# 本科生《计算机视觉》
# 基于深度学习的视觉理解与生成

黄雷

人工智能研究院

huangleiAI@buaa.edu.cn

2023年10月12日

# 主要内容

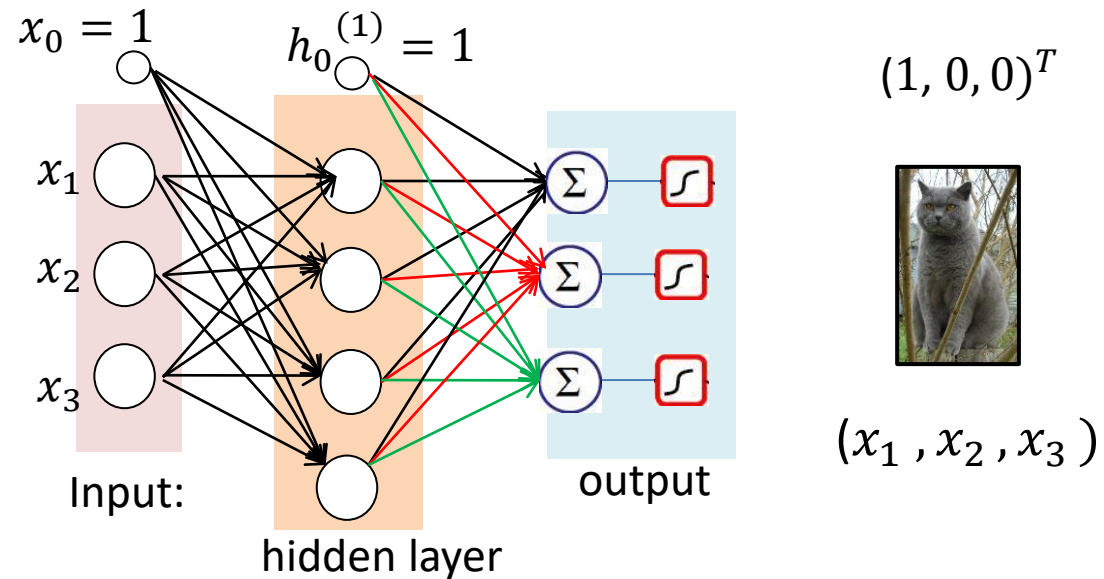- 深度学习基础
  - 神经网络及反向传播算法
  - <span style="color:red">卷积神经网络中的视觉表示思想</span>
- 视觉理解任务
  - 目标检测
  - 分割
- 视觉生成
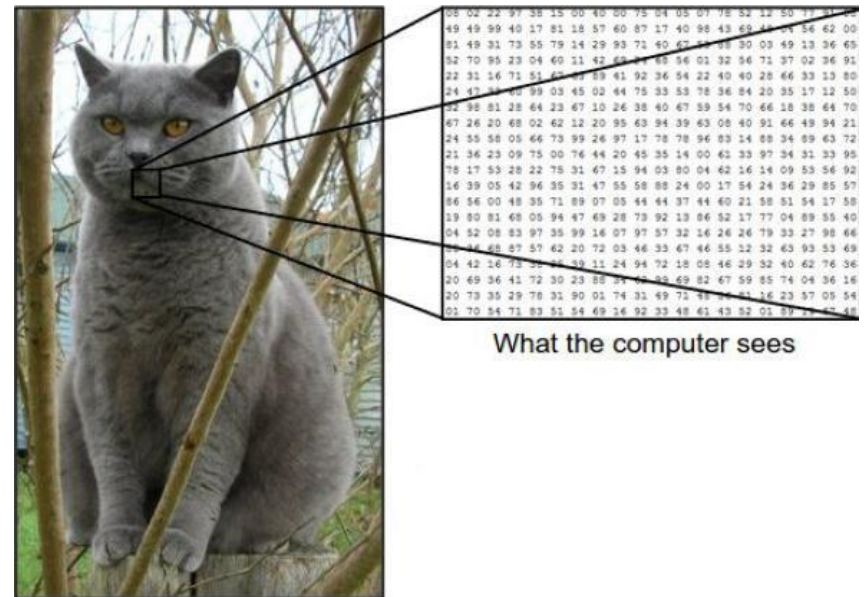  - 深度生成模型
  - 图像翻译任务详解
- 深度神经网络训练技巧

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution like template
  - Filters
  - Convolution module
- Pooling layer (module)

# Feature extraction

- Feature extraction
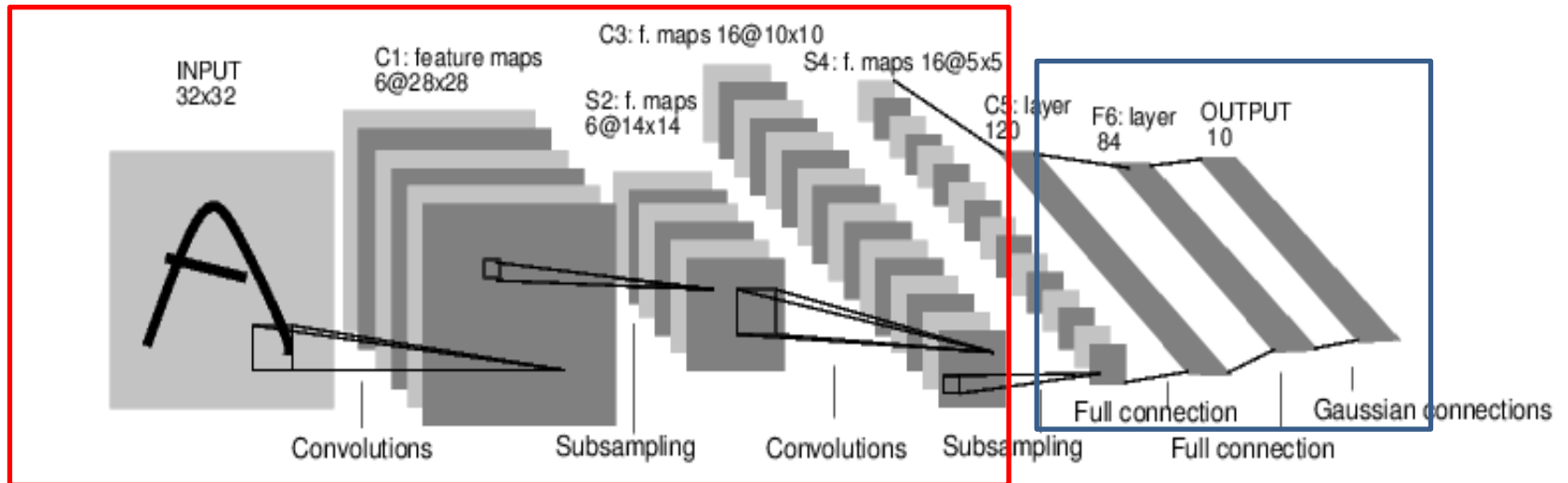  - Pixel-wise input
  - Correlation between features

$x_0 = 1$

$h_0^{(1)} = 1$

$(1, 0, 0)^T$

$(x_1, x_2, x_3)$

$x_1$

$x_2$

$x_3$

Input:

hidden layer

output

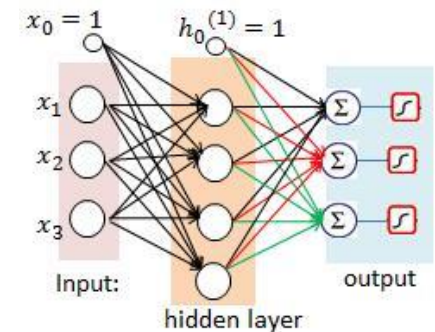Convolutional Neural Network(CNN),卷积神经网络

What the computer sees

# Convolution Neural Network

- Lenet-5



Convolution related layers

全连接层

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution like template
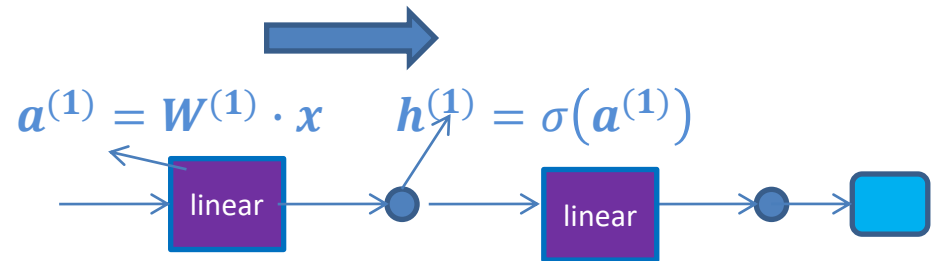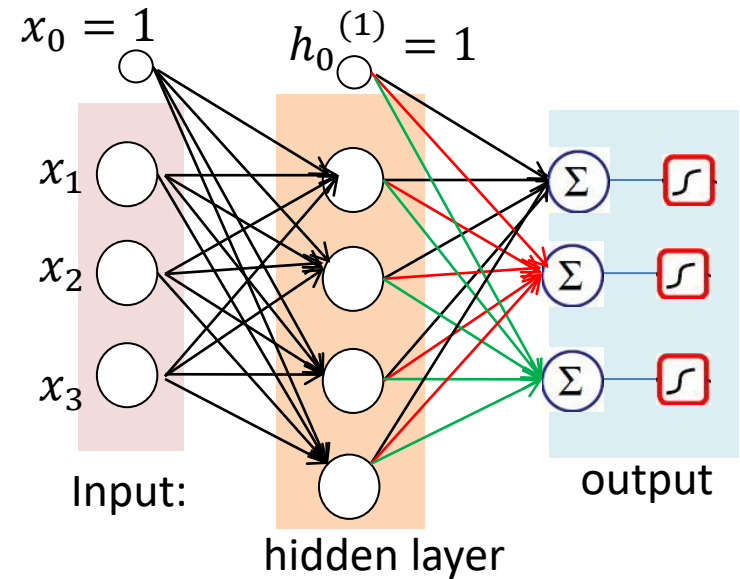  - Filters
  - Convolution module
- Pooling layer (module)

