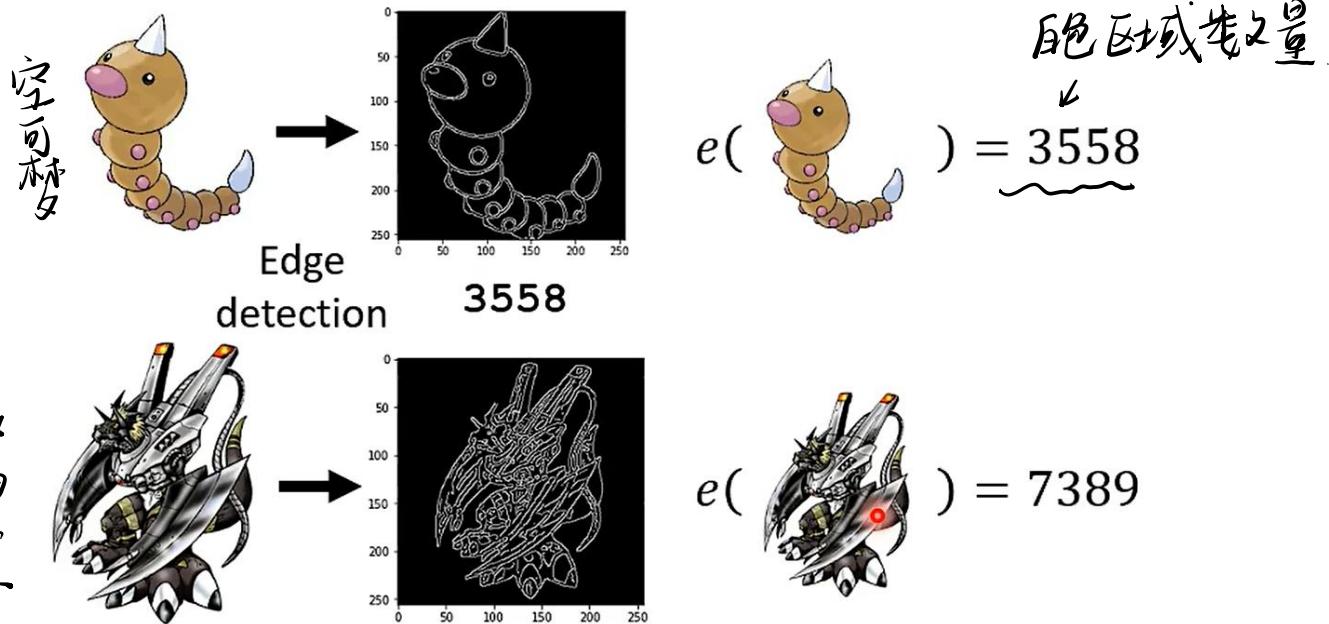


更多的parameter, easier to overfit ! why ?

对于数码宝贝与宝可梦：

$$f(\text{数码宝贝}) = \text{宝可梦 or 数码宝贝}$$

由图可知：数码宝贝比较复杂，宝可梦比较简单



可以用 Edge detection 技术 将图变成一个黑白图，白也是线边
7389

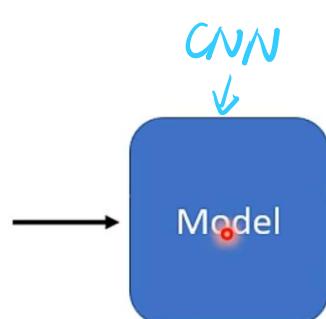
模型复杂度高，代表有很多的未知参数 (feature)

模型的好坏取决于取到何 train，若不是独立同分布的话，train 的模型会
较差、则在 model 遇到一个 new test 时，效果会更差。

具体的公式推导没看良，在 P23 的视频中讲的比较详细。

CNN的第一种解释：

Image Classification



y'

$$\begin{bmatrix} \vdots \\ 0.2 \\ 0.7 \\ 0.1 \\ \vdots \end{bmatrix}$$

这有多少个数，就代表
可以辨识多少只动物

$$\begin{array}{l} \uparrow \\ \begin{bmatrix} \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \end{array}$$

$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \rightarrow \text{cat}$

可辨识4个，其中
输入的图片是 cat 的结果

(All the images to be classified have the same size.)

原本是概率， \Rightarrow 过一个 cross entropy 得到分类结果

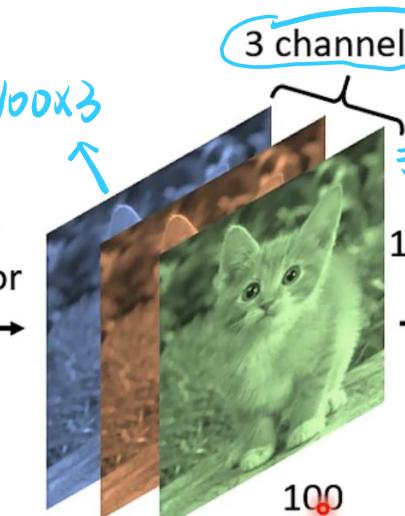
Image Classification



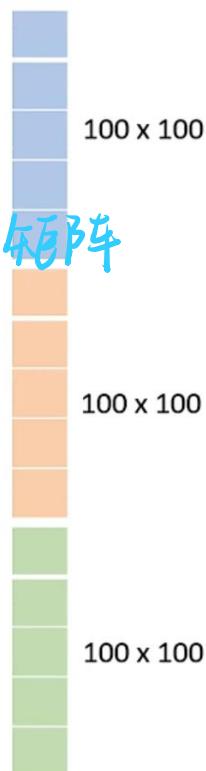
100 x 100 x 3

3-D tensor

100



100

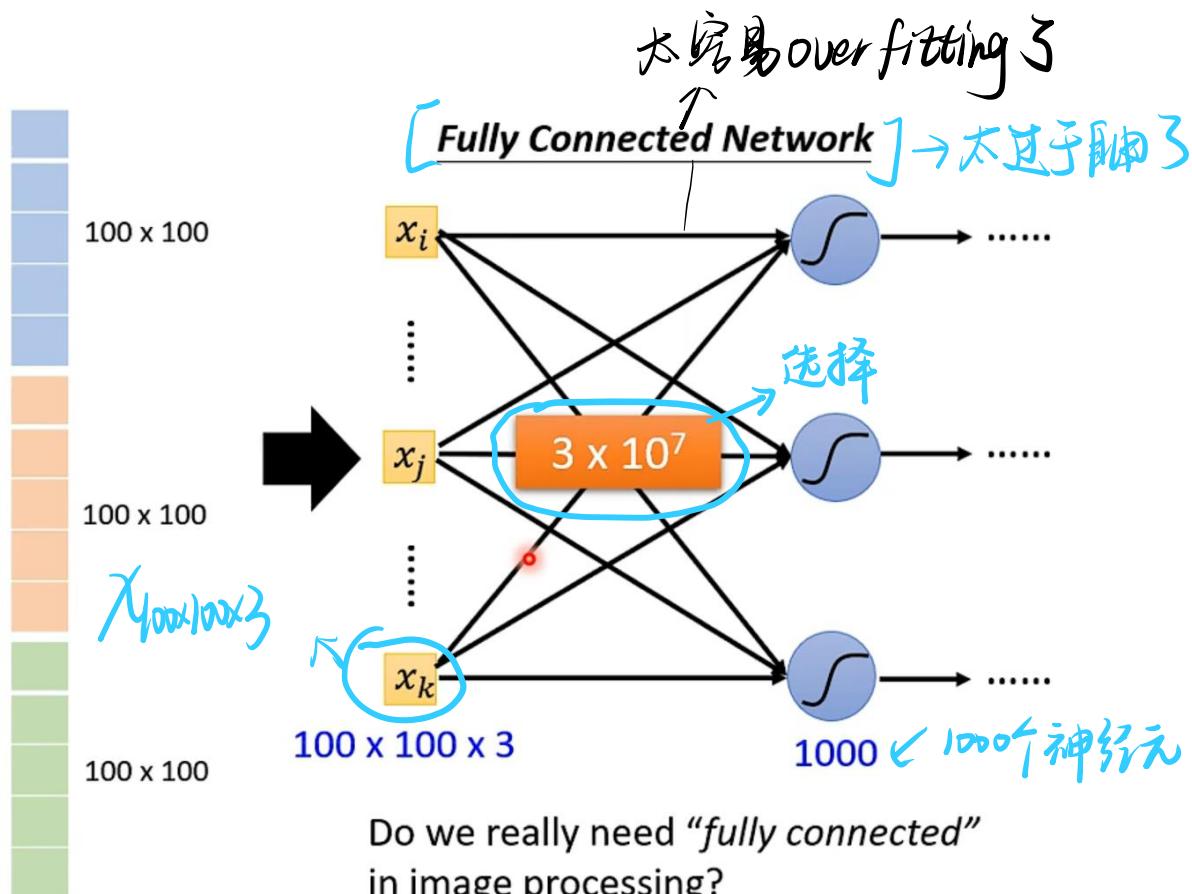


对于一个黑白图：channel = 1

一个图片是 RGB，即三色

故 channel = 3

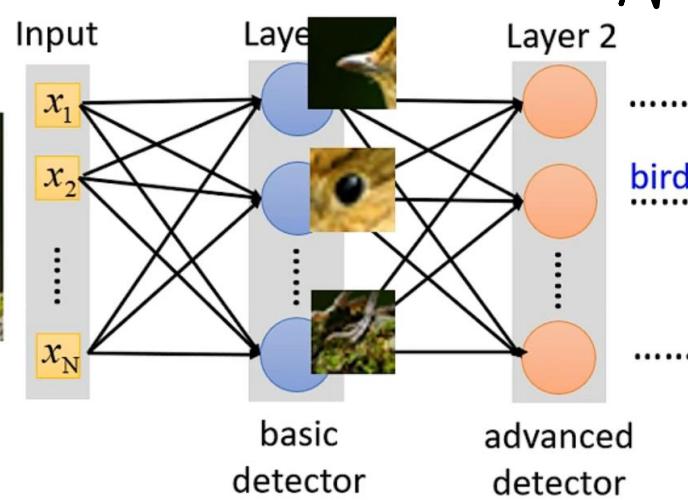
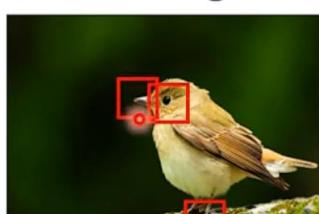
如果我们将这个 picture 的 $3 \times 100 \times 100$ 进行 Fully Connected (全连接)



Do we really need "fully connected" in image processing?

