

# 深度学习平台与应用

第六讲: 卷积神经网络架构

范琦

fanqi@nju.edu.cn

2024年10月16日





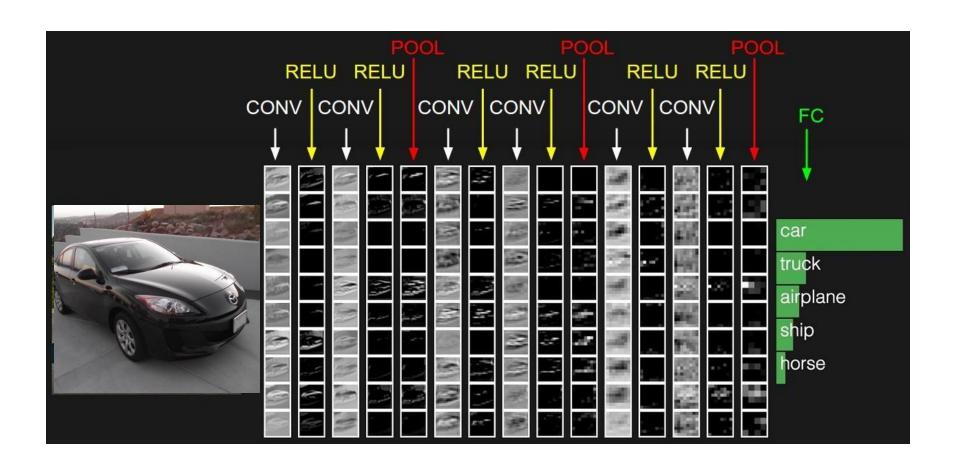


- ■归一化层
- ■经典卷积网络结构

# 卷积神经网络架构



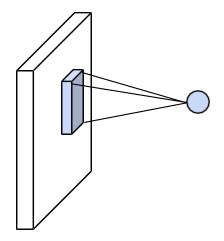
■ 卷积神经网络结构



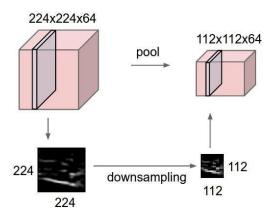
# 卷积神经网络组成部分



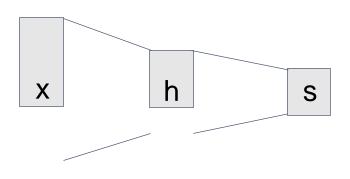
## 卷积层



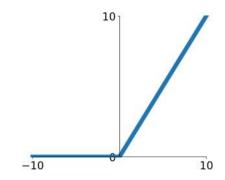
## 池化层



全连接层



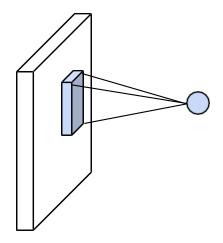
## 激活函数



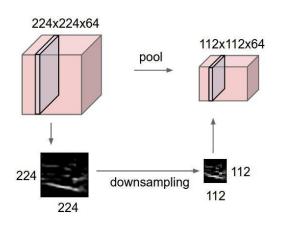
# 卷积神经网络组成部分



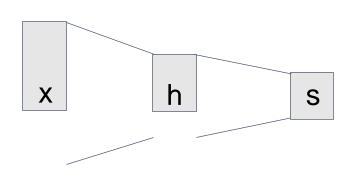
### 卷积层



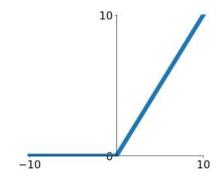
## 池化层



#### 全连接层



## 激活函数



### 批归一化层

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$



- 以下输入会导致网络难以优化:
  - 输入不以零为中心(有较大的偏置)
  - 输入的每个元素具有不同的缩放比例



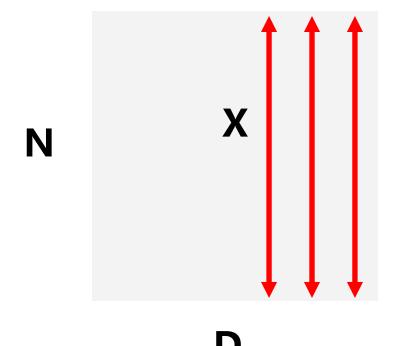
- 以下输入会导致网络难以优化:
  - 输入不以零为中心(有较大的偏置)
  - 输入的每个元素具有不同的缩放比例
- 解决方案:
  - 对输入进行缩放(归一化)



逐通道均值,

大小为D

输入:  $x: N \times D$ 



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差,  
大小为D  $\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{2}}$  归一化后的输入,

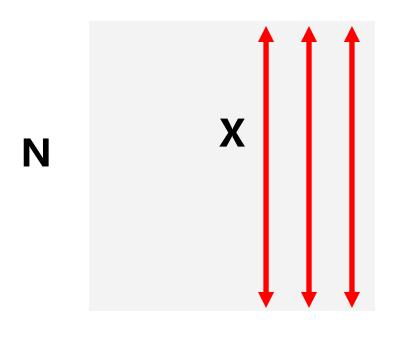
$$\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + arepsilon}}$$
 归一化后的输入, 大小为N x D



逐通道均值,

大小为D

输入:  $x:N\times D$ 



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

$$\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差, 大小为D

$$\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + arepsilon}}$$
 归一化后的输入, 大小为N  $_{\mathbf{X}}$  D

问题: 归一化后的输入丢失了 大量数据信息(均值、方差)