# Module-wise architecture

➢PyTorch 平台

- **Model construction**

```python
# Define model
class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10)
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```



$x_0 = 1$

$h_0^{(1)} = 1$

Input:

hidden layer

output

$a^{(1)} = W^{(1)} \cdot x$

$h^{(1)} = \sigma(a^{(1)})$

linear
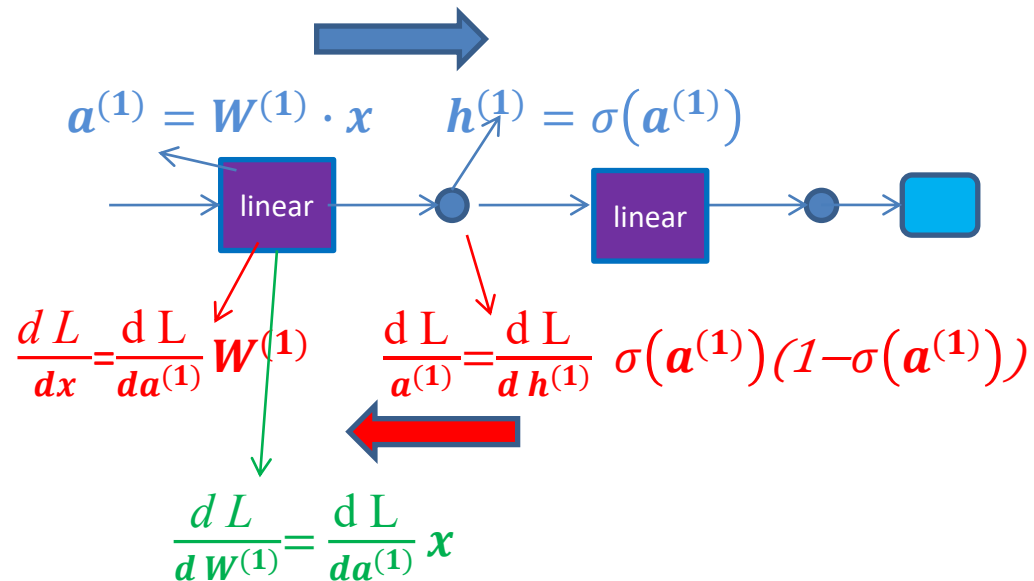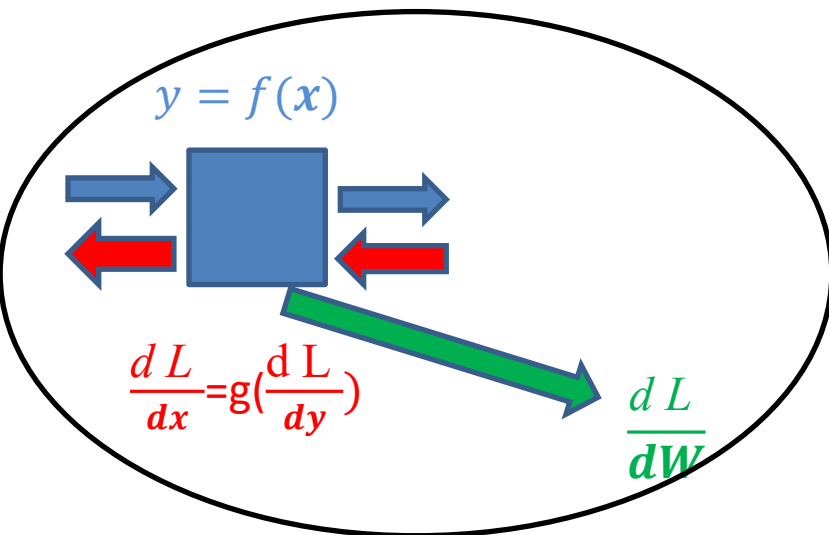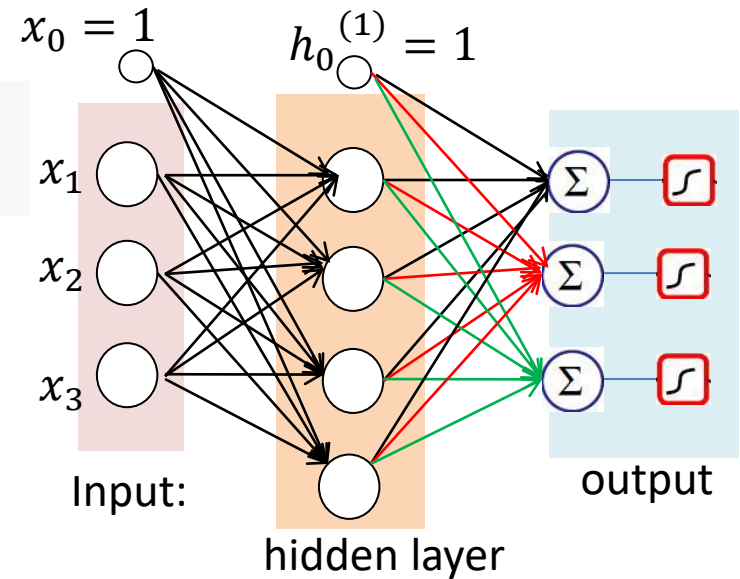
linear

# Module-wise architecture

➤ PyTorch 平台

```
loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
```

- **Training per iteration:**

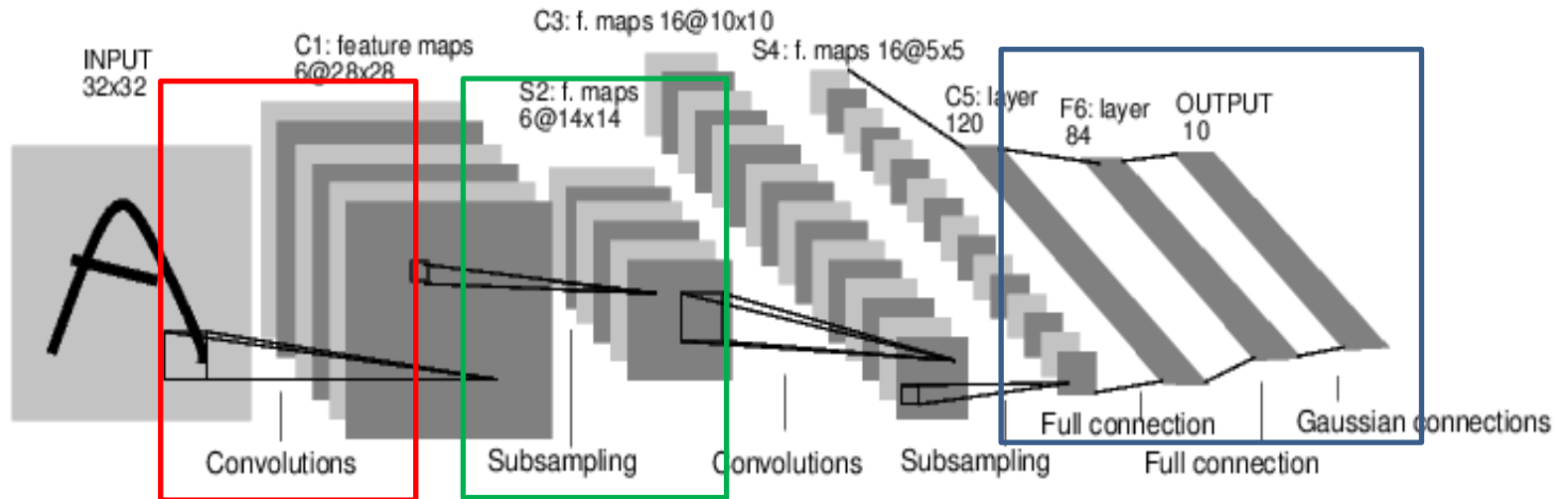```
# Compute prediction error
pred = model(X)
loss = loss_fn(pred, y)

# Backpropagation
optimizer.zero_grad()
loss.backward()
optimizer.step()
```



$x_0 = 1$

$h_0^{(1)} = 1$

$x_1$

$x_2$

$x_3$

Input:

hidden layer

output

$y = f(x)$

$$\frac{dL}{dx} = g\left(\frac{dL}{dy}\right)$$

$$\frac{dL}{dW}$$

$$a^{(1)} = W^{(1)} \cdot x \qquad h^{(1)} = \sigma(a^{(1)})$$

linear

linear

$$\frac{dL}{dx} = \frac{dL}{da^{(1)}} W^{(1)}$$

$$\frac{dL}{a^{(1)}} = \frac{dL}{d h^{(1)}} \sigma(a^{(1)})(1-\sigma(a^{(1)}))$$

$$\frac{dL}{d W^{(1)}} = \frac{dL}{da^{(1)}} x$$
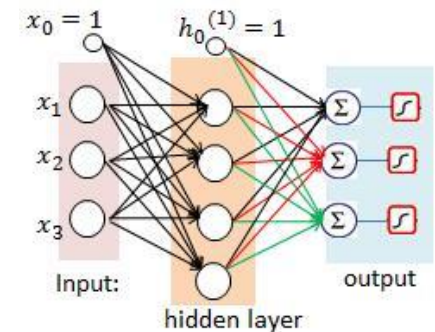
# Convolution Neural Network

