



Computer Vision

第九周 图像识别

庞彦

yanpang@gzhu.edu.cn



01

Object Recognition

图像识别初探究

Object Recognition

目标识别: 分类问题

Class 1: 笔记本电脑;

Class 2: 台式机电脑;

Class 3: 平板电脑。



Object Recognition

目标识别： 分类问题

对于人类：

共同点：
都是电脑…

不同点：
大小？ 形状？ 材质？ 便携性…

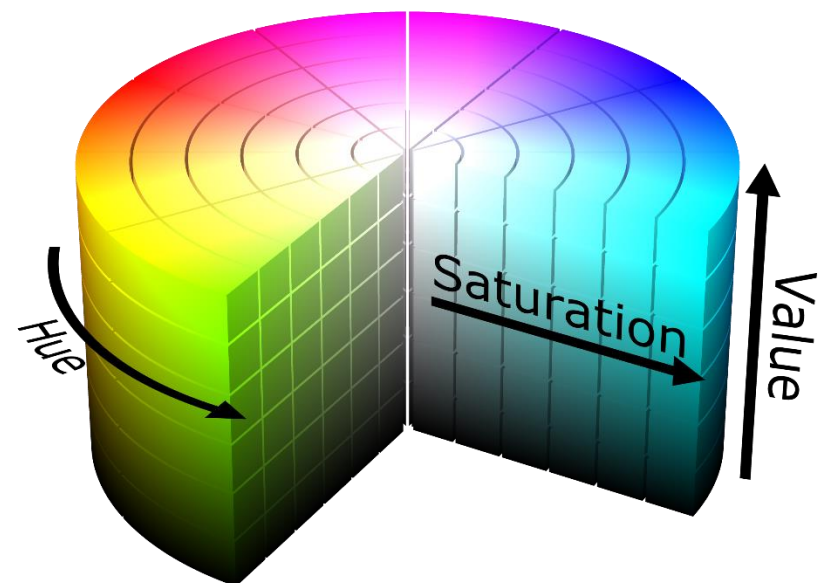
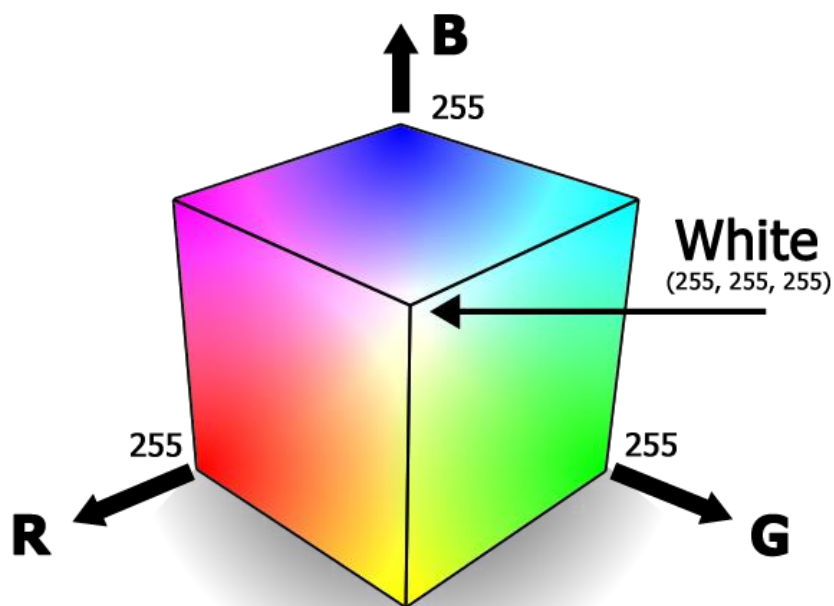


Features

计算机通过学习并挖掘目标属性的相关**特征**（Features）来对目标进行识别与探究。

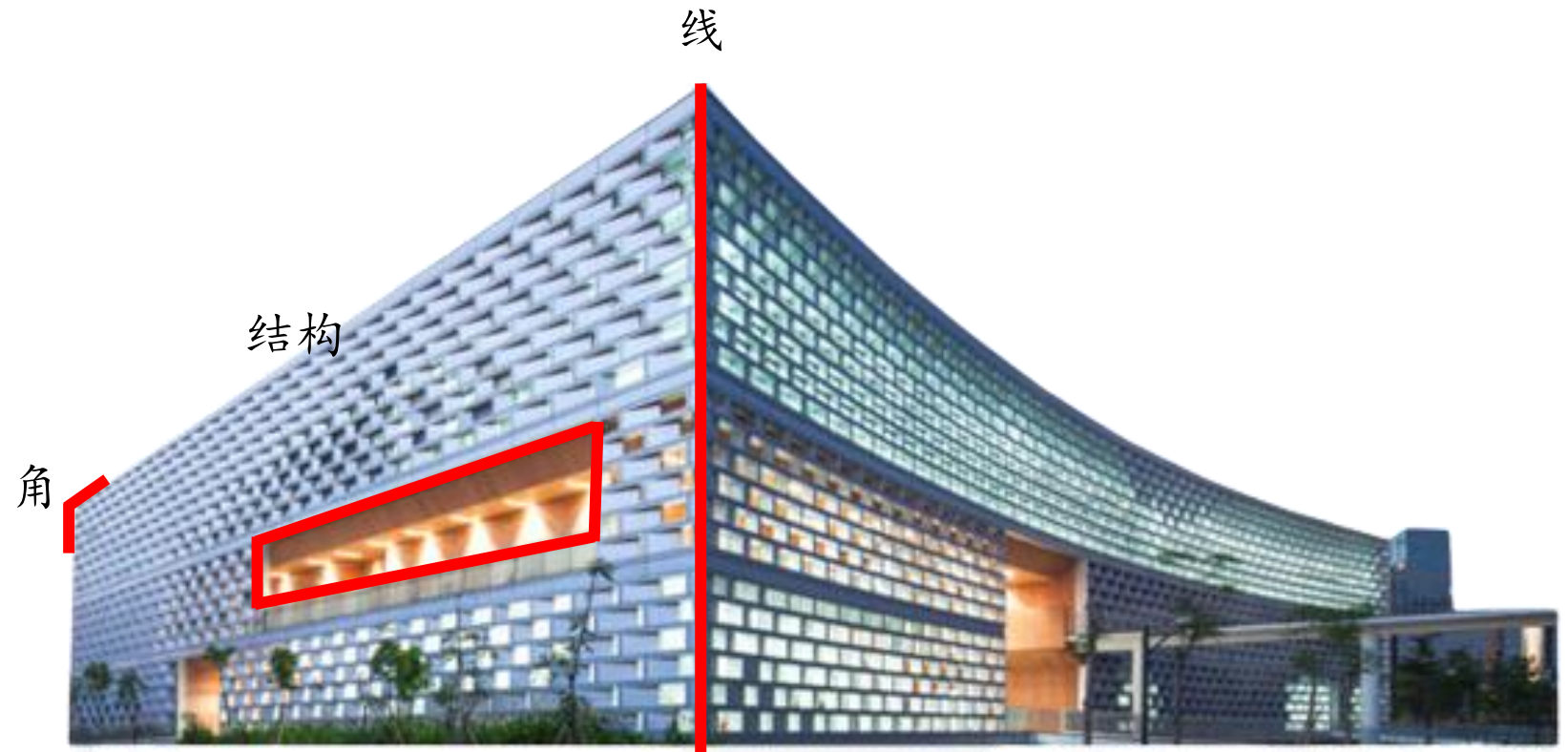
Features

颜色特征



Features

颜色特征
形状特征



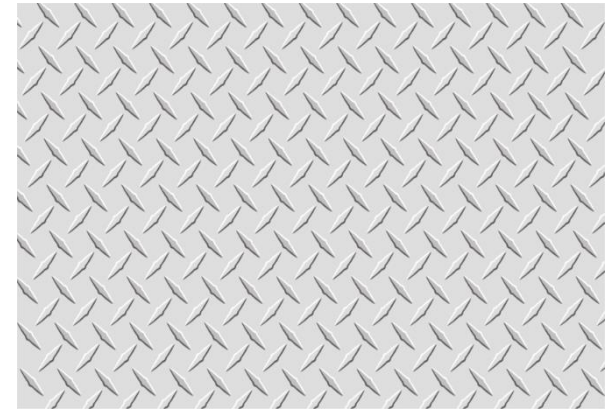
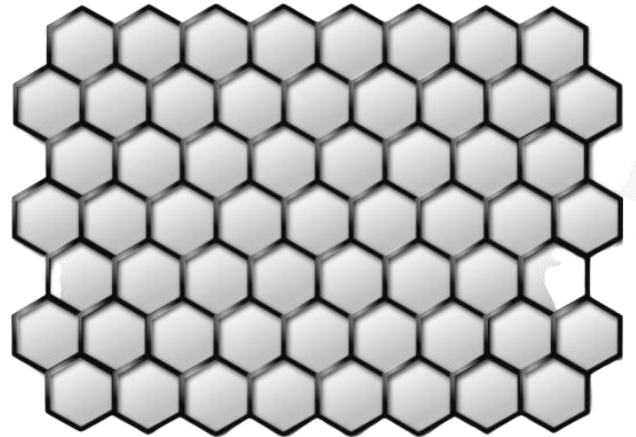
Features