但是：我们观察一张图片，要把全部都放进去吗？ → NO, 一般只用看局部便能识别

Need to see the whole image?

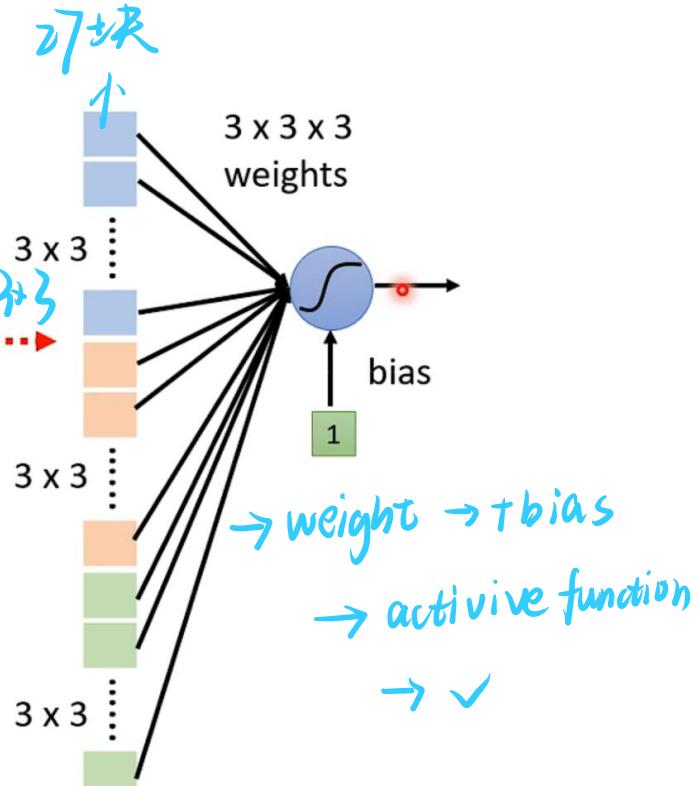
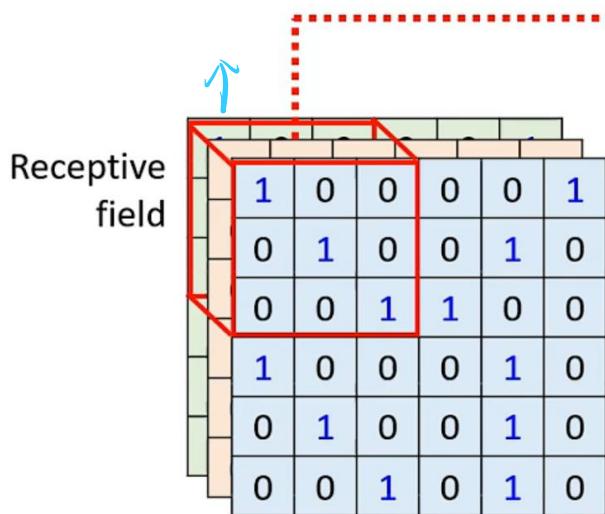


Some patterns are much smaller than the whole image.

Receptive field

Simplification 1

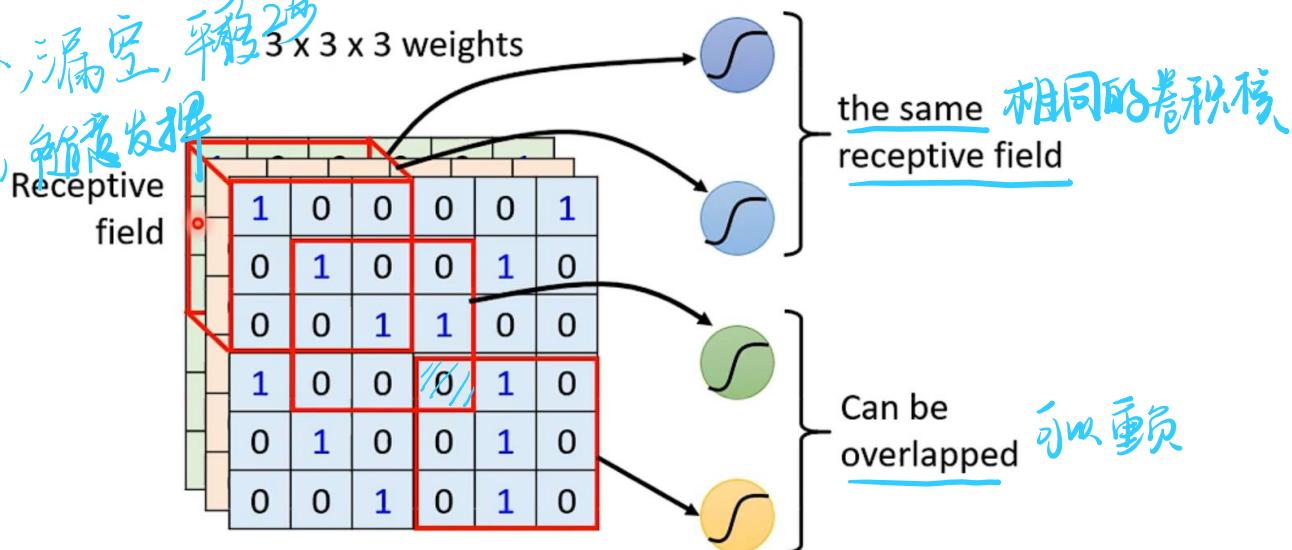
选择3x3的区域
→将其拉成3x3



- Can different neurons have different sizes of receptive field?
- Cover only some channels?
- Not square receptive field?

Simplification 1

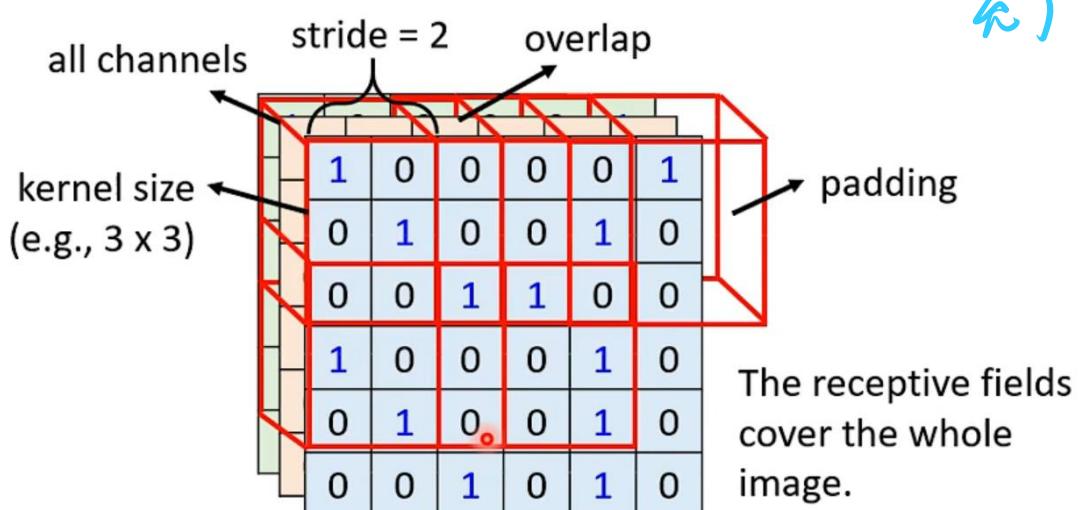
放大，缩小，漏空，平面 3x
不同形状，能发挥



经典1: Kernel size = 3x3 (卷积核)

存在部分重叠，多余部用padding(可用全0填充, 可用平均填充, 也可用补逆填充)

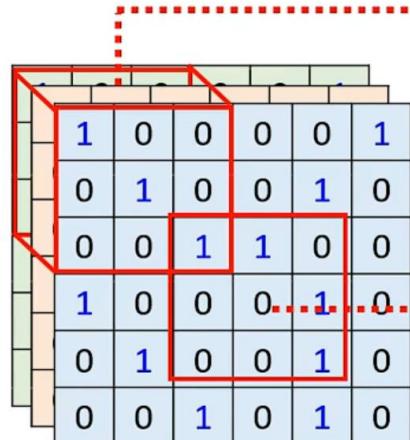
Each receptive field has a set of neurons (e.g., 64 neurons).