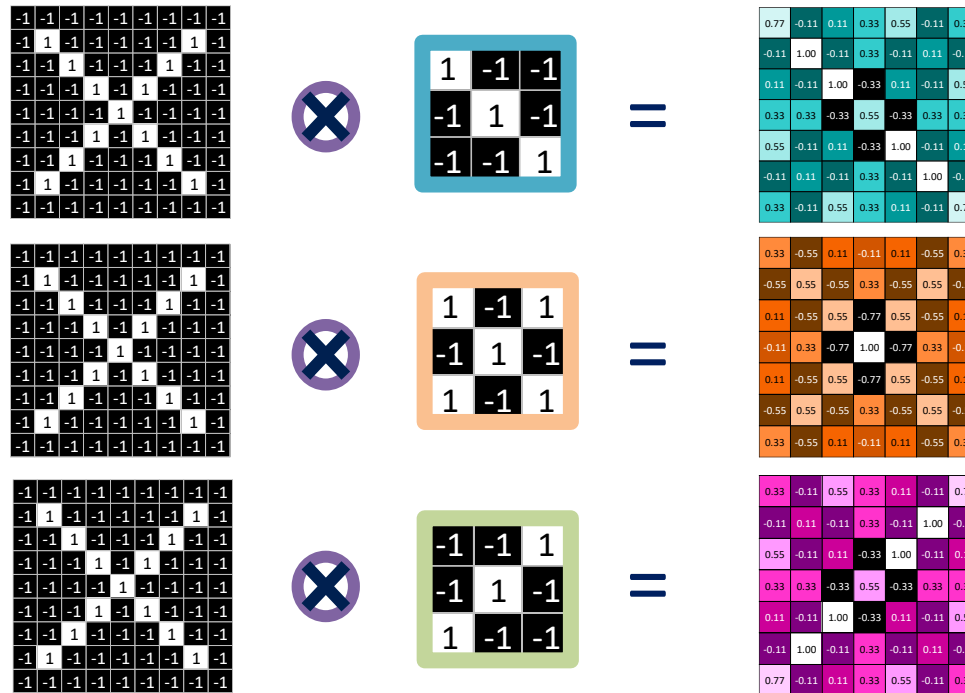- Lenet-5


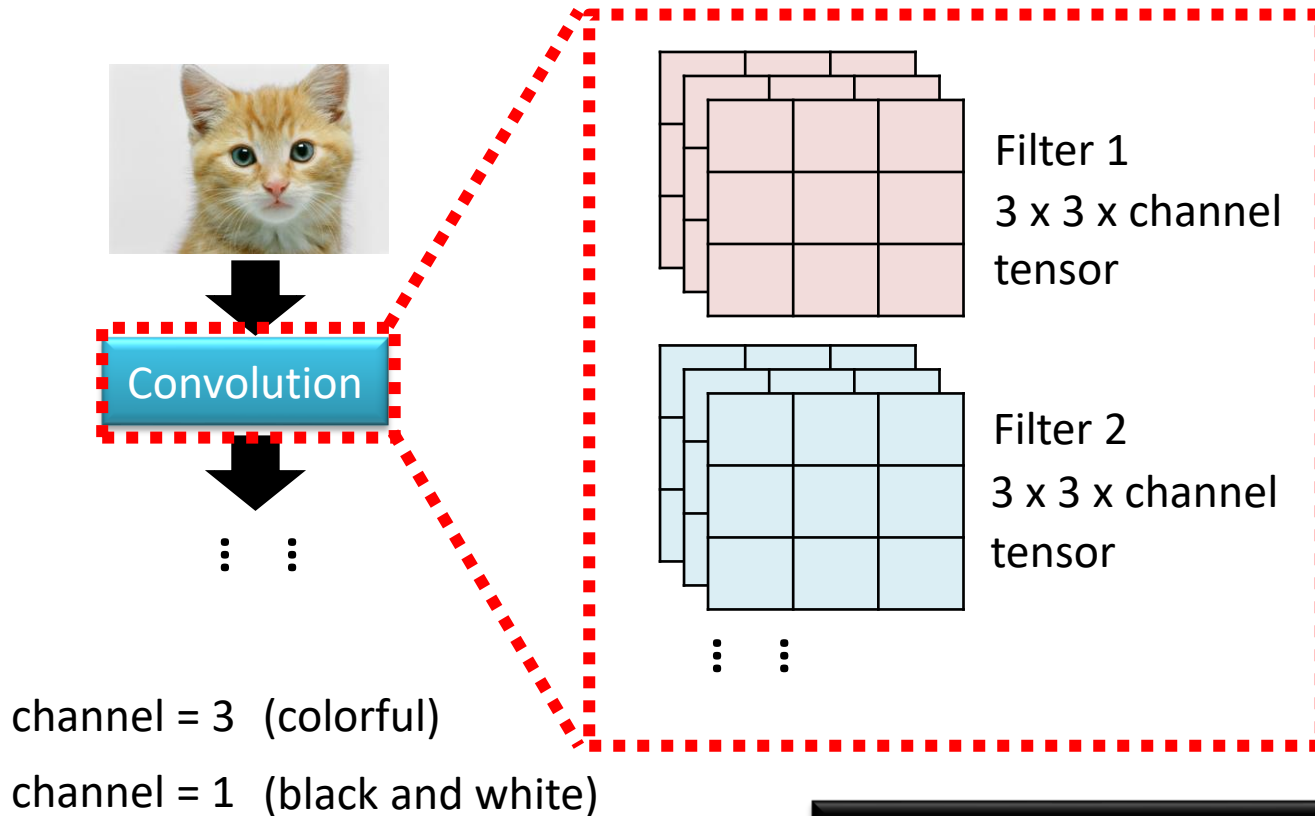
Convolution (卷积)   Pooling (池化)   全连接层

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution like template
  - Filters
  - Why convolution for vision
- Pooling layer (module)

# Convolution

- **Match like template**

# Convolutional Layer

Convolution

Filter 1
3 x 3 x channel
tensor

Filter 2
3 x 3 x channel
tensor

channel = 3   (colorful)

channel = 1   (black and white)

Each filter detects a small
pattern (3 x 3 x channel).

# Convolutional Layer

Consider channel = 1
(black and white image)

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

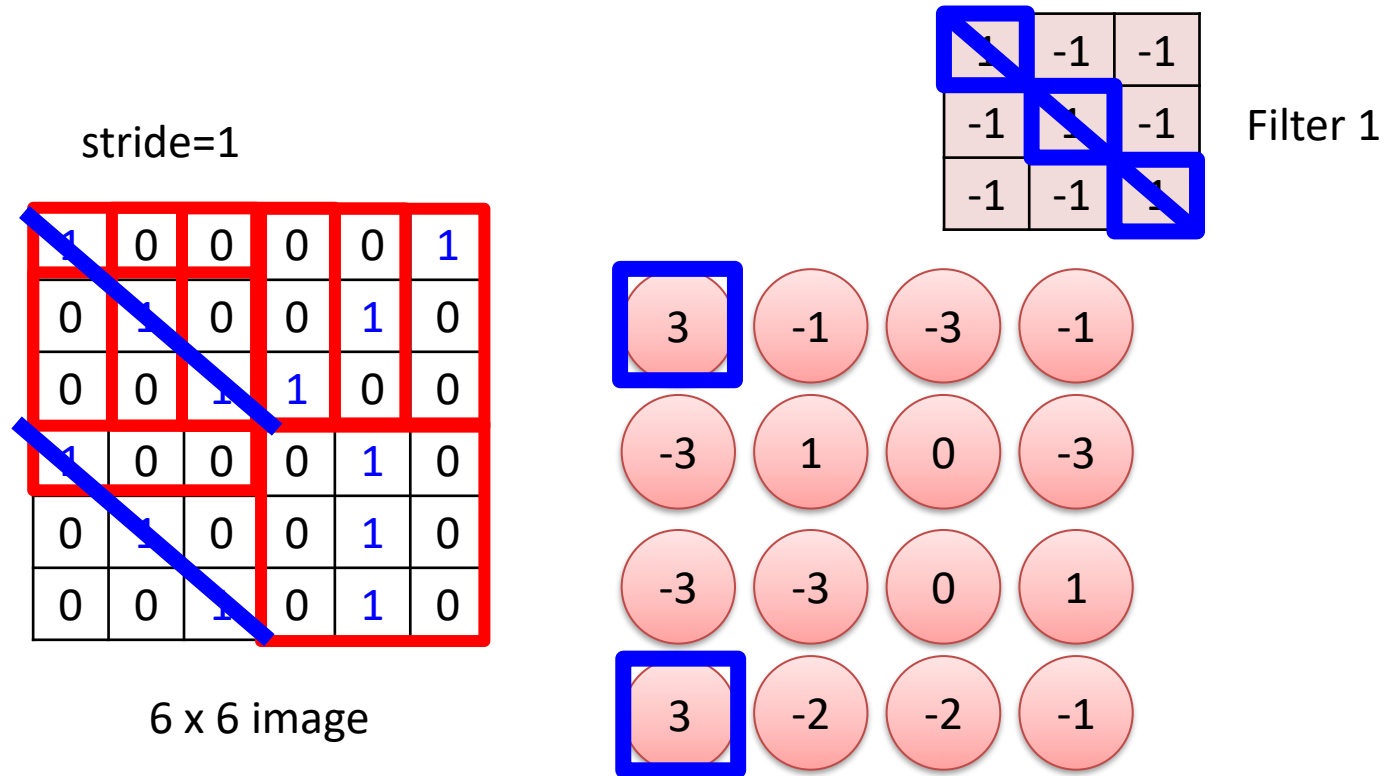| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

⋮  ⋮

(The values in the filters
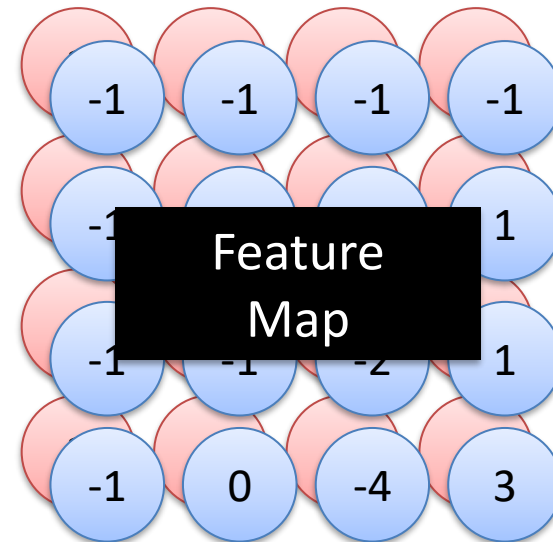are unknown parameters.)

13

# Convolutional Layer

stride=1

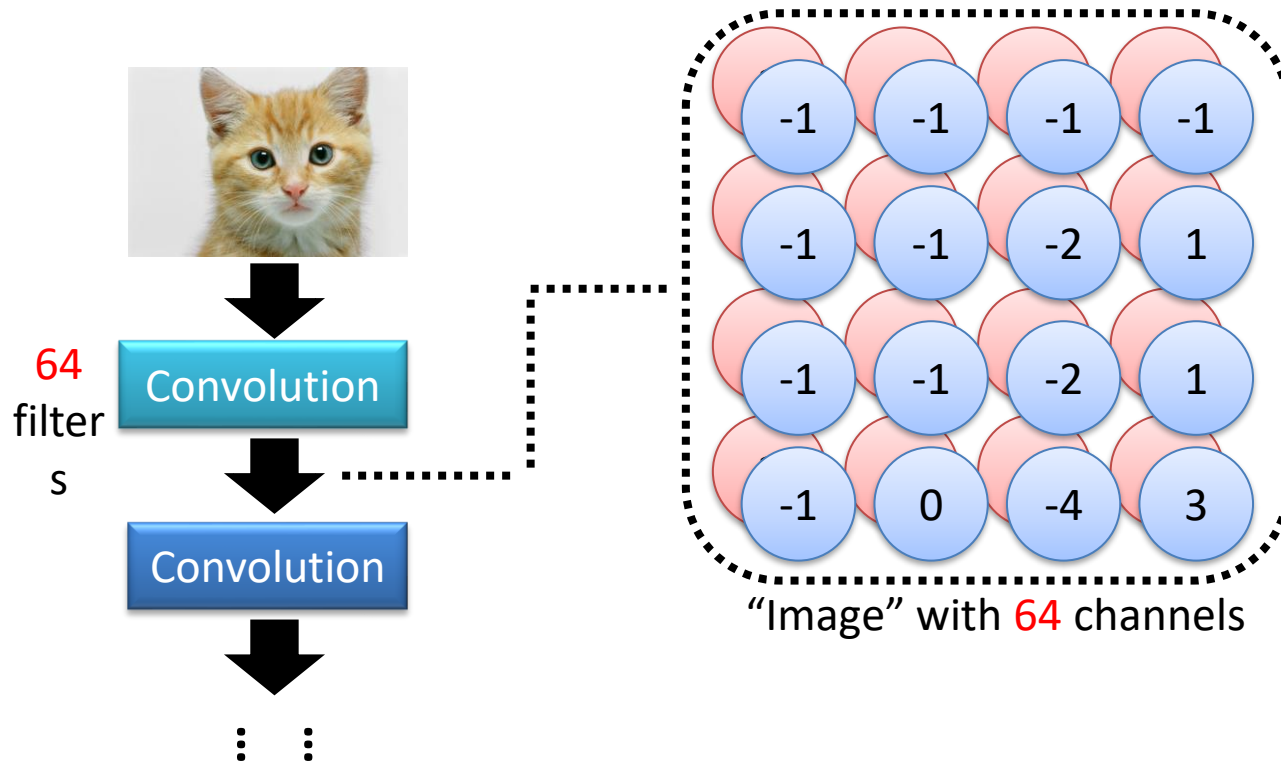| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 3 | -1 | -3 | -1 |
|---|---|---|---|
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