输入:  $x: N \times D$ 

可学习的缩放和偏 移参数:

 $\gamma, \beta: D$ 

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差,  
大小为D

$$\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + arepsilon}}$$
 归一化后的输入, 大小为N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

逐通道均值, 大小为D

缩放后的输出, 大小为N x D



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

逐通道均值, 大小为D

$$\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差,大小为D

$$\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + arepsilon}}$$
 归一化后的输入, 大小为N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

缩放后的输出, 大小为N x D

训练: 从minibatch中估计



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

逐通道均值, 大小为D

$$\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差,大小为D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

归一化后的输入, 大小为N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

缩放后的输出, 大小为N x D

训练: 从minibatch中估计

测试: 使用从训练集中估计

的参数 (移动平均)



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

逐通道均值, 大小为D 训练: 从minibatch中估计

 $\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$  逐通道方差, 大小为D

测试: 使用从训练集中估计

的参数 (移动平均)

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

归一化后的输入, 大小为N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

缩放后的输出, 大小为N x D 训练: 从数据集中学习



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

逐通道均值, 大小为D

训练: 从minibatch中估计

 $\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$  逐通道方差, 大小为D

测试: 使用从训练集中估计

的参数 (移动平均)

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

归一化后的输入, 大小为N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

缩放后的输出, 大小为N x D

训练:从数据集中学习

测试: 使用训练得到的参数



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

逐通道均值, 大小为D

训练: 从minibatch中估计

$$\sigma_j^2 = rac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
 逐通道方差,大小为D

归一化后的输入, 大小为N x D

测试: 使用从训练集中估计

的参数 (移动平均)

 $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_i^2 + \varepsilon}}$ 

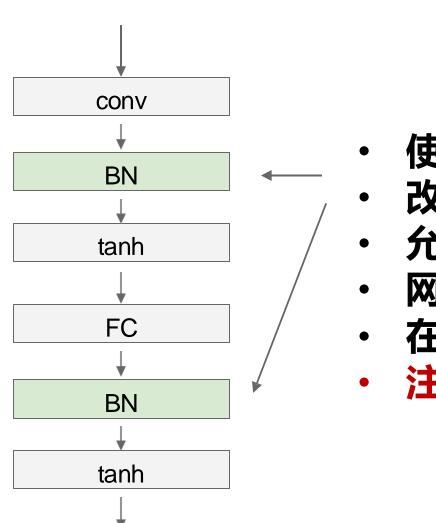
在测试时, batchnorm 变 成了一个线性运算符,可以 与之前的FC或conv层融合

 $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ 缩放后的输出, 大小为N x D

训练: 从数据集中学习

测试: 使用训练得到的参数

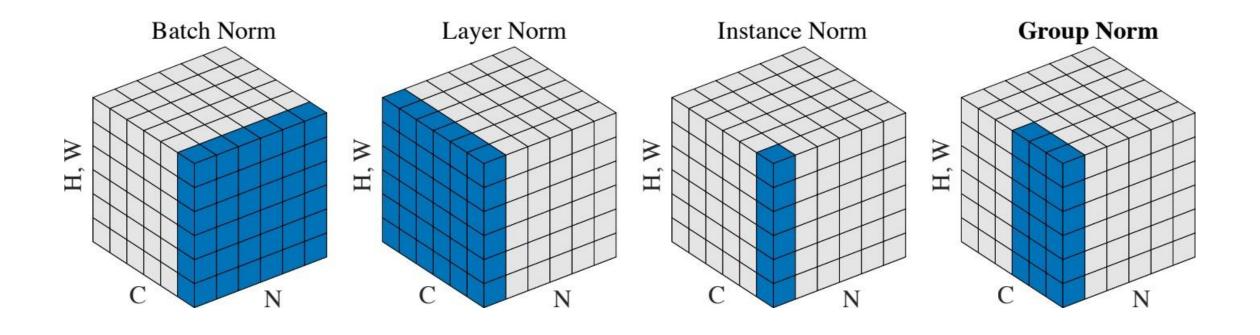




- 使深度网络更容易训练!
- 改善梯度流
- 允许更高的学习率,更快的收敛
- 网络对初始化变得更加稳健
- 在训练过程中起到正则化的作用
- 注意:训练和测试是不同的!