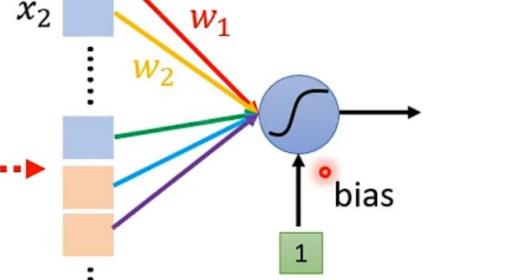
经典2: sharing parameters.

来源是因为有些图片是位置不一, 而使用的参数数量不多, 才想到
3 sharing parameters.

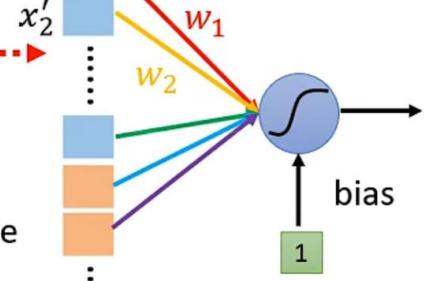
Simplification 2



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



$$\sigma(w_1 x'_1 + w_2 x'_2 + \dots)$$



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

对于不同的卷积层，导致相同的结果吗？

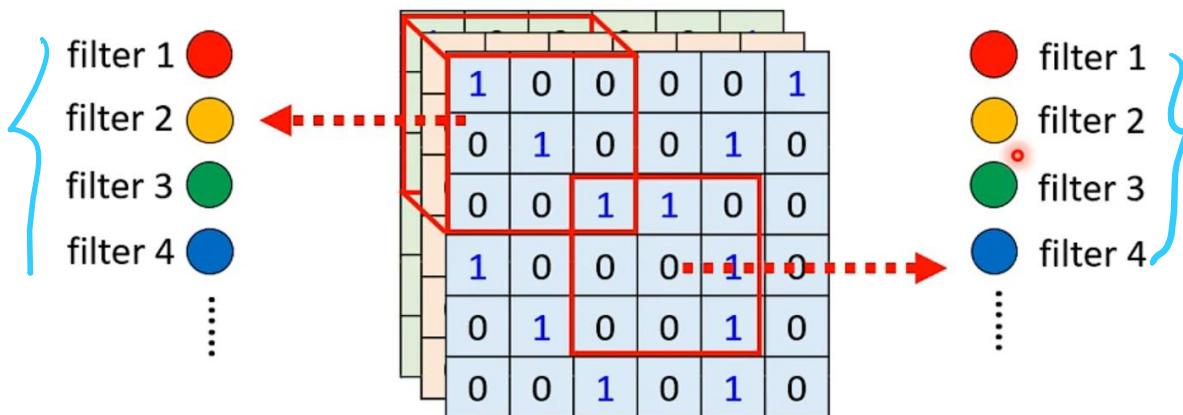
problem? weight & bias相同，会



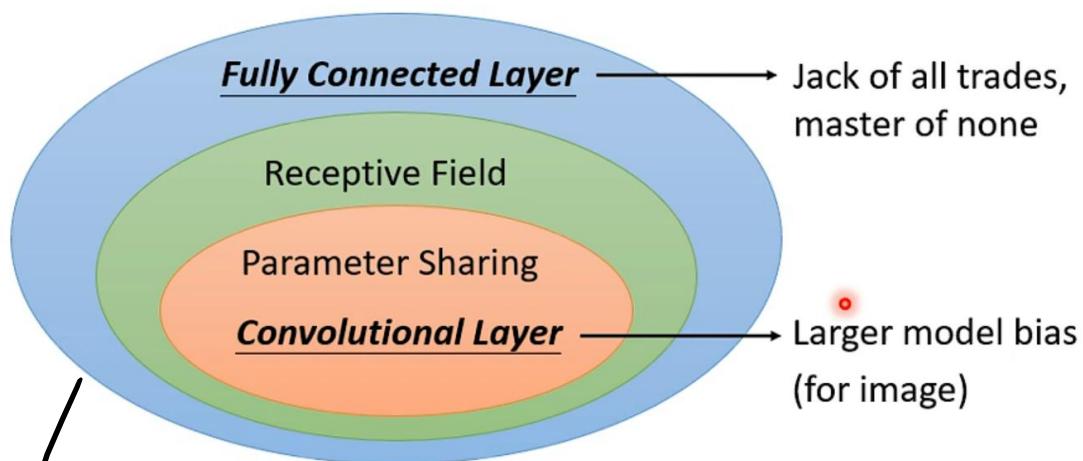
Answer = 不会，因为输入不同，其输出

Each receptive field has a set of neurons (e.g., 64 neurons). 也会不同！

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



Benefit of Convolutional Layer



- Some patterns are much smaller than the whole image.
- The same patterns appear in different regions.

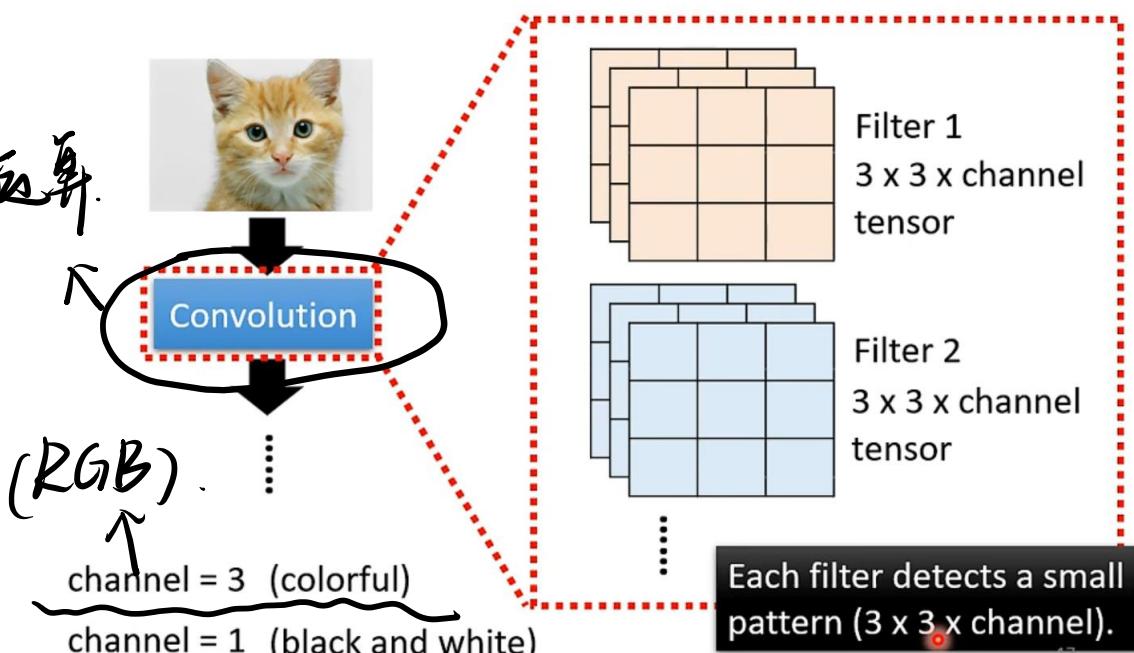
- 层-层的限制，使 fully Connected 并没有那么 freestyle