颜色特征

形状特征

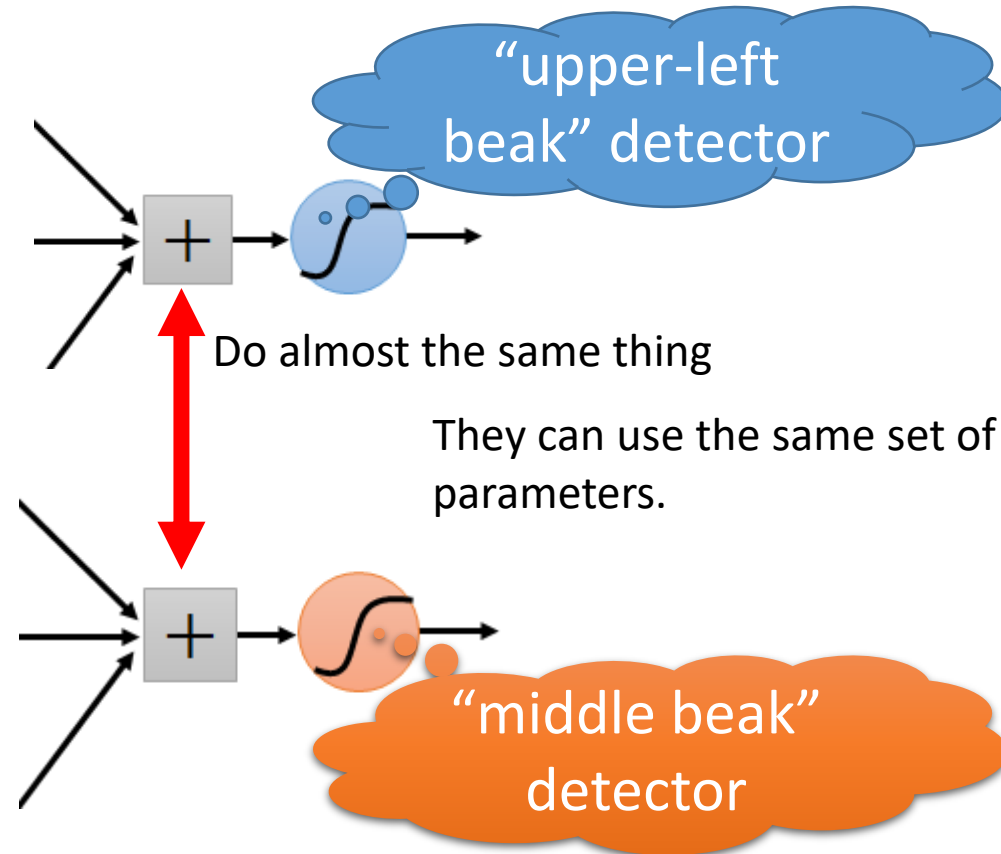
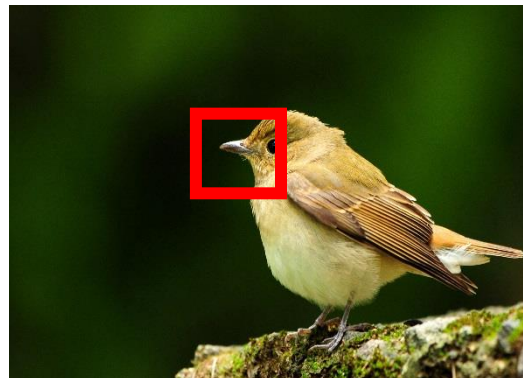
纹理特征

... ..



Why CNN for Image

- The same patterns appear in different regions.



Why CNN for Image

- ✓ **Subsampling** the pixels will not change the object;
- ✓ It just make the image smaller.

parrot

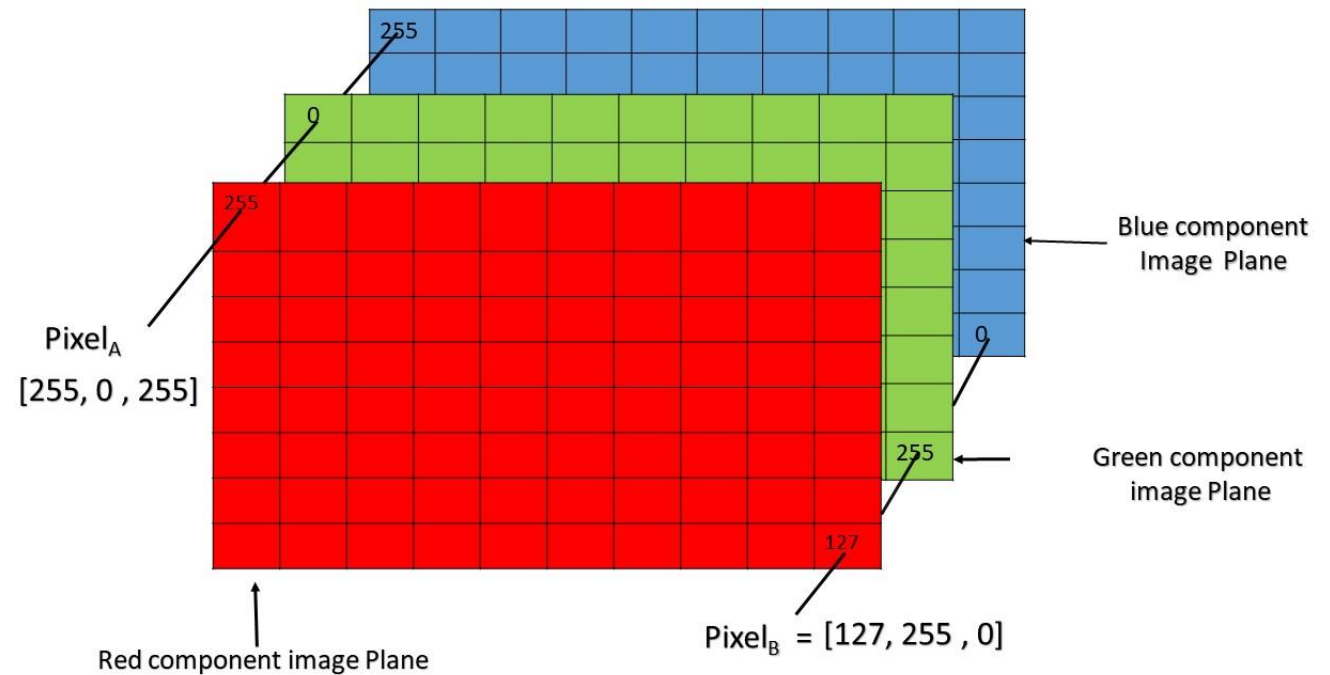


parrot



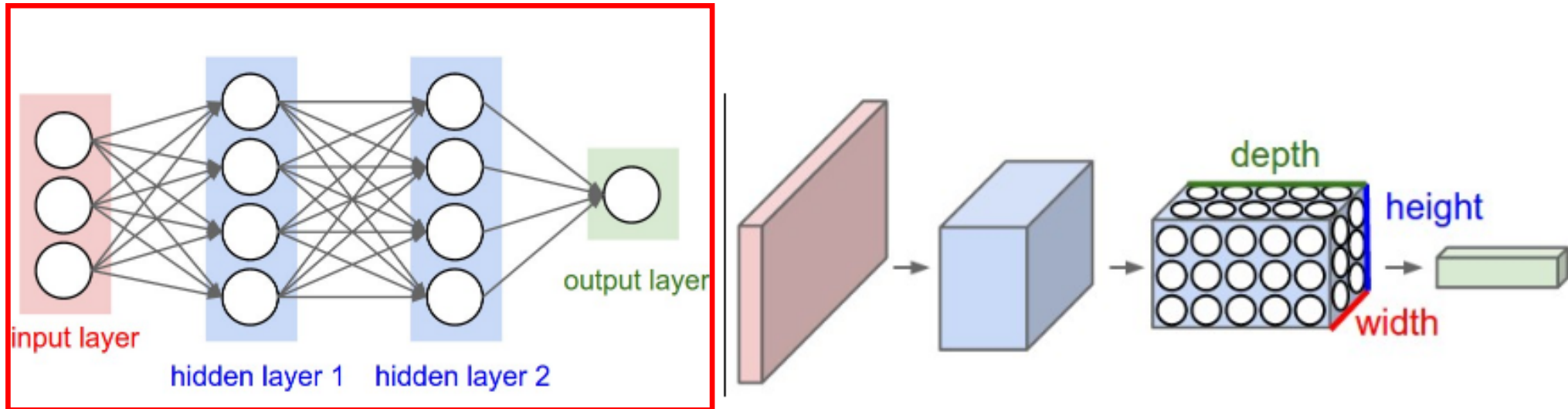
Less parameters for the network to process the image

Color Image: 3 channels



Pixel of an RGB image are formed from the corresponding pixel of the three component images

Convolutional Neural Networks (CNNs)



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

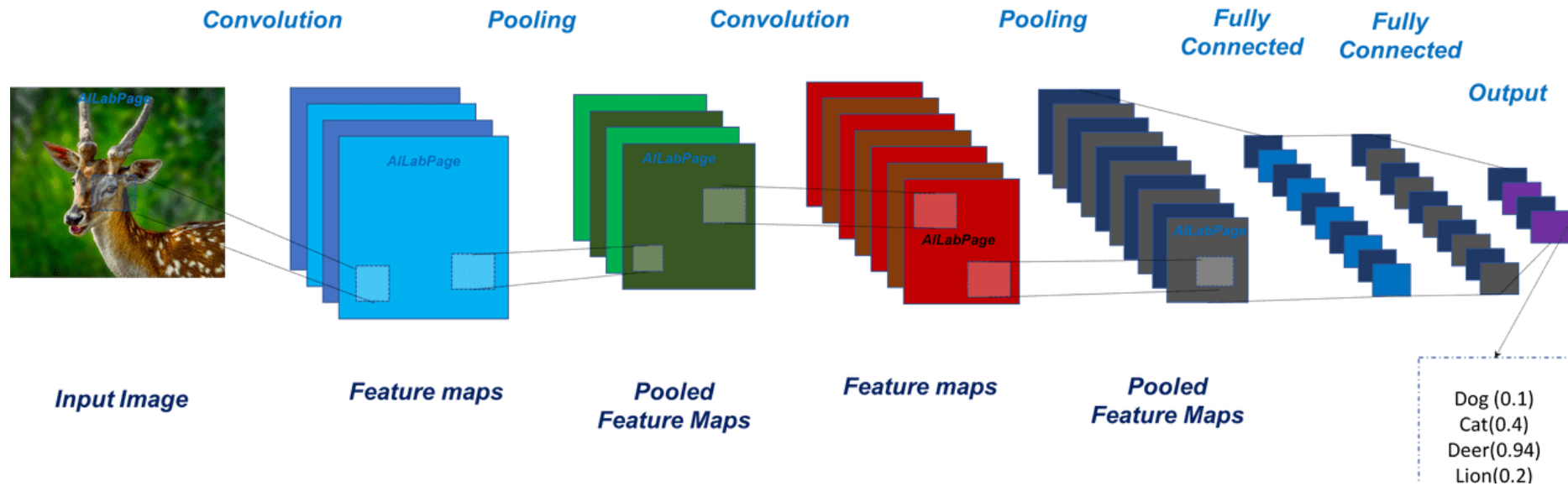
Convolutional Neural Networks (CNNs)

