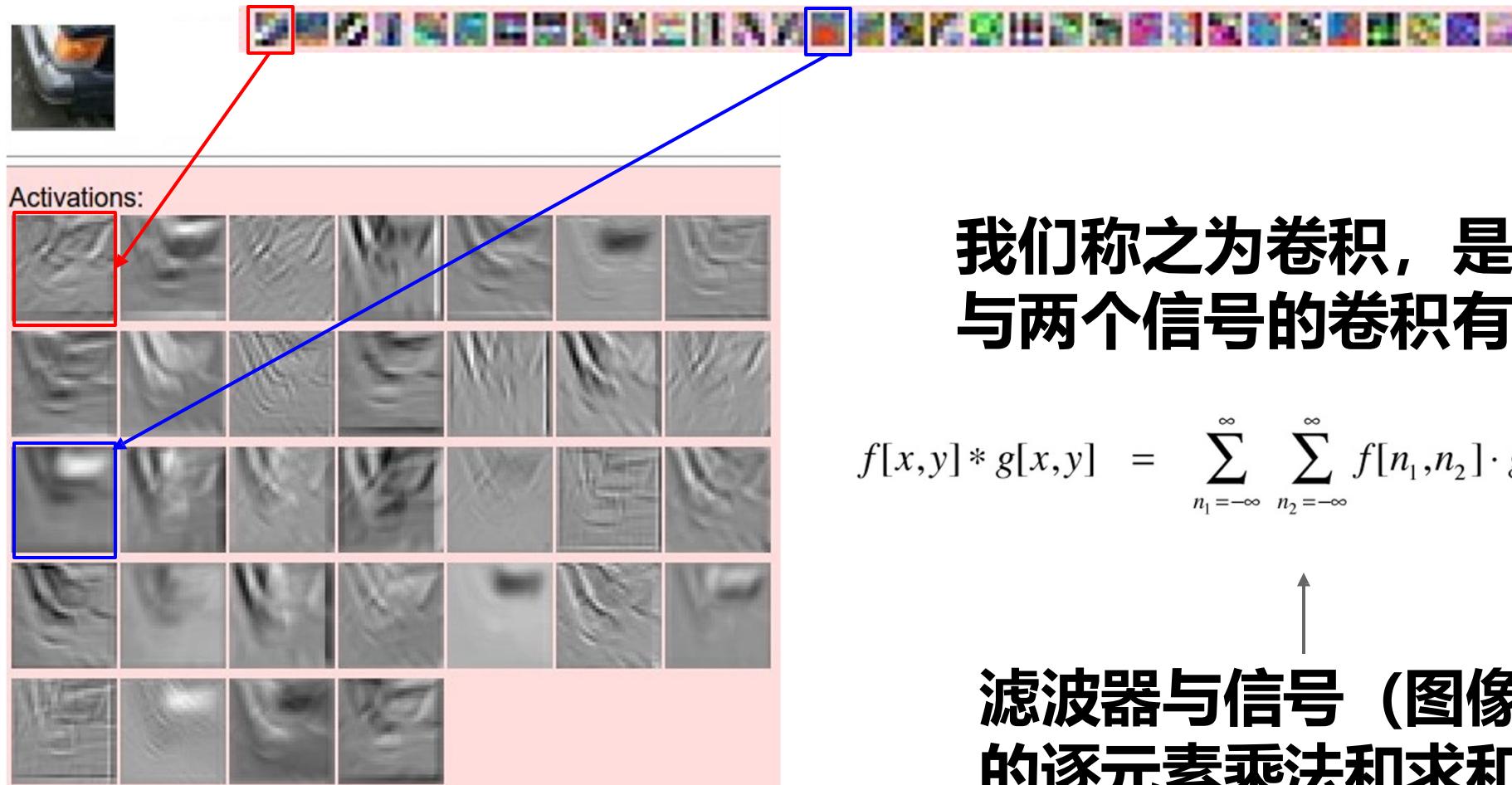


# 卷积神经网络



## ■ 卷积神经网络

一个滤波器-->一个激活图



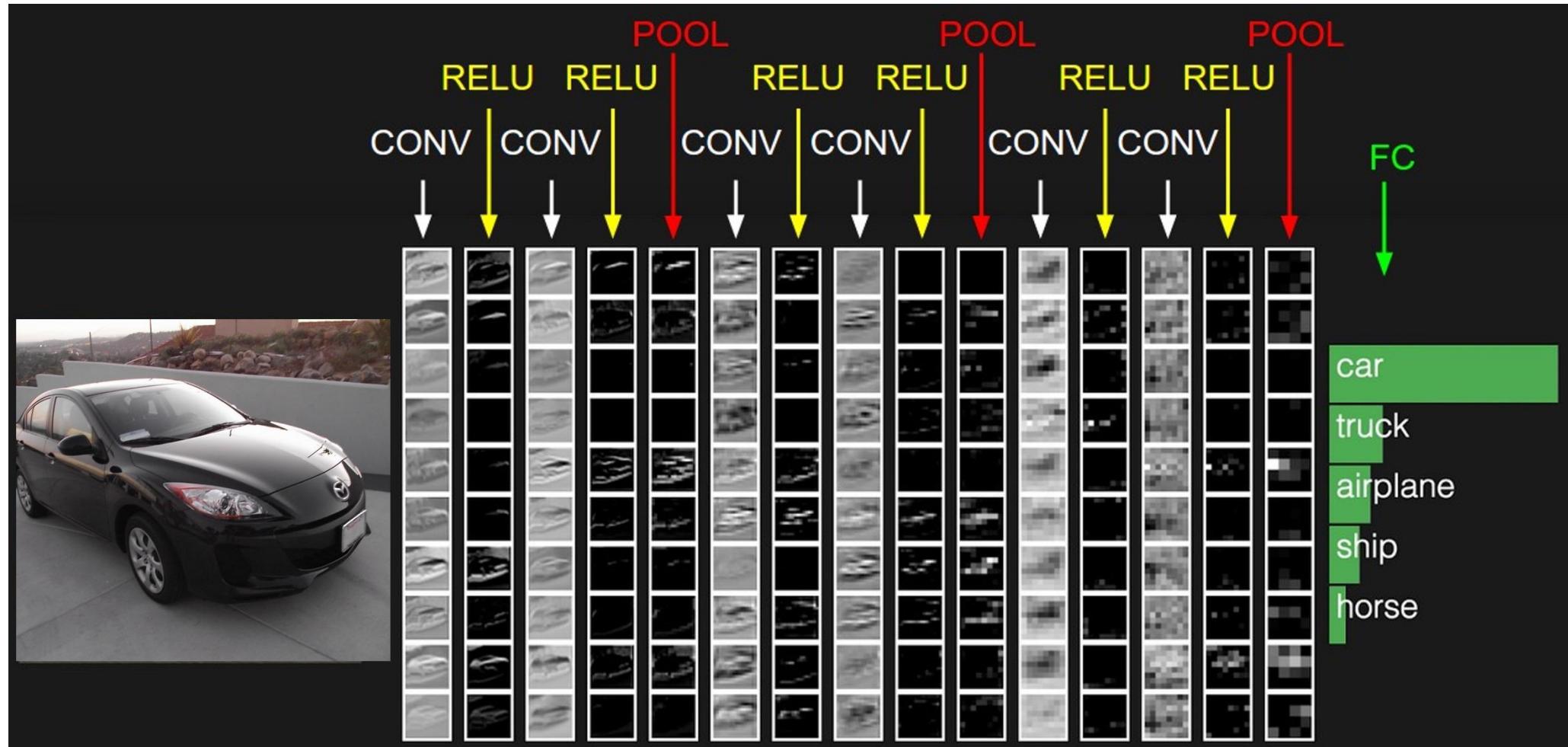
我们称之为卷积，是因为它与两个信号的卷积有关：

$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x - n_1, y - n_2]$$

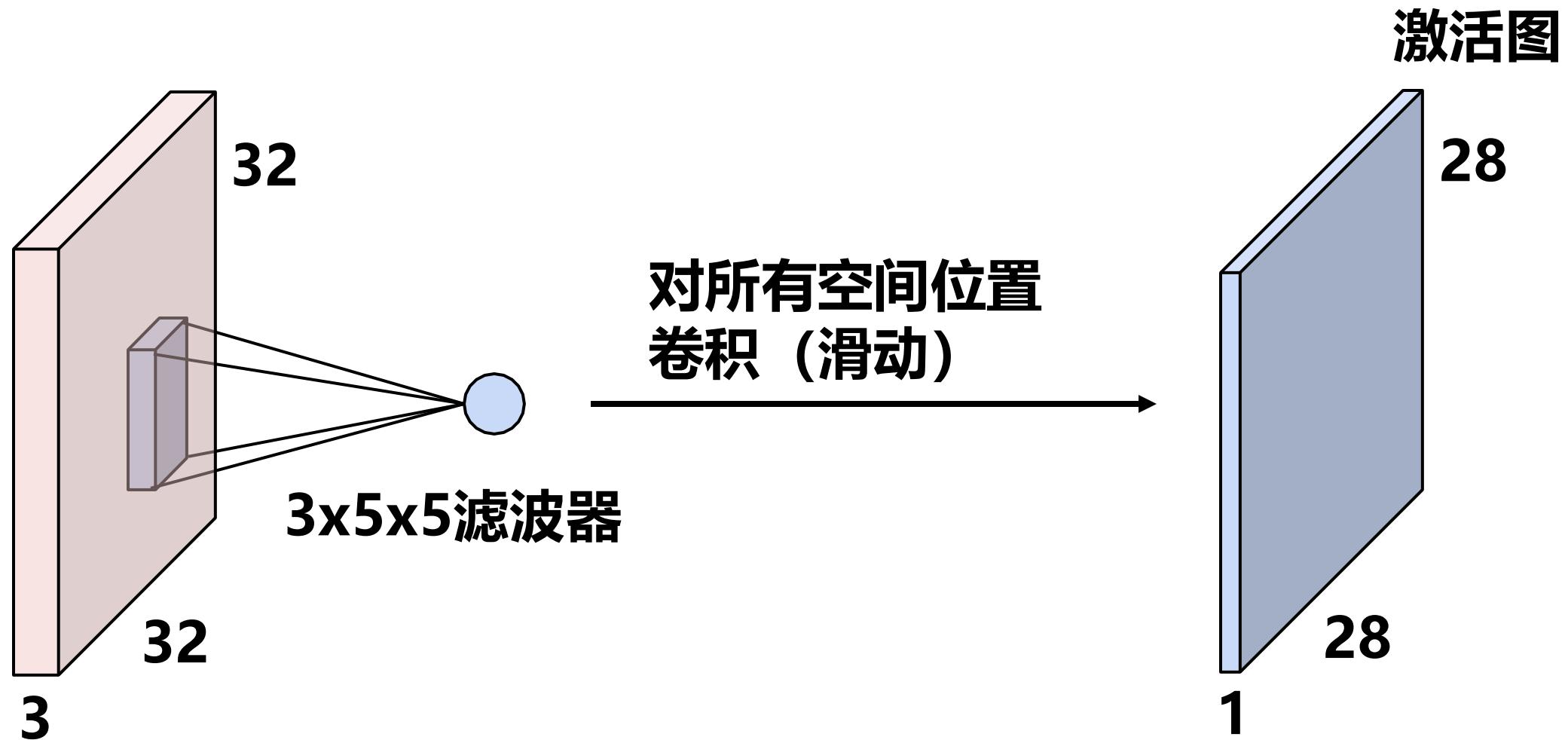


滤波器与信号（图像）的逐元素乘法和求和

## 卷积神经网络

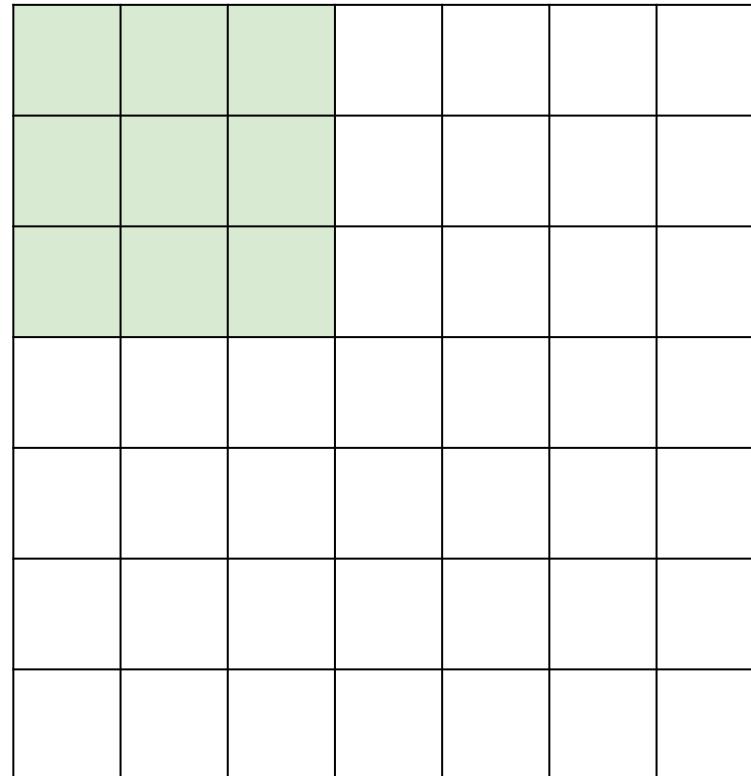


## ■ 卷积输出的空间尺寸



## ■ 卷积输出的空间尺寸

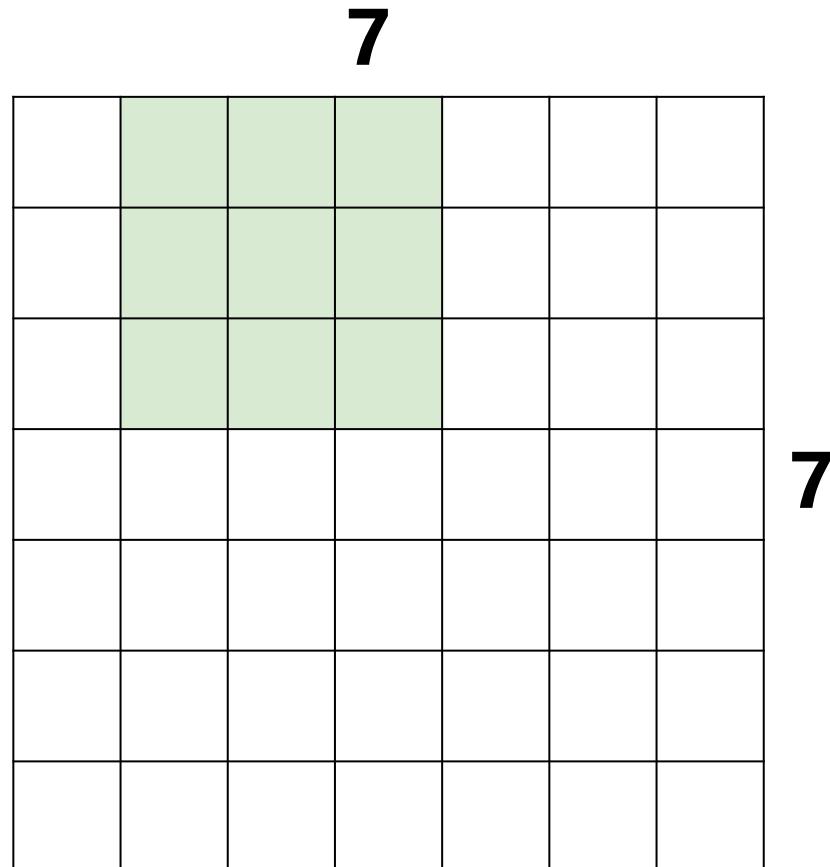
7



7x7输入  
3x3滤波器

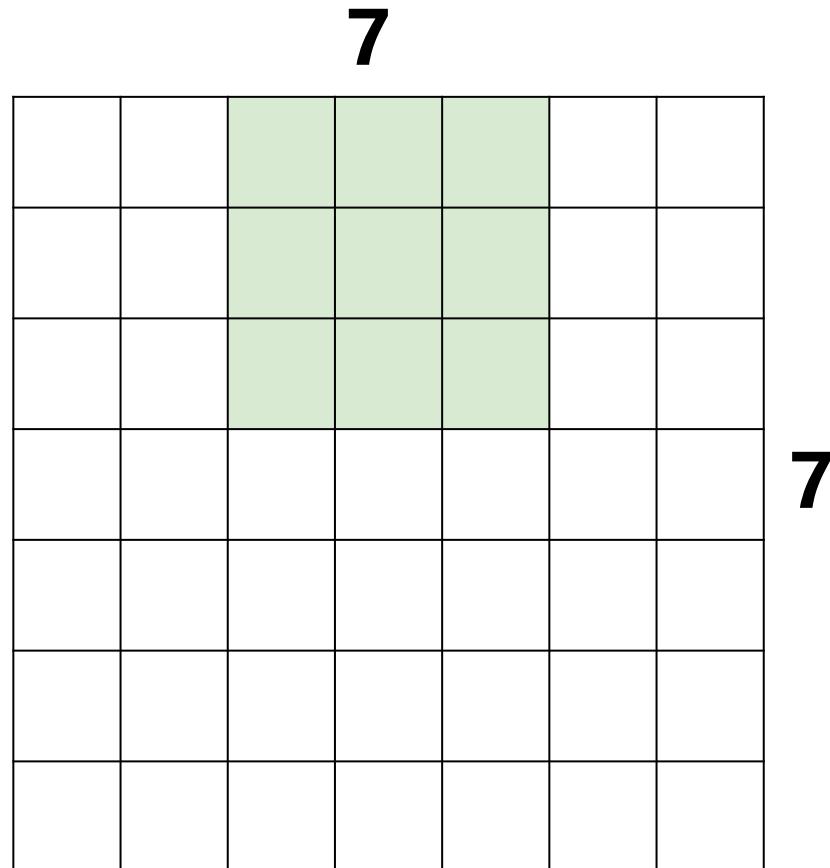
7

## ■ 卷积输出的空间尺寸



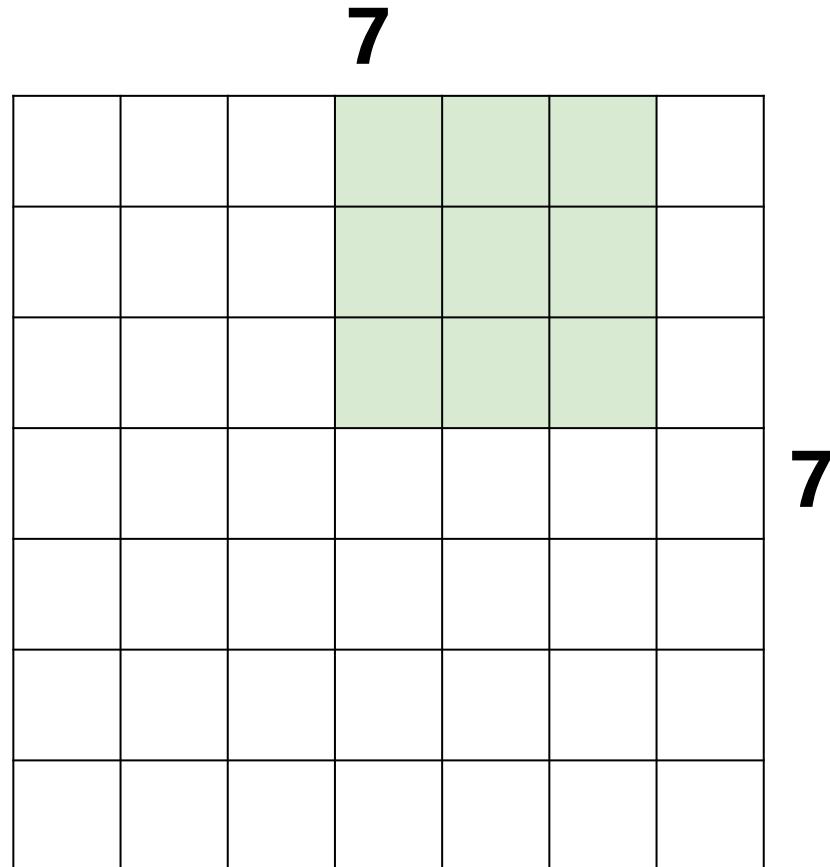
7x7输入  
3x3濾波器

## ■ 卷积输出的空间尺寸



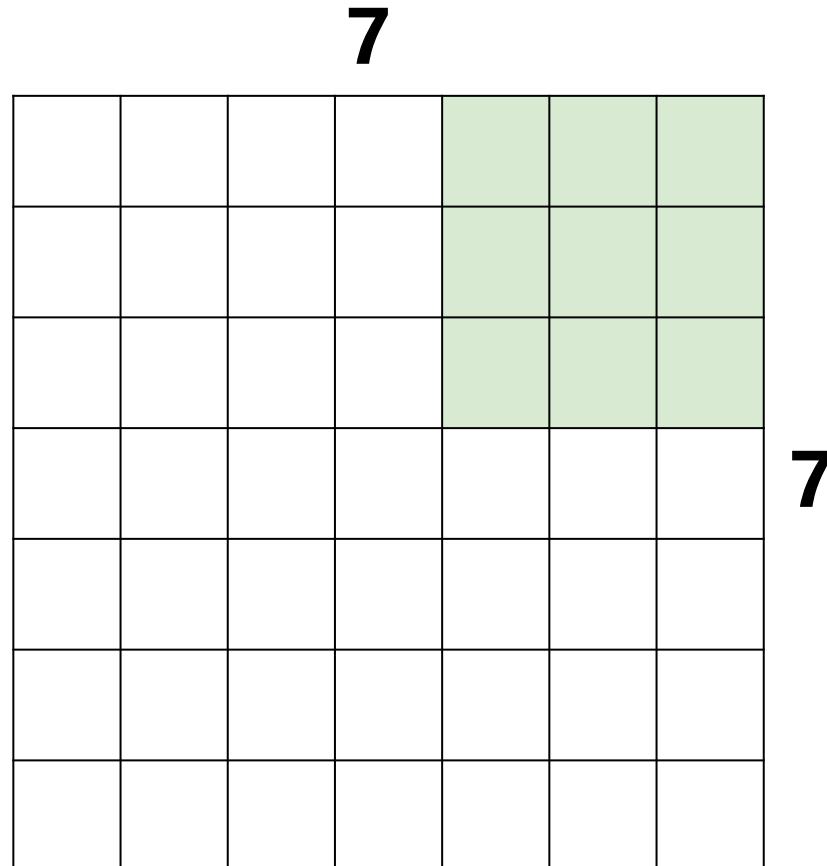
7x7输入  
3x3濾波器

## ■ 卷积输出的空间尺寸



7x7输入  
3x3滤波器

## ■ 卷积输出的空间尺寸

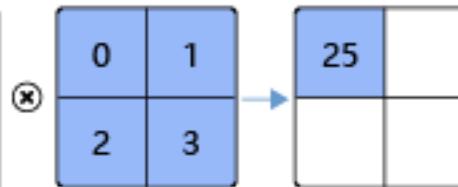


7x7输入  
3x3濾波器

5x5输出

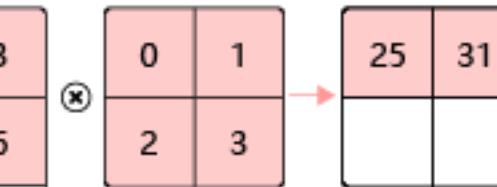
## ■ 卷积层

1	2	3
4	5	6
7	8	9



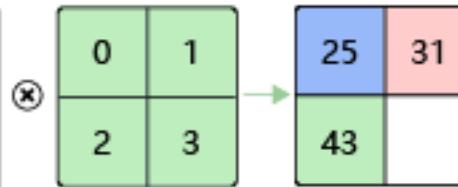