14

# Convolutional Layer

stride=1

| | | | | | |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

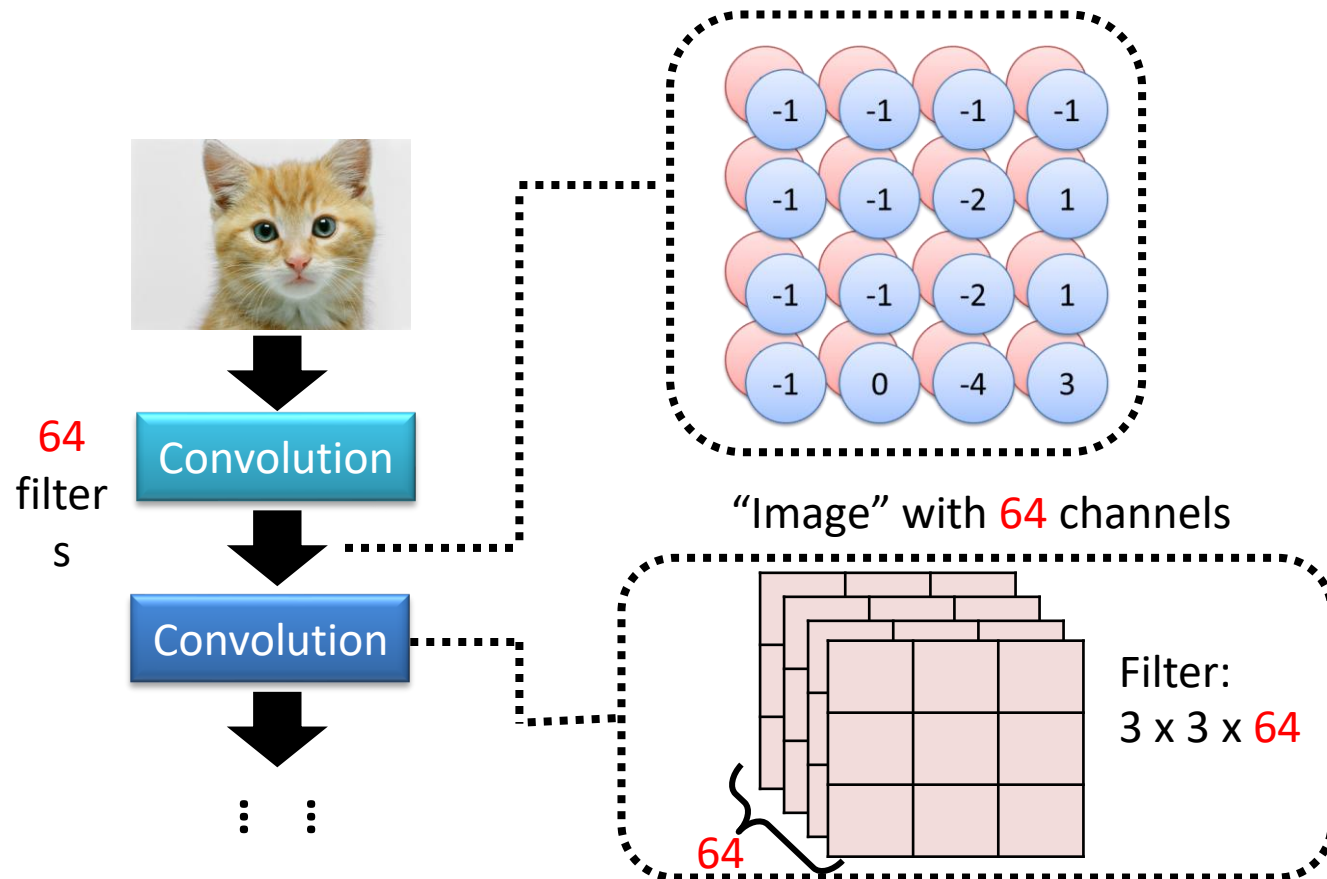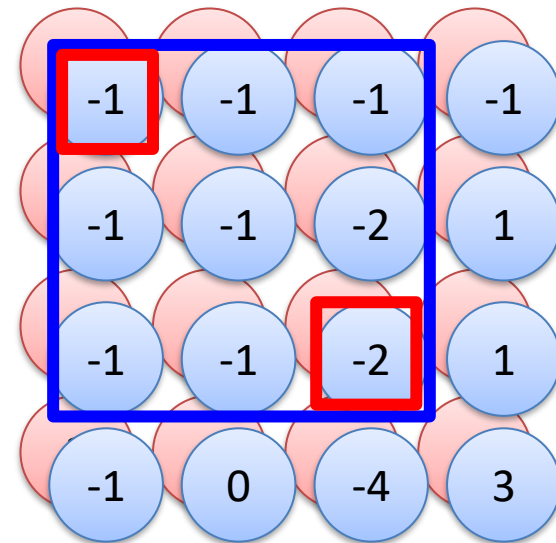| -1 | 1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

Do the same process for every filter

| -1 | -1 | -1 | -1 |
|---|---|---|---|
| -1 | | | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

Feature Map

# Convolutional Layer



64 filters

Convolution

Convolution

⋮ ⋮

|  |  |  |  |
|---|---|---|---|
| -1 | -1 | -1 | -1 |
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

"Image" with 64 channels

# Convolutional Layer



64 filters

Convolution

Convolution

⋮ ⋮

| -1 | -1 | -1 | -1 |
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

"Image" with 64 channels

Filter:
3 x 3 x 64

64

# Receptive Field（感受野）



64 filters

Convolution

Convolution

⋮ ⋮

| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

| -1 | -1 | -1 | -1 |
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

# Convolution Layer (卷积层)

Input: $X \in R^{d_{in} \times h \times w}$

weight: $W \in R^{d_{out} \times d_{in} \times F_h \times F_w}$

output: $Y \in R^{d_{out} \times h \times w}$

channels

Feature maps

Input: $x \in R^{d_{in}}$

weight: $W \in R^{d_{out} \times d_{in}}$

output: $y \in R^{d_{out}}$

$x_0 = 1$    $h_0^{(1)} = 1$

$x_1$

$x_2$

$x_3$

$\Sigma$

$\Sigma$

$\Sigma$

Input:

hidden layer

output

# Feature detection (特征检测)

- Learning filters (weights)

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution like template
  - <span style="color:red">Filters</span>
  - Why convolution for vision
- Pooling layer (module)

# Practice with linear filters(线性滤波器)



Original

| 0 | 0 | 0 |
|---|---|---|
| 0 | 1 | 0 |
| 0 | 0 | 0 |

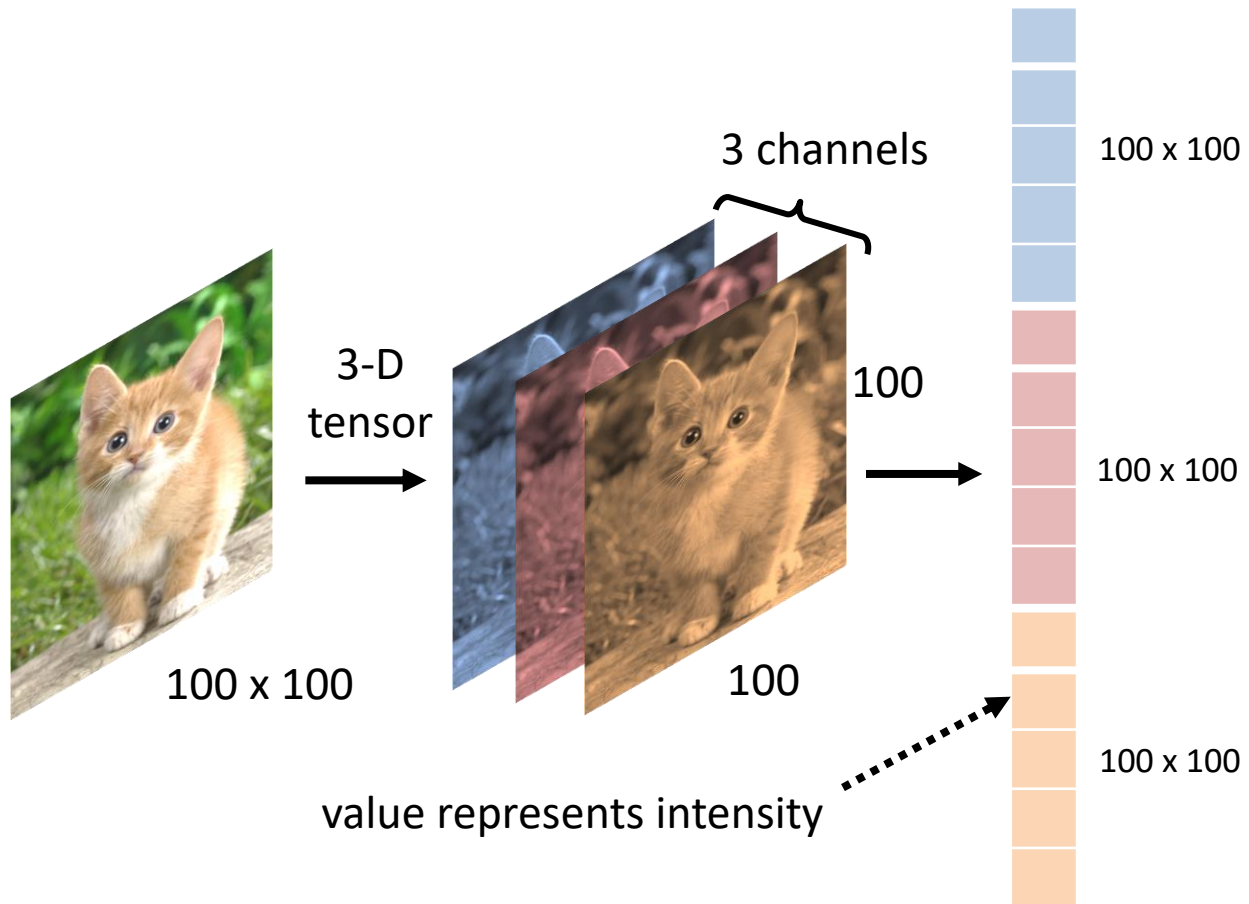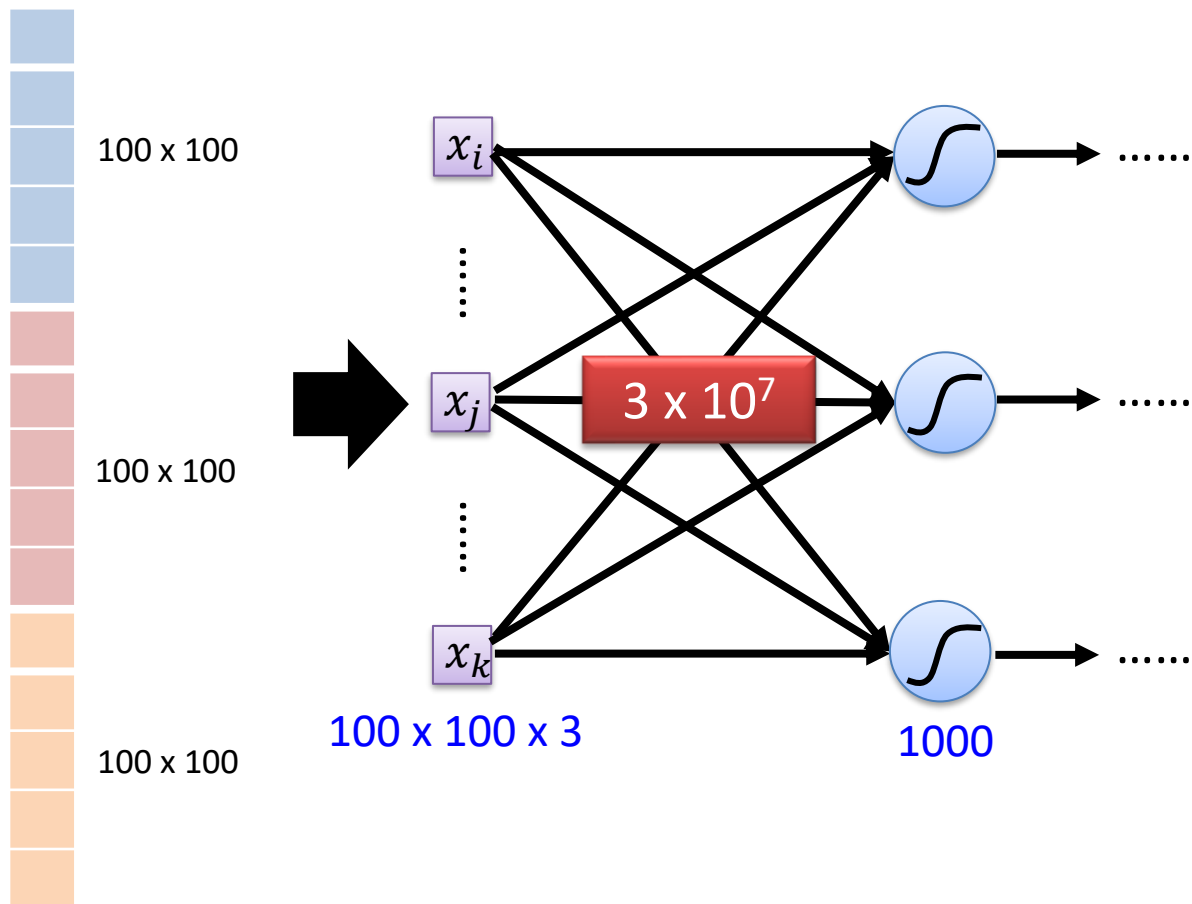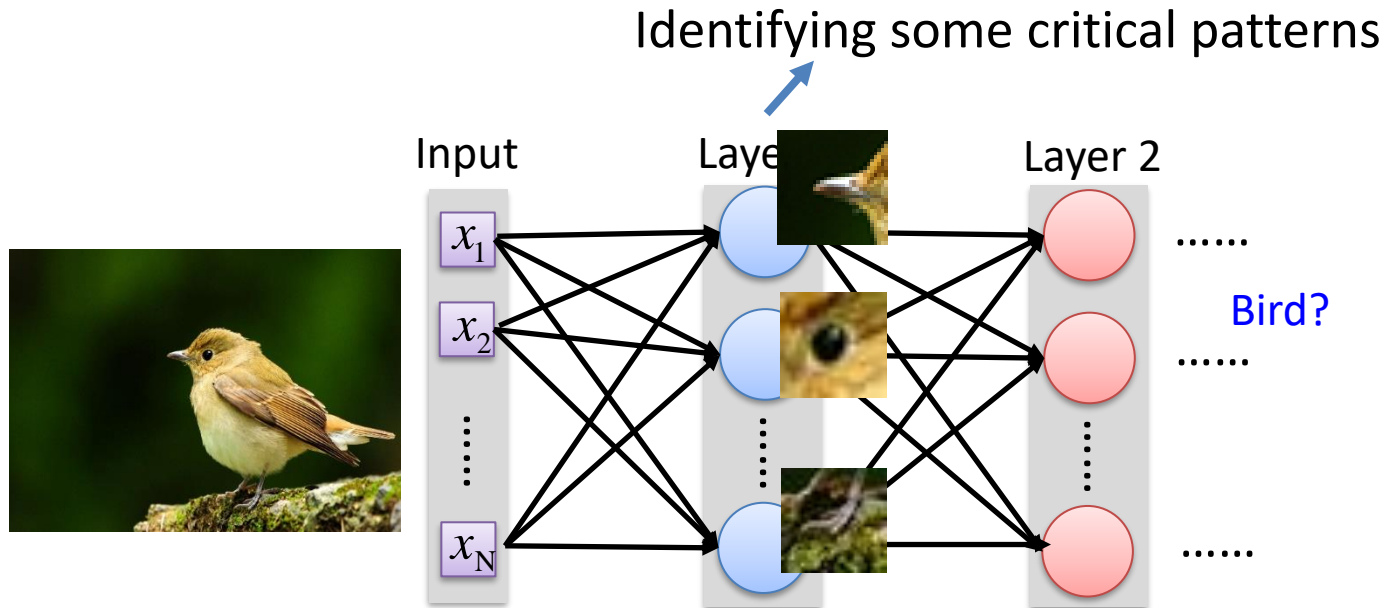Filter