卷积层	Convolutional Layer
池化层	Pooling Layer
全连接层	Fully-Connected Layer

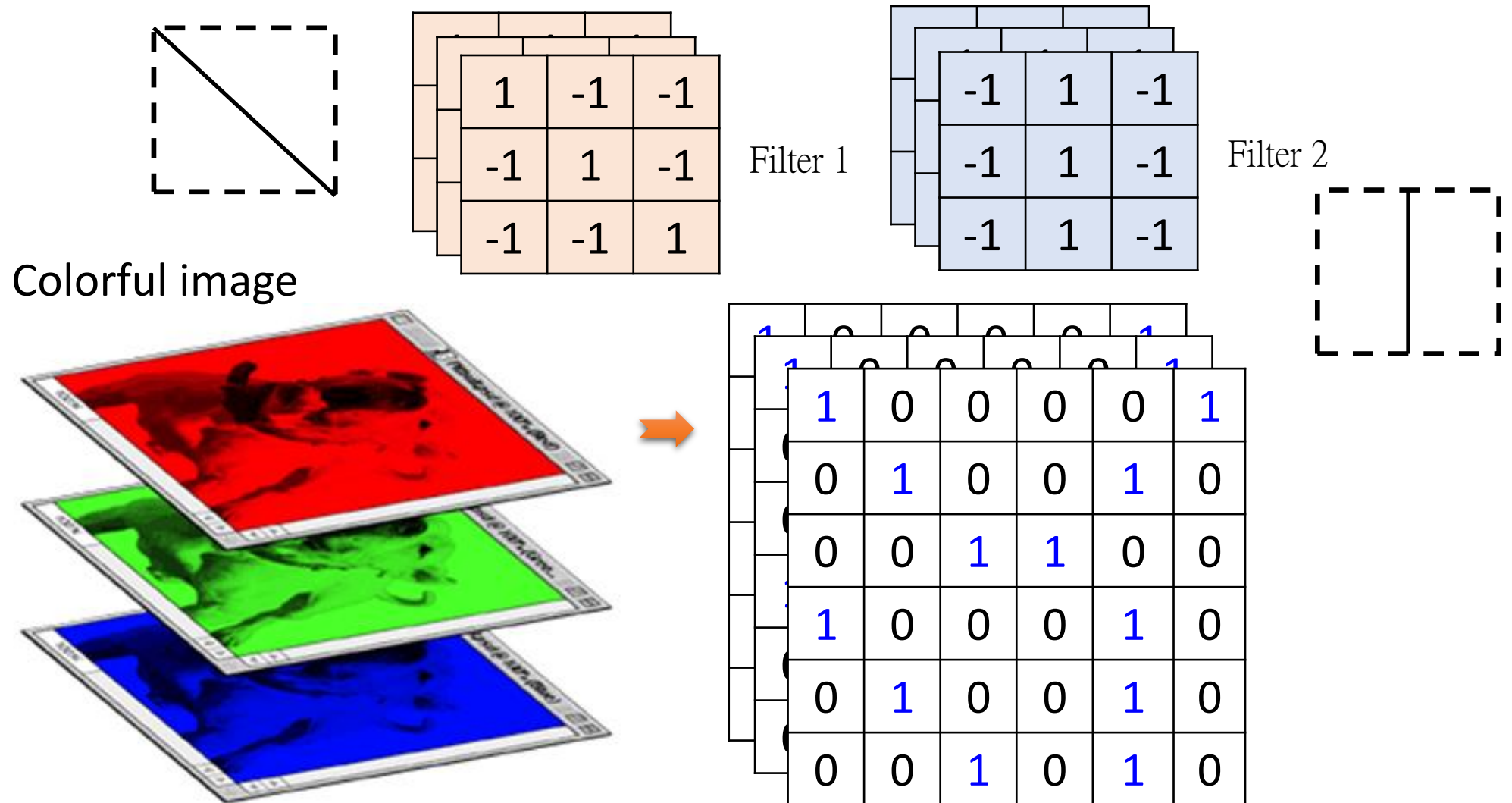
We will stack these layers to form a full **ConvNet** architecture.

[**INPUT - CONV - POOL – FC - OUTPUT**]

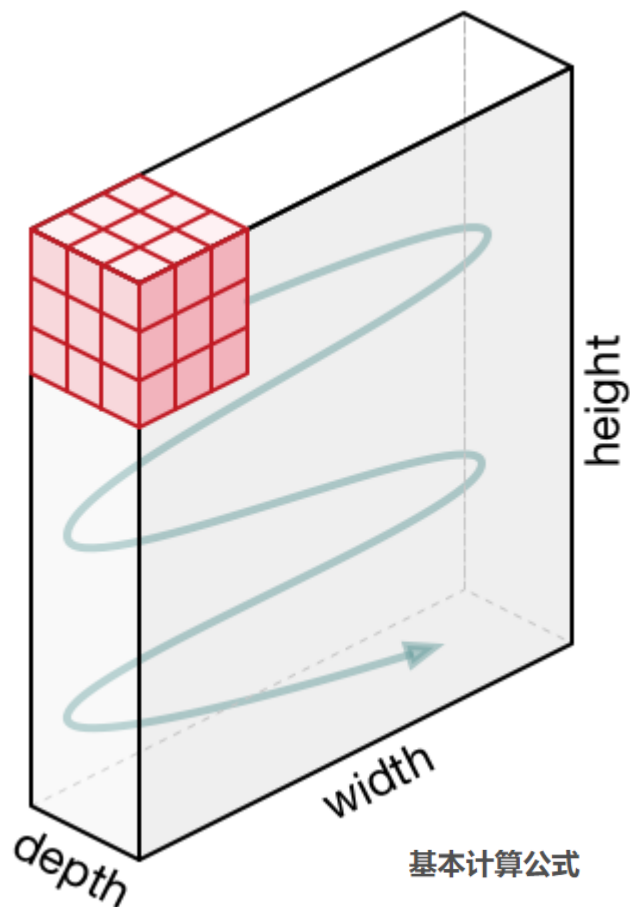
Convolutional Neural Networks (CNNs)



Convolution Layer: The Kernel



Convolution Layer: The Kernel



0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1

308

+

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2

-498

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

164

+

Bias = 1

+ 1 = -25

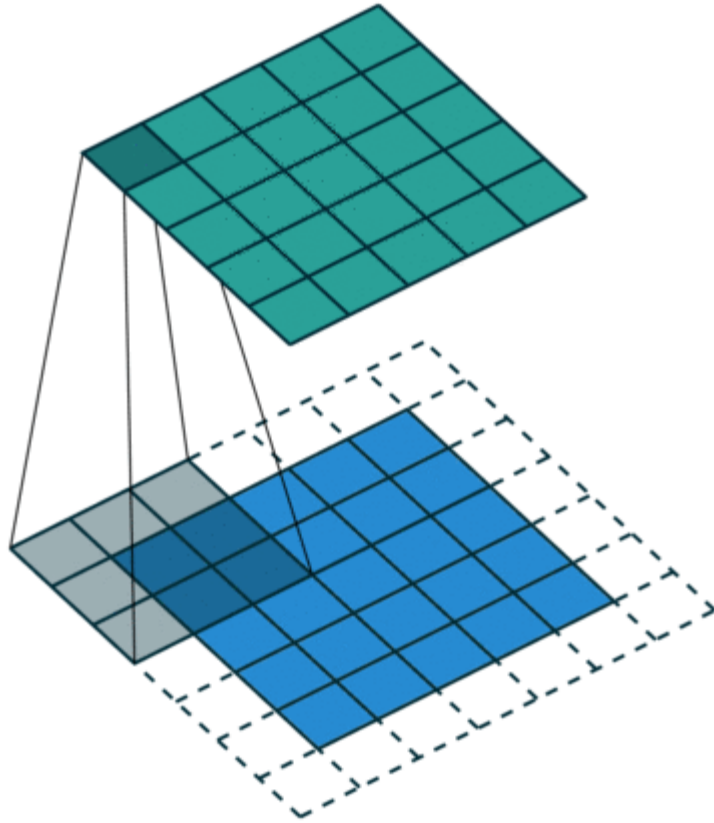
-25				...
				...
				...
				...
...

基本计算公式

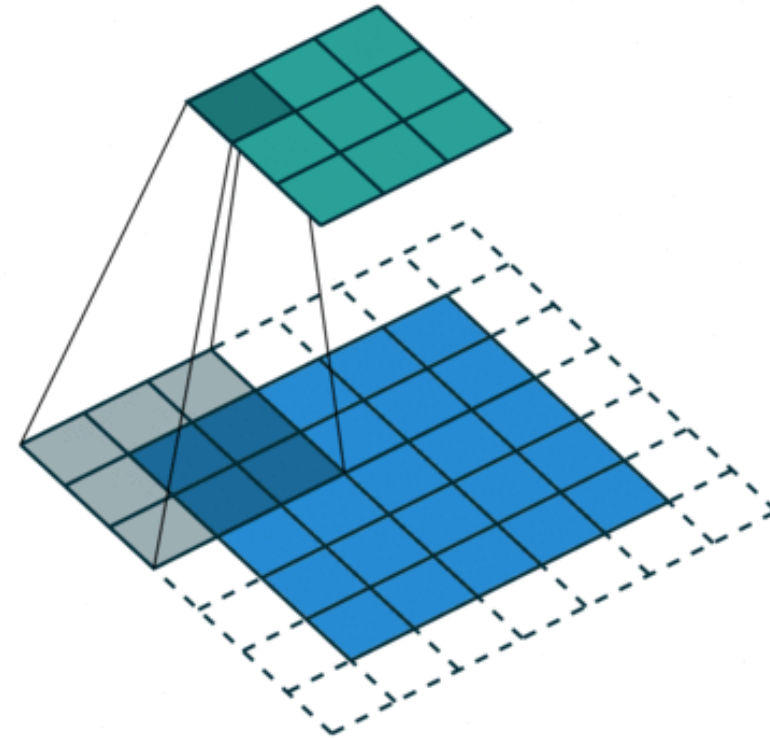
W 为输入大小, F 为卷积核大小, P 为填充大小(padding), S 为步长(stride), N 为输出大小。有如下计算公式:

$$N = \frac{(W - F + 2P)}{S} + 1$$

Convolution Layer: Padding and Stride



padding = 1
stride = 1



padding = 1
stride = 2

Pooling Layer

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

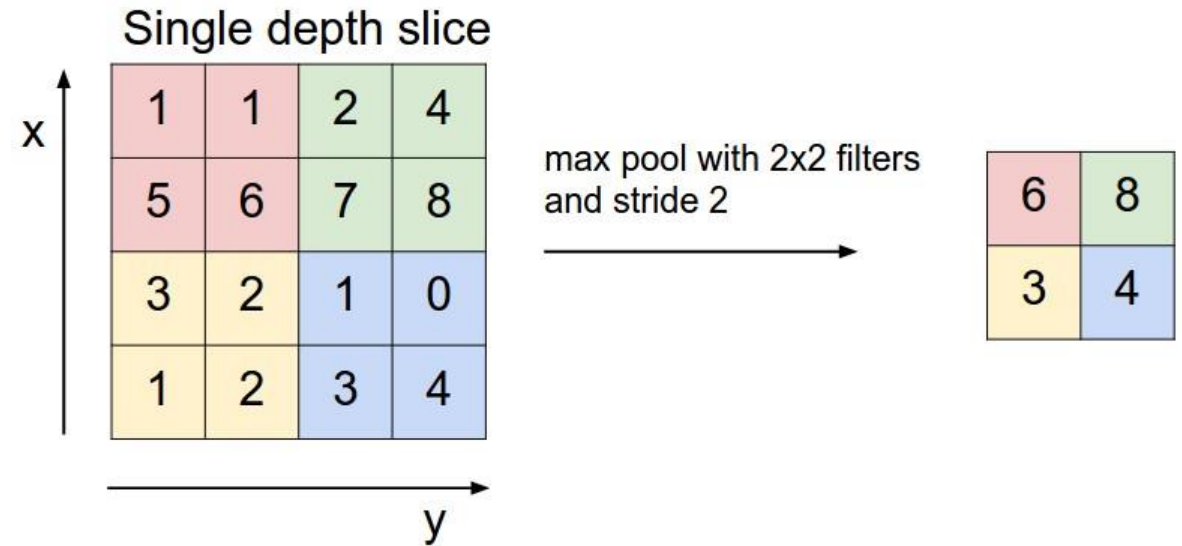
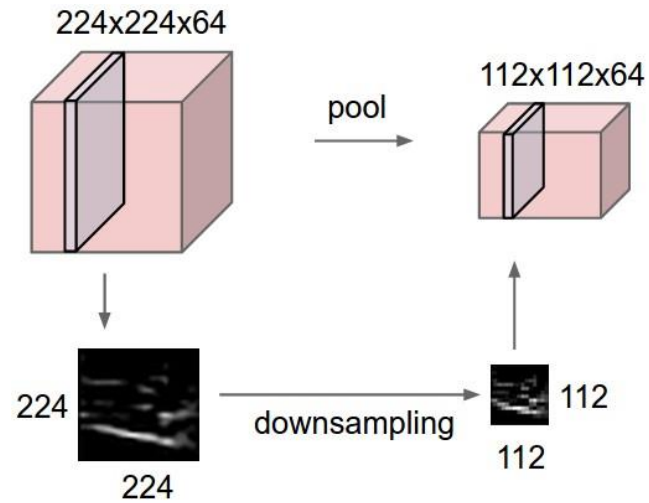
max pooling

20	30
112	37

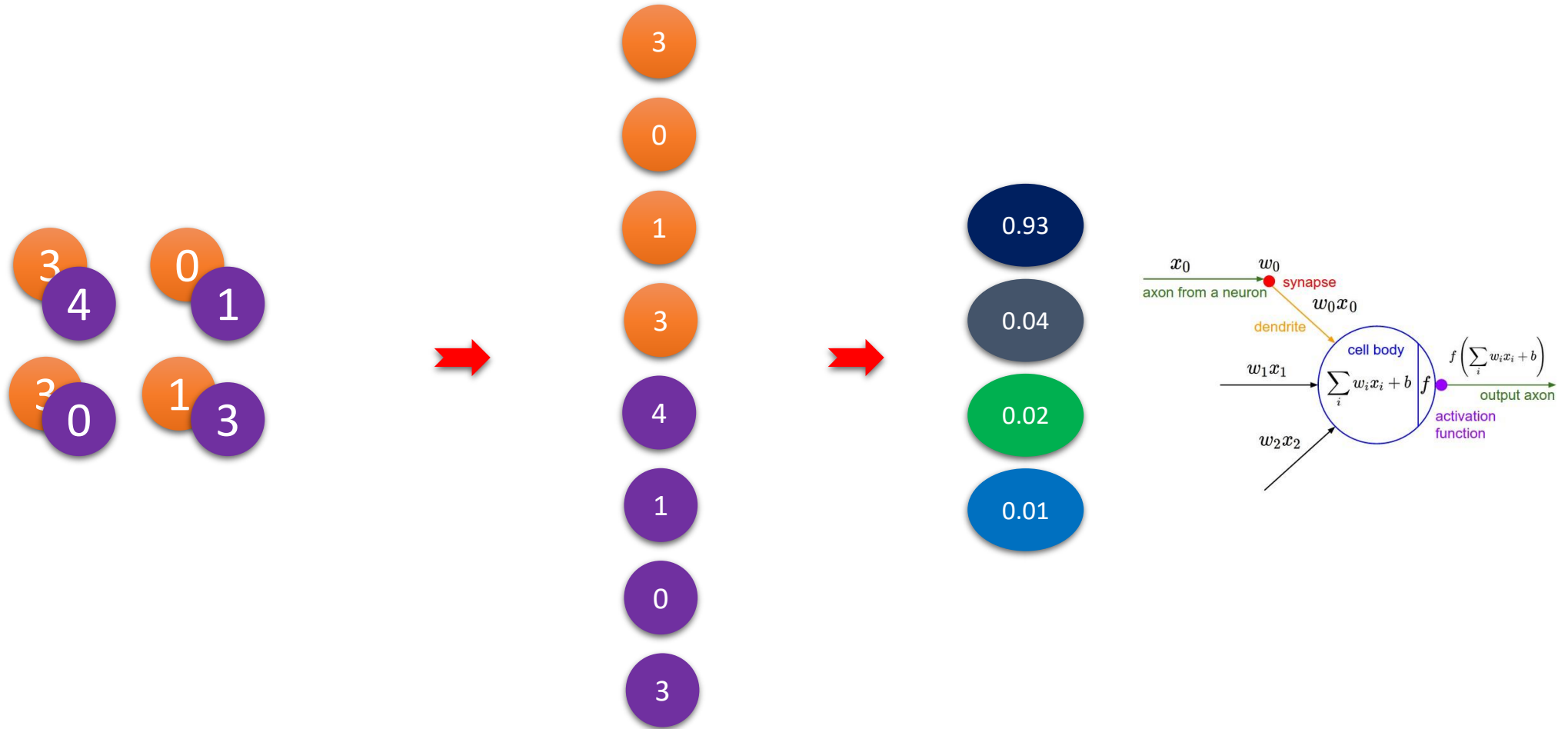
average pooling

13	8
79	20

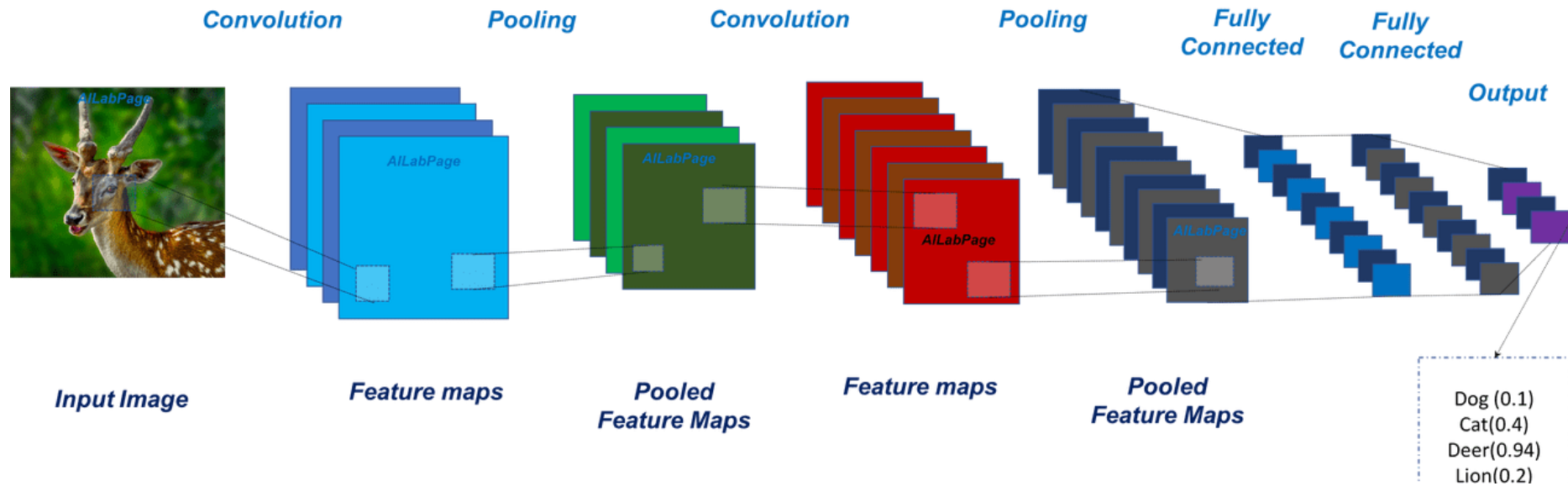
Pooling Layer



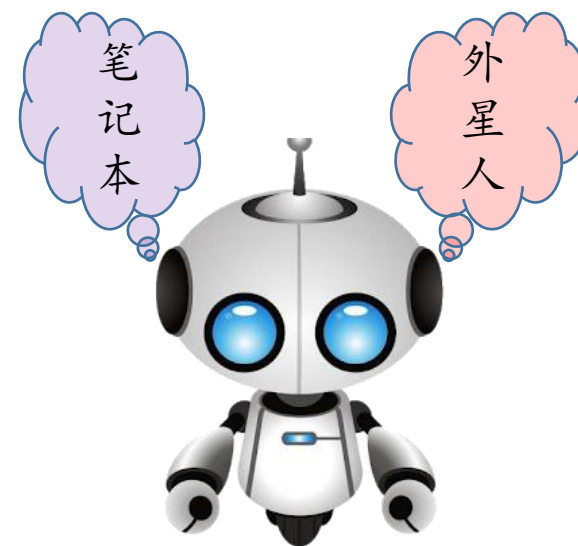
Flatten



Review: Convolutional Neural Networks (CNNs)