$$(a) 0 \times 1 + 1 \times 2 + 2 \times 4 + 3 \times 5 = 25$$

1	2	3
4	5	6
7	8	9



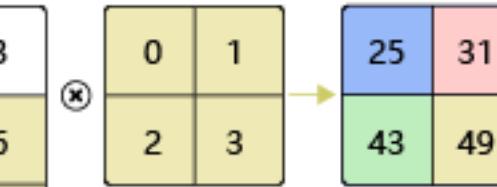
$$(b) 0 \times 2 + 1 \times 3 + 2 \times 5 + 3 \times 6 = 31$$

1	2	3
4	5	6
7	8	9



$$(c) 0 \times 4 + 1 \times 5 + 2 \times 7 + 3 \times 8 = 43$$

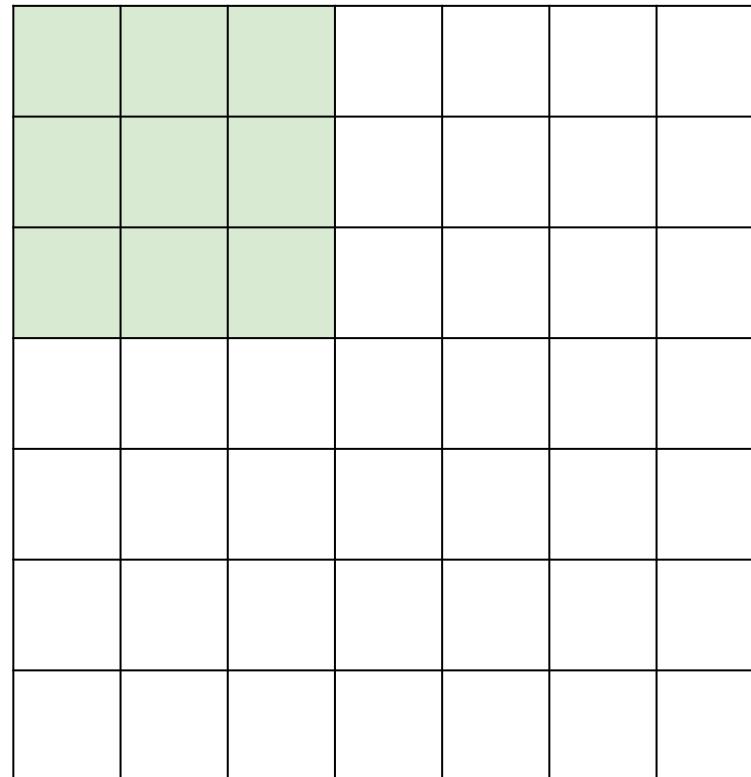
1	2	3
4	5	6
7	8	9



$$(d) 0 \times 5 + 1 \times 6 + 2 \times 8 + 3 \times 9 = 49$$

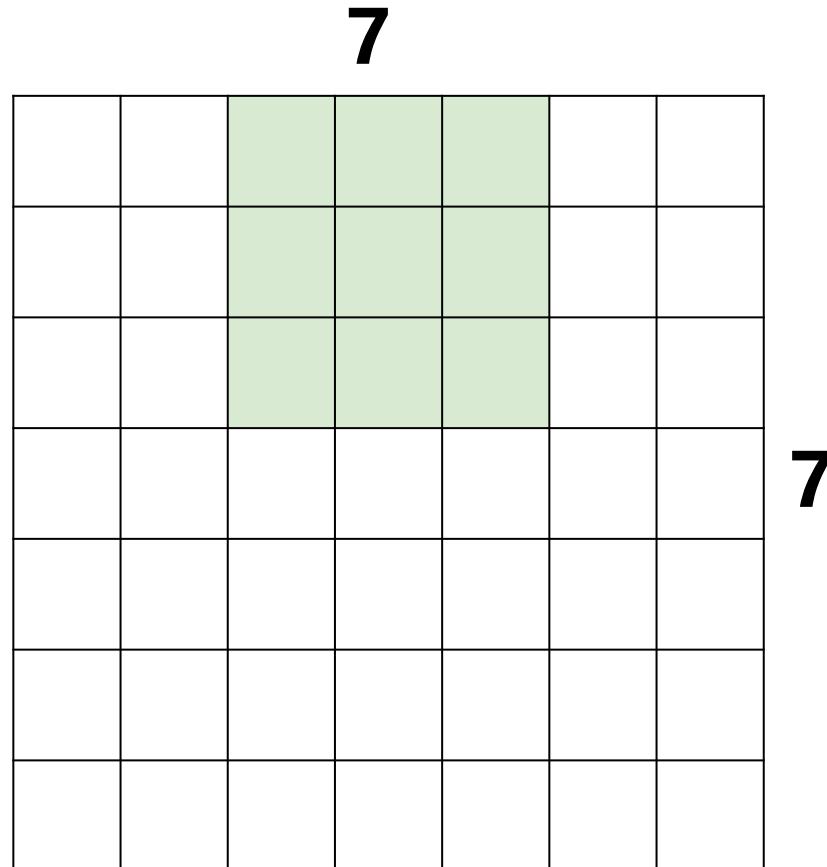
## ■ 卷积输出的空间尺寸

7



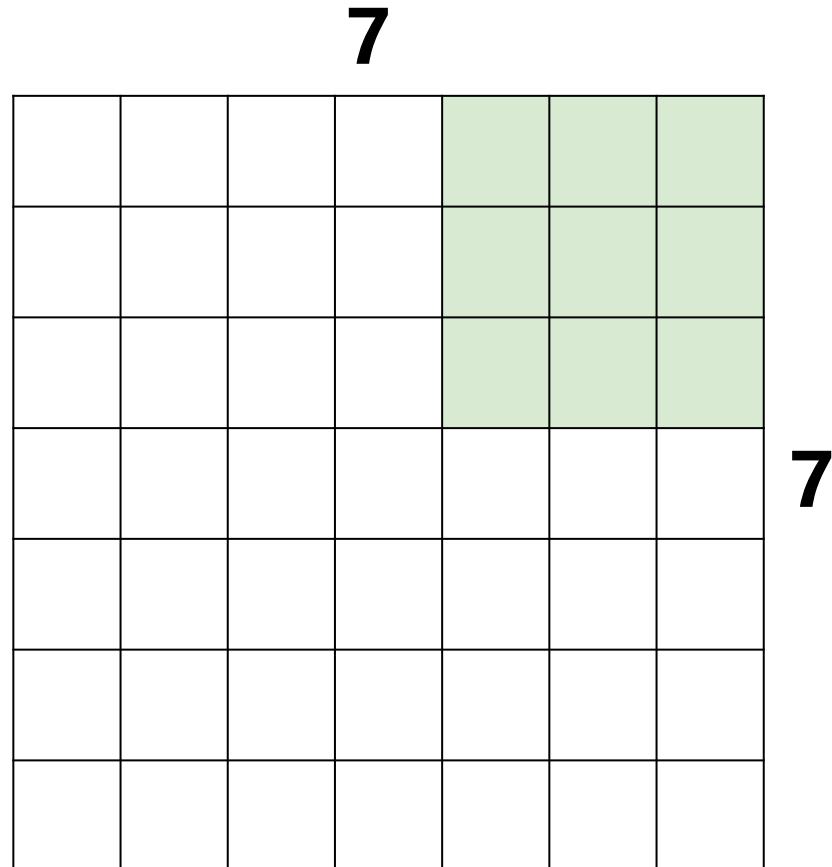
7x7输入  
3x3滤波器  
2步长

## ■ 卷积输出的空间尺寸



**7x7输入  
3x3滤波器  
2步长**

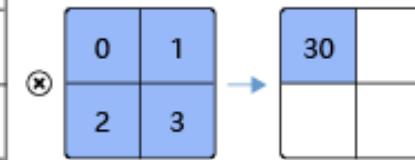
## ■ 卷积输出的空间尺寸



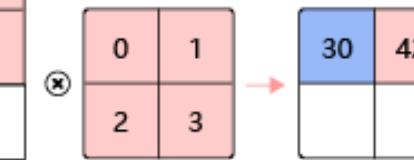
7x7输入  
3x3濾波器  
2步长  
3x3输出

## ■ 卷积层

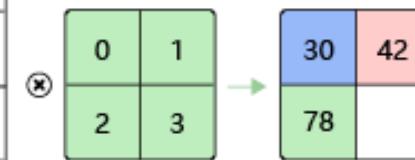
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



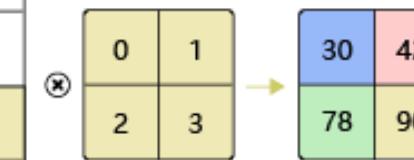
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