# 不同的归一化层









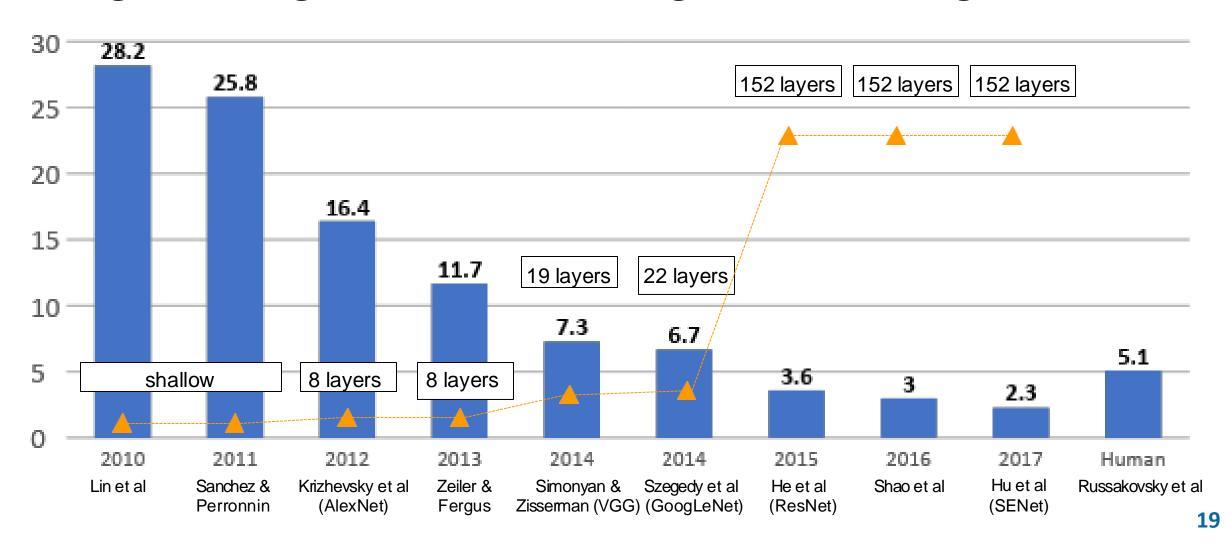


- ■归一化层
- 经典卷积网络结构

# 历年 ILSVRC 比赛冠军



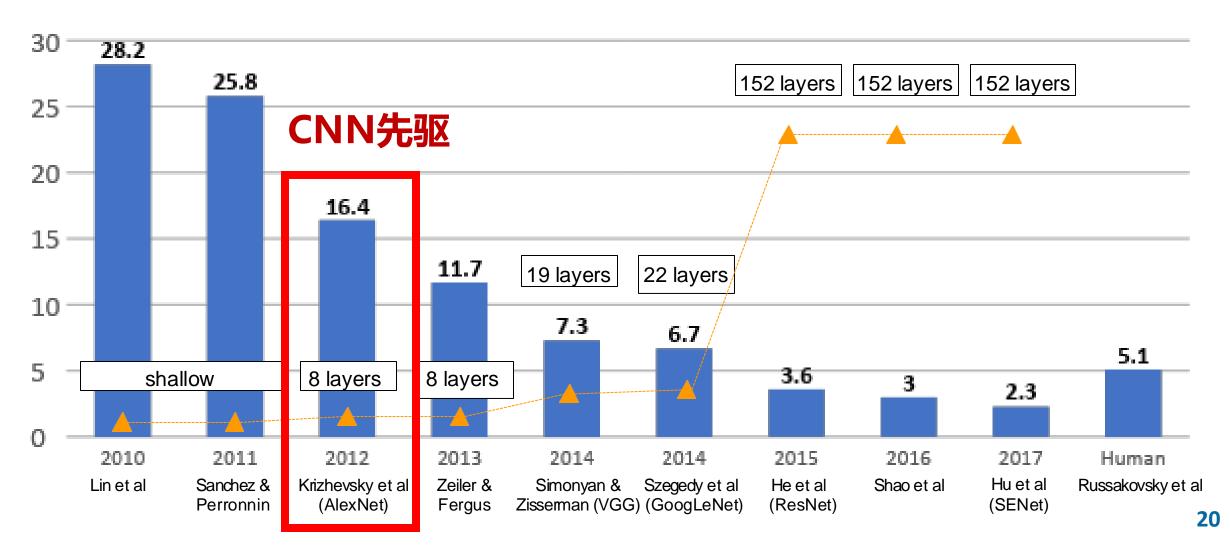
## ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



# 历年 ILSVRC 比赛冠军



## ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



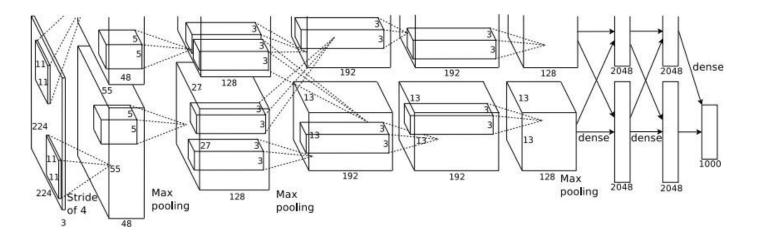


结构:
CONV1
MAX POOL1
NORM1
CONV2
MAX POOL2
NORM2
CONV3
CONV4
CONV5
Max POOL3

FC6

FC7

FC8



100 miles

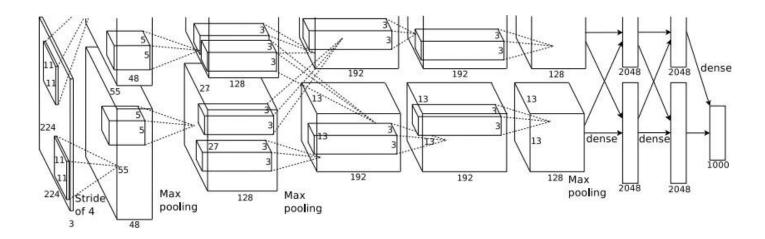
#### NIPS papers

https://papers.nips.cc > paper > 4824-imagenet-classificat... :

#### ImageNet Classification with Deep Convolutional Neural ...

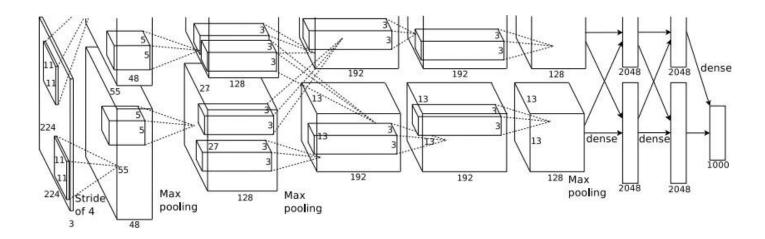
by A Krizhevsky · 2012 · Cited by 134307 — The neural network, which has 60 million parameters and 500,000 neurons, consists of five convolutional layers, some of which are...