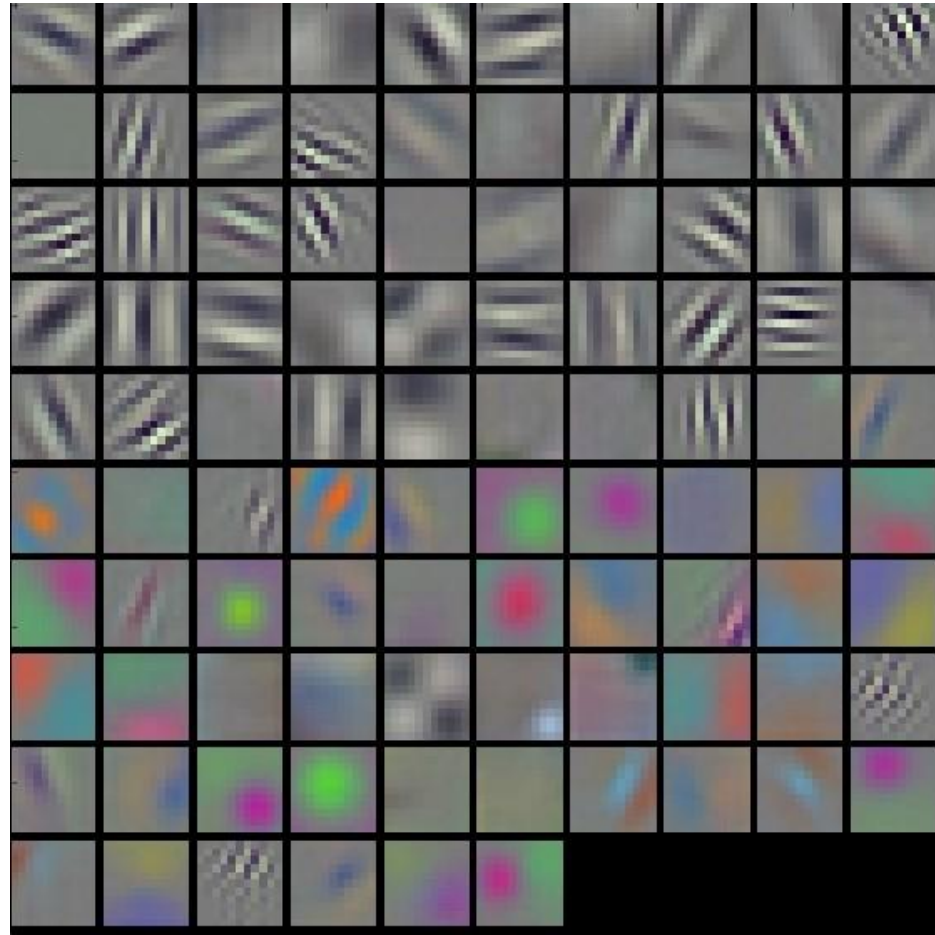
What does the machine learn?



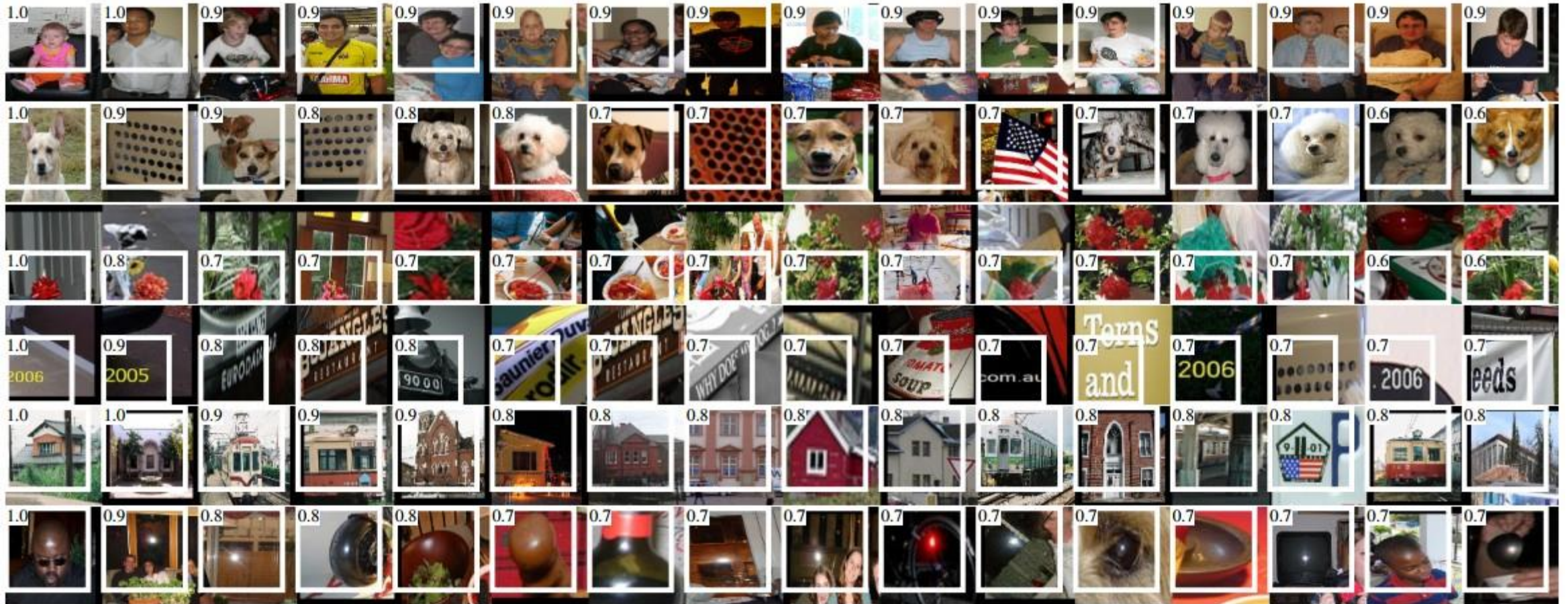
Features

- Typical-looking filters on the trained first layer

11 x 11
(AlexNet)



Features



Maximally activating images for some POOL5 (5th pool layer) neurons of an AlexNet.

More Applications: Deep Dream



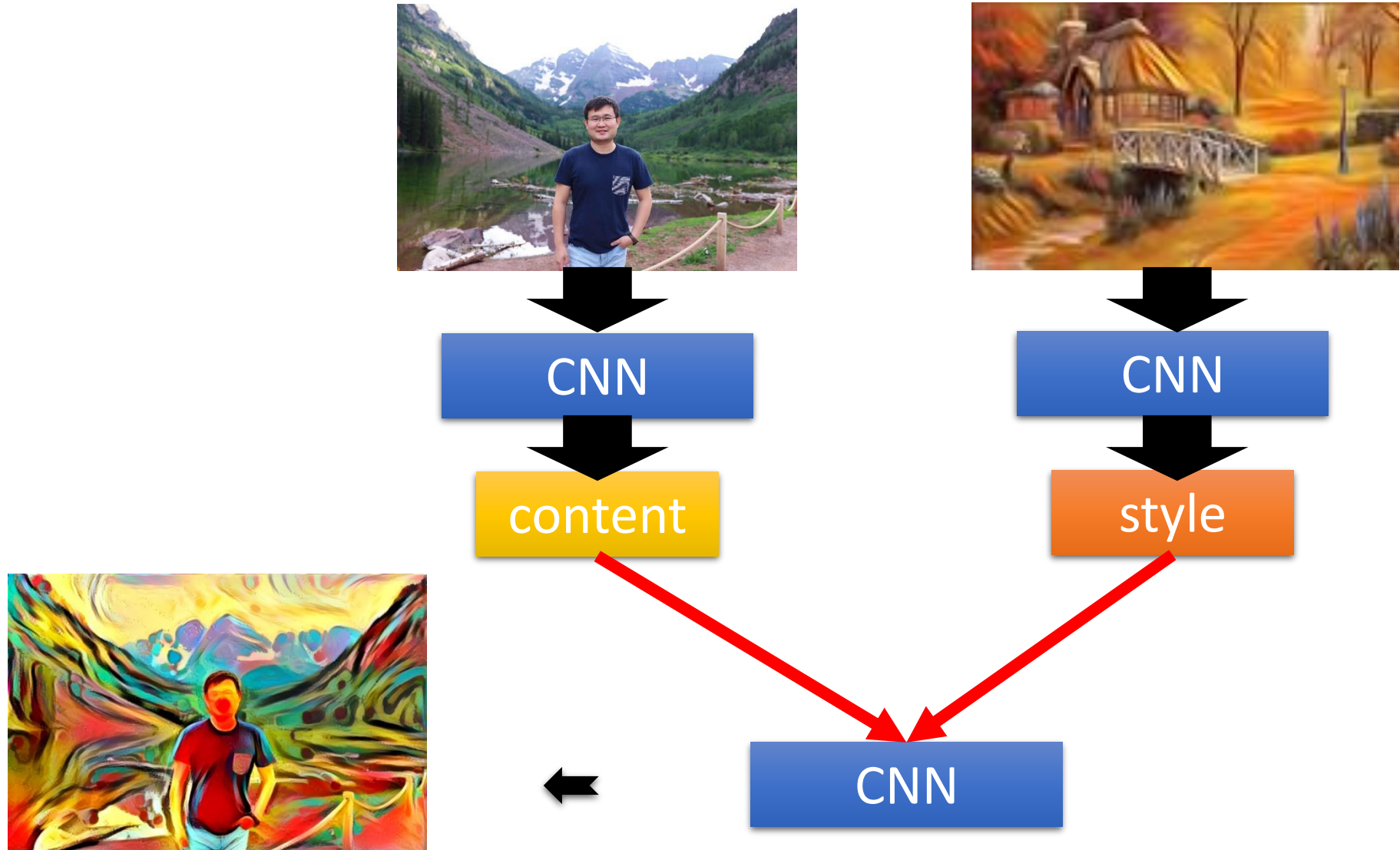
+



=



Given a photo, make its style like famous paintings





02

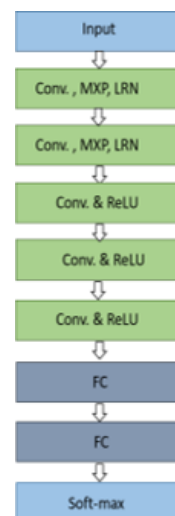
VGG → CSPNet

从VGG到CSPN

Deep Convolutional Neural Network

Inaccuracy: 16.4%

8 layers



AlexNet (2012)

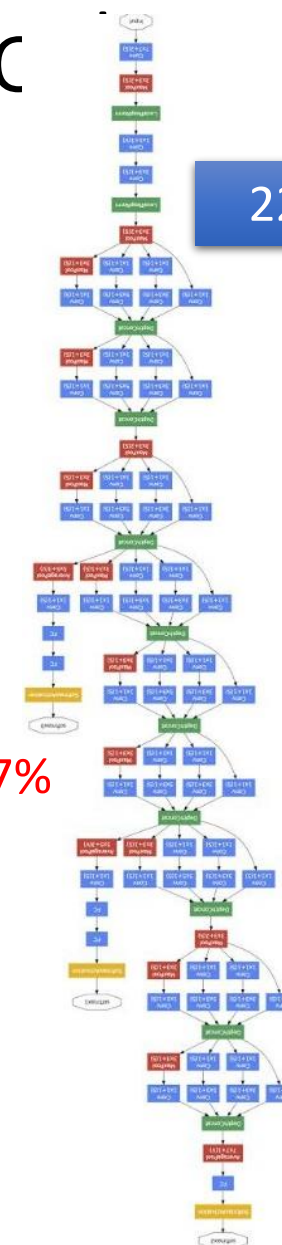
19 layers



VGG (2014)