而 CNN 却是在共享参数中的先验得来的

第二种常见的关于CNN的解释：

Convolutional Layer

卷积运算



那 convolution 如何运算 =

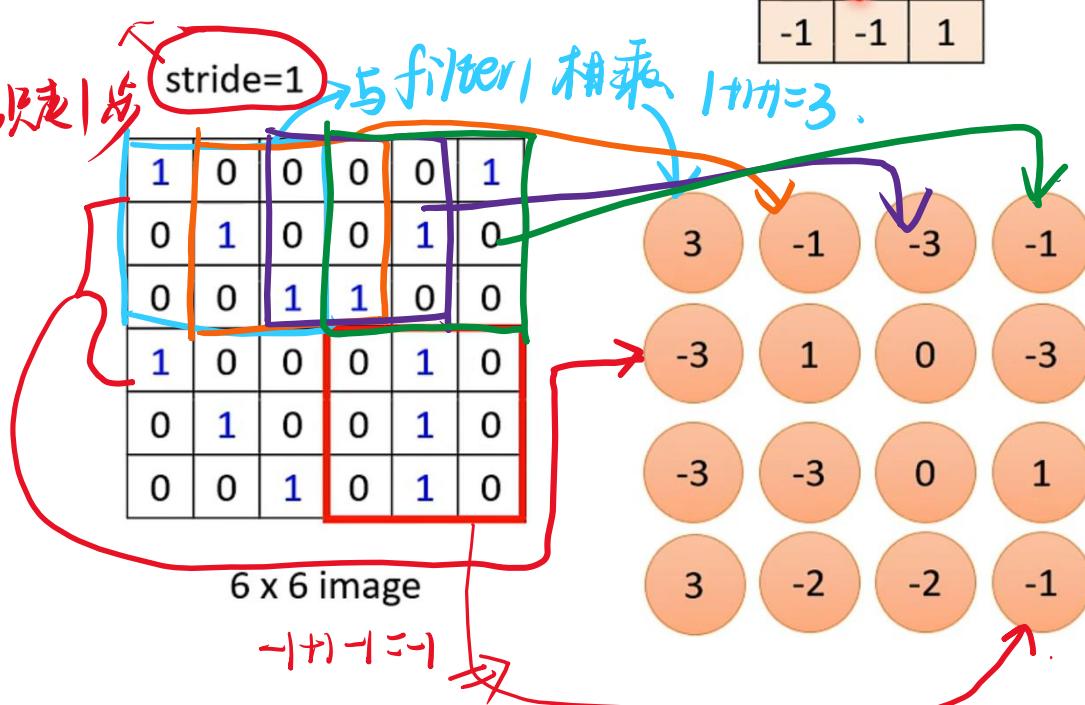
→ convolution 是 Part of filter

step=1
↓
每次平移只走1步

Convolutional Layer

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



Convolutional Layer

stride=1

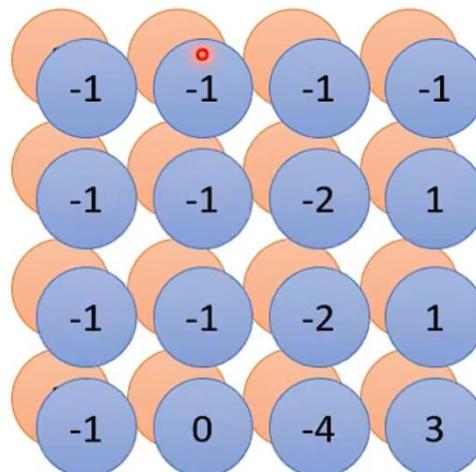
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

用 filter2



Do the same process for every filter



Convolutional Layer

故而有多少个 Filter, 就有多少层结果
这些结果的总和被称为什么

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

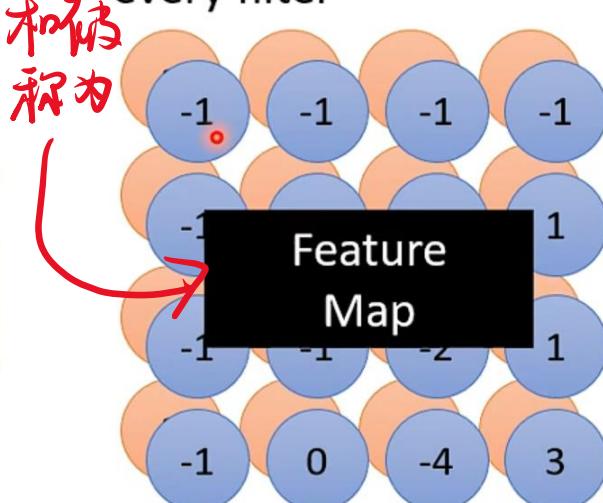
6 x 6 image

若有 5 个 filter, 故而有 5 x 4 x 4

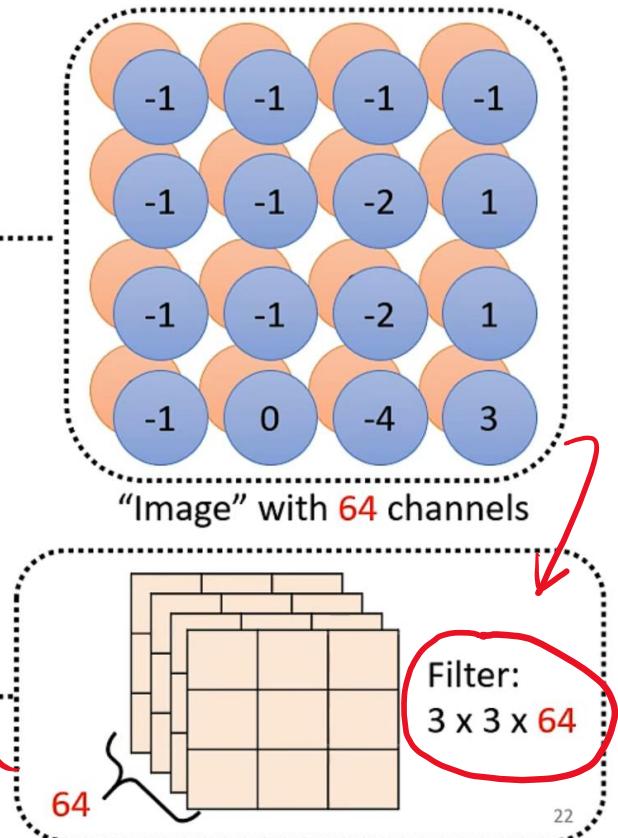
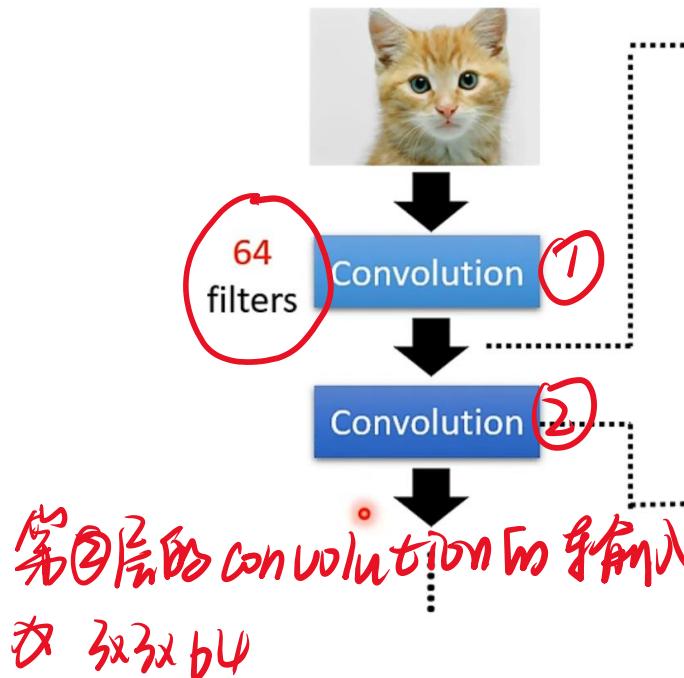
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

Do the same process for every filter

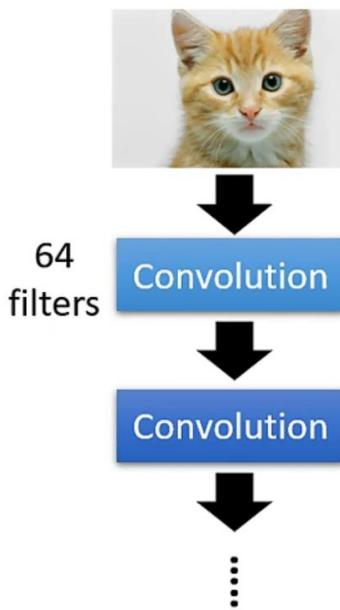


Multiple Convolutional Layers

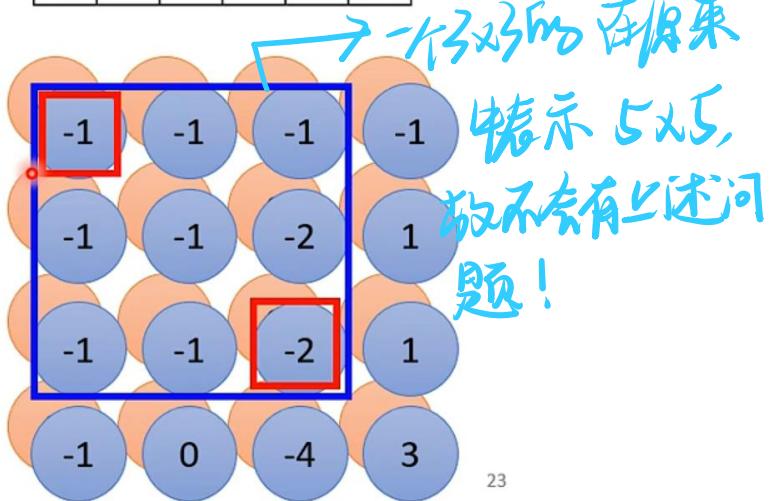


problem: 很明显从626变成424,会不会有点消失或者扭曲呢?