Filtered
(no change)

# Practice with linear filters



Original

| 0 | 0 | 0 |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 0 | 0 |

Filter

Shifted *left*
By 1 pixel

# Practice with linear filters



Original

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

Filter

Blur (with a box filter)

# Practice with linear filters



Original

| 0 | 1 | 0 |
| 1 | -4 | 1 |
| 0 | 1 | 0 |

Filter

Output Image

**Edge detect** （边缘检测）

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution like template
  - Filters
  - Why convolution for vision
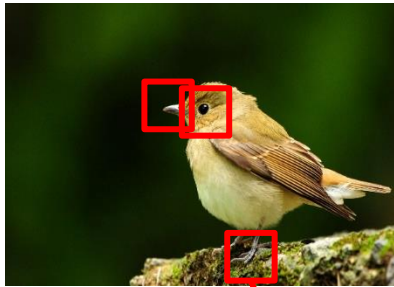- Pooling layer (module)

# Image Classification

3 channels

3-D tensor

100 x 100

100

100

value represents intensity

100 x 100

100 x 100

100 x 100

100 x 100

100 x 100

100 x 100

$x_i$

$x_j$

$x_k$

3 x 10$^7$

100 x 100 x 3

1000

......

......

......

# Vision model

## How MLP to CNN?

➢ Model locality
➢ Parameter Efficiency

# Observation 1

# Observation 1



Need to see the whole image?

A neuron does not have to see the whole image.

Some patterns are much smaller than the whole image.

Input   Layer 1   Layer 2

$x_1$

$x_2$

$x_N$

...... bird ......

basic detector

advanced detector

31

# Simplification 1

# Simplification 1



3 x 3 x 3 weights

Receptive field

the same receptive field

Can be overlapped

# Simplification 1 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).



The receptive fields cover the whole image.
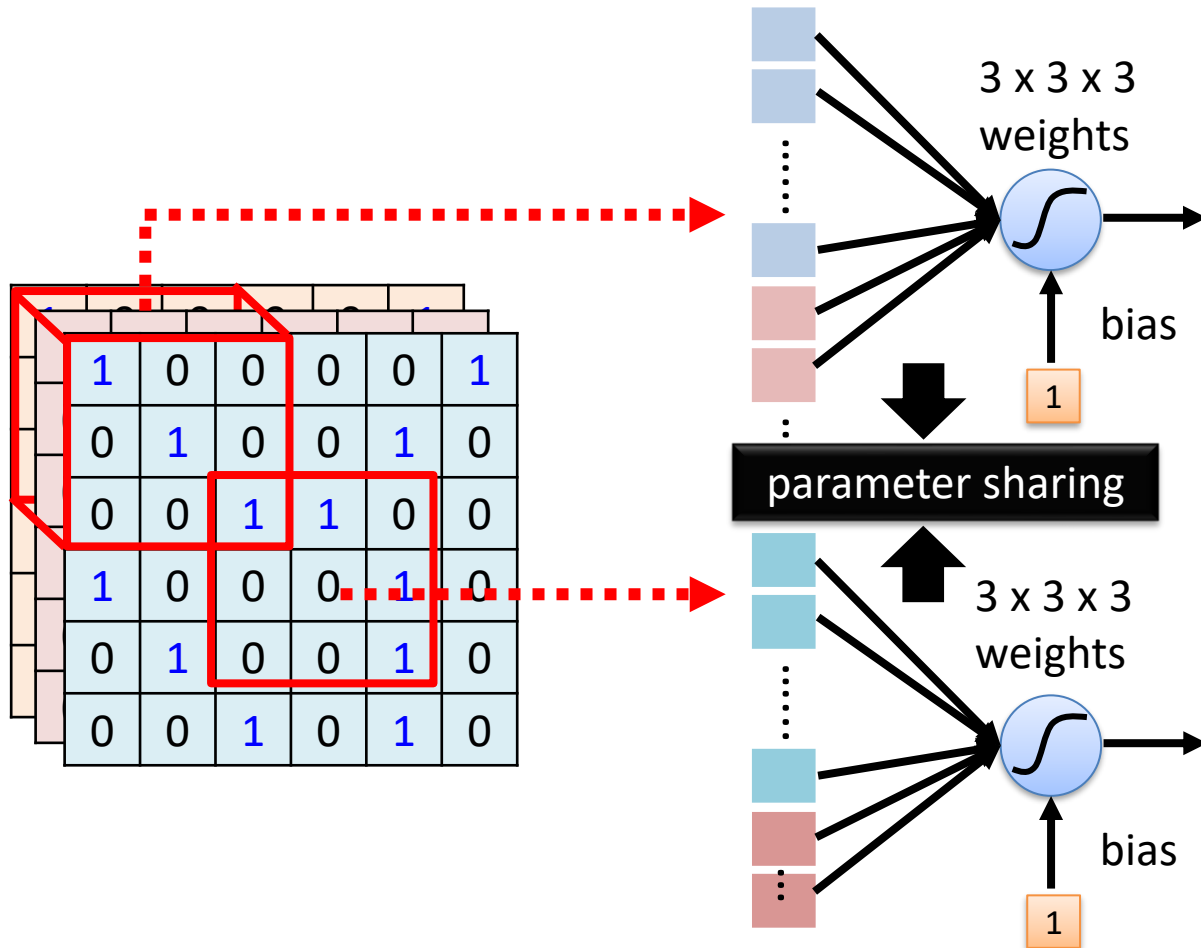
# Vision model

## How MLP to CNN?

➢ Model locality
➢ Parameter Efficiency

# Observation 2

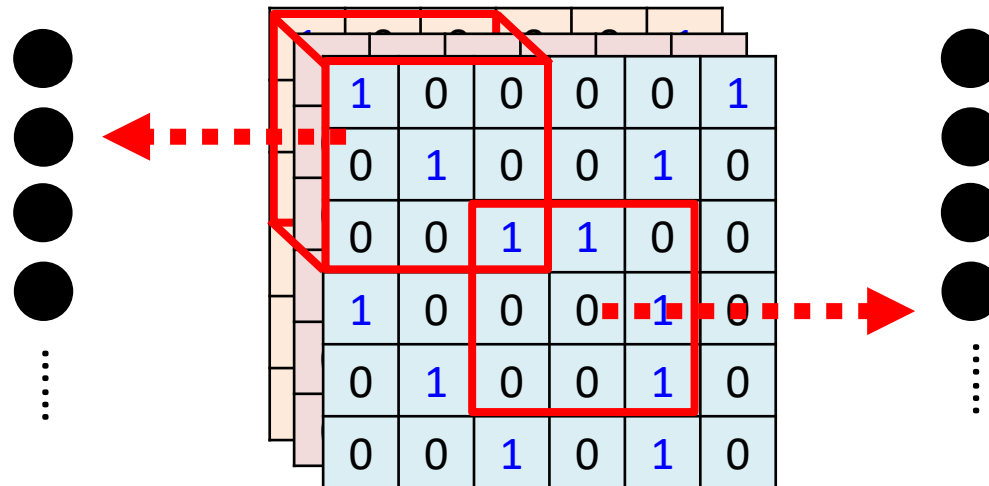- The same patterns appear in different regions.
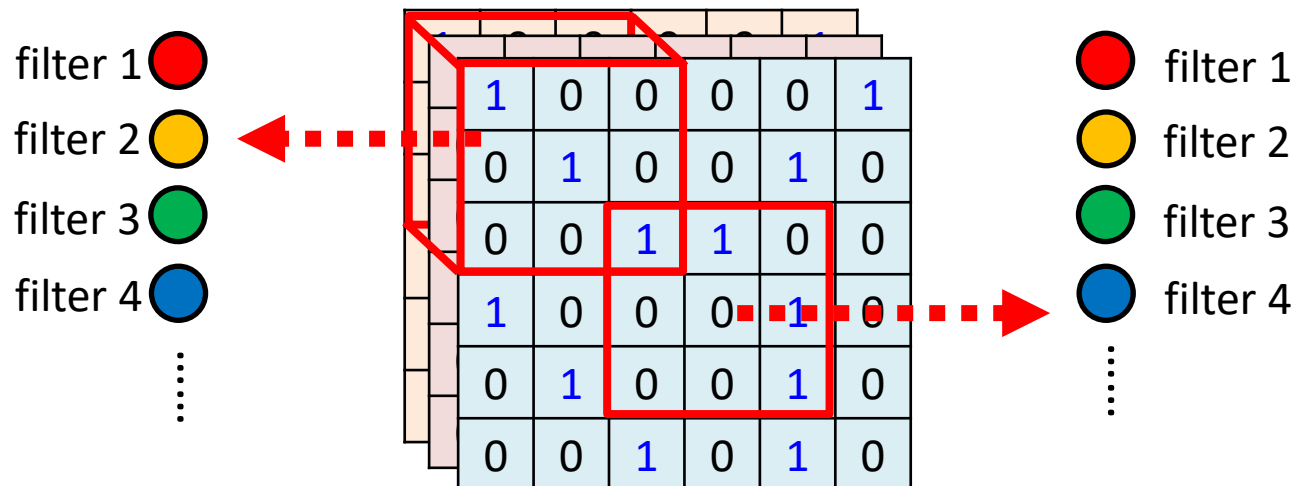
# Simplification 2

# Simplification 2



$x_1$    $\sigma(w_1 x_1 + w_2 x_2 + \cdots)$

$x_2$

$w_1$

$w_2$

bias

1

$x_1'$    $\sigma(w_1 x_1' + w_2 x_2' + \cdots)$

$x_2'$

$w_1$

$w_2$

bias

1

Two neurons with the same receptive field would not share parameters.