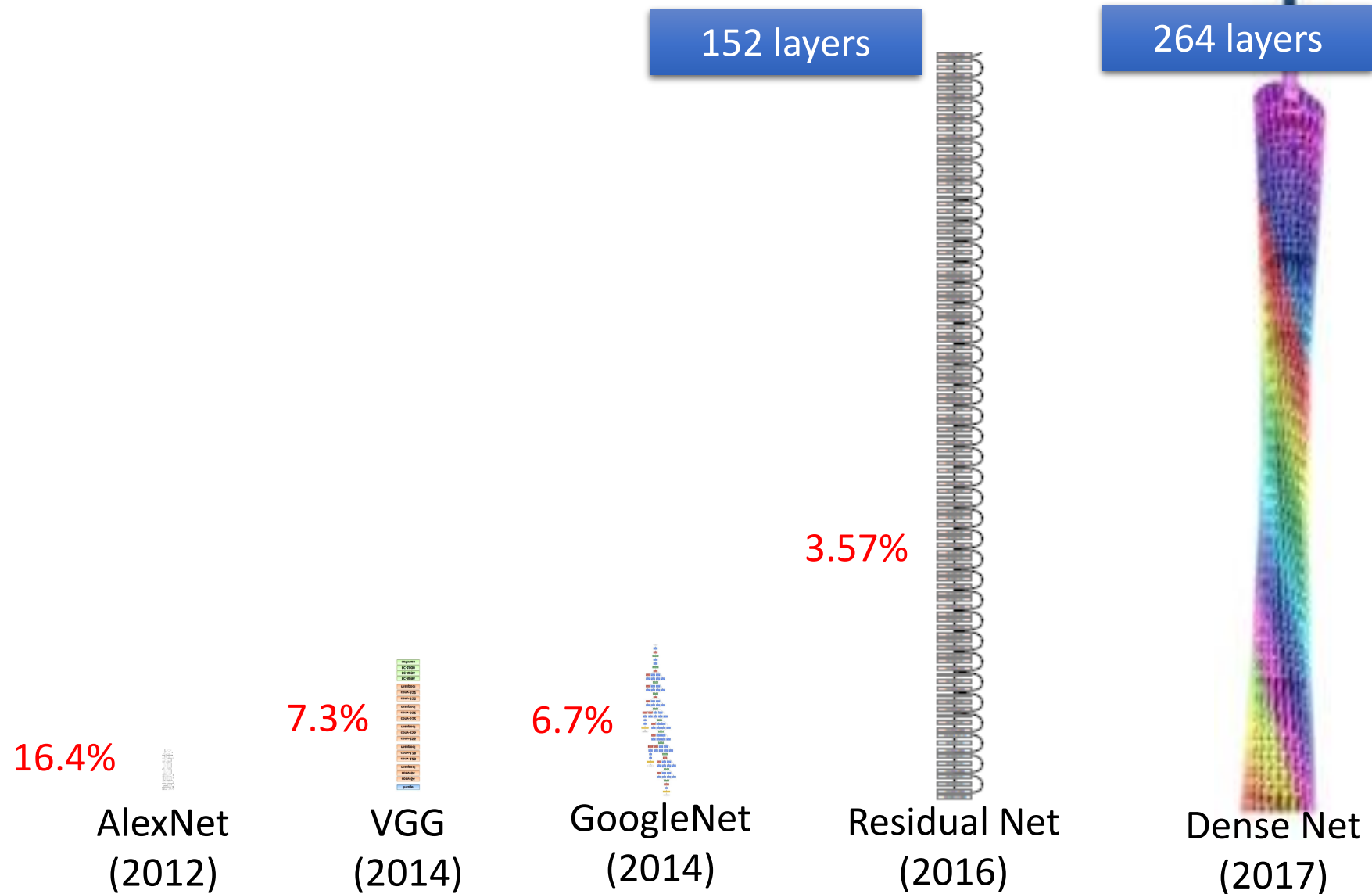
22 layers

6.7%

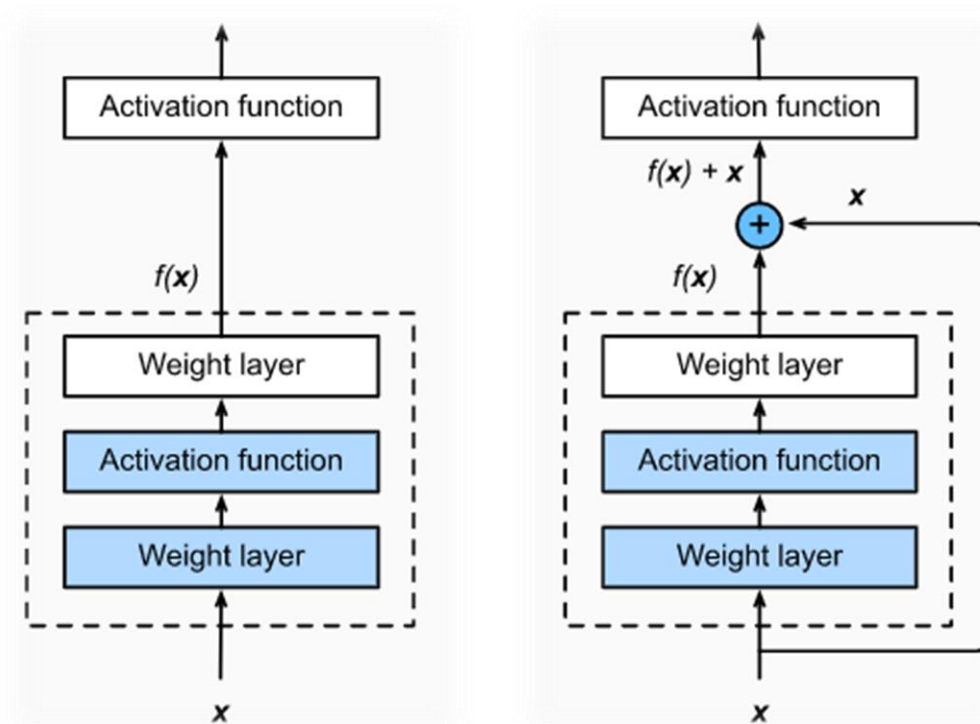


GoogleNet (2014)

Deep Convolutional Neural Network

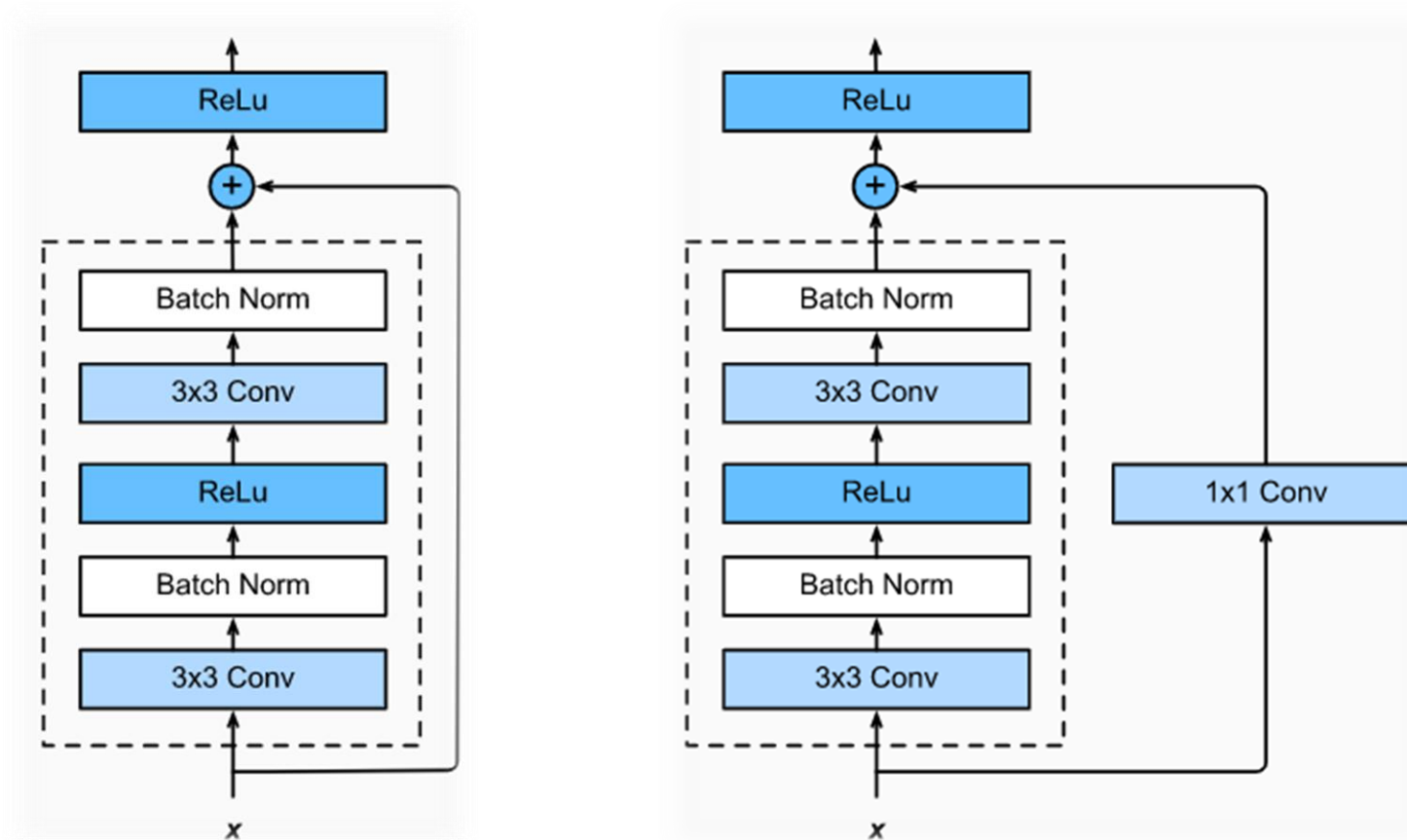


ResNet



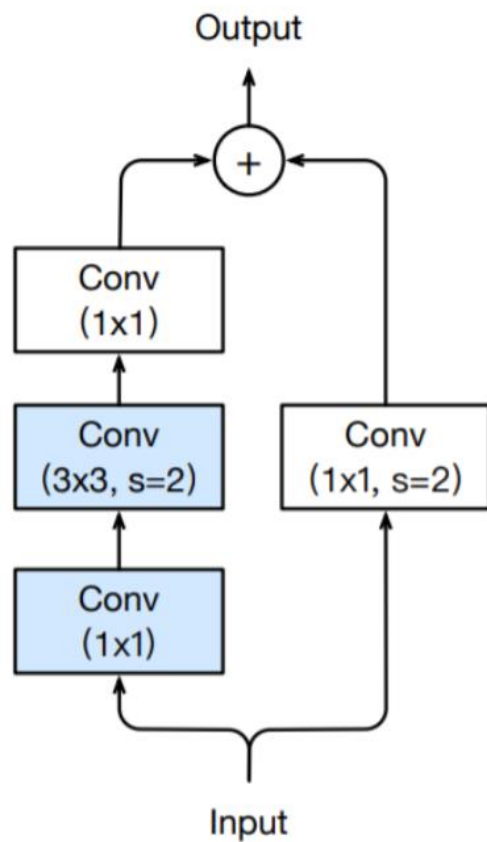
The difference between a regular block (left) and a residual block (right).