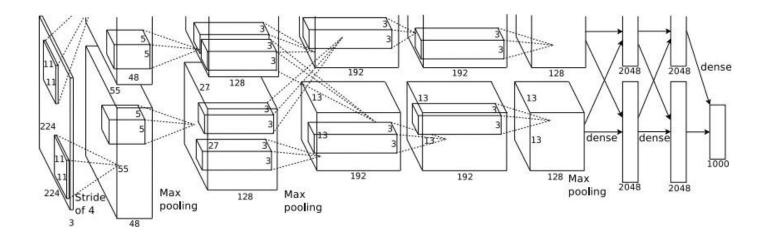
- 输入: 227x227x3
- 第一层: CONV1, 96个11x11滤波器, 步长为4





- 输入: 227x227x3
- 第一层: CONV1, 96个11x11滤波器, 步长为4
- 输出: 55x55x96 W' = (W F + 2P) / S + 1





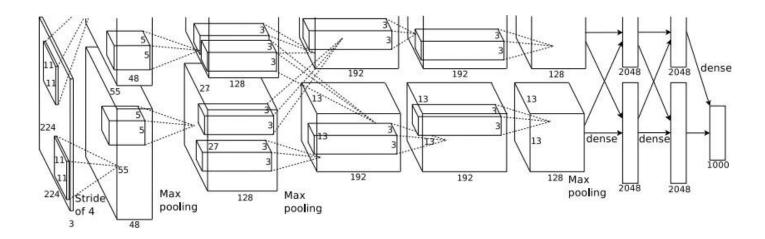
■ 输入: 227x227x3

■ 第一层: CONV1, 96个11x11滤波器, 步长为4

■ 输出: 55x55x96

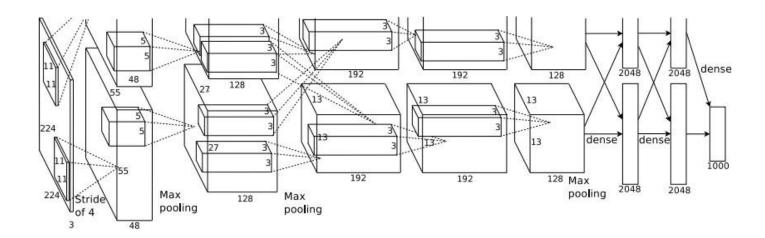
■ 参数量: (11\*11\*3 + 1)\*96 = 35K





- 输入: 227x227x3
- 第二层: POOL1, 3x3大小, 步长为2





■ 输入: 227x227x3

■ 第二层: POOL1, 3x3大小, 步长为2

■ 输出: 27x27x96

■ 参数量: 0



#### AlexNet 结构:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

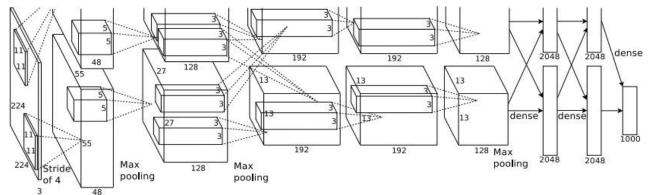
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)





#### AlexNet 结构:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

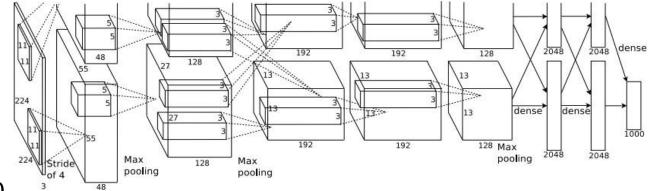
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



- 首次使用 ReLU
- 使用 LRN 层 (现在不常用)
- 大量使用数据增强
- dropout 0.5
- -batch size 128
- SGD Momentum 0.9
- 学习率 1e-2, 在验证集准确率达到高原时学习率乘以0.1
- -L2 weight decay 5e-4
- 7个模型融合: 18.2% -> 15.4%

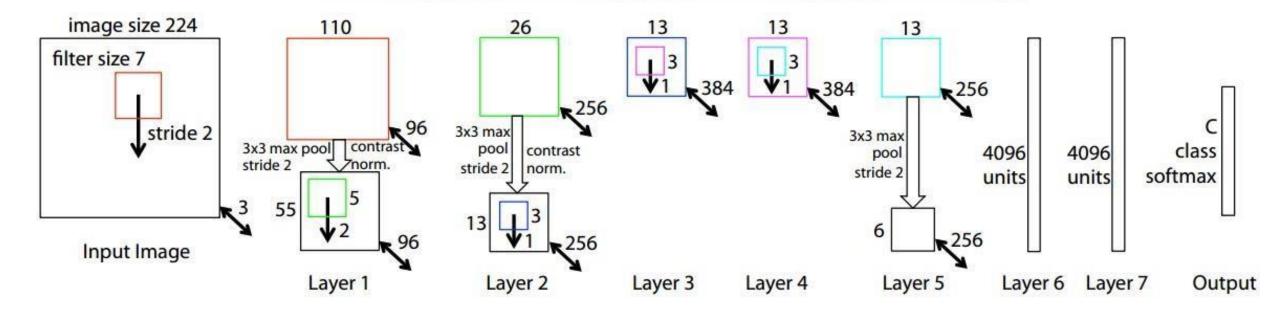






#### Visualizing and Understanding Convolutional Networks

by MD Zeiler  $\cdot$  2013  $\cdot$  Cited by 23189 — We introduce a novel visualization technique that gives insight into the function of intermediate feature layers and the operation of the classifier.