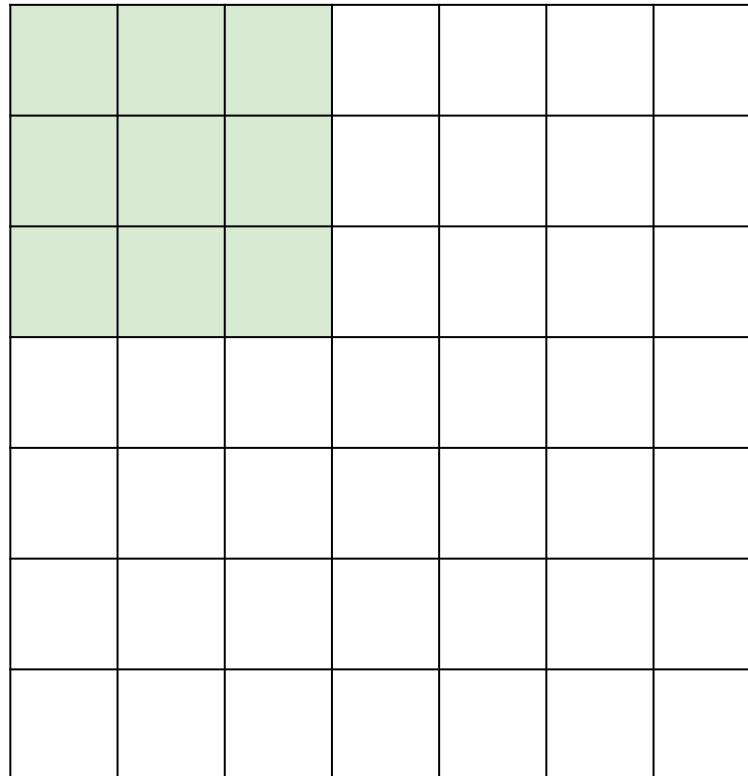


1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



## ■ 卷积输出的空间尺寸

7

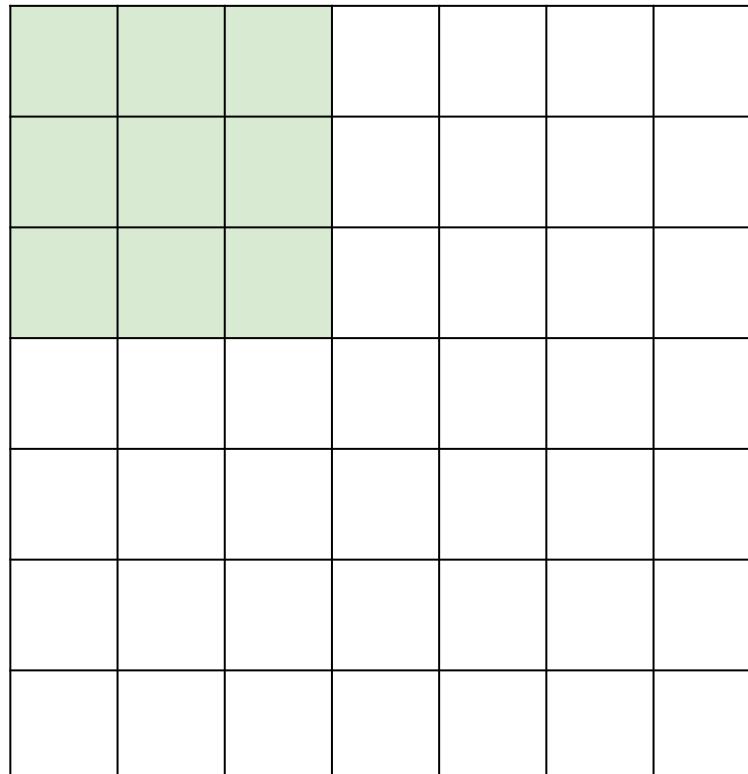


7

7x7输入  
3x3滤波器  
3步长?

不可以!  
尺寸不匹配

## ■ 卷积输出的空间尺寸



**NxN 输入  
FxF 滤波器  
S 步长  
输出大小:  $(N-F)/S+1$**

**例如:  $N = 7, F = 3:$**   
**步长 1 =>  $(7 - 3)/1 + 1 = 5$**   
**步长 2 =>  $(7 - 3)/2 + 1 = 3$**   
**步长 3 =>  $(7 - 3)/3 + 1 = 2.33$**

## ■ 卷积输出的空间尺寸

0	0	0	0	0	0			
0								
0								
0								
0								

7x7 输入 NxN  
3x3 滤波器 FxF  
1 步长 S  
1 Padding P  
输出尺寸?

## ■ 卷积输出的空间尺寸

0	0	0	0	0	0
0	13	14	15	16	0
0	9	10	11	12	0
0	5	6	7	8	0
0	1	2	3	4	0
0	0	0	0	0	0

(a)padding=1

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	13	14	15	16	0	0
0	0	9	10	11	12	0	0
0	0	5	6	7	8	0	0
0	0	1	2	3	4	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

(b)padding=2

## ■ 卷积输出的空间尺寸

0	0	0	0	0	0			
0								
0								
0								
0								

7x7 输入 NxN  
3x3 滤波器 FxF  
1 步长 S  
1 Padding P  
7x7 输出

输出:  $(N+2P-F)/S+1$

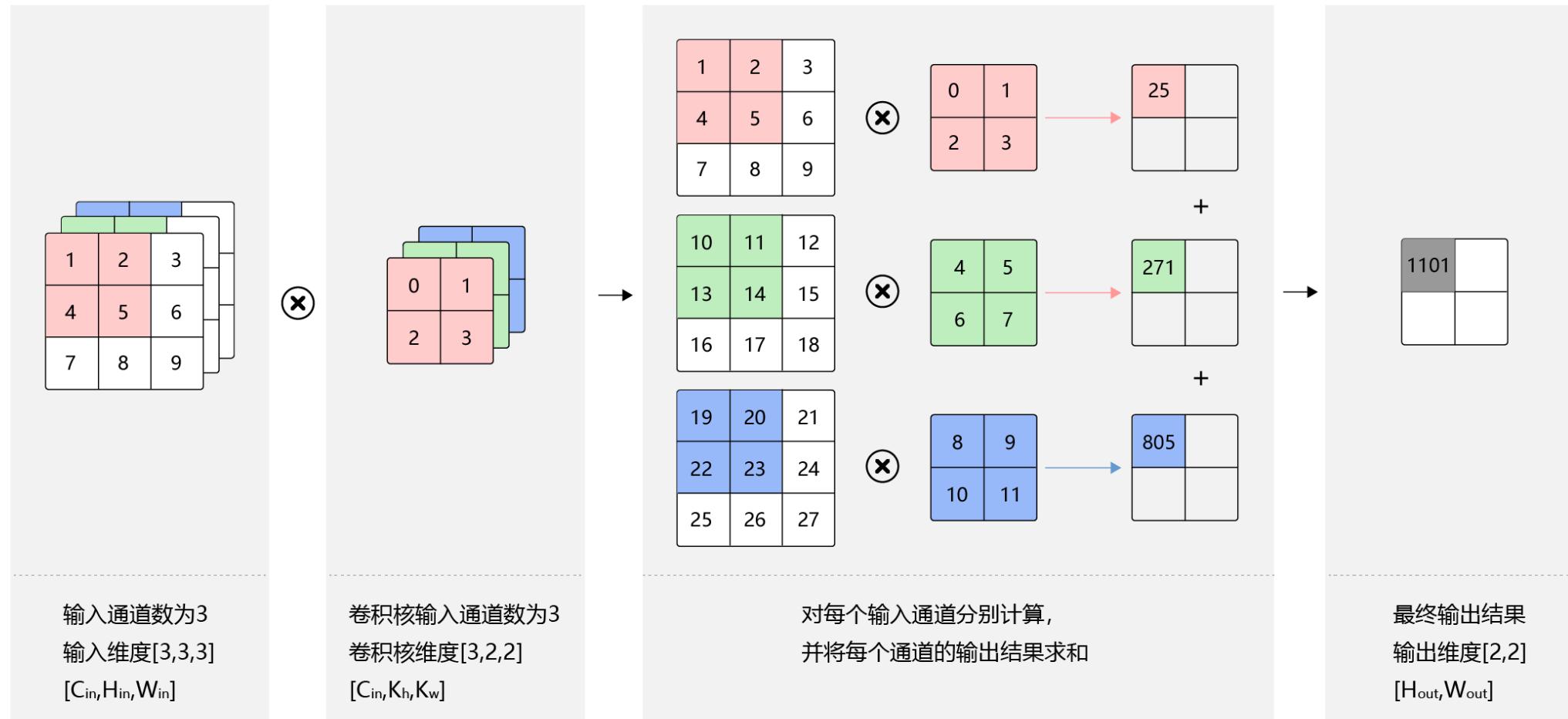
## ■ 卷积输出的空间尺寸

0	0	0	0	0	0			
0								
0								
0								
0								

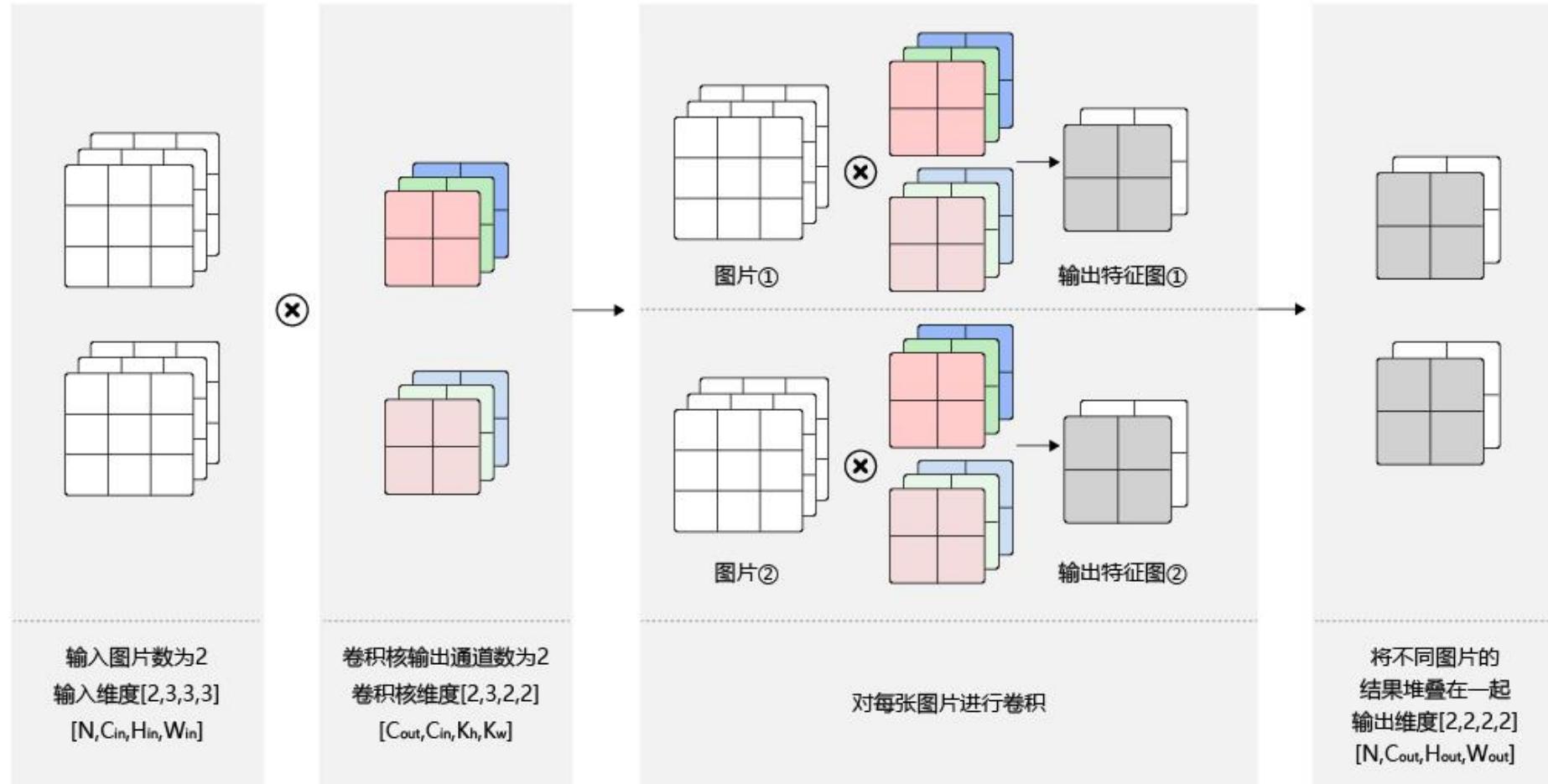
为了保持输出的空间尺寸大小：  
**CONV层步幅为1**  
**滤波器大小为 $F \times F$**   
**零填充为  $(F-1) / 2$**

**$F = 3 \Rightarrow$  zero pad with 1**  
 **$F = 5 \Rightarrow$  zero pad with 2**  
 **$F = 7 \Rightarrow$  zero pad with 3**

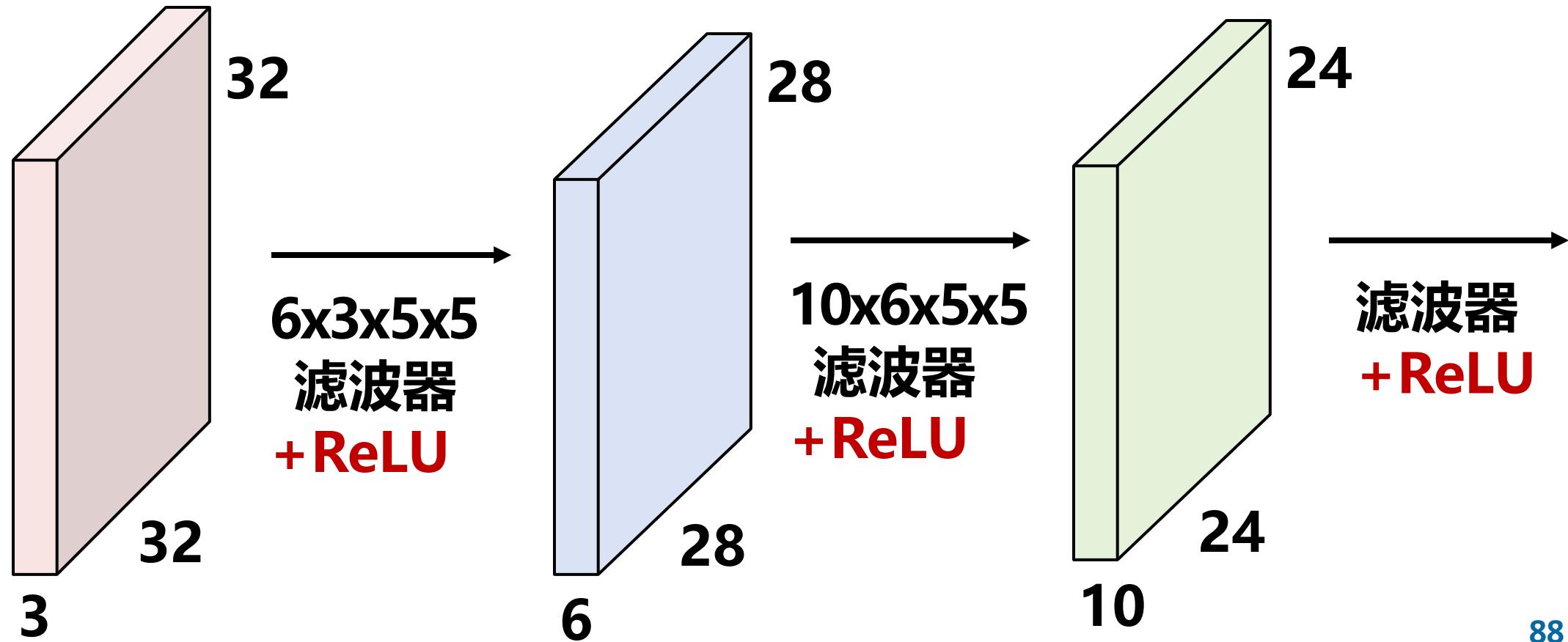
## ■ 多输入通道卷积



## ■ 多输出通道卷积



■ 激活图尺寸:  $32 \times 32 \rightarrow 28 \times 28 \rightarrow 24 \times 24$



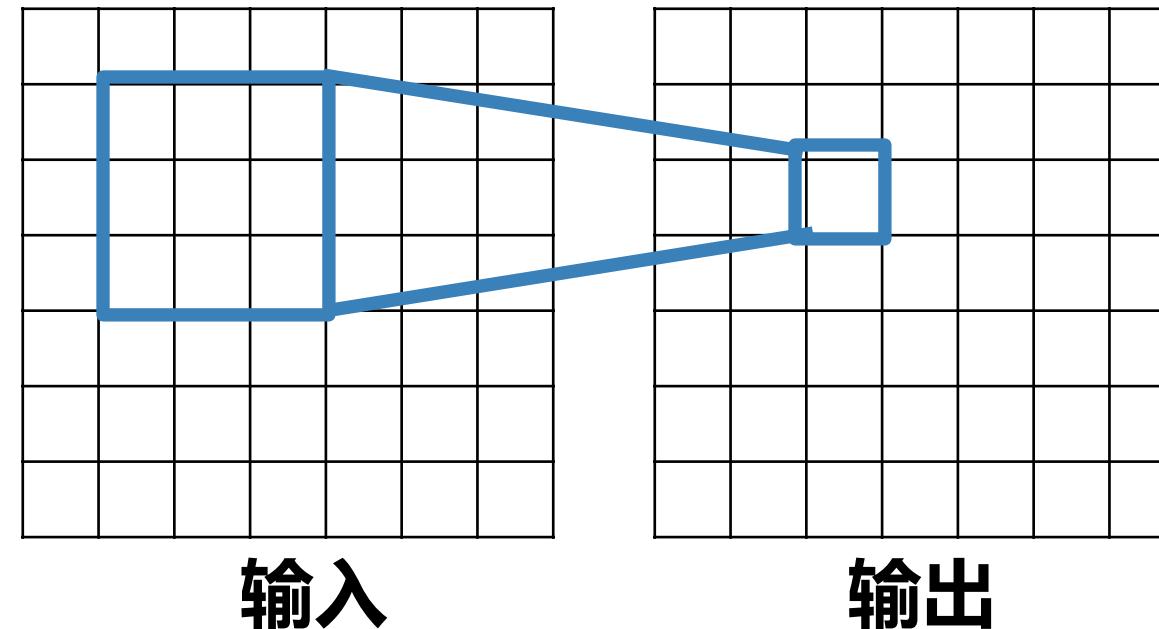
- 卷积输出的空间尺寸
- 假设：
  - 输入：  $3 \times 32 \times 32$
  - 滤波器：  $10 \times 5 \times 5$
  - 步长： 1
  - Padding： 2
- 输出多少？