Multiple Convolutional Layers

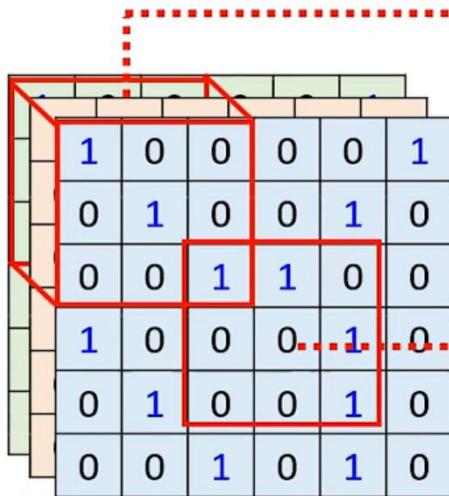


1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

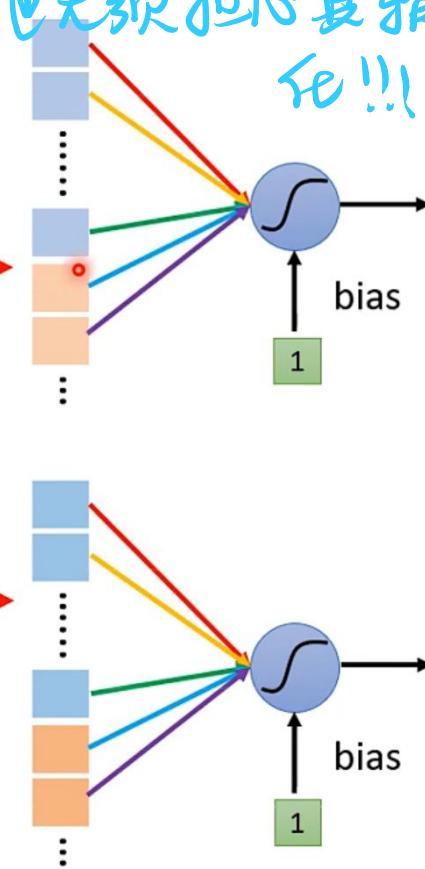


故而同样对此进行 share parameters 处理, 对于一些 weight 和 bias 设置也相同, 因输入不同也无须担心其输出的变 Fe!!!

The neurons with different receptive fields **share the parameters.**



Each filter convolves over the input image.



池化 / 下采样 / pooling. 可以使用 max pooling 也可

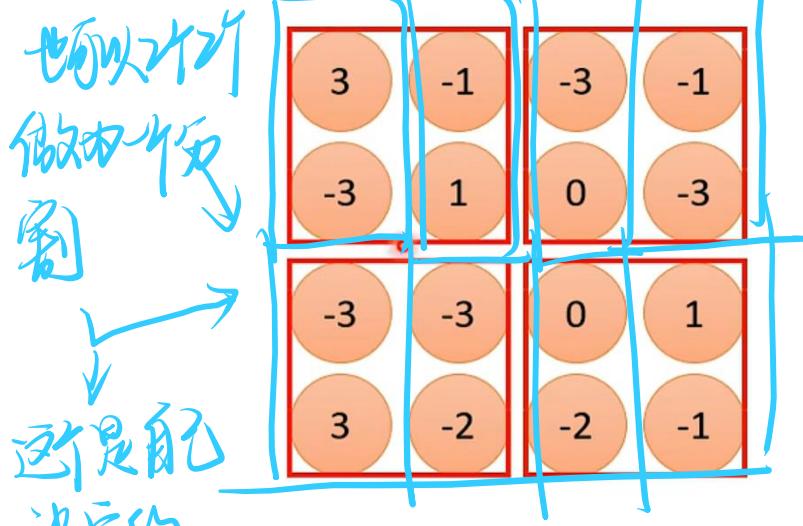
Pooling – Max Pooling 采用 min pooling 或 mean pooling 等

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2



-1	-1	-1	-1
-1	-1	-2	1
-1	-1	-1	1

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

Pooling – Max Pooling

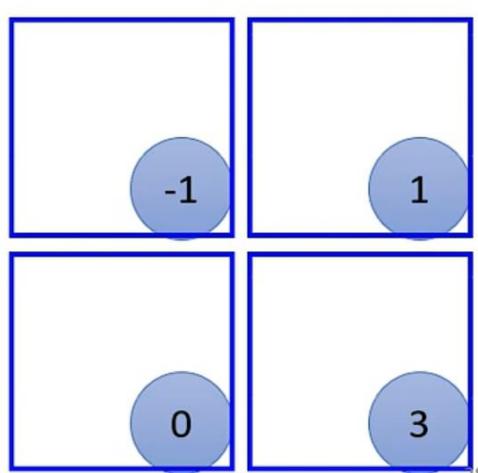
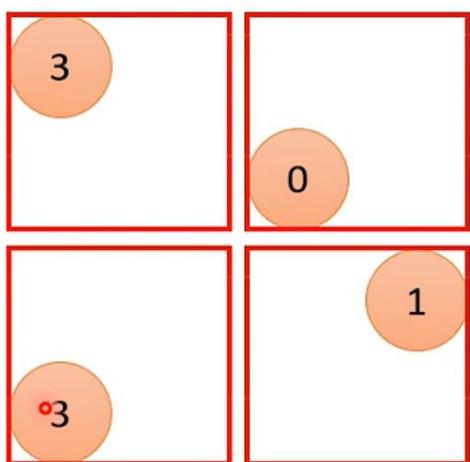
提取最大值

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

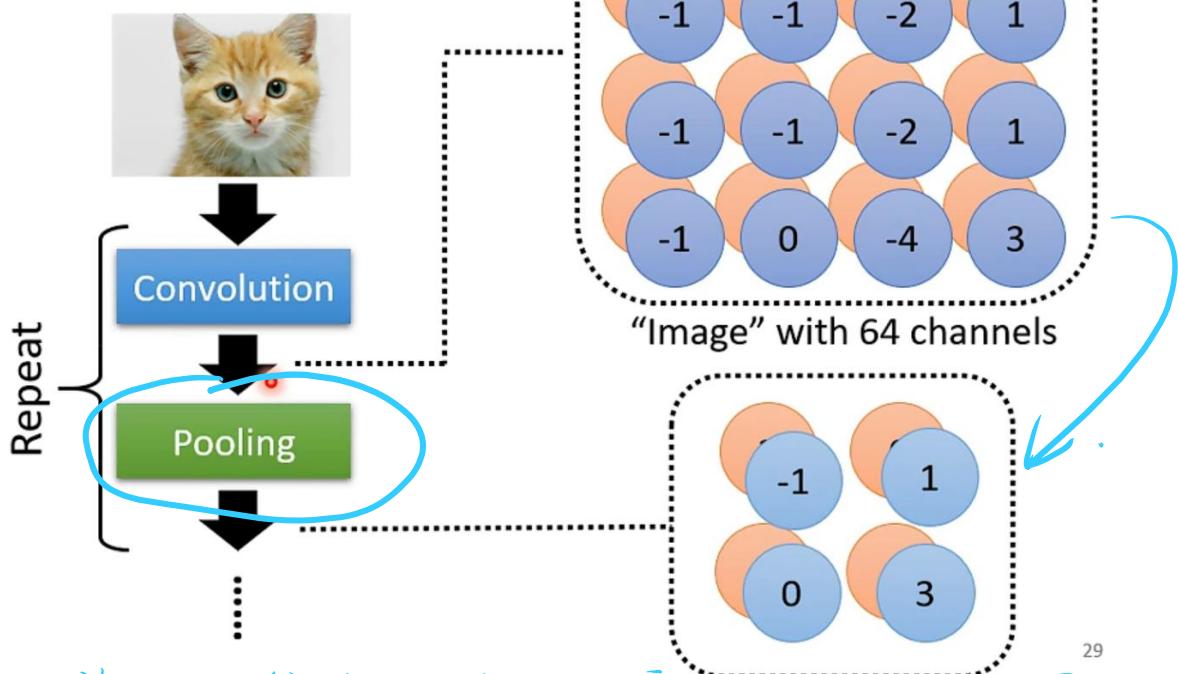
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2



在一些大型数据集中, pooling 可以减少运算量,使其简化(若 GPU 够,也可不进行 pooling)

Convolutional Layers + Pooling



这个 pooling 可以使 $4 \times 4 \times 6$ 的输入减少为 $2 \times 2 \times 64$ 但确实消耗了 CPU 和 GPU 资源, 但对于一些小数据集作

精细处理，一般采用 pooling 效果会很差！

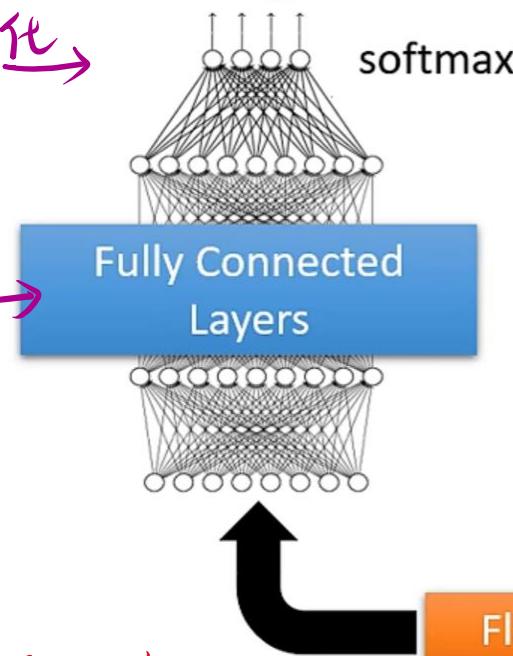
↓ 一整个 CNN 的步骤图

The whole CNN

$y = [0, 1, 0, 0, \dots, 0] \rightarrow$ cat dog

softmax 量化

全连接层



Convolution

Pooling

Convolution

Pooling

repeat
可选操作
Pn

但是，真的需要 pooling 吗???

Why CNN for Go playing?

- Subsampling the pixels will not change the object



Pooling

How to explain this???

围棋大师选取的
最优策略 48

(48 表示)

输入: 19x19x48

①. 23x23, 用 0

填充

② 用 5x5, 步 1
进行 平滑.

..