### 在AlexNet的基础上进行超参数优化:

CONV1: 由 (11x11 stride 4) 改为 (7x7 stride 2)

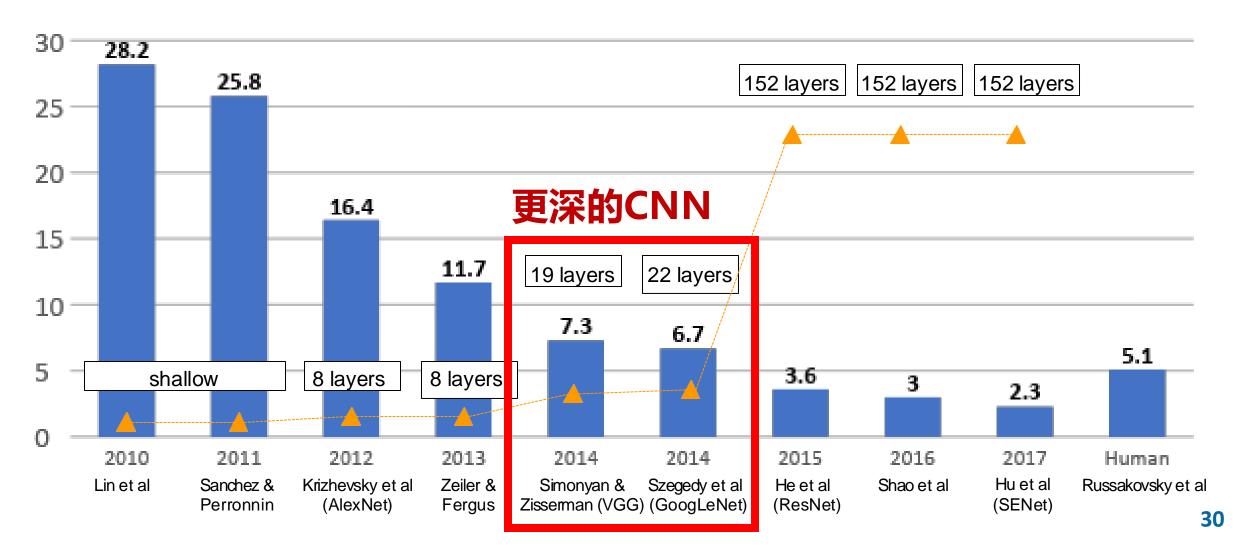
CONV3,4,5: 通道数由 384, 384, 256 改为 512, 1024, 512

**ImageNet top5 error: 16.4%** → 11.7%

# 更深的卷积神经网络



## ImageNet Large Scale Visual Recognition Challenge (ILSVRC)





- 更小的滤波器,更深的网络
- 只使用
  - 3x3 Conv, 歩长1
  - 2x2 Max Pooling, 歩长2

■ 深度: 8→16-19

■ 精度: 11.7%→7.3%

arXiv https://arxiv.org > cs

Very Deep Convolutional Networks for Large-Scale Image ...

by K Simonyan · 2014 · Cited by 131122 — In this work we investigate the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting.

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 256
3x3 conv, 384
Pool
3x3 conv, 384
Pool
5x5 conv, 256
11x11 conv, 96
Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
Pool
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

Table 1: **ConvNet configurations** (shown in columns). The depth of the configurations increases from the left (A) to the right (E), as more layers are added (the added layers are shown in bold). The convolutional layer parameters are denoted as "conv $\langle$ receptive field size $\rangle$ - $\langle$ number of channels $\rangle$ ". The ReLU activation function is not shown for brevity.



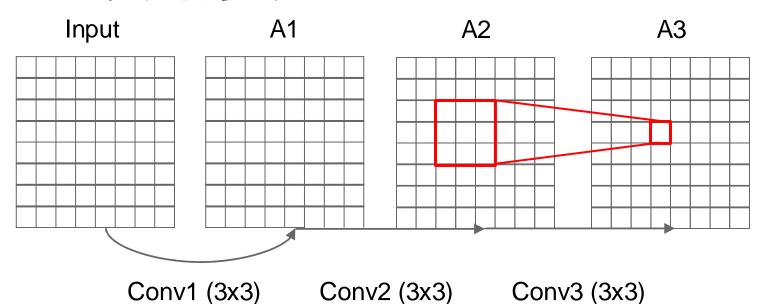
		ConvNet C	onfiguration	26.30	7,70,7
A	A-LRN	В	С	D	Е
11 weight	11 weight	13 weight	16 weight	16 weight	19 weight
layers	layers	layers	layers	layers	layers
	i	nput ( $224 \times 2$	24 RGB image	e)	76.
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64
	LRN	conv3-64	conv3-64	conv3-64	conv3-64
	16,97,76,5	max	pool		_ %, %
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
	37, 2, 4.	conv3-128	conv3-128	conv3-128	conv3-128
			pool	60	2,70,7
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
			conv1-256	conv3-256	conv3-256
	3, 3, 3	0	77, 25	3.0	conv3-256
	5 92 72 05 T		pool	2,70	5 V2 V2
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
			conv1-512	conv3-512	conv3-512
	8, 6, 7,	6	72, 70, 7	2.70	conv3-512
	_6,95,76,5		pool		_6,0%
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
		2	conv1-512	conv3-512	conv3-512
		6	142/1762	2007	conv3-512
	1, 0, 7, 7, 7		pool	Te.	15 0
			4096	9,5	₹0.19/1
	7/2/25:25		4096	3,0	7/2/3
	5° 45.775.765		1000	2.50	5/2/2
		soft	-max	27	· · · · · · · · · · · · · · · · · · ·

Table 2: Number of parameters (in millions).