- 卷积输出的空间尺寸
- 假设：
  - 输入：  $3 \times 32 \times 32$
  - 滤波器：  $10 \times 5 \times 5$
  - 步长： 1
  - Padding： 2
  - 输出：  $32 \times 32 \times 10$ , 即  $(32 + 2 * 2 - 5) / 1 + 1 = 32$

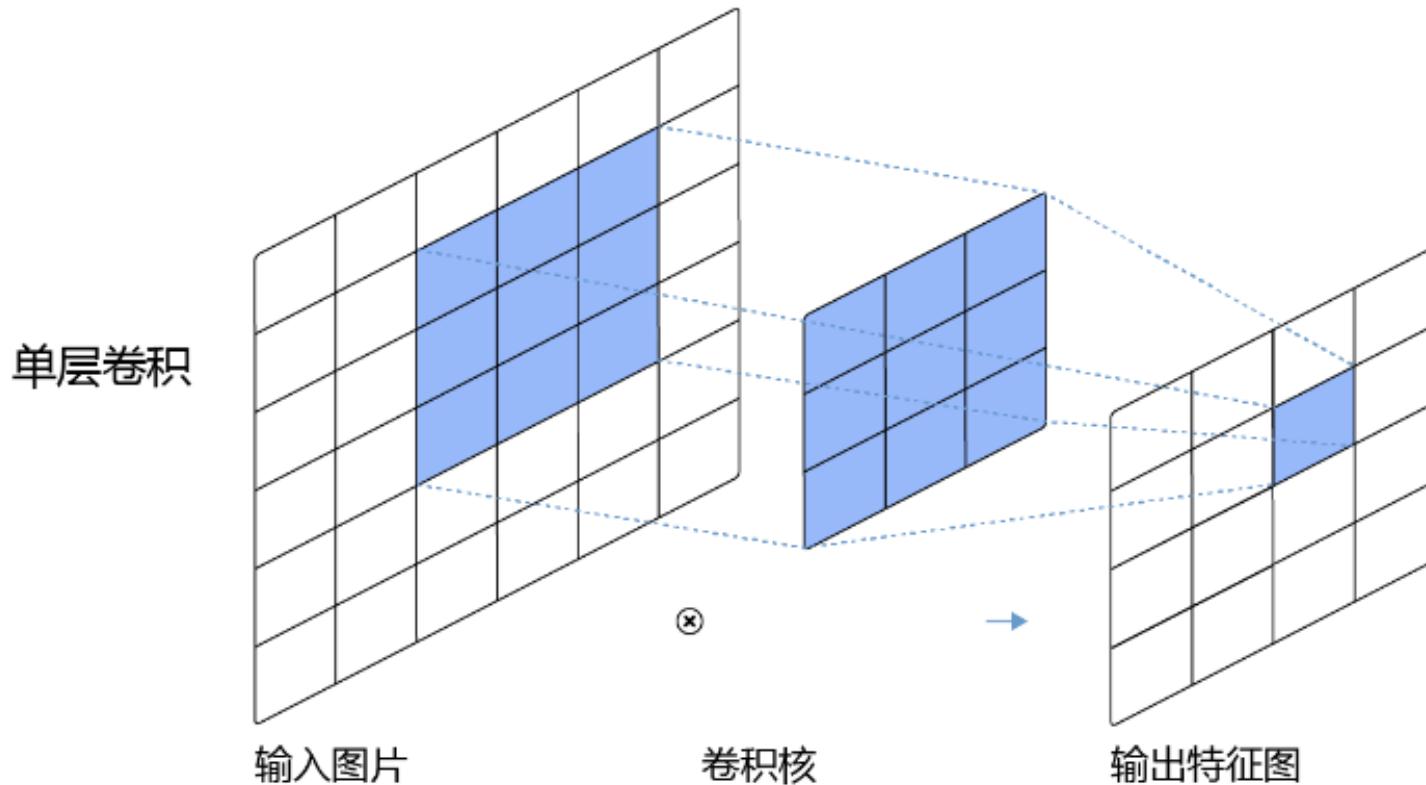
- 卷积输出的空间尺寸
- 假设：
  - 输入：  $3 \times 32 \times 32$
  - 滤波器：  $10 \times 5 \times 5$
  - 步长： 1
  - Padding： 2
- 参数量？

- 卷积输出的空间尺寸
- 假设：
  - 输入：  $3 \times 32 \times 32$
  - 滤波器：  $10 \times 5 \times 5$
  - 步长： 1
  - Padding： 2
- 参数量：  $(5 \times 5 \times 3 + 1) \times 10$

- 卷积层的感受野
- 对于核大小为K的卷积，输出中的每个元素都取决于输入中的 $K \times K$ 感受野

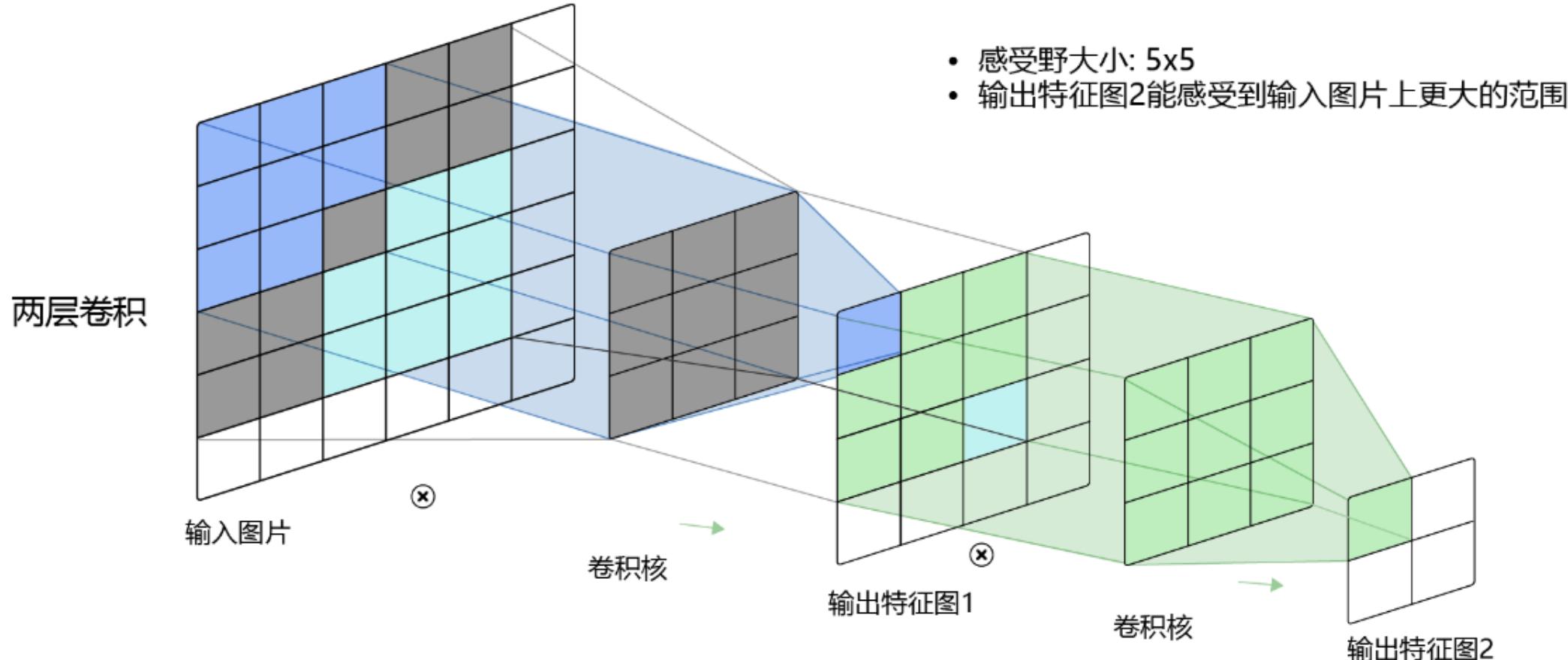


## ■ 卷积层的感受野



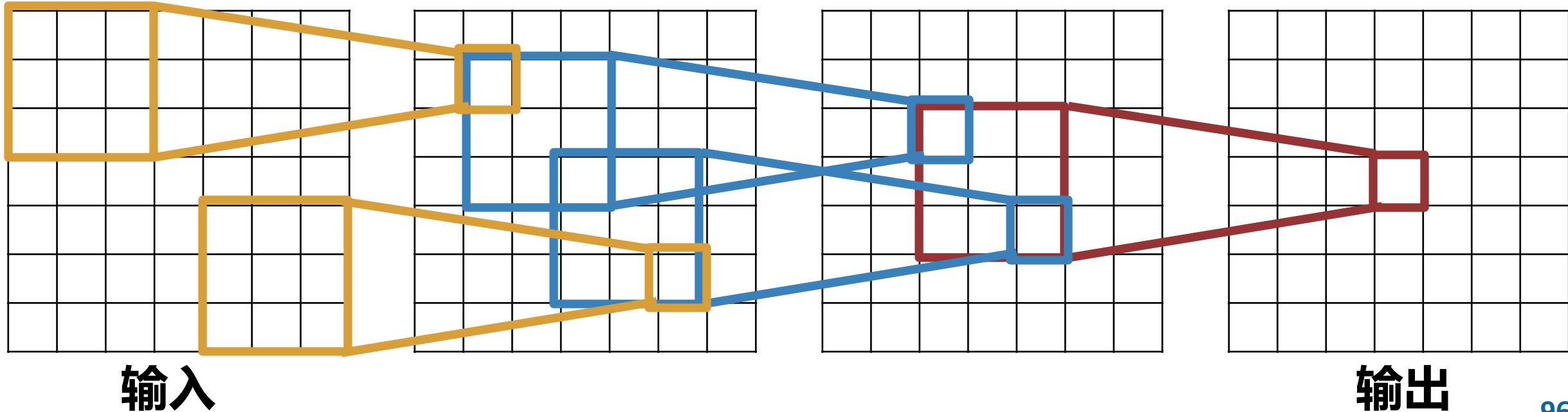
- 感受野大小:  $3 \times 3$
- 输出特征图上的像素点所能感受到的输入数据的范围

## ■ 卷积层的感受野



- 卷积层的感受野
- 每次连续的卷积都会使感受野大小增加 $K-1$ 。  
对于L层，感受野大小为 $1 + L * (K-1)$

注意区分：  
在前一层上的感受野  
在输入图像上的感受野



- 卷积层的感受野
- 问题：对于大图像，我们需要使用许多卷积层，才能“看到”整个图像
- 解决方案：下采样
  - 大步长
  - 池化层

## ■ 卷积层总结

■ 输入:  $W_1 \times H_1 \times C$

■ 卷积超参:

■ 滤波器数量  $K$

■ 滤波器大小  $F$

■ 步长  $S$

■ Padding  $P$

■ 输出:  $W_2 \times H_2 \times C$

■  $W_2 = (W_1 - F + 2P)/S + 1$

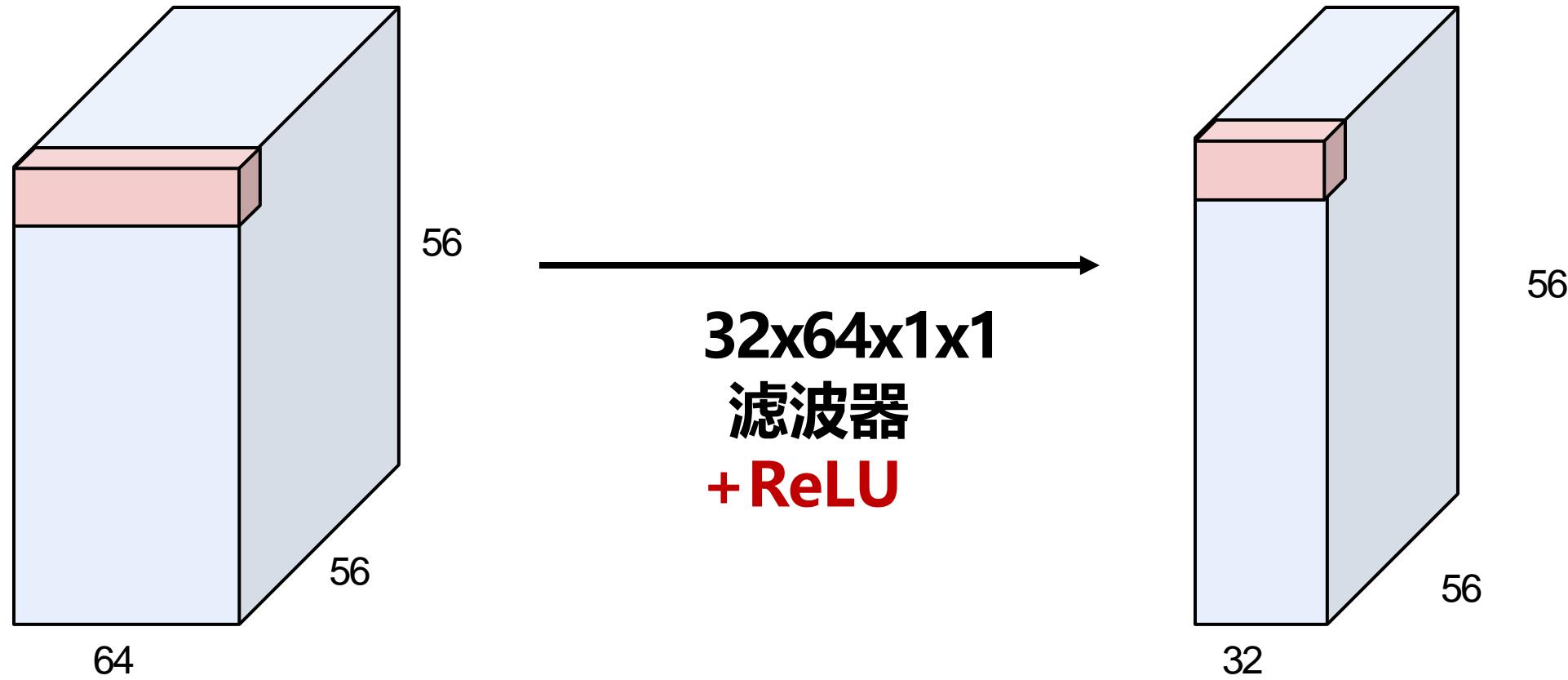
■  $H_2 = (H_1 - F + 2P)/S + 1$

■ 参数量:

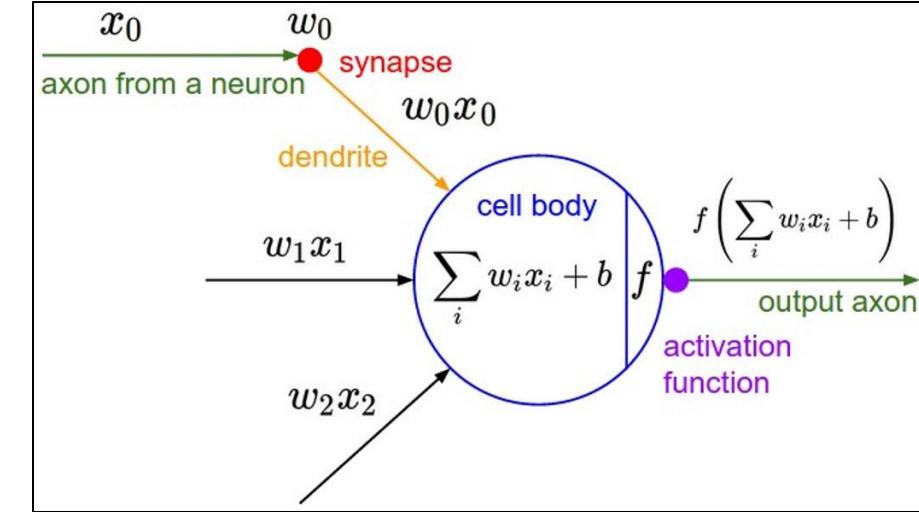
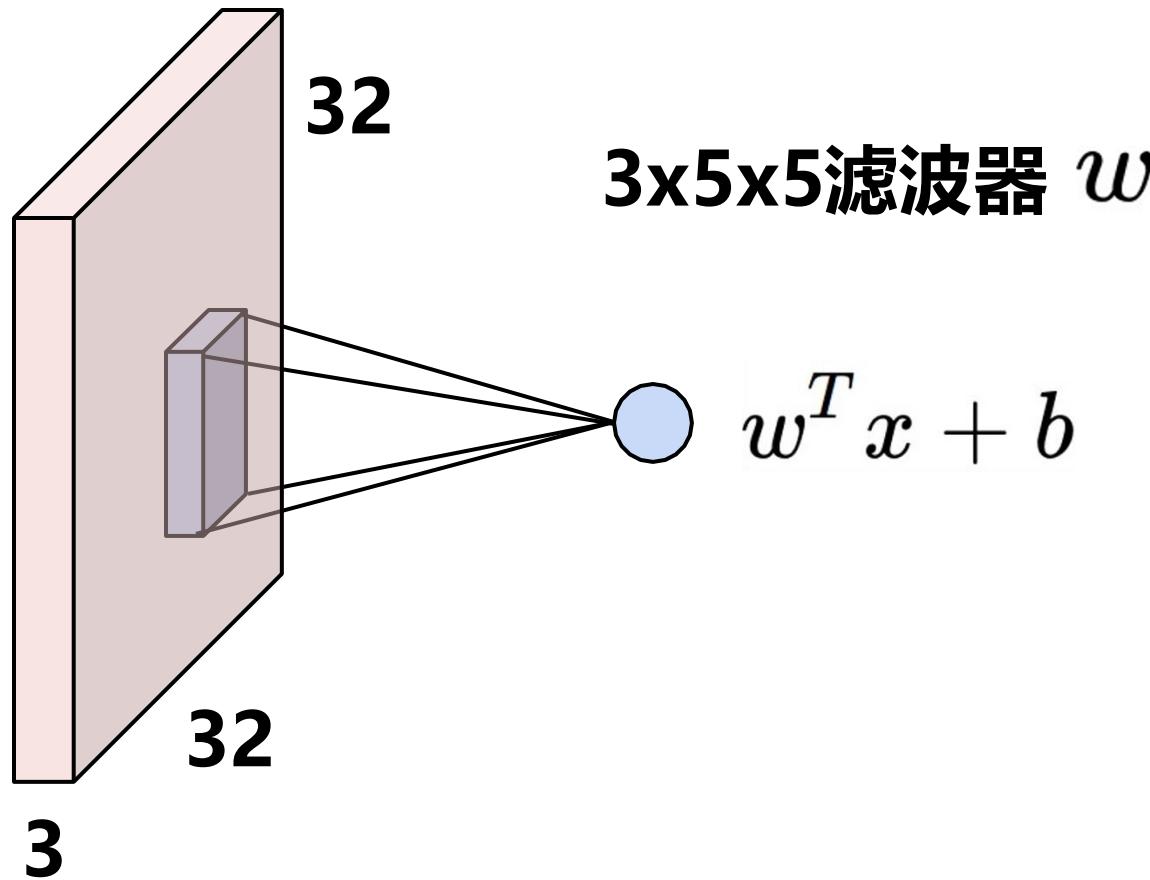
■  $F^2CK + K$ , 其中  $K$  是 bias 数量

- 卷积层总结
- 输入： $W_1 \times H_1 \times C$       ■  $K = (2\text{的指数, e.g. } 32, 64, 128, 512)$
- 卷积超参：
  - 滤波器数量  $K$
  - 滤波器大小  $F$
  - 步长  $S$
  - Padding  $P$
- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$  (匹配尺寸)
- $F = 1, S = 1, P = 0$

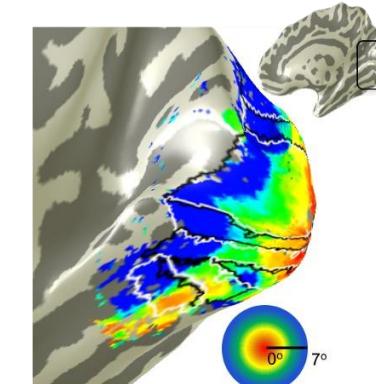
## ■ 1x1卷积层



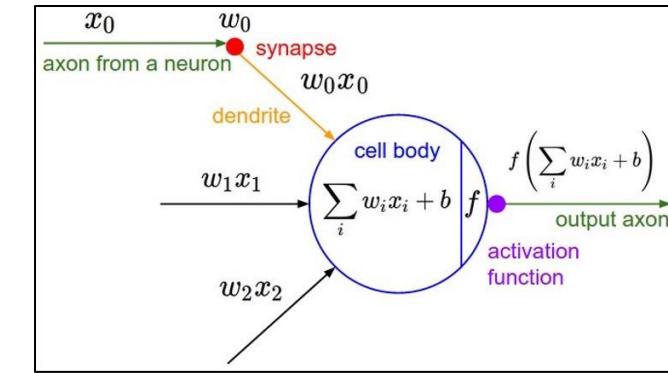
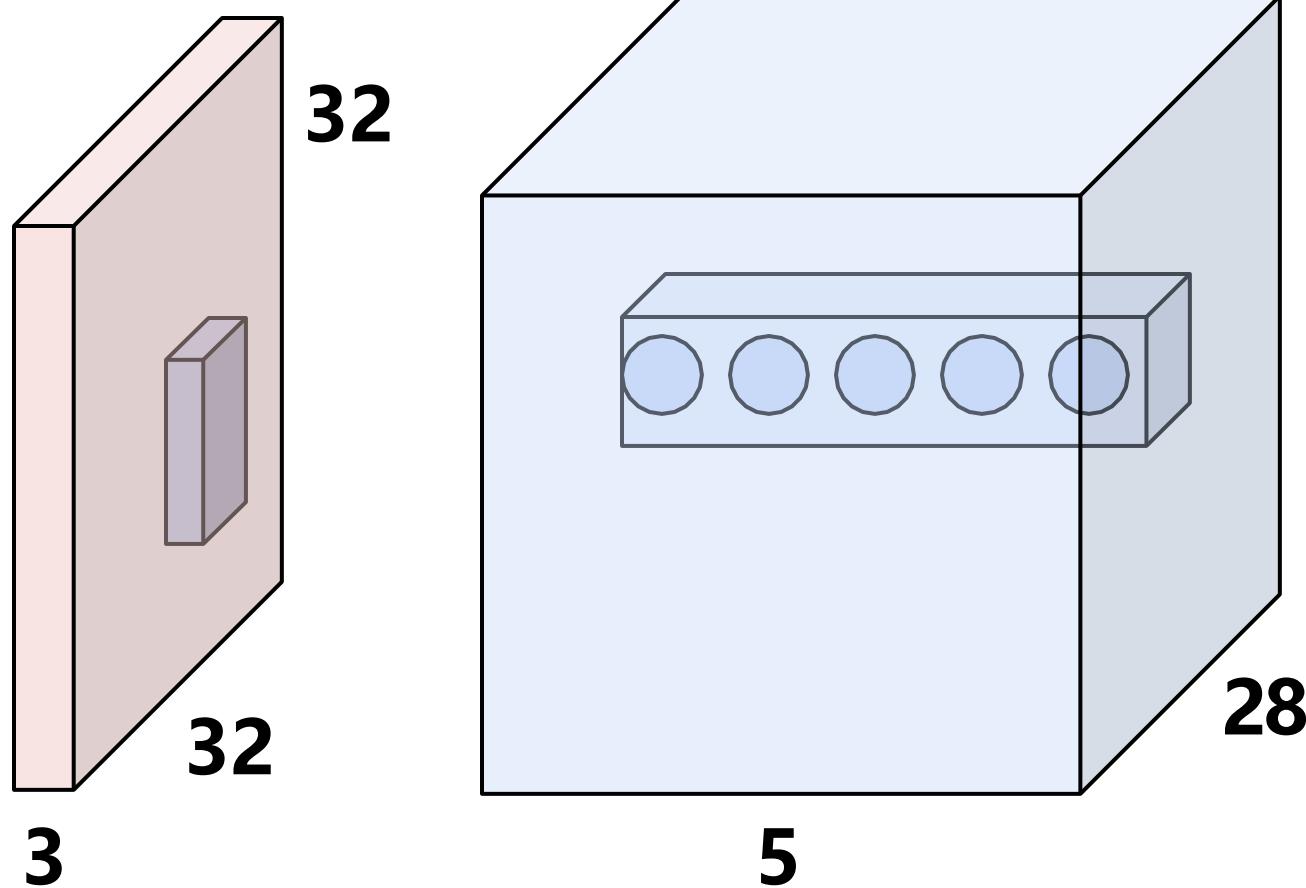
## ■ 从大脑/神经元角度理解卷积层