Neural network architecture. The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23×23 image, then convolves k filters of kernel size 5×5 with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used $k = 192$ filters; Fig. 2b and Extended Data Table 1 show the architecture. Alpha Go does not use Pooling

33

其中 Alpha Go 没有采用一个 pooling，因为本身棋盘下棋要放到很精确，若 pooling 的话，误差会很大，从而结果很垃圾。故

不用 pooling 在 Alpha Go 上。

但是对于 CNN 而言，若将 picture 放大或缩小，可能无法识别 picture 是狗了。这时便

引入一个新层技术：Spatial Transformer

- CNN is not invariant to scaling and rotation (we need data augmentation 😊).



Dog ✓

Spatial Transformer Layer



Cat X

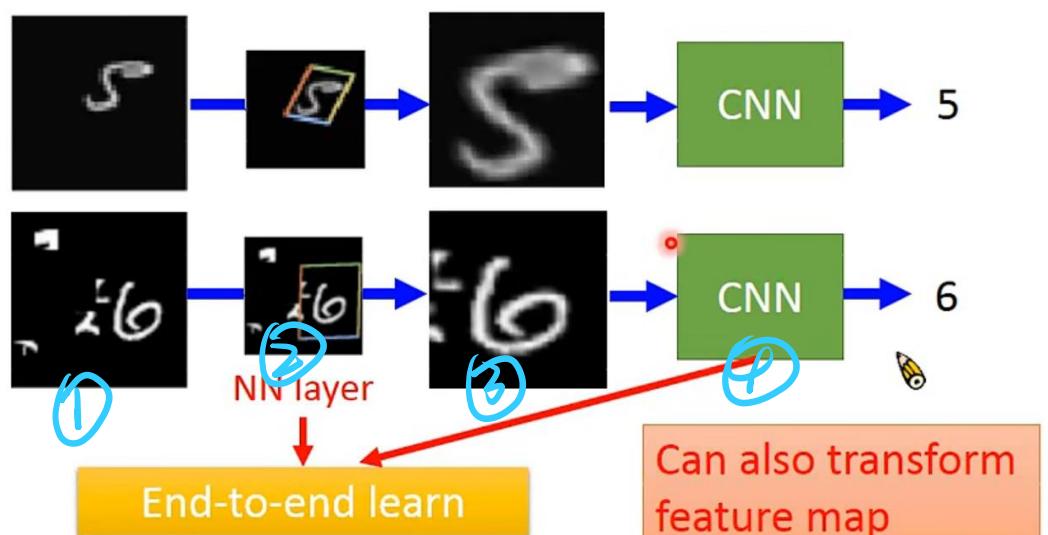


<https://youtu.be/SoCyz1hZak>
(in Mandarin)

Spatial Transformer Layer

CNN 无法识别 缩放与旋转者。已

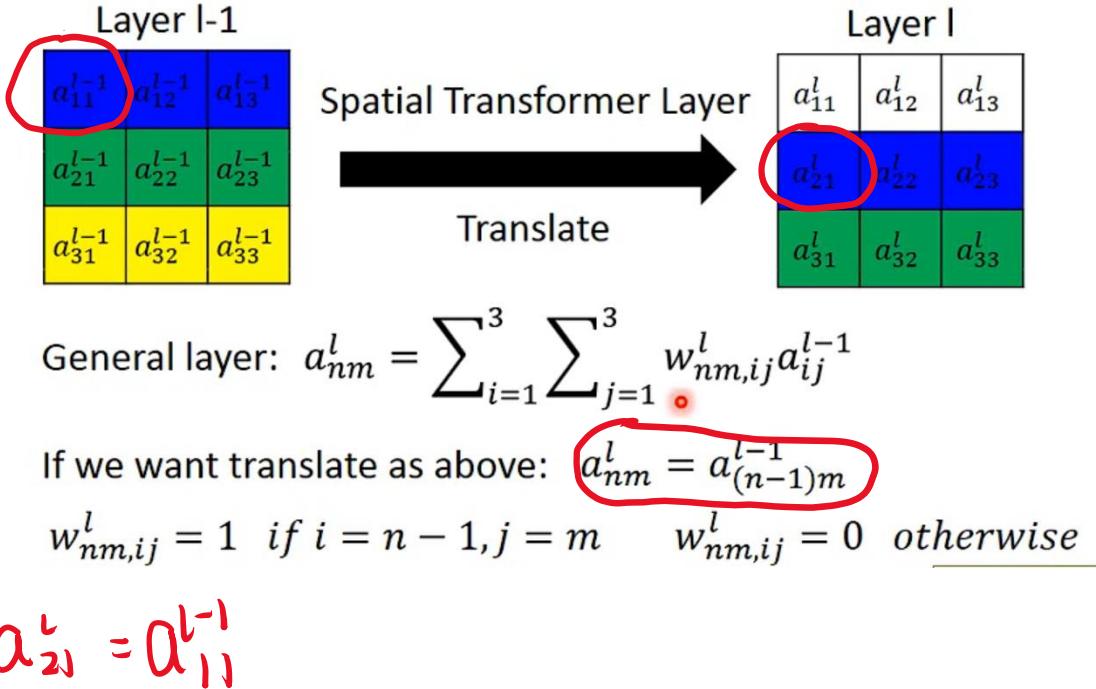
- CNN is not invariant to scaling and rotation