# Simplification 2 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).

# Simplification 2 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).

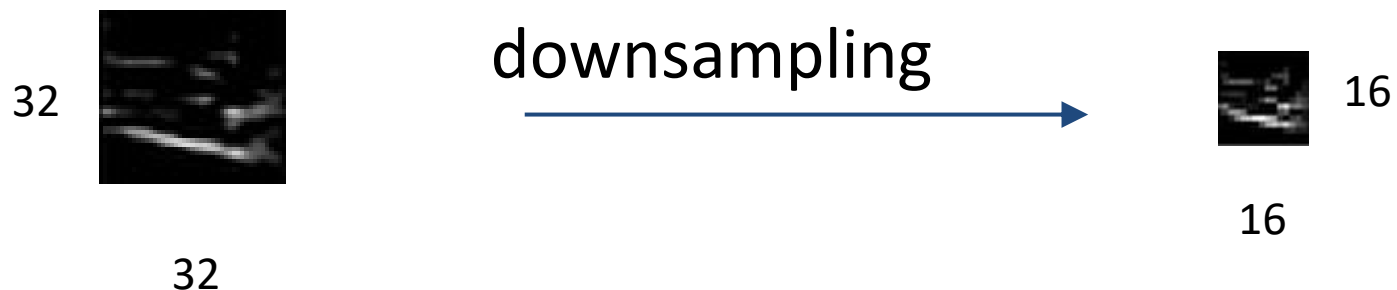Each receptive field has the neurons with the same set of parameters.

# outline

- Modeling of CNN
  - Module-wise architecture模块化结构
- Convolutional layer (module)
  - Convolution in general
  - Filters
  - Convolution module
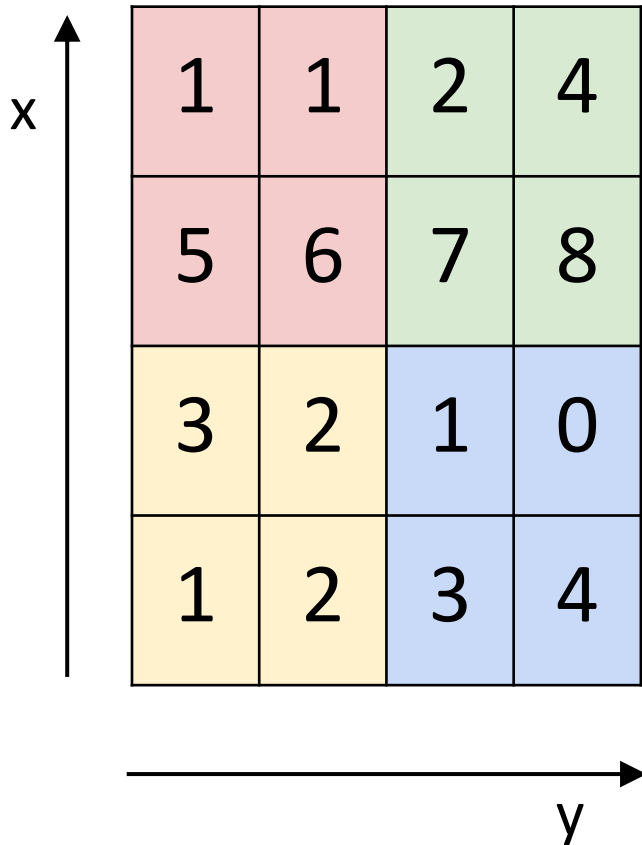- <span style="color:red">Pooling layer (module)</span>

# POOLING Layer

- In ConvNet architectures, **Conv** layers are often followed by **Pooling** layers
  - makes the representations smaller and more manageable without losing too much information.
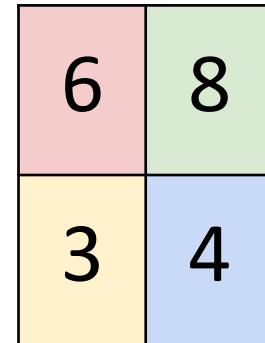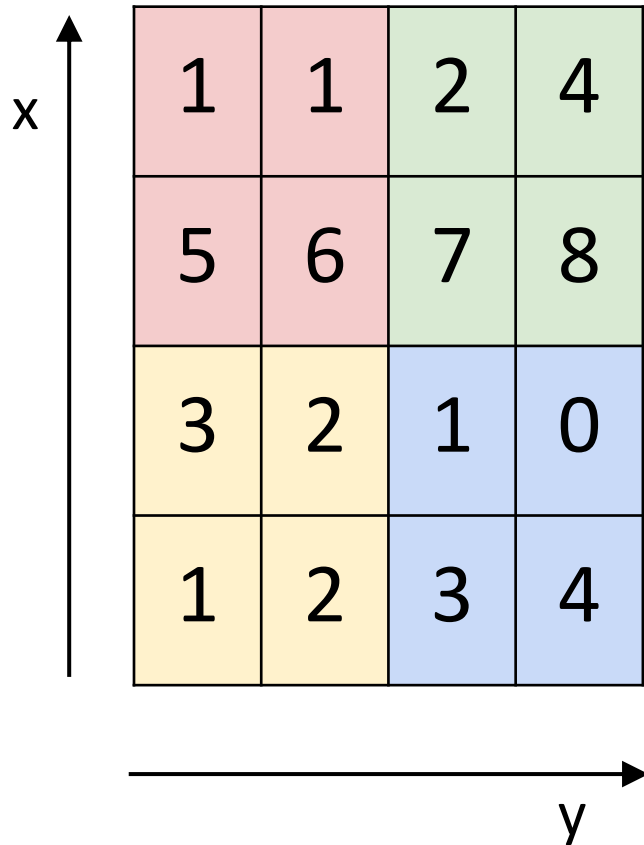  - Invariant in region.

32

32

downsampling

16

16

# MAX POOLING

Single depth slice

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

x

y

max pool with 2x2 filters and stride 2

| | |
|---|---|
| 6 | 8 |
| 3 | 4 |

# Average POOLING

Single depth slice

x

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

y

average pool with 2x2 filters and stride 2

| 4.25 | 5.25 |
|------|------|
| 2 | 2 |

# Intuitive example



CONV ReLU CONV ReLU POOL CONV ReLU CONV ReLU POOL CONV ReLU CONV ReLU POOL FC (Fully-connected)
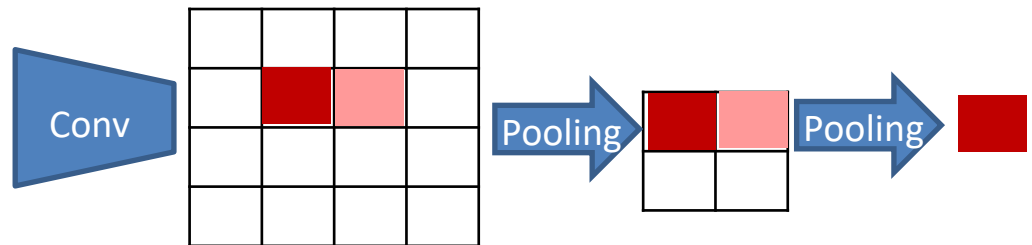
truck
car
airplane
ship
horse

Source: Andrej Karpathy & Fei-Fei Li

# The characteristics of CNN

- Equivalent (等变性)----object detection
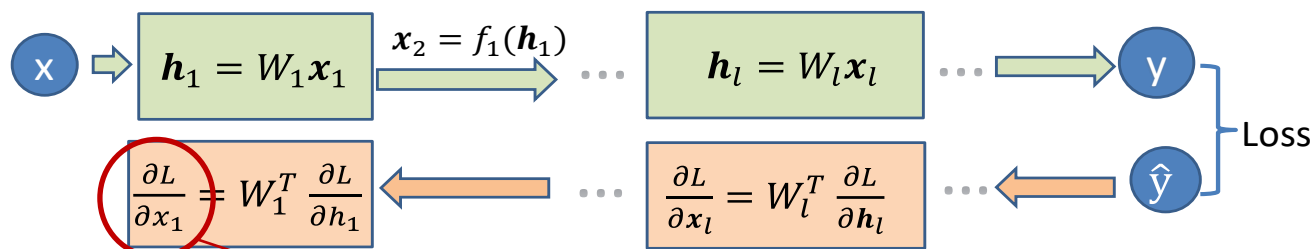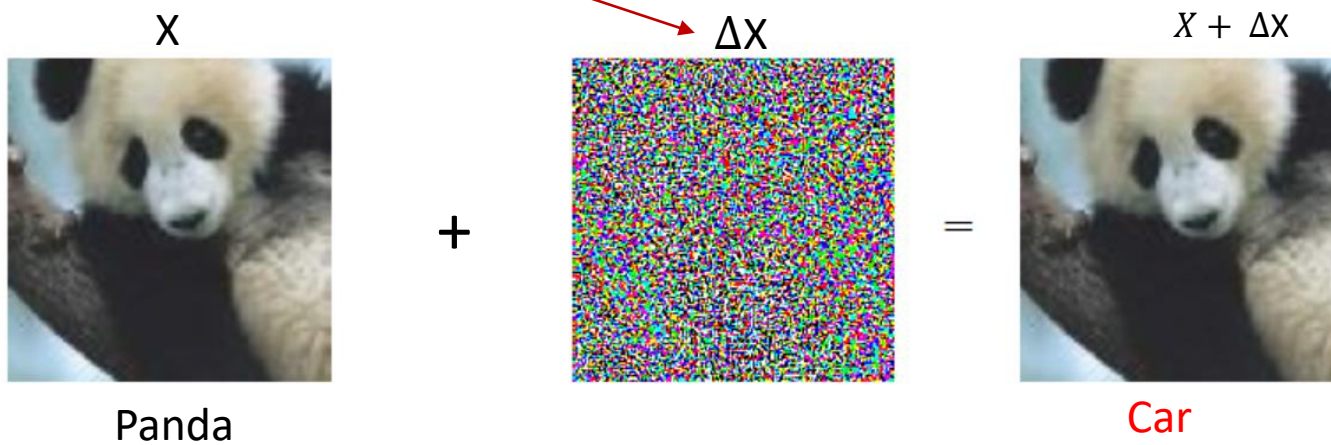- Invariant (不变性)----Image Classification