具有局部连接  
的神经元

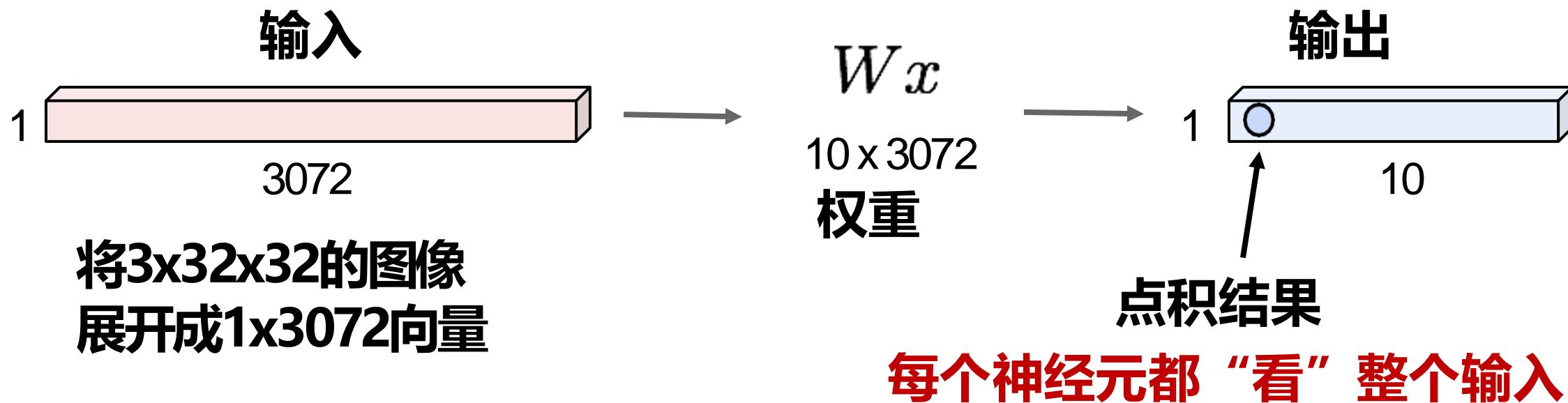


## ■ 从大脑/神经元角度理解卷积层



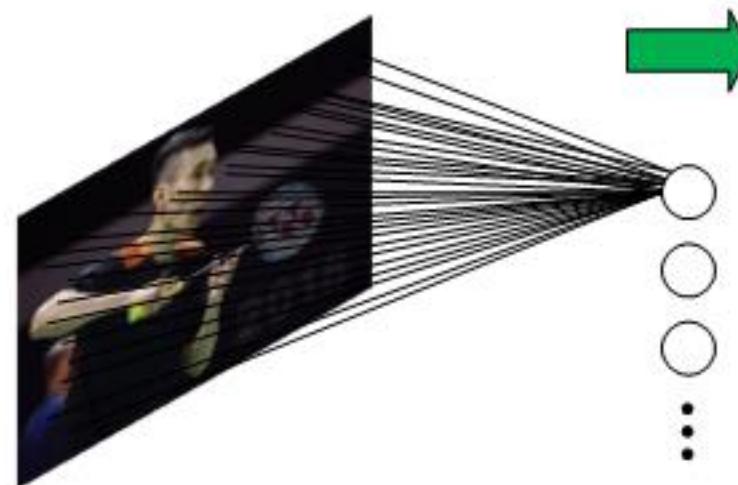
5个不同的神经元都  
在观察输入体积中的  
同一区域

## ■ 从大脑/神经元角度理解全连接层

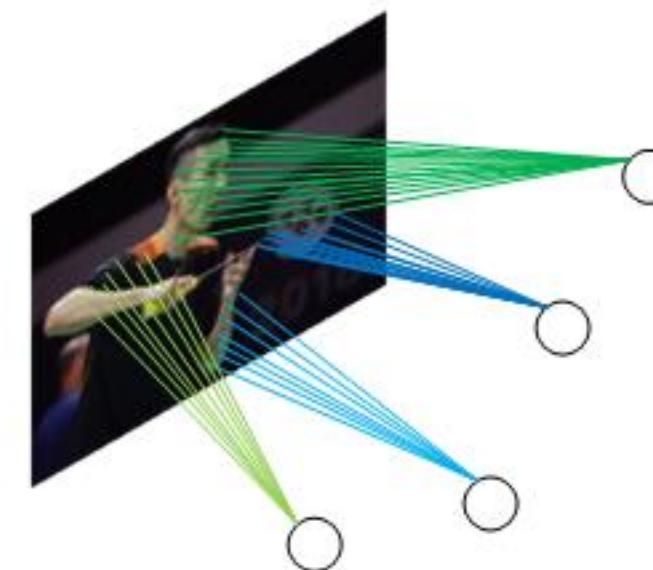


## ■ 卷积层 v.s. 全连接层

FULLY CONNECTED NEURAL NET

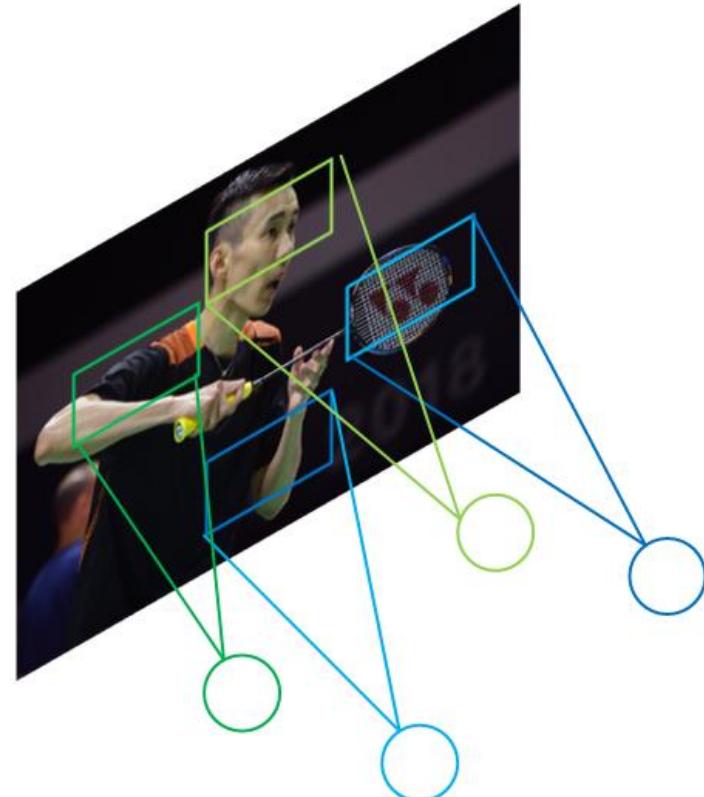


LOCALLY CONNECTED NEURAL NET

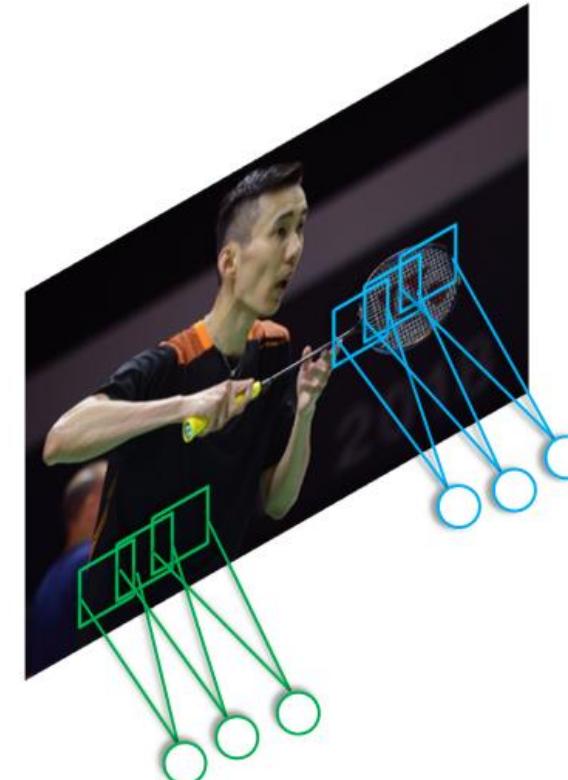


## ■ 卷积层 v.s. 局部全连接层

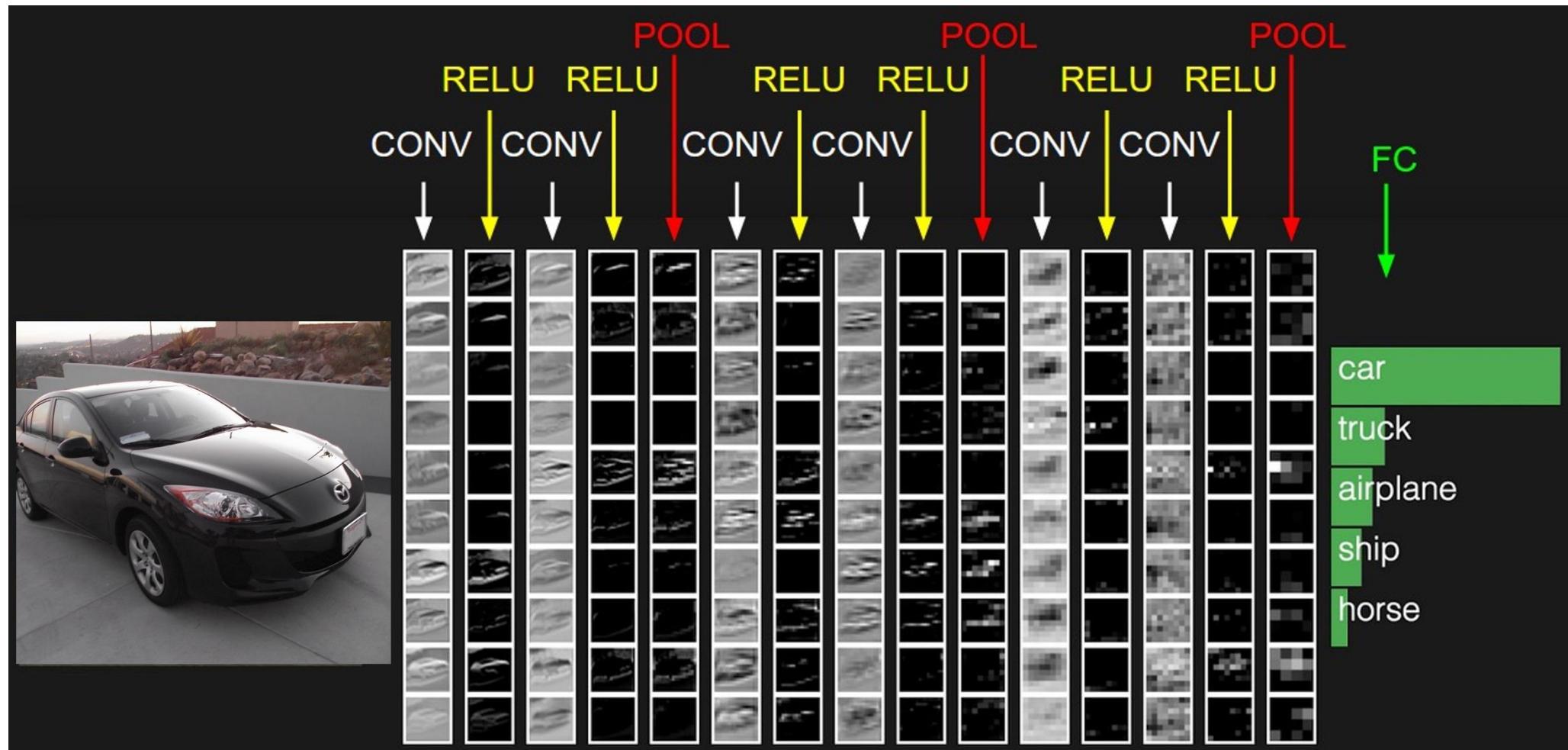
LOCALLY CONNECTED NEURAL NET



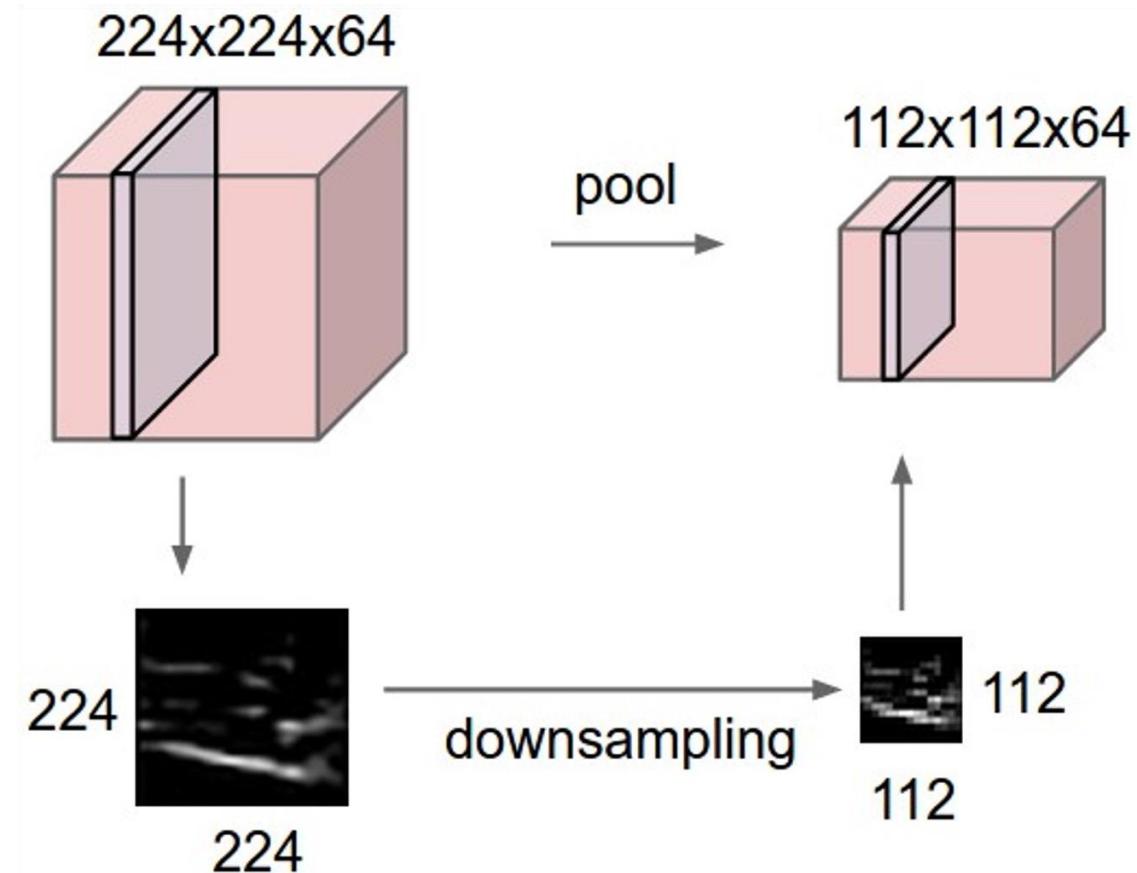
CONVOLUTIONAL NET



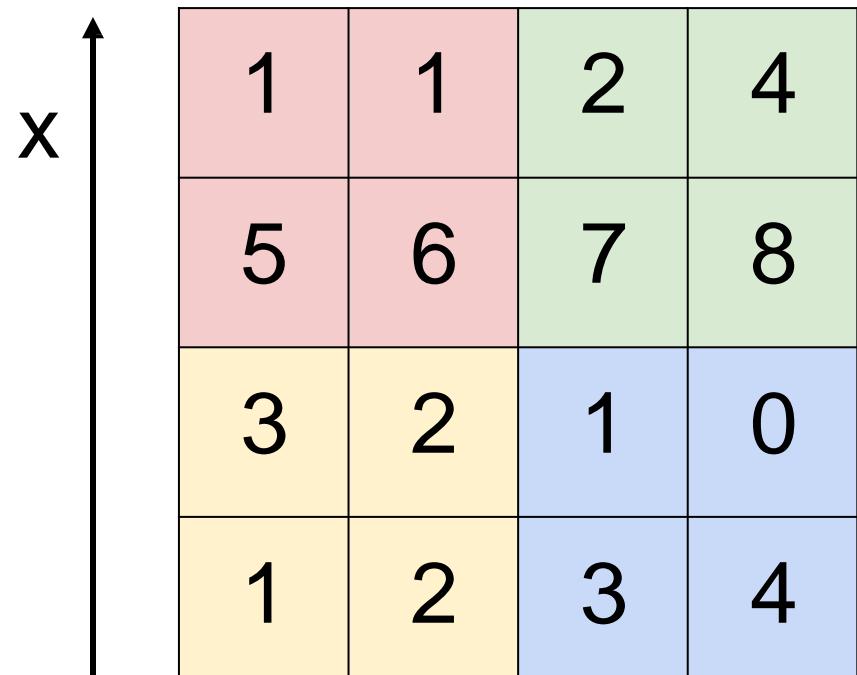
## ■ 池化层



- 池化层
- 使特征更小，更易于处理
- 独立操作每个激活图



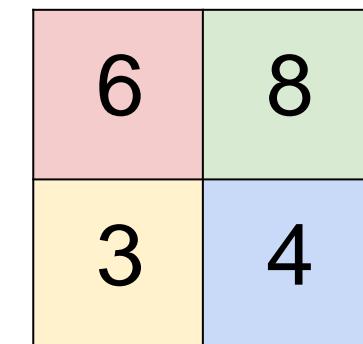
## ■ 最大池化 (Max Pooling)



A 4x4 input feature map with values ranging from 0 to 8. The values are arranged in a 4x4 grid:

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

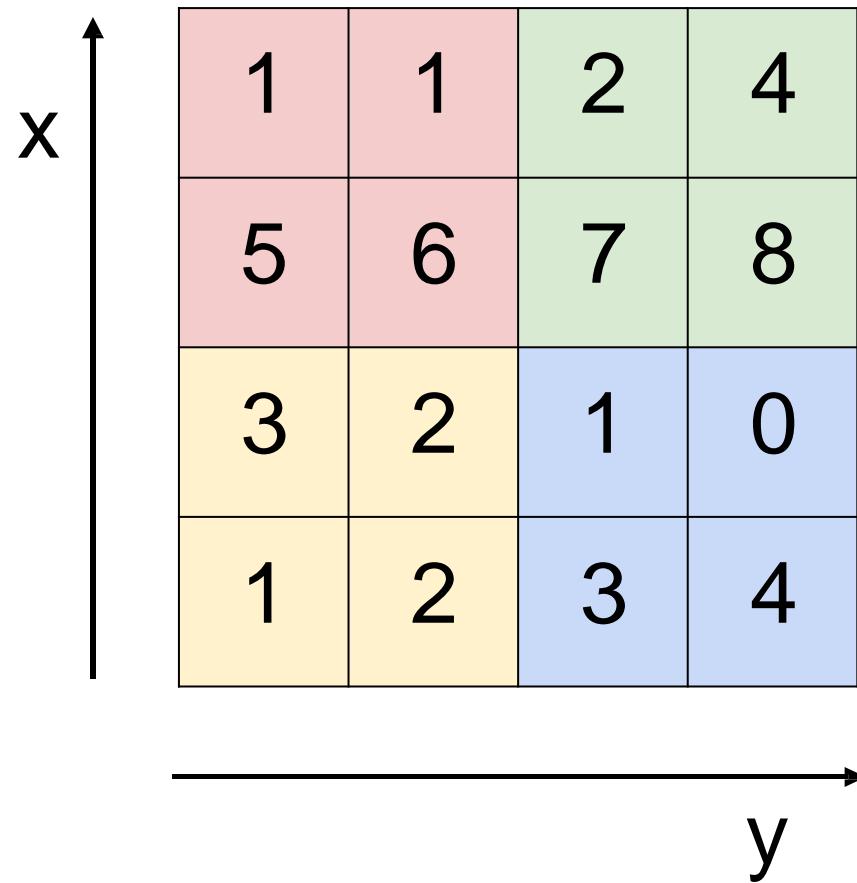
**Max Pooling  
2x2 大小和 2 步长**



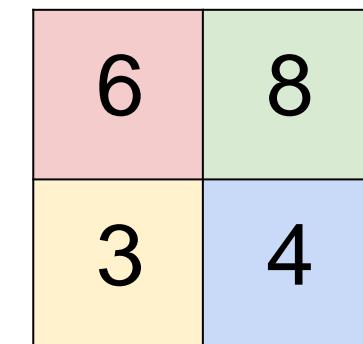
The output feature map after applying Max Pooling with a 2x2 kernel and a stride of 2. It has a size of 2x2. The values are:

6	8
3	4

## ■ 最大池化 (Max Pooling)



**Max Pooling  
2x2 大小和 2 步长**



The result of applying Max Pooling with a 2x2 kernel and a stride of 2 to the input feature map. The resulting output feature map has a 2x2 dimension with values 6, 8, 3, and 4.

6	8
3	4

**没有可学习参数  
引入空间不变性**

## ■ 池化层总结

■ 输入:  $W_1 \times H_1 \times C$

■ 卷积超参:

■ 滤波器大小  $F$

■ 步长  $S$

■ 输出:  $W_2 \times H_2 \times C$

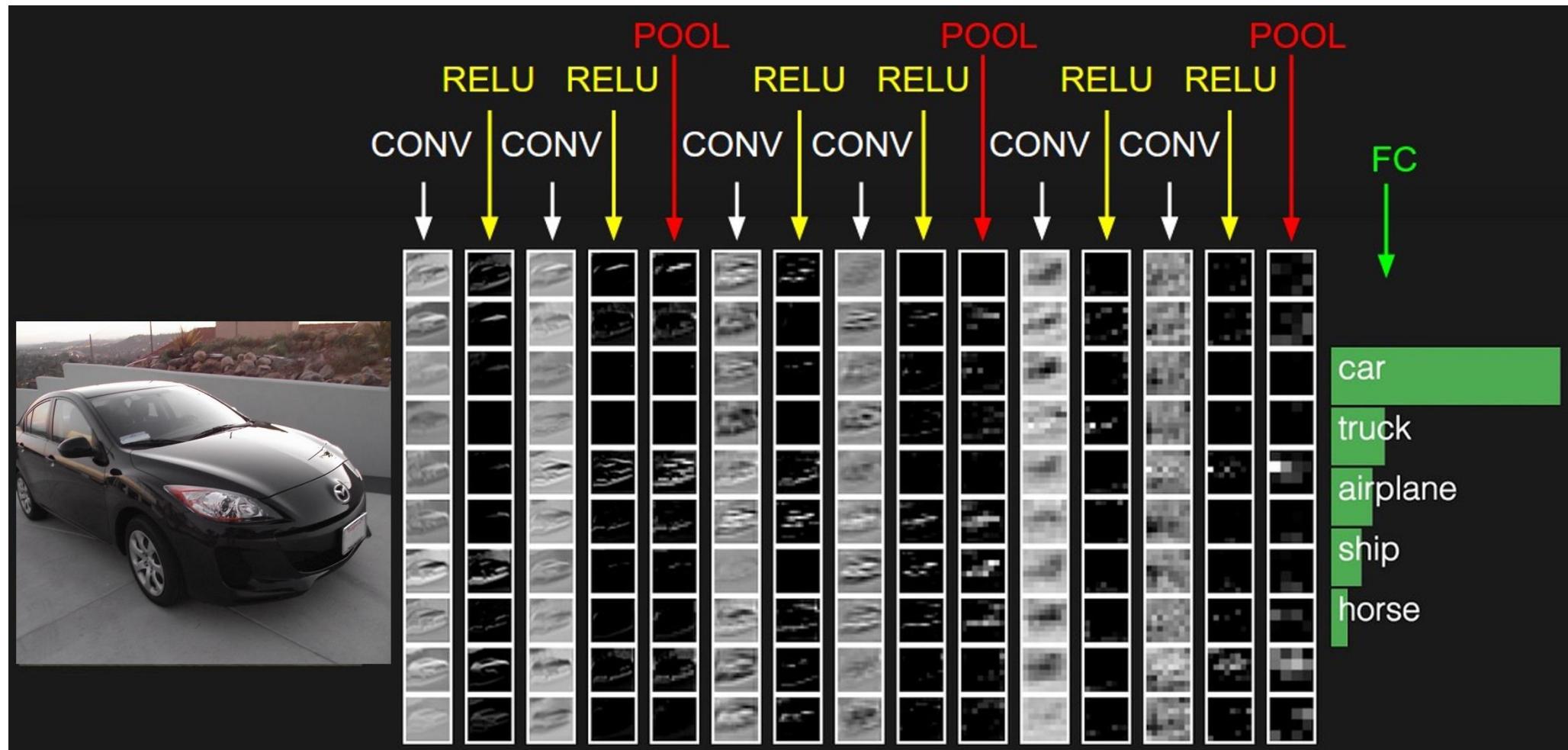
■  $W_2 = (W_1 - F)/S + 1$

■  $H_2 = (H_1 - F)/S + 1$

■ 参数数量:

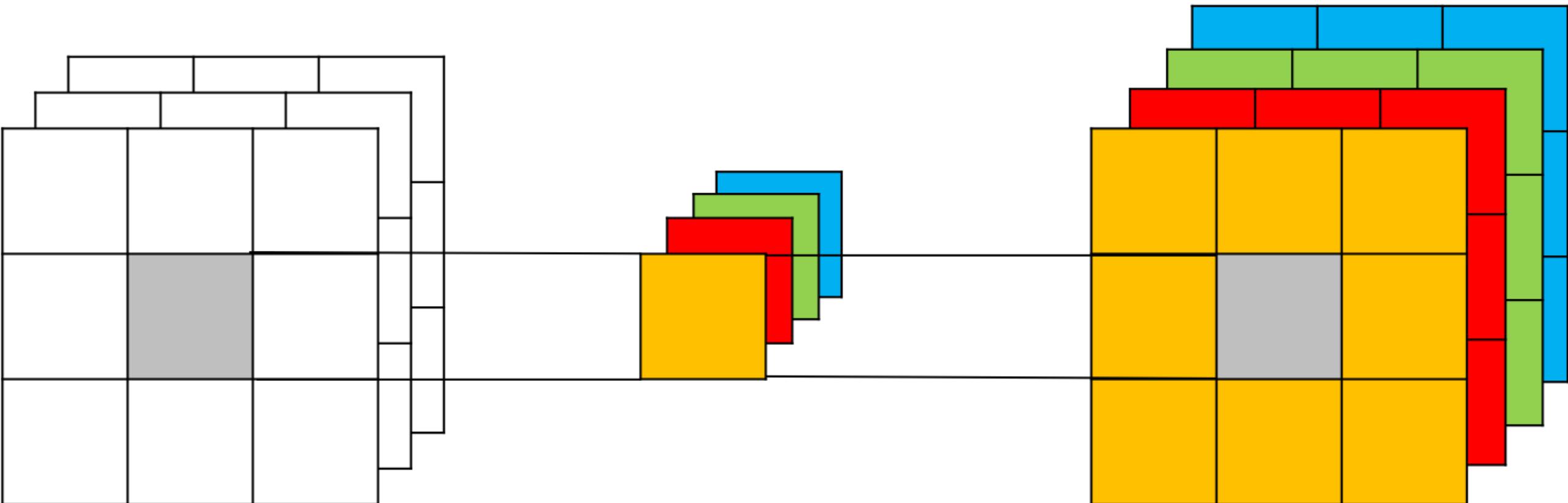
■ 0

## 全连接层

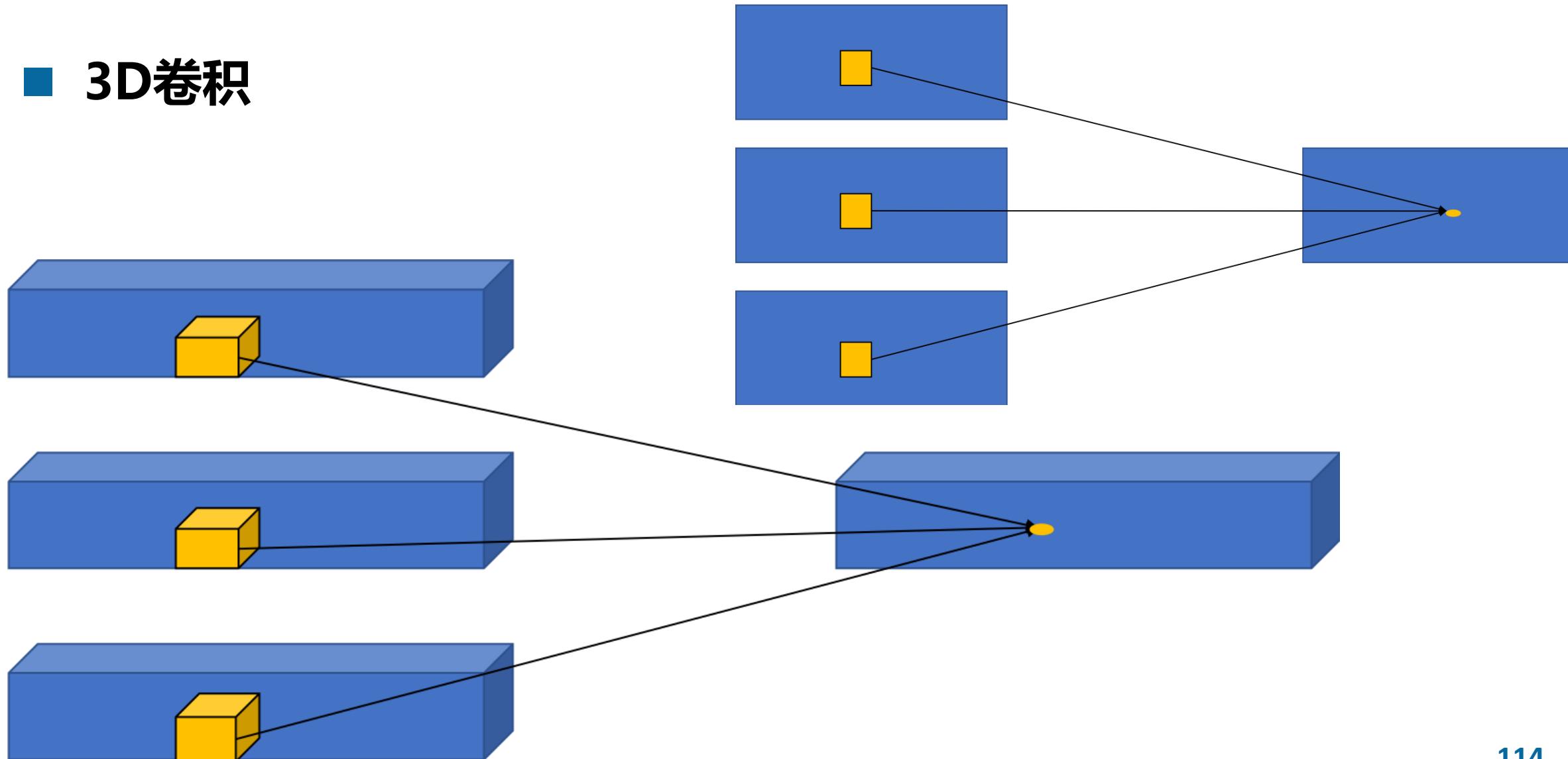


- 卷积层+池化层+全连层+激活函数
- 趋向于更小的滤波器和更深的层
- 趋向于舍弃池化层和全连接层
- 以往常见的卷积网络架构：
  - $[(\text{Conv-ReLU})^*N\text{-Pooling}]^*M - (\text{FC-ReLU})^*K\text{-Softmax}$
  - N一般为~5, M较大, K一般在0~2
- 新型卷积网络 (ResNet, GoogLeNet) 提出了新的架构

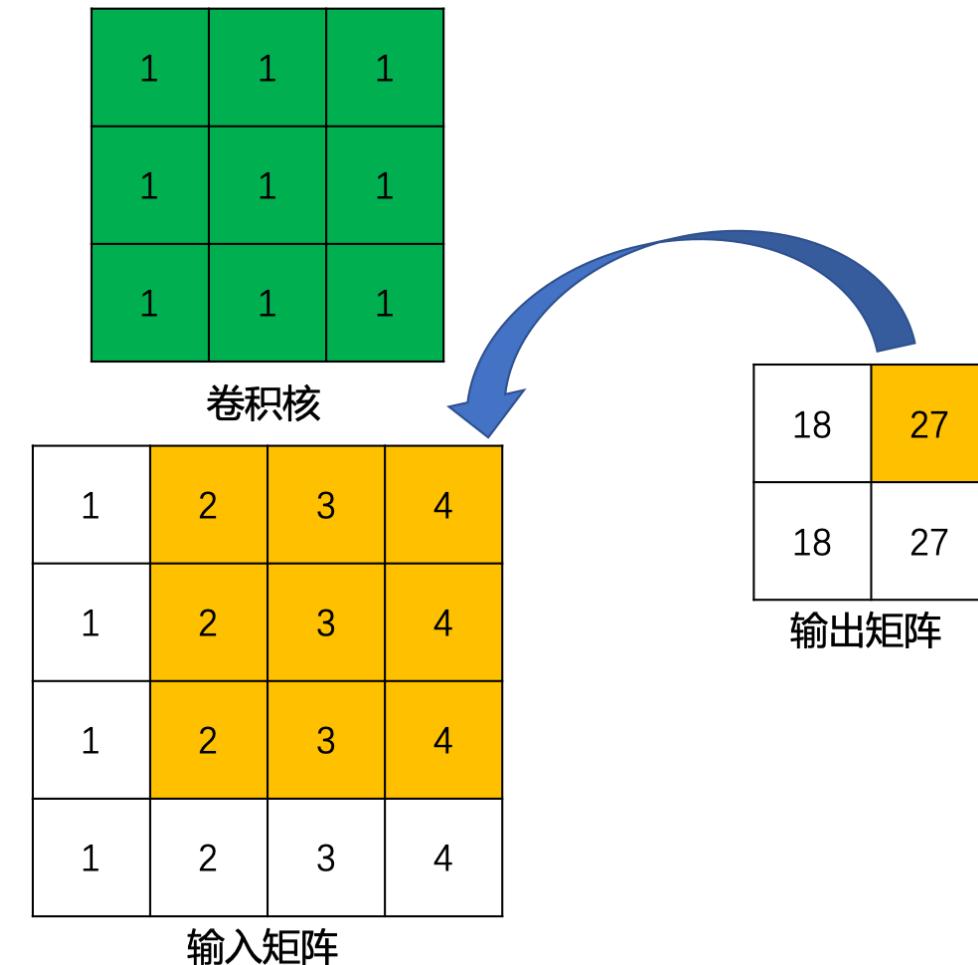
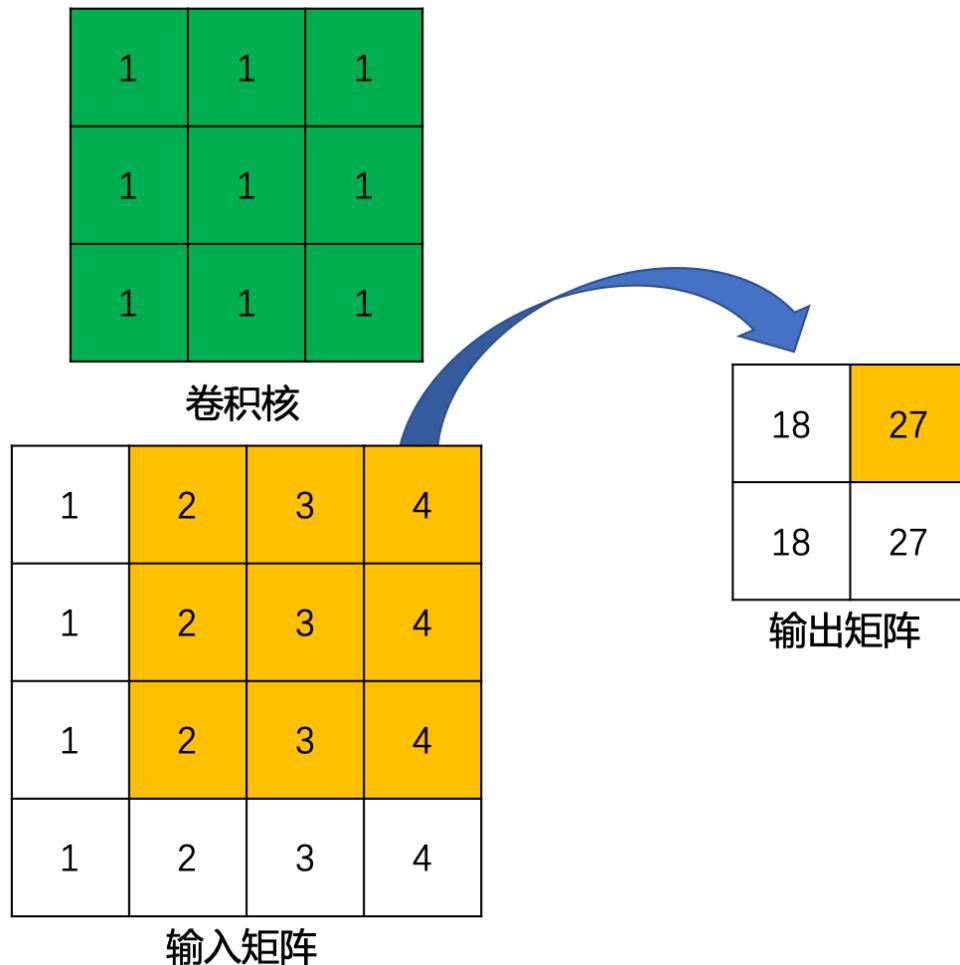
## ■ 1x1 卷积