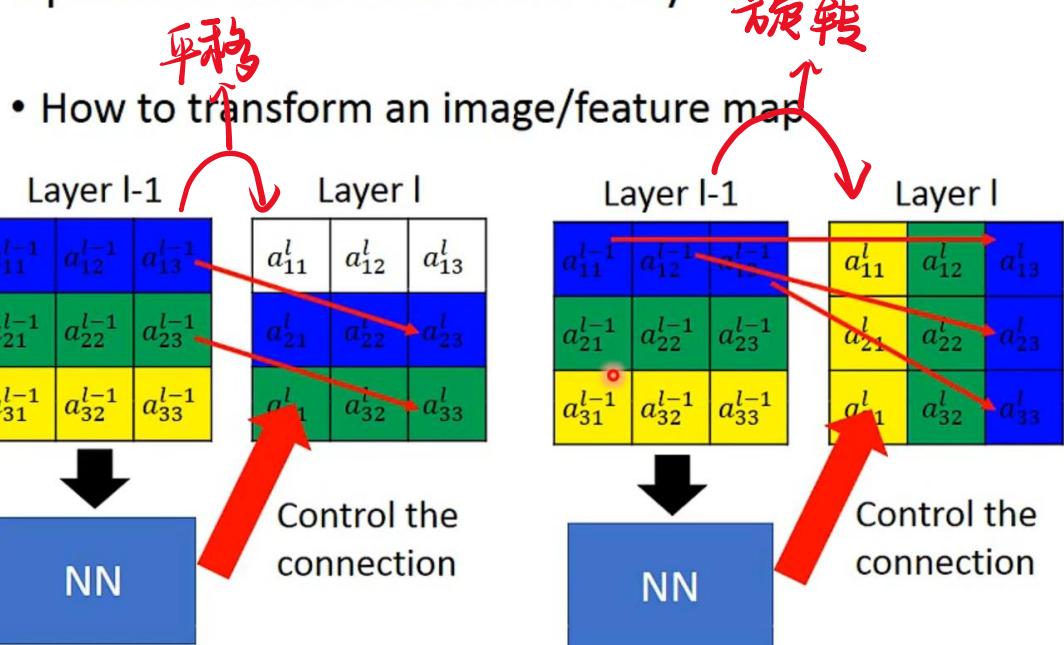
Can also transform
feature map

Spatial Transformer Layer

- How to transform an image/feature map



Spatial Transformer Layer

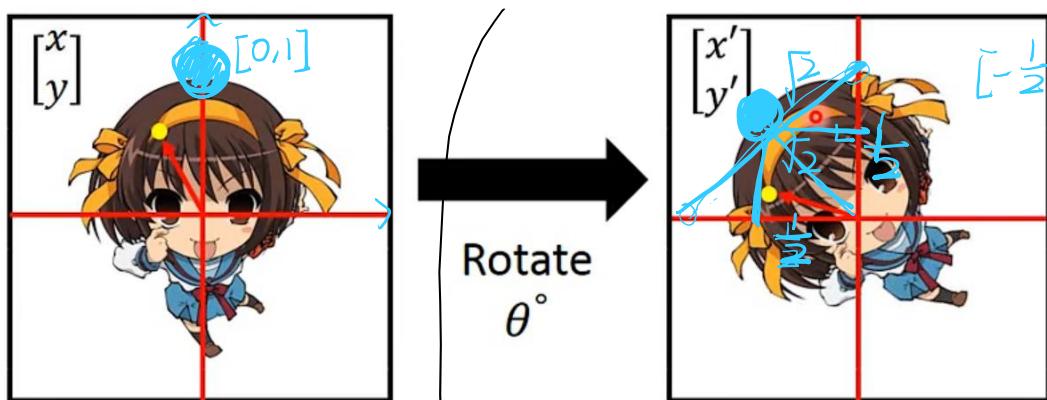
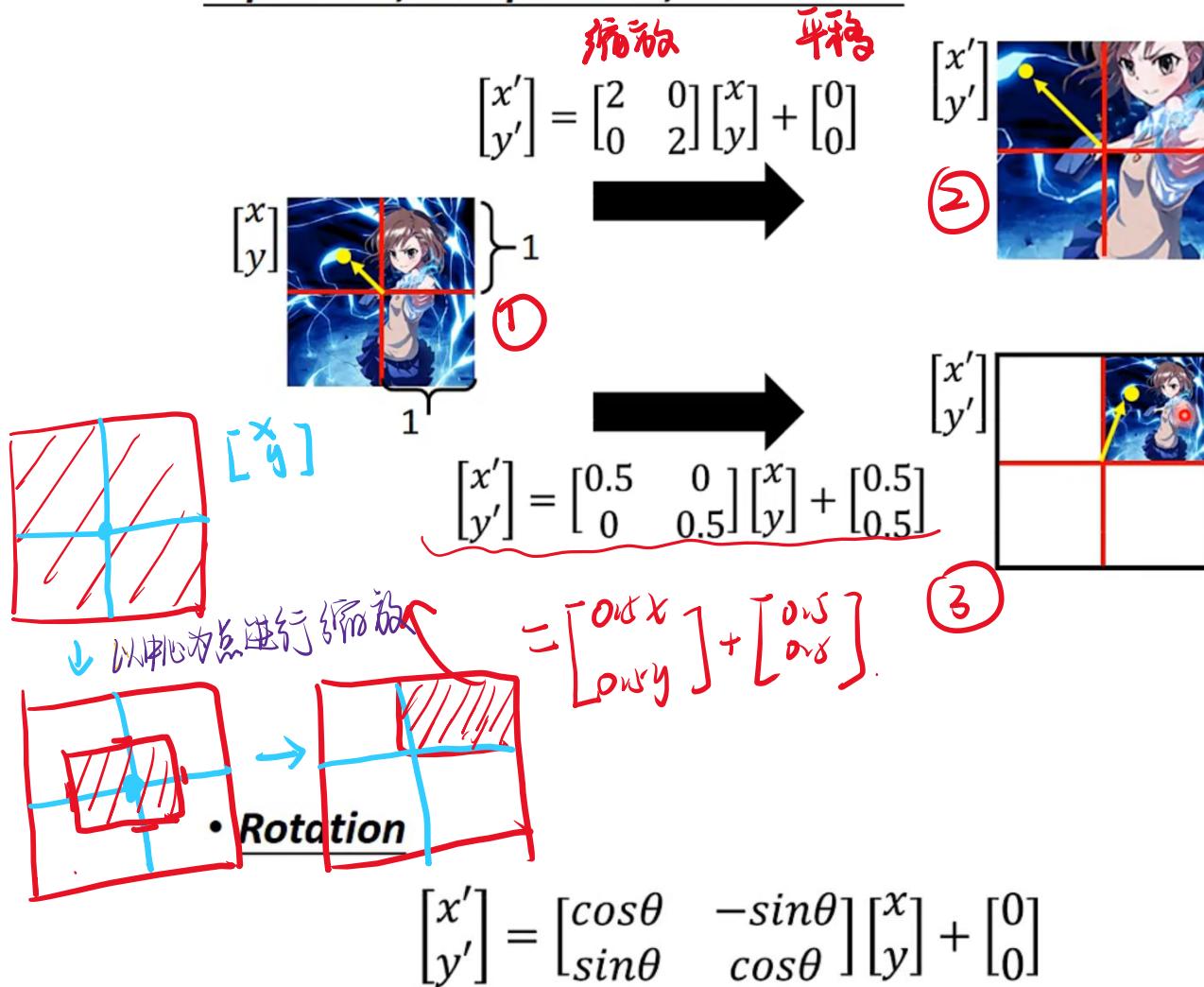


$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix} = 2 \begin{bmatrix} x \\ y \end{bmatrix}$$

即②是由①在中心进行扩大2倍得到的结果

Image Transformation

Expansion, Compression, Translation



$$\theta = 45^\circ, \quad \begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \frac{1}{2}x - \frac{1}{2}y \\ \frac{1}{2}x + \frac{1}{2}y \end{bmatrix}$$

由 $[x] = [9]$ 故 $[x'] = \left[\begin{array}{c} -\frac{1}{2} \\ \frac{1}{2} \end{array} \right]$

结果一致 OK

若

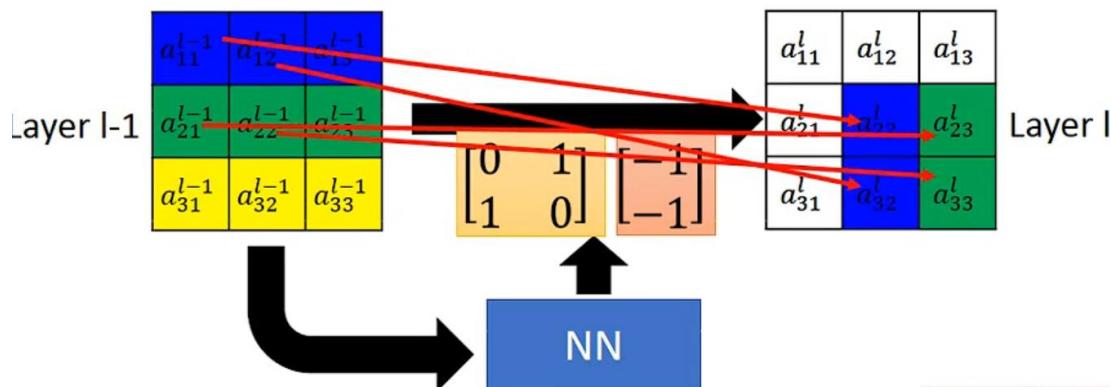
均勻整數，此時，一塊對應一塊，完全 NO problem!

Spatial Transformer Layer

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

Index of layer l-1 Index of layer l

6 parameters to describe the affine transformation



Created with EverCam

可是，若參照小數部分呢？

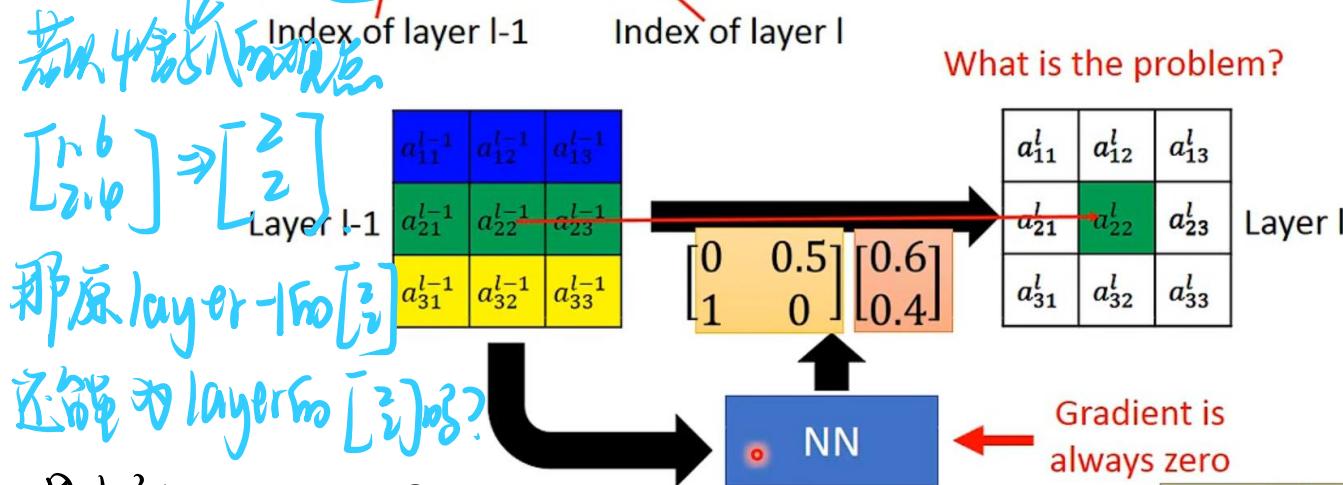
Spatial Transformer Layer

$$\begin{bmatrix} 1.6 \\ 2.4 \end{bmatrix} = \begin{bmatrix} 0 & 0.5 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$$

Index of layer l-1 Index of layer l

6 parameters to describe the affine transformation

What is the problem?



Gradient is always zero

顯然是不正確的，因為如果小數部分全為一個樣子，此時會由一個連續變動分段：



那 Gradient = 0，无法再继续

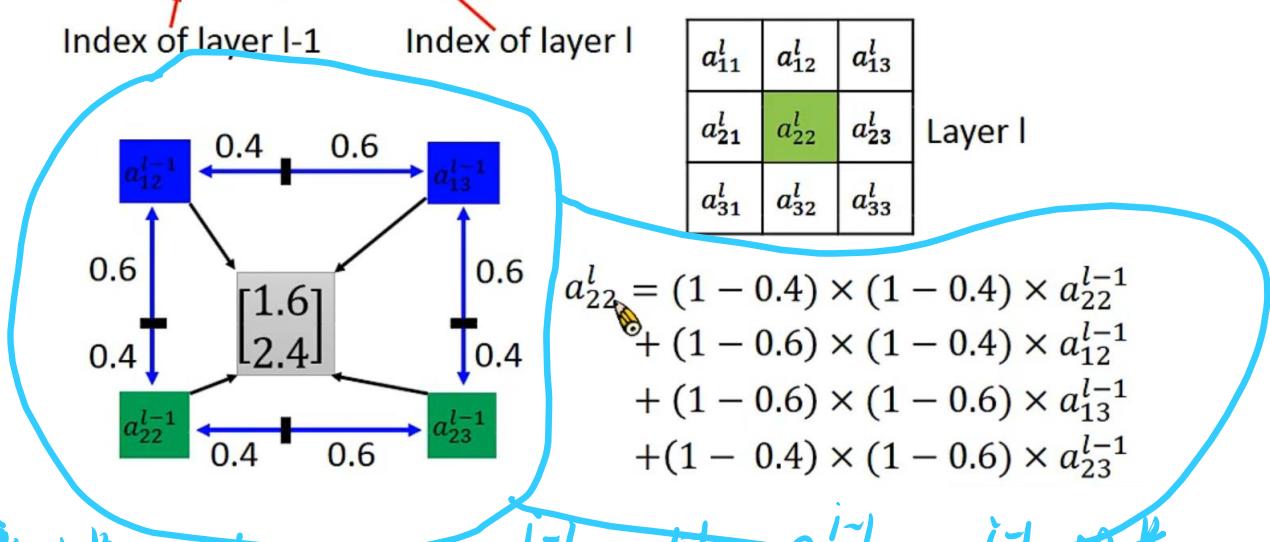
故而：不可采用简单的回传至入 new way!