ResNet

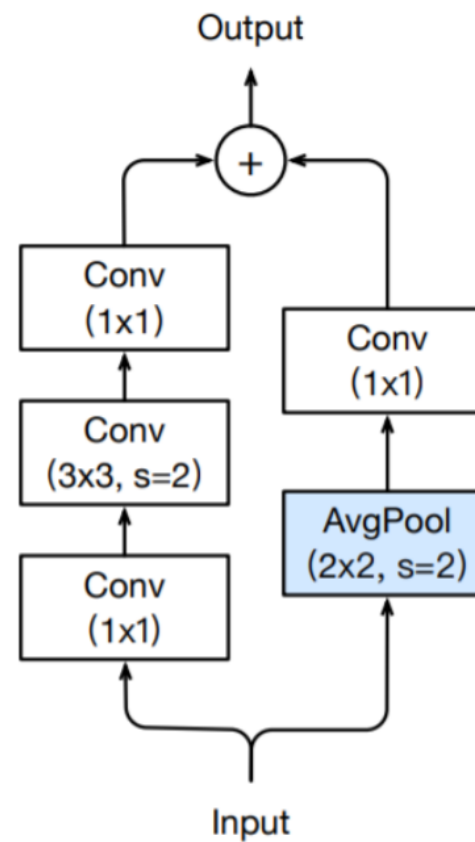


Left: regular ResNet block; Right: ResNet block with 1x1 convolution

ResNet

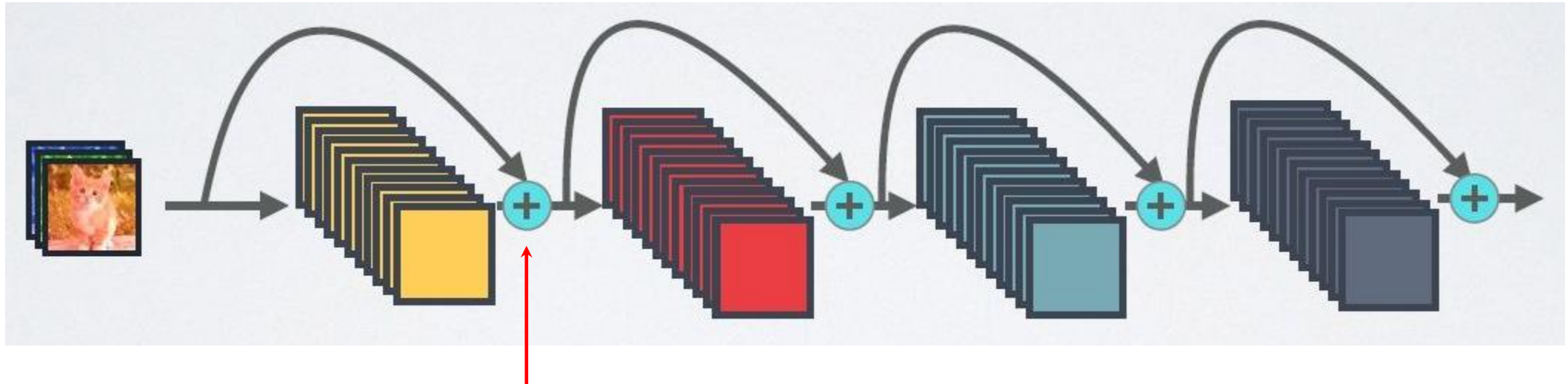


(a) ResNet-B



(c) ResNet-D

ResNet

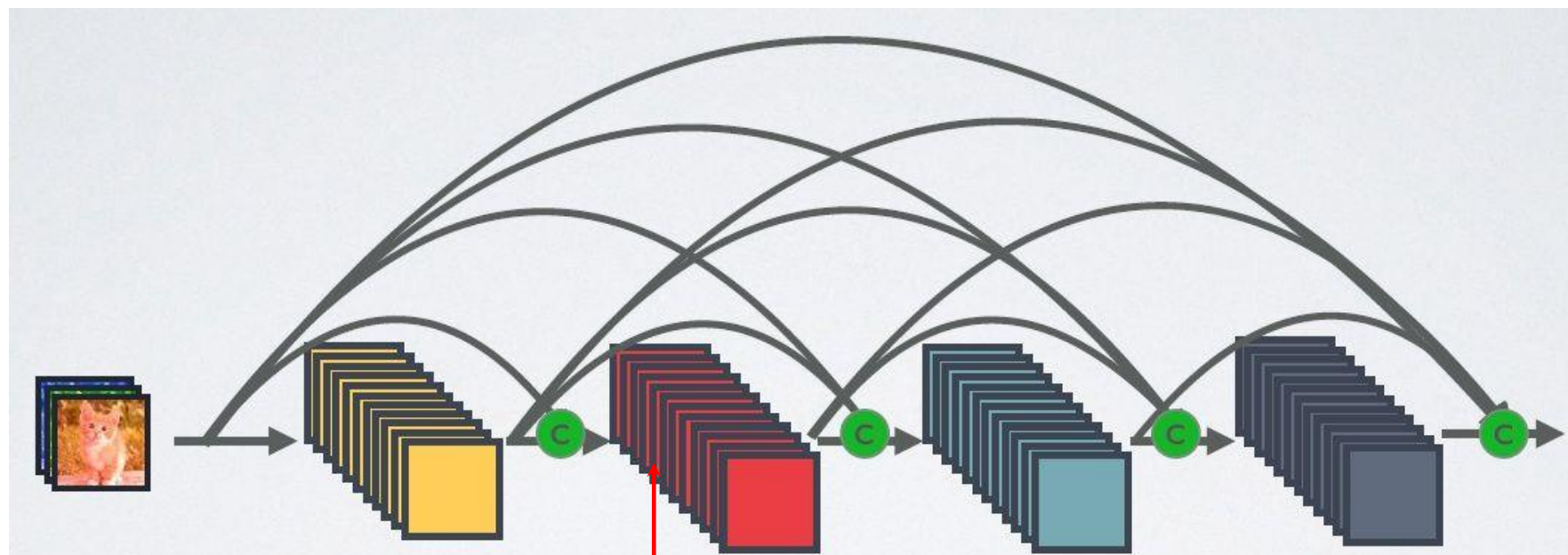


Element-wise Addition

- 1、减轻了vanishing-gradient
- 2、加强了feature的传递

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016.

DenseNet



Channel-wise Concatenation

- 1、减轻了vanishing-gradient
- 2、加强了feature的传递
- 3、更有效地利用了feature
- 4、一定程度上减少了参数数量

Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 4700–4708, 2017

DenseNet: Dense Block

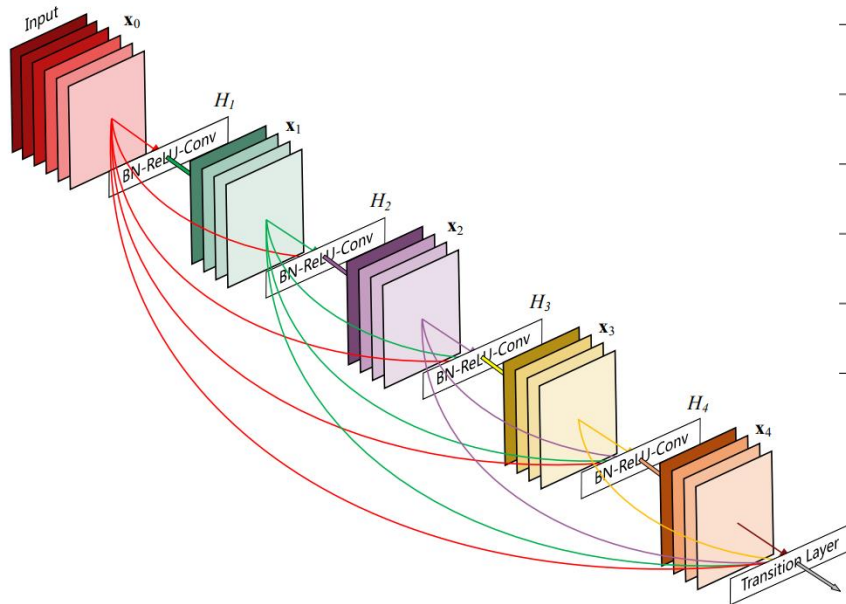


Figure 1: A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.