Network	A,A-LRN	В	С	D	Е	
Number of parameters	133	133	134	138	144	2



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量

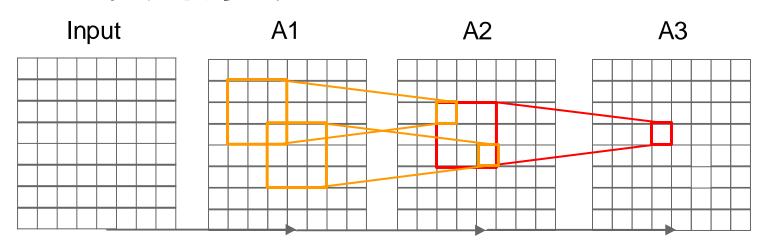


Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
Pool
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量



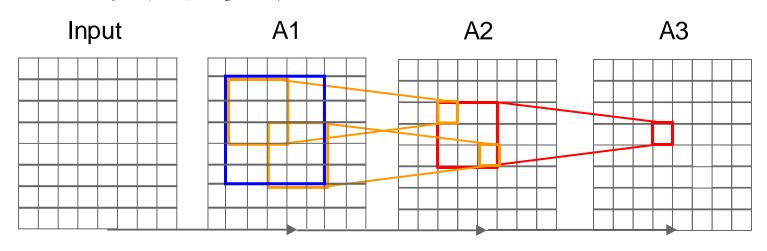
Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
Pool
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

34



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量



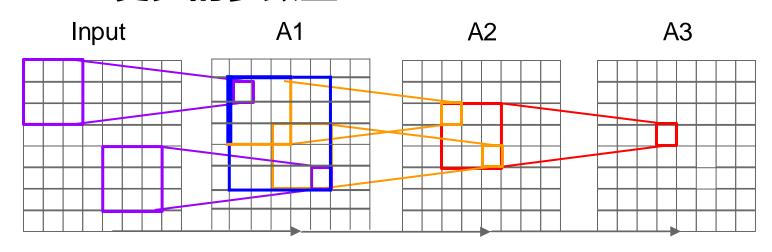
Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64 Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
Pool
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

**35** 



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量



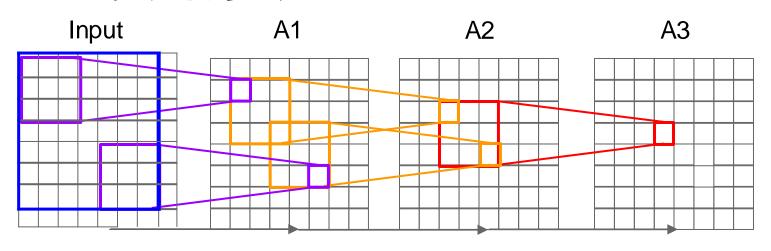
Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
Pool
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

36



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量



Softmax							
FC 1000							
FC 4096							
FC 4096							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 256							
3x3 conv, 256							
Pool							
3x3 conv, 128							
3x3 conv, 128							
Pool							
3x3 conv, 64							
3x3 conv, 64 3x3 conv, 64 Input							

Softmax							
FC 1000							
FC 4096							
FC 4096							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 256							
3x3 conv, 256							
Pool							
3x3 conv, 128							
3x3 conv, 128							
Pool							
3x3 conv, 64							
3x3 conv, 64							
Input							

**37** 



- 为什么只使用3x3 Conv?
- 多个 3x3 Conv 相对一个 7x7 Conv:
  - 相同的感受野
  - 更少的参数量
  - 更深, 更多的非线性变换

Softmax						
FC 1000						
FC 4096						
FC 4096						
Pool						
3x3 conv, 512						
3x3 conv, 512						
3x3 conv, 512						
Pool						
3x3 conv, 512						
3x3 conv, 512						
3x3 conv, 512						
Pool						
3x3 conv, 256						
3x3 conv, 256						
Pool						
3x3 conv, 128						
3x3 conv, 128						
Pool						
3x3 conv, 64						
3x3 conv, 64						
Input						

-							
Softmax							
FC 1000							
FC 4096							
FC 4096							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
3x3 conv, 512							
Pool							
3x3 conv, 256							
3x3 conv, 256							
Pool							
3x3 conv, 128							
3x3 conv, 128							
Pool							
3x3 conv, 64							
3x3 conv, 64							
Input							



INPUT: [224x224x3] memory: 224*224*3=150K params: 0 不计算bias CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728								
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864								
POOL2: [112x112x64] memory: 112*112*64=800K params: 0								
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728								
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456								
POOL2: [56x56x128] memory: 56*56*128=400K params: 0								
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912								
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824								
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824								
POOL2: [28x28x256] memory: 28*28*256=200K params: 0								
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648								
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296								
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296								
POOL2: [14x14x512] memory: 14*14*512=100K params: 0								
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296								
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296								
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296								
POOL2: [7x7x512] memory: 7*7*512=25K params: 0 FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = 102,760,448 FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216								
FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000								