Interpolation

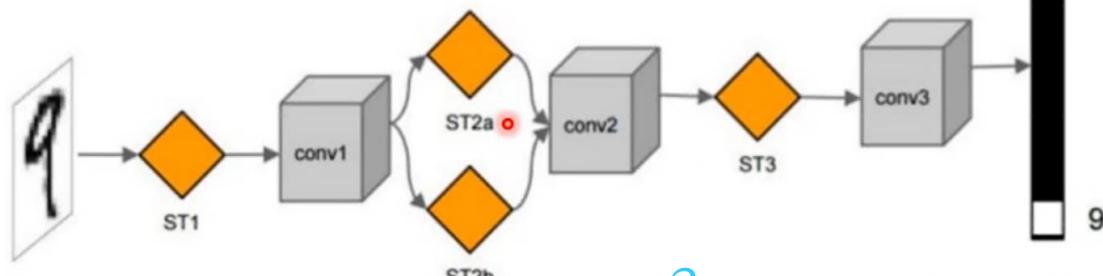
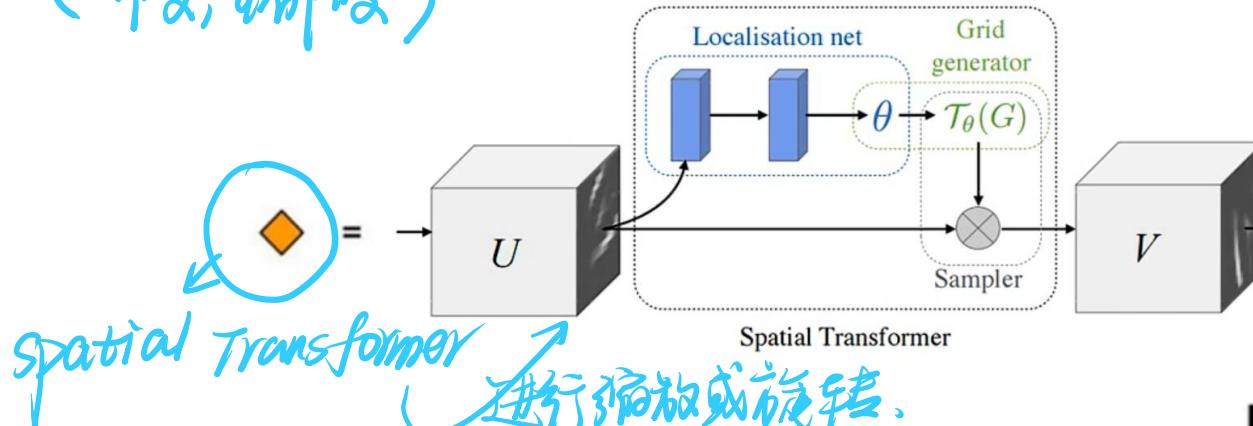
Now we can use gradient descent

$$\begin{bmatrix} 1.6 \\ 2.4 \end{bmatrix} = \begin{bmatrix} 0 & 0.5 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$$

6 parameters to describe the affine transformation



也就意味着 a_{22}^l 的值与 $a_{12}^{l-1}, a_{13}^{l-1}, a_{22}^{l-1}, a_{23}^{l-1}$ 相关
(一像, 二像)

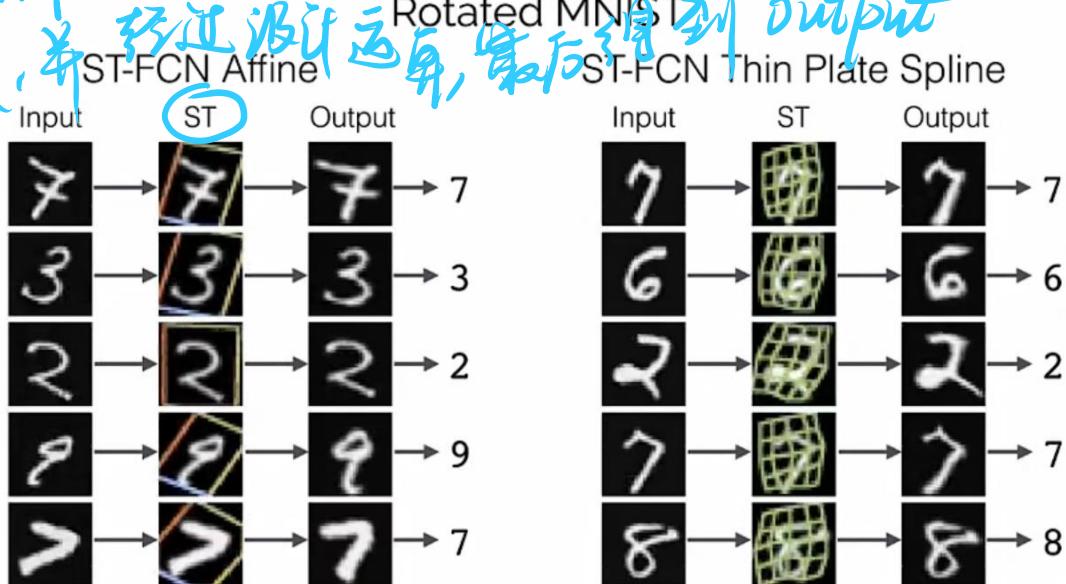


而 spatial transformer 不仅可以用于最开始的 convolution, 也可以在后面使用, 也可以 2个 spatial transformer 同时使用!

对于手写 dataset = MNIST 为例

ST \Rightarrow spatial transformer

ST \Rightarrow 把 input 的 picture 从 重要区域固定起来，再放入 CNN 处理，并经过逆运算，最后得到 output

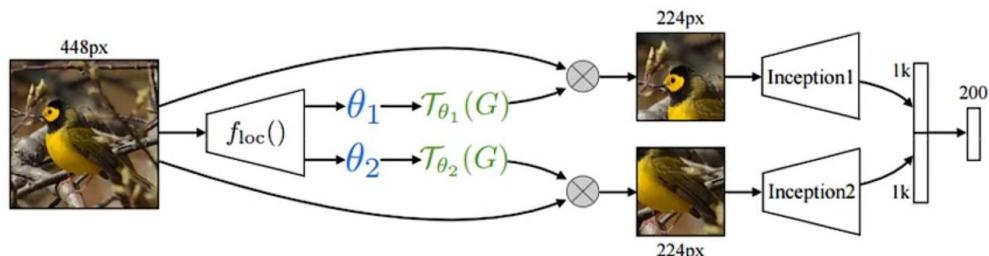


对于 Bird 图分类：

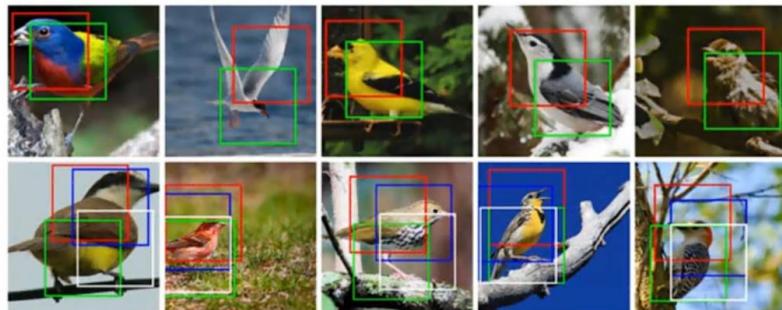
在正确分类图片最终 CNN 的前面采用了 ST 技术（对于本题：红色区域更多的是收集鸟的嘴，绿色部分更多是收集翅膀）

Bird Recognition

$$\begin{bmatrix} a & 0 \\ 0 & d \end{bmatrix} \begin{bmatrix} e \\ f \end{bmatrix}$$



Model	
Cimpoi '15 [4]	66.7
Zhang '14 [30]	74.9
Branson '14 [2]	75.7
Lin '15 [20]	80.9
Simon '15 [24]	81.0
CNN (ours)	82.3
2×ST-CNN 224px	83.1
2×ST-CNN 448px	83.9
4×ST-CNN 448px	84.1



最终，可以识别 放大缩小，旋转的 picture 是哪种 呢？

更好的