Layers	Output Size	DenseNet-121	DenseNet-169	DenseNet-201	DenseNet-264
Convolution	112×112	7×7 conv, stride 2			
Pooling	56×56	3×3 max pool, stride 2			
Dense Block (1)	56×56	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$
Transition Layer (1)	56×56	1×1 conv			
	28×28	2×2 average pool, stride 2			
Dense Block (2)	28×28	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$
Transition Layer (2)	28×28	1×1 conv			
	14×14	2×2 average pool, stride 2			
Dense Block (3)	14×14	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 24$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 64$
Transition Layer (3)	14×14	1×1 conv			
	7×7	2×2 average pool, stride 2			
Dense Block (4)	7×7	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 16$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$
Classification Layer	1×1	7×7 global average pool			
		1000D fully-connected, softmax			

CSPNet

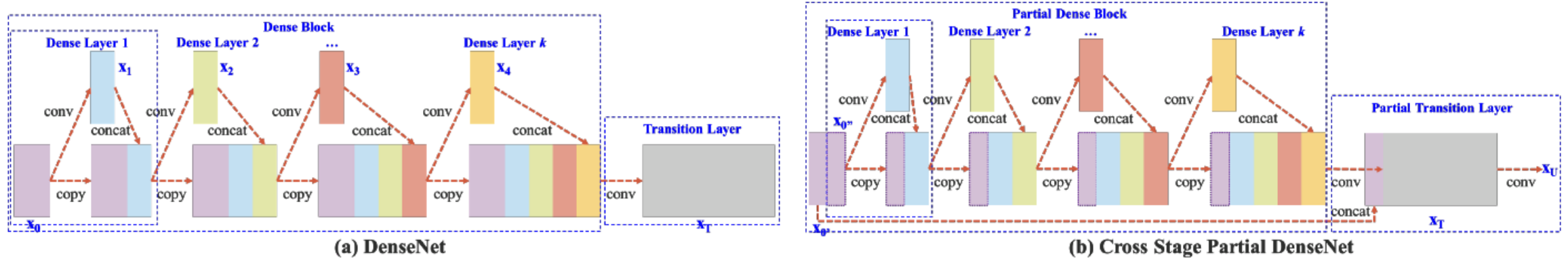
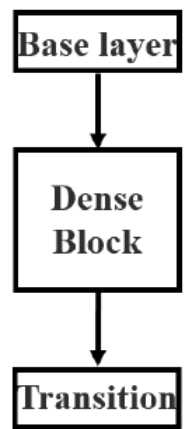


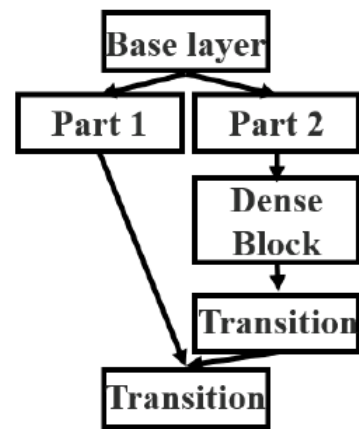
Figure 2: Illustrations of (a) DenseNet and (b) our proposed Cross Stage Partial DenseNet (CSPDenseNet). CSPNet separates feature map of the base layer into two part, one part will go through a dense block and a transition layer; the other one part is then combined with transmitted feature map to the next stage.

Wang, Chien-Yao, et al. "CSPNet: A New Backbone that can Enhance Learning Capability of CNN." arXiv preprint arXiv:1911.11929 (2019).

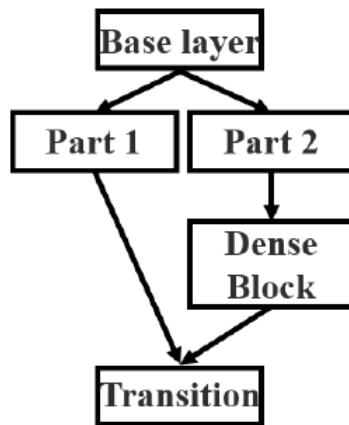
CSPNet



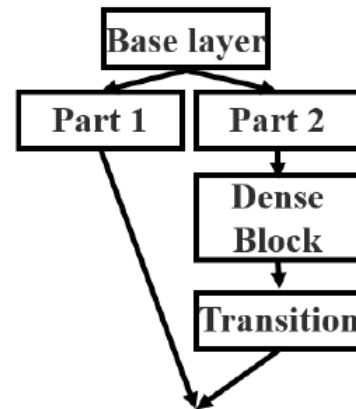
(a) DenseNet



(b) CSPDenseNet



(c) Fusion First



(d) Fusion Last

CSPNet

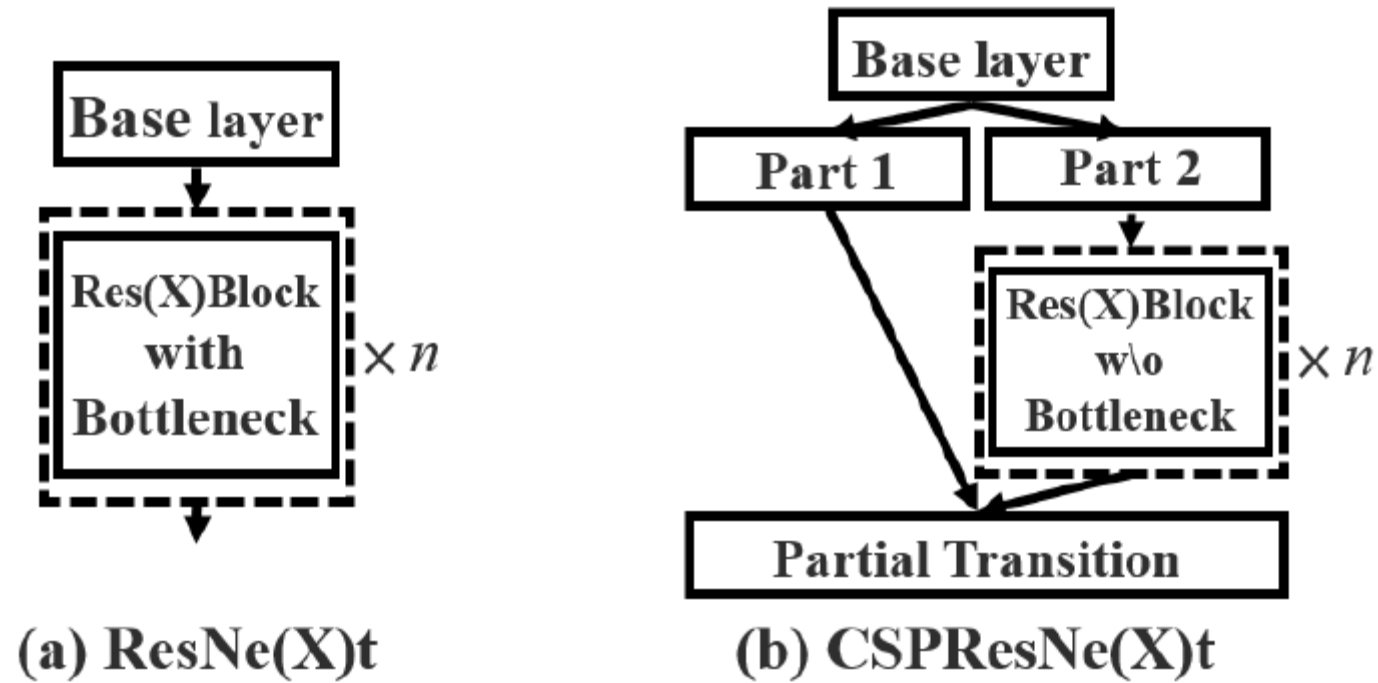


Figure 5: Applying CSPNet to ResNe(X)t.

CSPNet

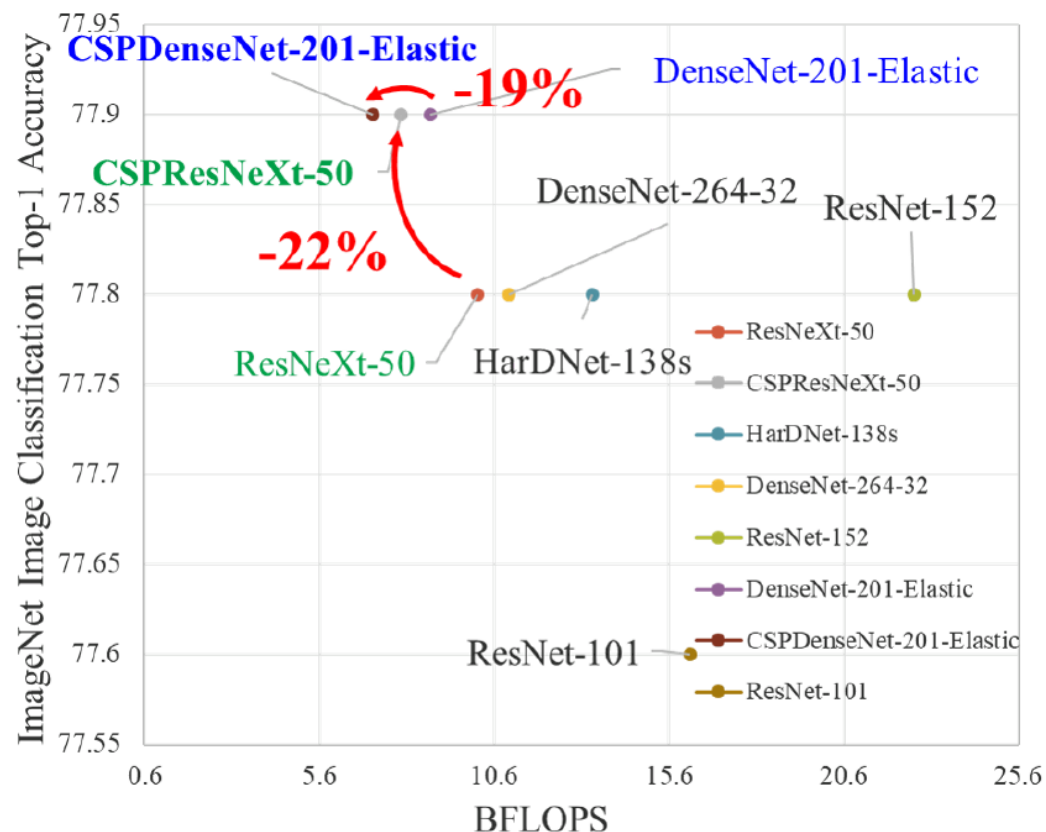


Table 3: Compare with state-of-the-art methods on ImageNet.

Model	#Parameter	BFLOPs	Top-1	Top-5
ResNet-10 [7]	5.24M	2.273	63.5%	85.0%
CSPResNet-10	2.73M	1.905 (-16%)	65.3%	86.5%
ResNeXt-50 [39]	22.19M	10.11	77.8%	94.2%
CSPResNeXt-50	20.50M	7.93 (-22%)	77.9%	94.0%
HarDNet-138s [1]	35.5M	13.4	77.8%	-
DenseNet-264-32 [11]	27.21M	11.03	77.8%	93.9%
ResNet-152 [7]	60.2M	22.6	77.8%	93.6%
DenseNet-201-Elastic [36]	19.48M	8.77	77.9%	94.0%
CSPDenseNet-201-Elastic	20.17M	7.13 (-19%)	77.9%	94.0%

Q&A



Spring 2023