Softmax FC 1000 FC 4096 FC 4096 Pool 3x3 conv, 512 3x3 conv, 512 3x3 conv, 512 Pool 3x3 conv, 512 3x3 conv, 512

3x3 conv, 512 Pool 3x3 conv, 256 3x3 conv, 256 Pool 3x3 conv, 128 3x3 conv, 128 Pool 3x3 conv, 64 3x3 conv, 64 Input



INPUT: [224x224x3]

CONV3-64: [224x224

CONV3-64: [224x224

POOL2: [112x112x64

CONV3-128: [112x11

CONV3-128: [112x11

POOL2: [56x56x128]

CONV3-256: [56x56]

CONV3-256: [56x56]

CONV3-256: [56x56x

Memory for parameters

Memory for param gradients

Memory for momentum

Memory for layer outputs

Memory for layer error

POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*256)\*512 = 1,179,648

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512=2,359,296

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512 = 2,359,296

POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512=2,359,296

POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0

FC: [1x1x4096] memory: 4096 params: 7\*7\*512\*4096 = 102,760,448

FC: [1x1x4096] memory: 4096 params: 4096\*4096 = 16,777,216

FC: [1x1x1000] memory: 1000 params: 4096\*1000 = 4,096,000

总显存占用: 24M \* 4 bytes

~= 96MB / image

总参数量: 138M parameters



不计算bias INPUT: [224x224x3] memory: 224\*224\*3=150K params: 0 CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M (3\*3\*3)\*64 = 1,728 CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params: (3\*3\*64)\*64=3.36POOL2: [112x112x64] memory: 112\*112\*64=800K params: 0 CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params: (3\*3\*64)\*128 = 73,728 CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params: (3\*3\*128)\*128 = 147,456POOL2: [56x56x128] memory: 56\*56\*128=400K params: 0 CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*128)\*256=294,912CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*256)\*256=589,824CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*256)\*256=589,824POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0 CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*256)\*512 = 1,179,648 CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512=2,359,296CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512=2,359,296POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0 CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,205CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*542 = 2,359,296 POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0 FC: [1x1x4096] memory: 4096 params: 7\*7\*512\*4096 = 102,760,448

FC: [1x1x4096] memory: 4096 params: 4096\*4096 = 16,777,216

FC: [1x1x1000] memory: 1000 params: 4096\*1000 = 4,096,000

显存占用瓶颈 浅层Conv

参数量瓶颈

FC层

总显存占用: 24M \* 4 bytes

~= 96MB / image

总参数量: 138M parameters

FC: [1x1x1000] memory: 1000 params: 4096\*1000 = 4,096,000



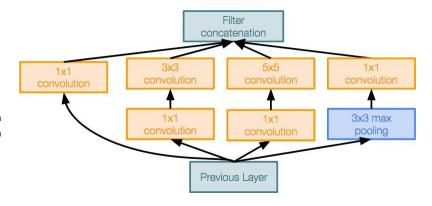
INPLIT: [224x224x3]		
1141 O 1. [22+x22+x0] Memory. 22+ 22+ 0=10010 params. 0	Softmax	
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728	FC 1000	fc8
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864	FC 4096	fc7
POOL2: [112x112x64] memory: 112*112*64=800K params: 0	FC 4096	fc6
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728	Pool	
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456	3x3 conv, 512	conv5-3
POOL2: [56x56x128] memory: 56*56*128=400K params: 0	3x3 conv, 512	conv5-2
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912	3x3 conv, 512	conv5-1
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824	Pool 3x3 conv, 512	conv4-3
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824	3x3 conv, 512	conv4-2
POOL2: [28x28x256] memory: 28*28*256=200K params: 0	3x3 conv, 512	conv4-1
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648	Pool	
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296	3x3 conv, 256	conv3-2
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296	3x3 conv, 256	conv3-1
POOL2: [14x14x512] memory: 14*14*512=100K params: 0	Pool	
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296	3x3 conv, 128	conv2-2
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296	3x3 conv, 128	conv2-1
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296	Pool 3x3 conv, 64	conv1-2
DOOL 0. [7.7.7.540] magazini 7*7*540 OFK magazini 0	3x3 conv, 64	conv1-1
	Input	
FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = 102,760,448 ~= 96MB / image		
FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216	VGG16	

总参数量: 138M parameters

## GoogLeNet



- 更深的网络,更高的计算效率
- ILSVRC′ 14 分类任务冠军 (6.7% top 5 error)
- 22层
- 只有 5 million 参数
  - 12x less than AlexNet
  - 27x less than VGG-16



Inception module

- 高效的 Inception 模块
- 没有 FC 隐藏层



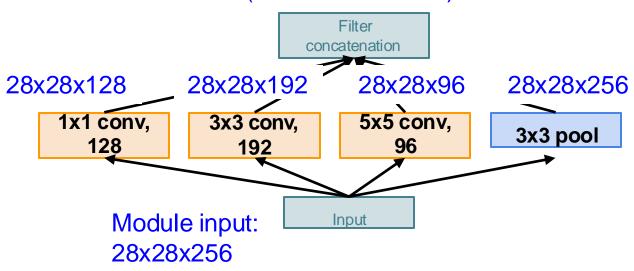
#### [1409.4842] Going Deeper with Convolutions

by C Szegedy · 2014 · Cited by 62440 — One particular incarnation used in our submission for ILSVRC 2014 is called **GoogLeNet**, a 22 layers deep network, the quality of which is ...