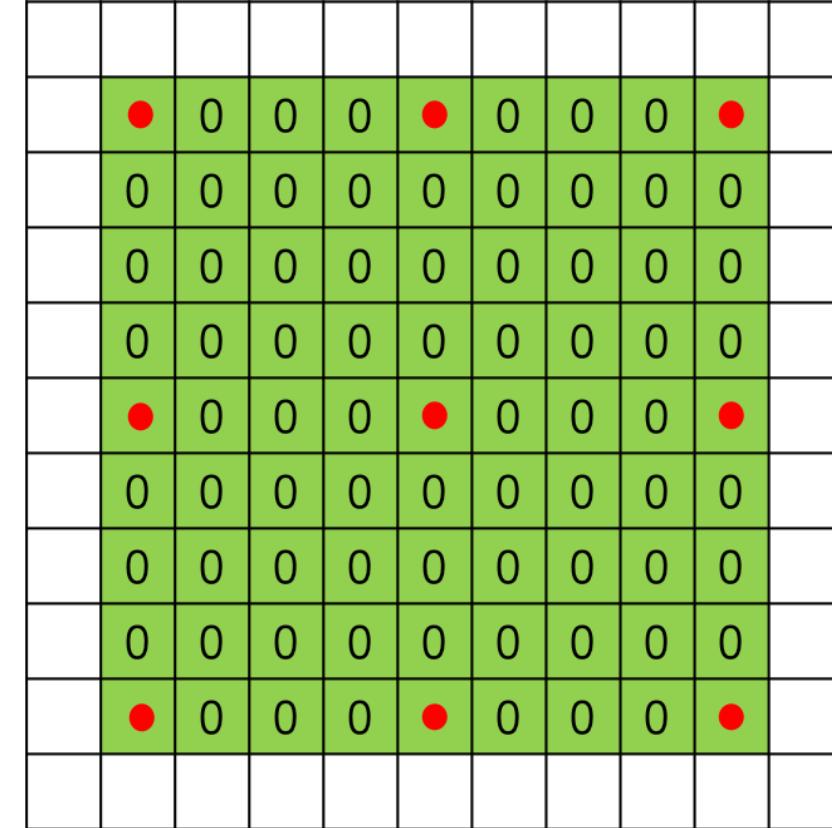
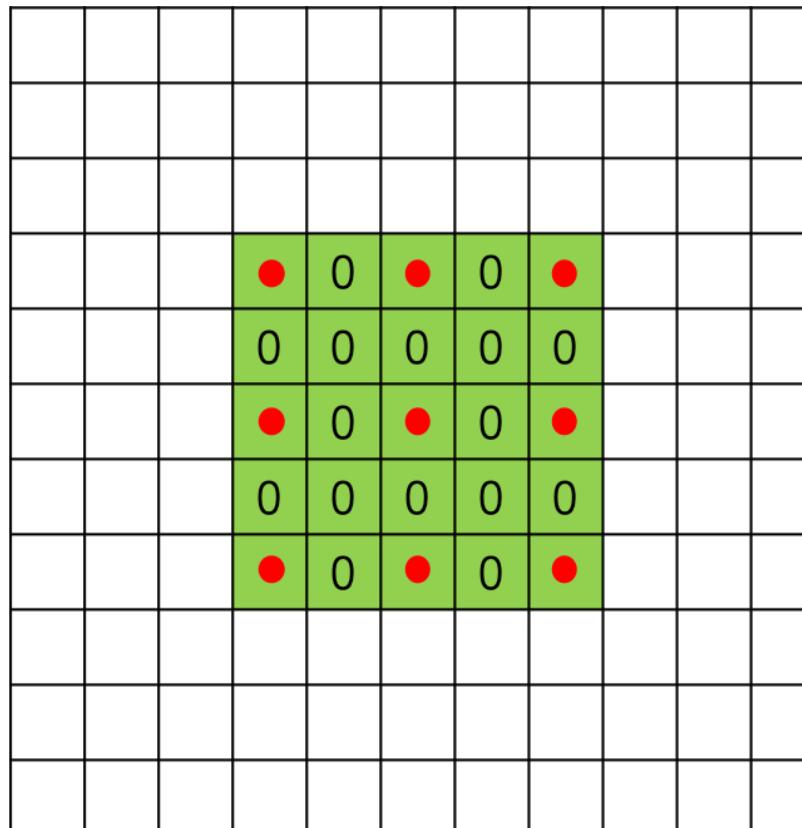
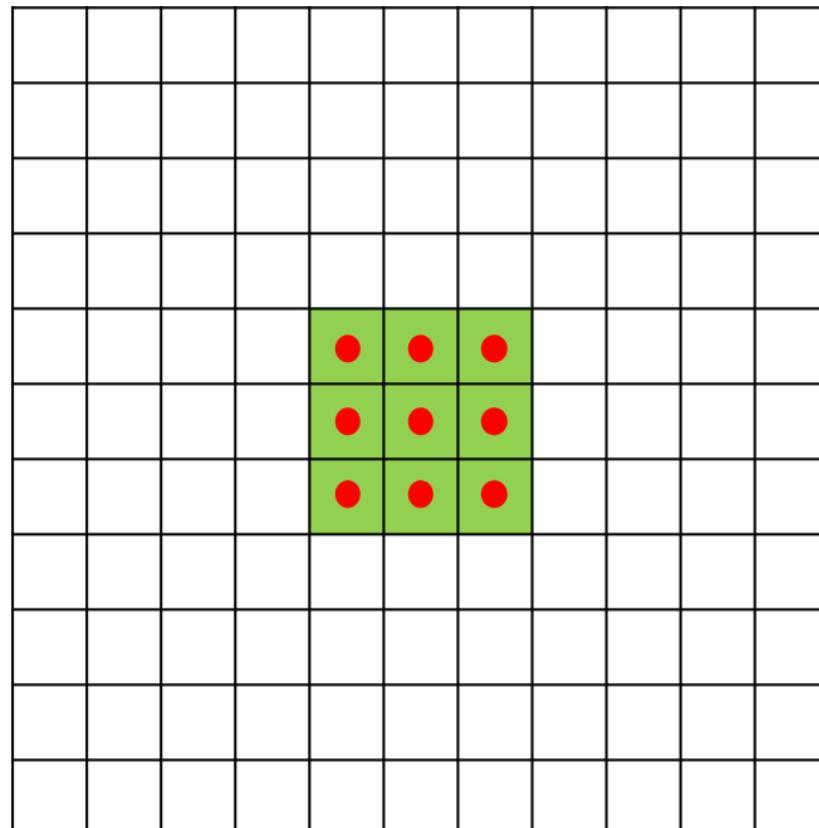
## ■ 3D 卷积



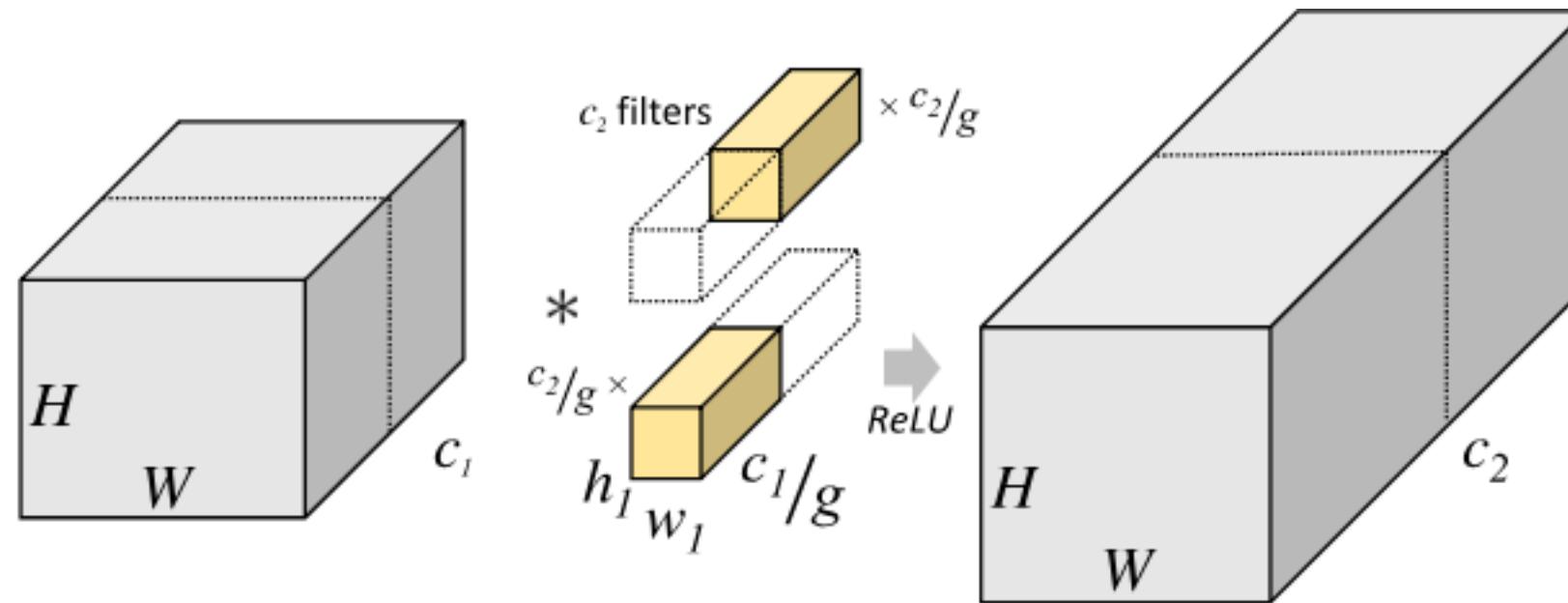
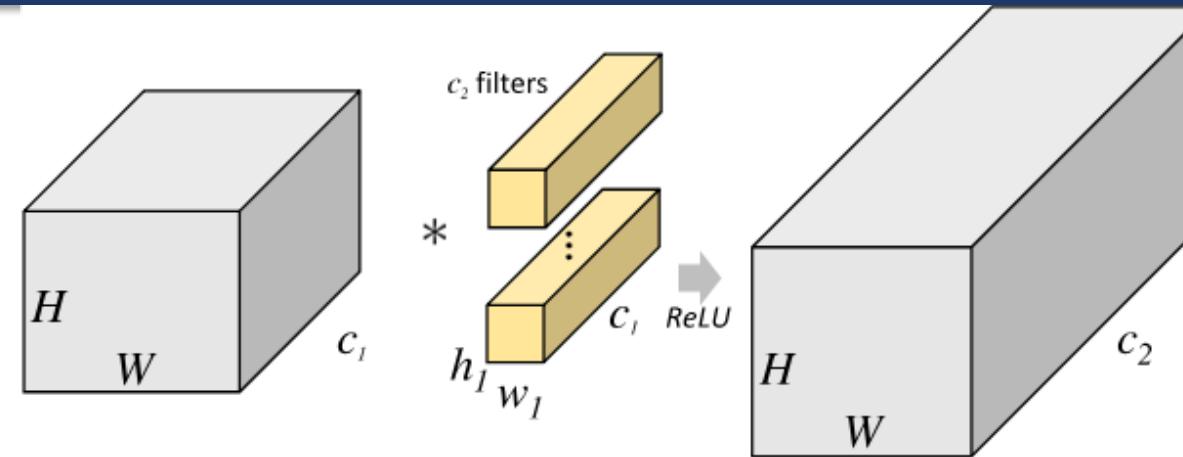
## ■ 转置卷积



## ■ 空洞卷积

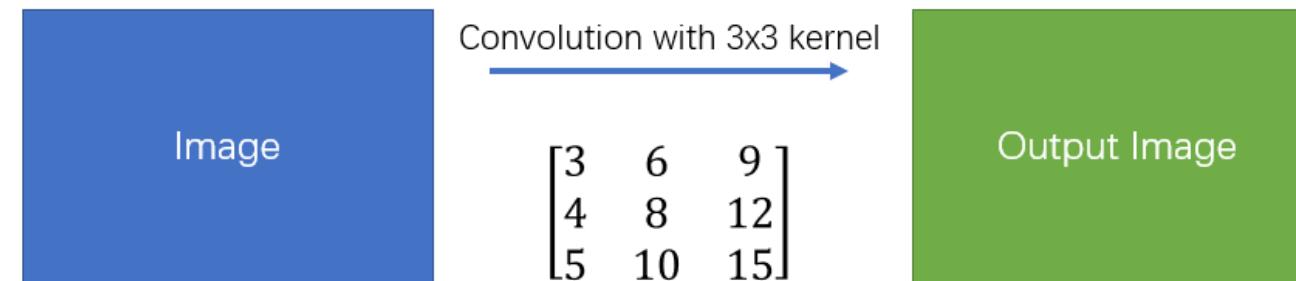


## ■ 分组卷积

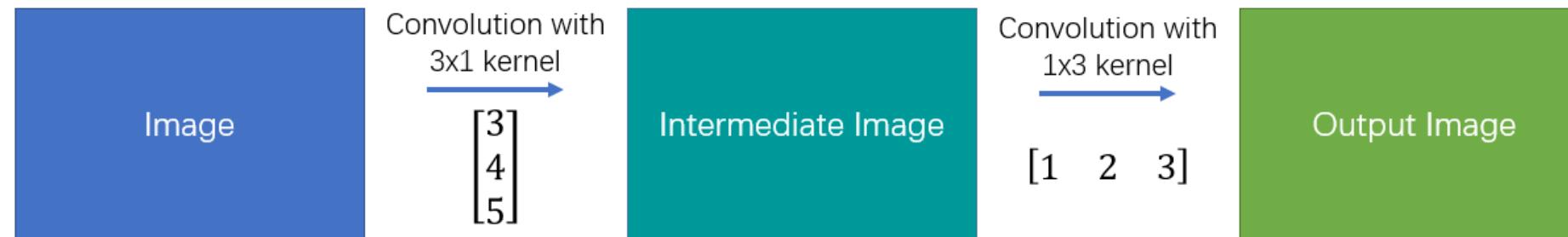


## ■ 可分离卷积

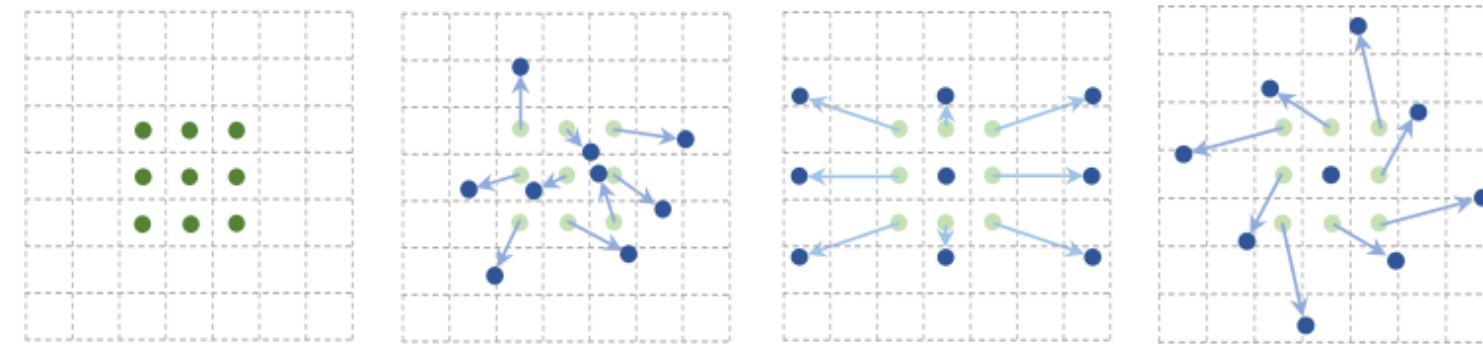
### Simple Convolution



### Spatial Separable Convolution



## ■ 可变形卷积

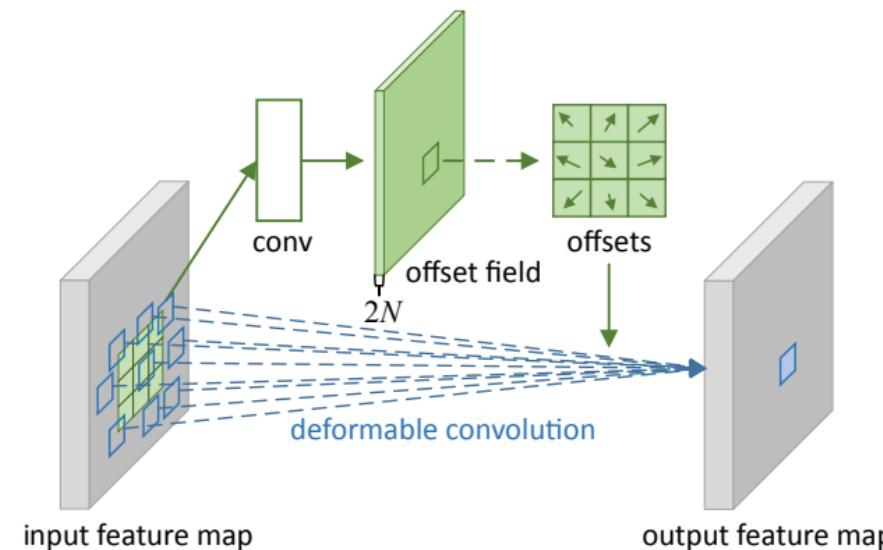
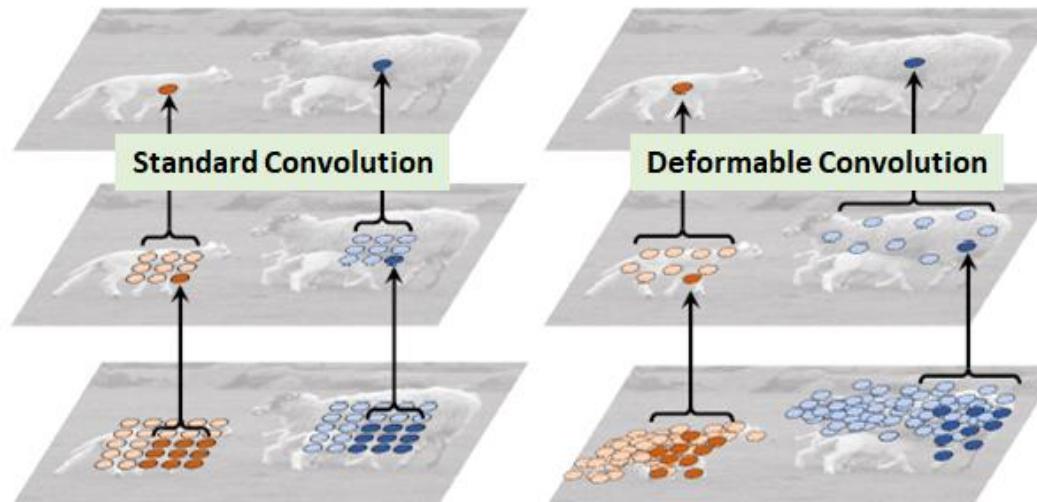


(a)

(b)

(c)

(d)



# 卷积神经网络

## ■ 下节课：卷积神经网络架构