receptive field （感受野）

Activation on Feature map
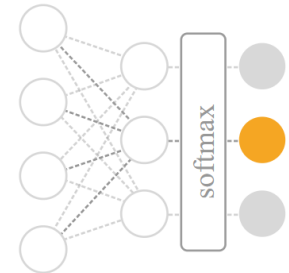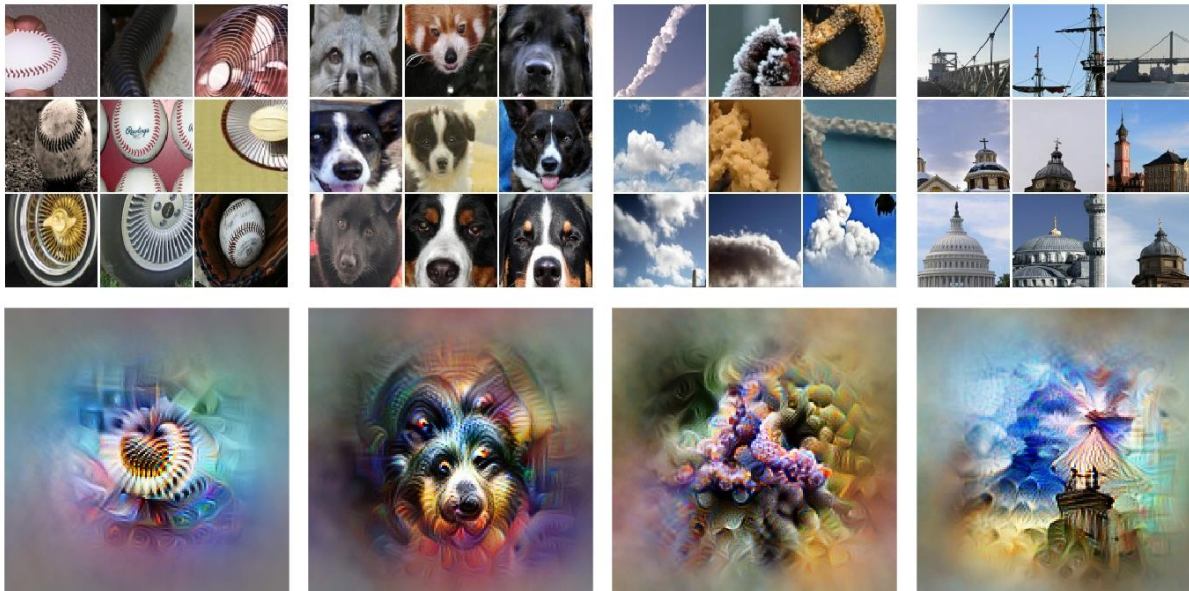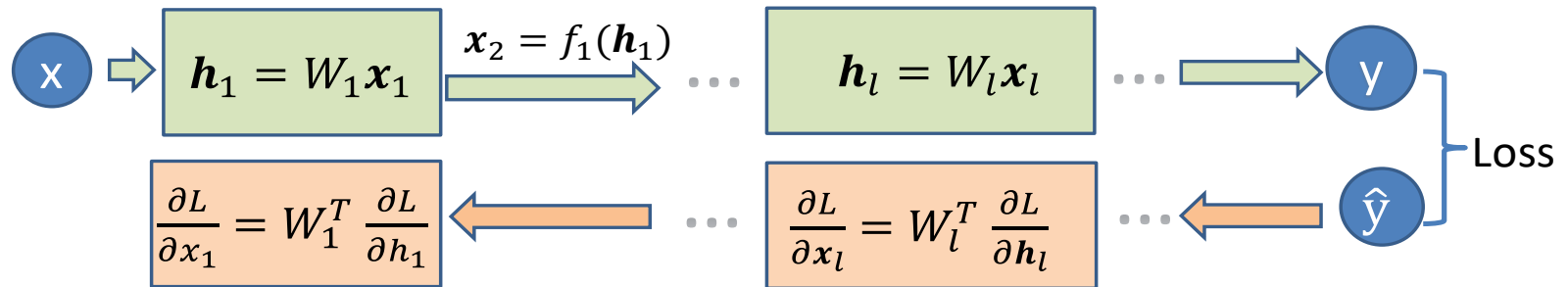
# Deep neural networks visualization
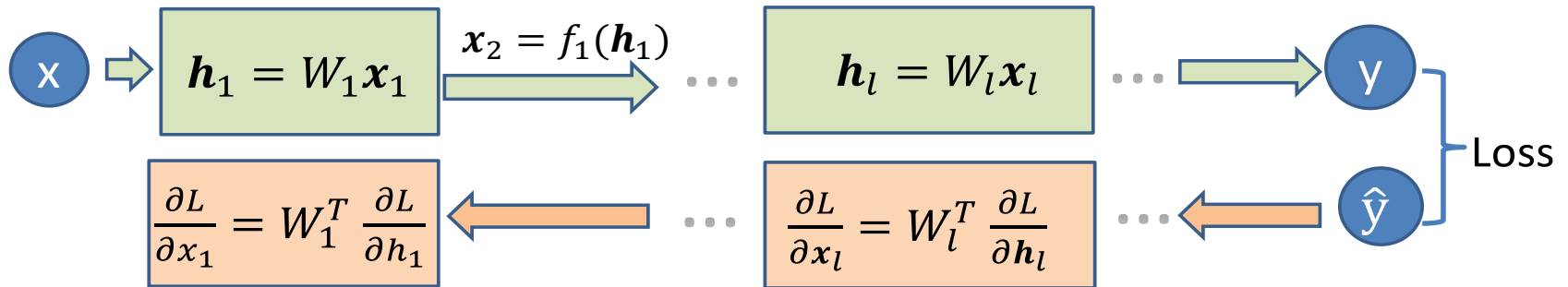
- 对抗样例（Adversarial example ）



$f(X + \Delta X) \neq f(X),$  $\Delta X$ should be imperceptible by human

X        $\Delta$X        $X + \Delta X$

+        =

Panda                                Car

# Deep neural networks visualization

# Neural Style Transfer



$$\boldsymbol{h}_1 = W_1 \boldsymbol{x}_1$$

$$\boldsymbol{x}_2 = f_1(\boldsymbol{h}_1)$$

$$\boldsymbol{h}_l = W_l \boldsymbol{x}_l$$

x

y

$$\frac{\partial L}{\partial \boldsymbol{x}_1} = W_1^T \frac{\partial L}{\partial \boldsymbol{h}_1}$$

$$\frac{\partial L}{\partial \boldsymbol{x}_l} = W_l^T \frac{\partial L}{\partial \boldsymbol{h}_l}$$
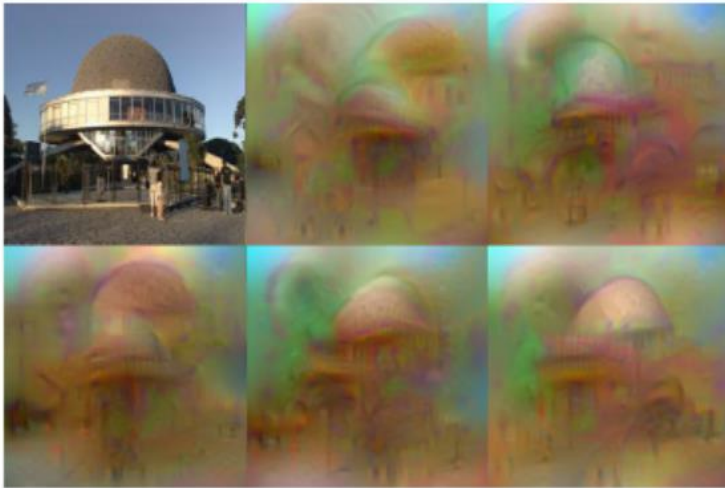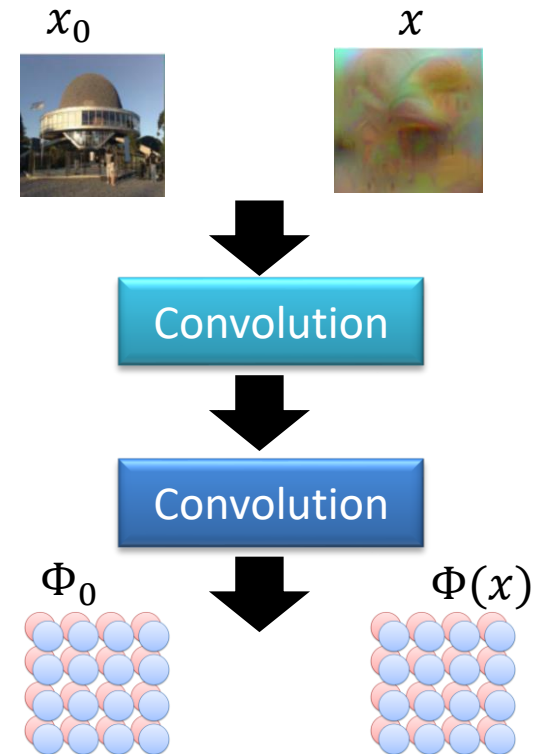
$\hat{y}$

Loss

➢ Style
➢ Content

# Reconstructing an image from a convolutional layer

- Representation function: $\Phi : \mathfrak{R}^{H \times W \times C} \to \mathfrak{R}^d$ (image space to feature space)

- Target Representation: $\Phi_0 = \Phi(x_0)$ (x0 is the original image)

- We need to find: $x \in \mathfrak{R}^{H \times W \times C}$ by minimizing:

$$x^* = \arg\min_{x \in \mathfrak{R}^{H \times W \times C}} l(\Phi(x), \Phi_0) + \lambda R(x)$$



"Understanding Deep Image Representations by Inverting Them", by Aravindh Mahendran and Andrea Vedaldi.
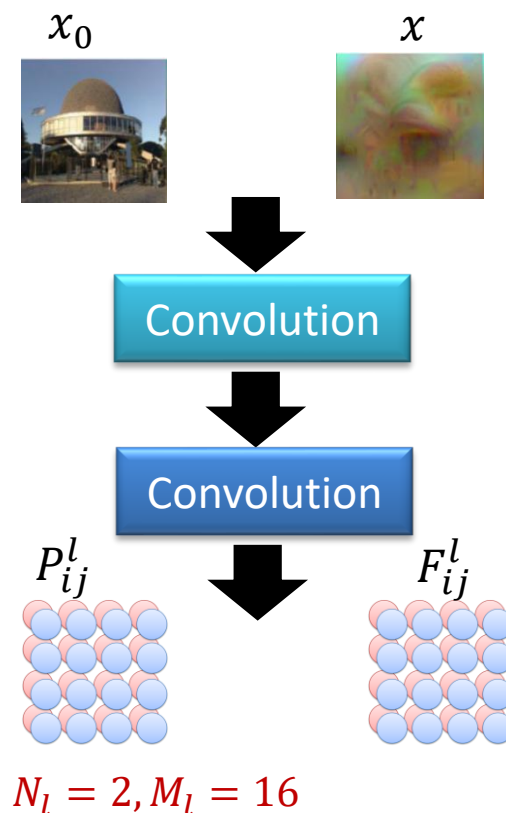
# Content Loss Function

- Filters (Depths) at layer l: $N_l$

- The height times the width of the feature map at layer l: $M_l$

- Response at layer l: $F_l \in \mathfrak{R}^{N_l \times M_l}$

  $F^l_{ij}$ represents the ith filter at position j in layer l

- Original image: $\vec{p}$

- We generate image: $\vec{x}$ (randomly initialized)

- Squared-error loss:

$$L_{content} = \frac{1}{2} \sum_{i,j} (F^l_{ij} - P^l_{ij})^2$$

$x_0$      $x$

Convolution

Convolution

$P^l_{ij}$      $F^l_{ij}$

$N_l = 2, M_l = 16$

# Style Loss Function

- Filter correlations are given by the Gram matrix:

$$G^l \in \Re^{N_l \times N_l}$$

- $G^l_{ij}$ is the inner product between the filters i and j in layer l:

$$G^l_{ij} = \sum_k F^l_{ik} F^l_{jk}$$

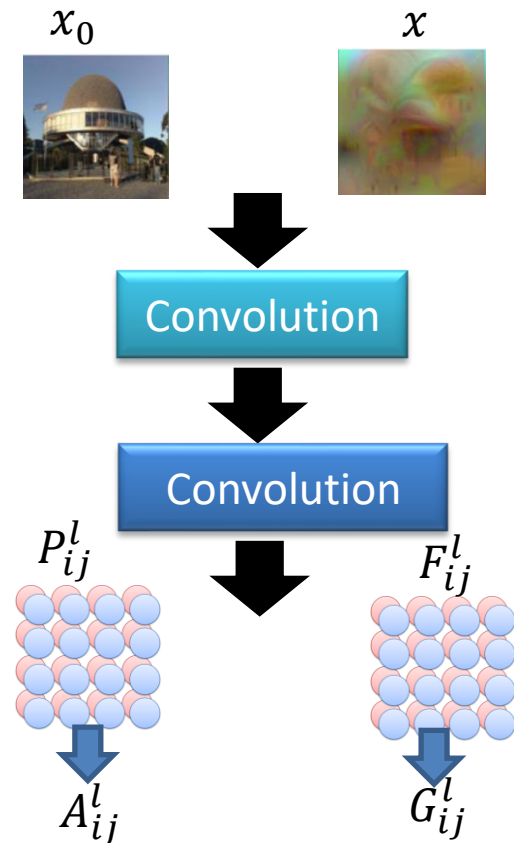- The loss at layer l:

$$E_l = \frac{1}{4N^2_l M^2_l} \sum_{i,j} (G^l_{ij} - A^l_{ij})^2$$

A <-> original image
G <-> generated image

- The total style loss:

$$L_{style} = \sum_{l=0}^{L} w_l E_l$$

$x_0$      $x$

Convolution

Convolution

$P^l_{ij}$      $F^l_{ij}$

$A^l_{ij}$      $G^l_{ij}$
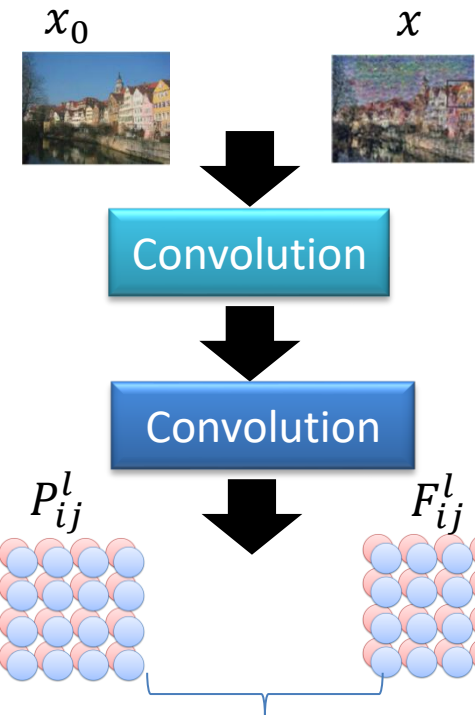
# Content Reconstruction



Image reconstructed from layers
(a)'conv1_1',
(b)'conv2_1',
(c)'conv3_1',
(d)'conv4_1' and
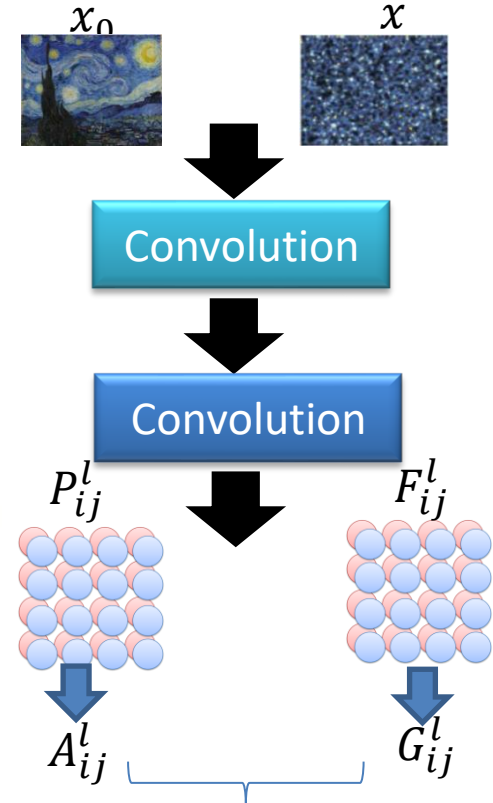(e)'conv5_1'
of the original VGG-Network

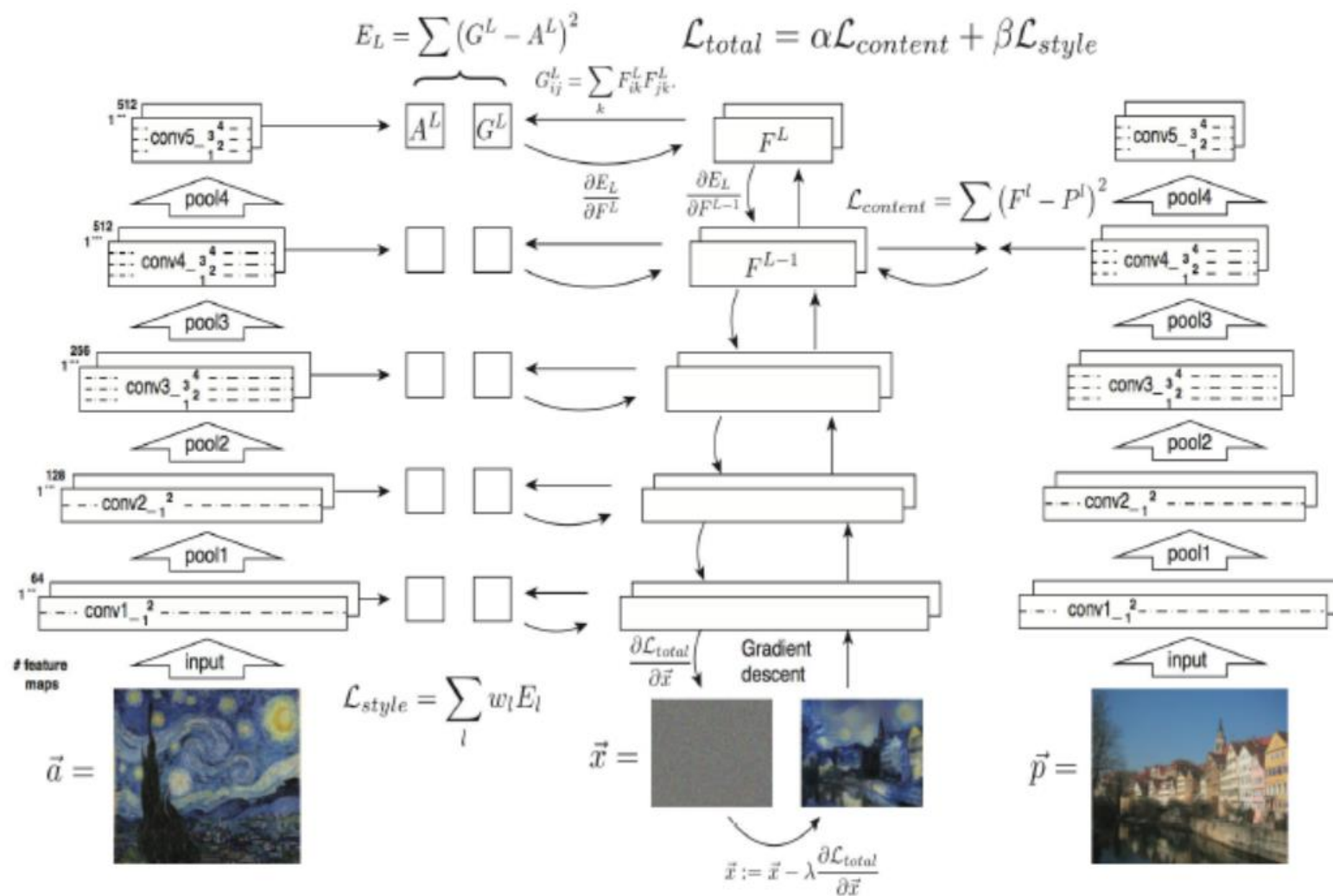$$L_{content} = \frac{1}{2} \sum_{i,j} (F^l_{ij} - P^l_{ij})^2$$

# Style Reconstruction



Style representations (filer correlations) from:
(a) 'conv1_1',
(b) 'conv1_1', 'conv2_1',
(c) 'conv1_1', 'conv2_1', 'conv3_1',
(d) 'conv1_1', 'conv2_1', 'conv3_1', 'conv4_1',
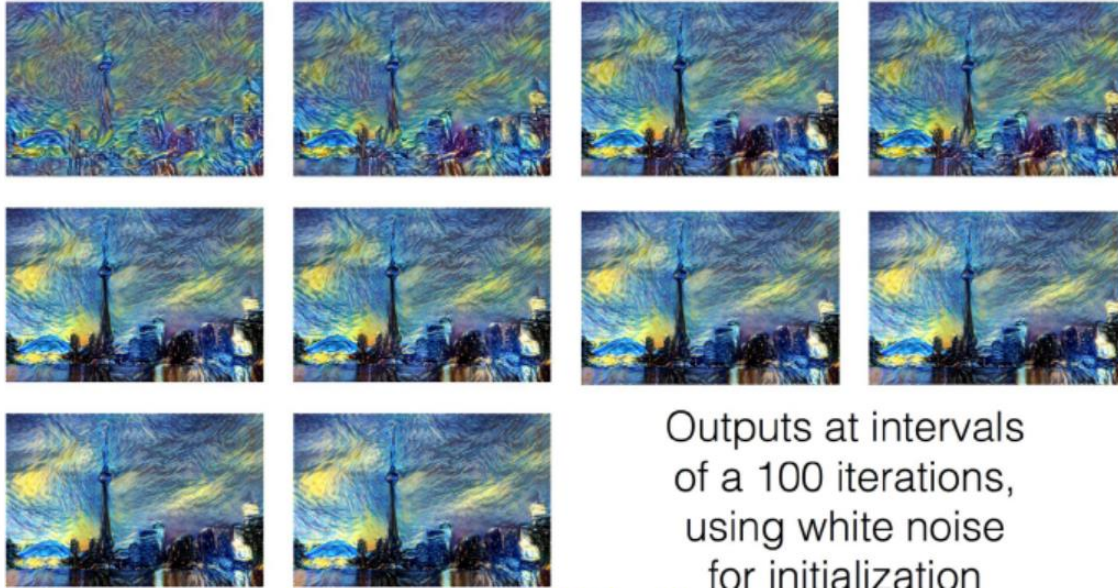(e) 'conv1_1', 'conv2_1', 'conv3_1', 'conv4_1', 'conv5_1'.

$x_0$     $x$

Convolution

Convolution

$P_{ij}^l$     $F_{ij}^l$

$A_{ij}^l$     $G_{ij}^l$

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

# The Total Loss Function

# Results



Outputs at intervals
of a 100 iterations,
using white noise
for initialization

show image every 10 iterations

# 谢谢！