■ 简单的 Inception: 并联多个不同的算子

28x28x(128+192+96+256) = 28x28x672





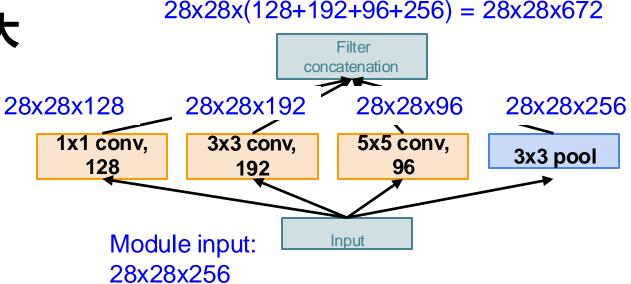
■ 简单的 Inception: 并联多个不同的算子

■ 问题一: 计算量和参数量太大

#### **Conv Ops:**

[1x1 conv, 128] 28x28x128x1x1x256 [3x3 conv, 192] 28x28x**192x3x3x256** [5x5 conv, 96] 28x28x**96x5x5x256** 

Total: 854M ops

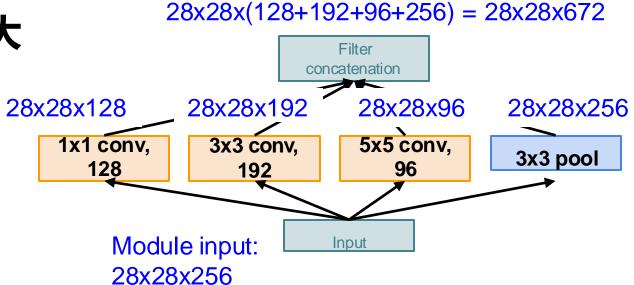




■ 简单的 Inception: 并联多个不同的算子

■ 问题一: 计算量和参数量太大

■ 问题二:通道数持续变大

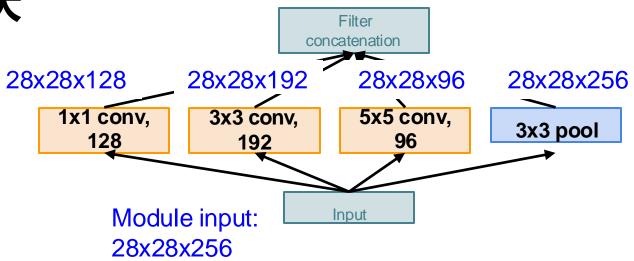




■ 简单的 Inception: 并联多个不同的算子

■ 问题一: 计算量和参数量太大

■ 问题二:通道数持续变大

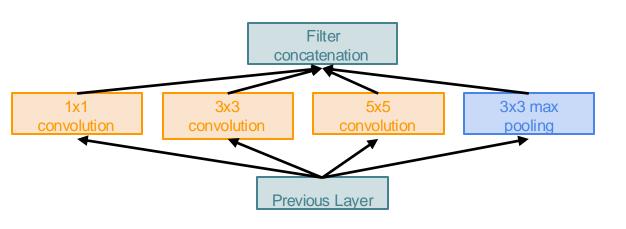


28x28x(128+192+96+256) = 28x28x672

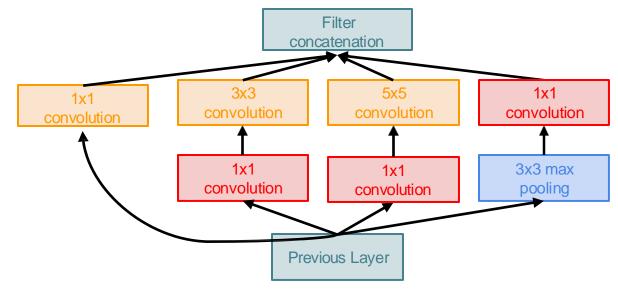
■ Solution: 使用 1x1 Conv 降低计算量、参数量、通道数



#### 1x1 conv 瓶颈层

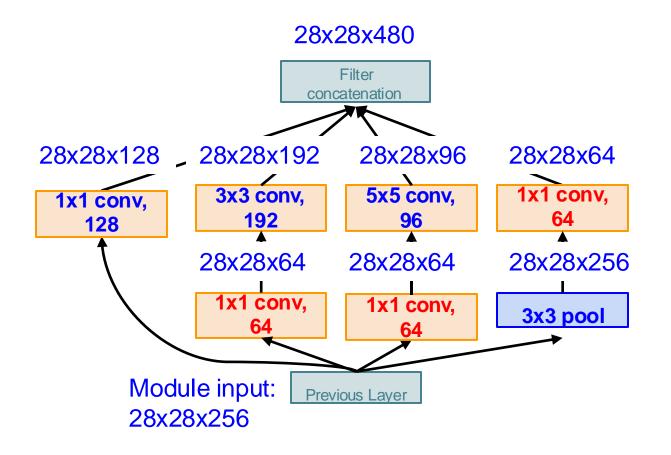


原始的 Inception 模块



GoogLeNet 中的 Inception 模块





#### **Conv Ops:**

[1x1 conv, 64] 28x28x64x1x1x256 [1x1 conv, 64] 28x28x64x1x1x256

[1x1 conv, 128] 28x28x128x1x1x256

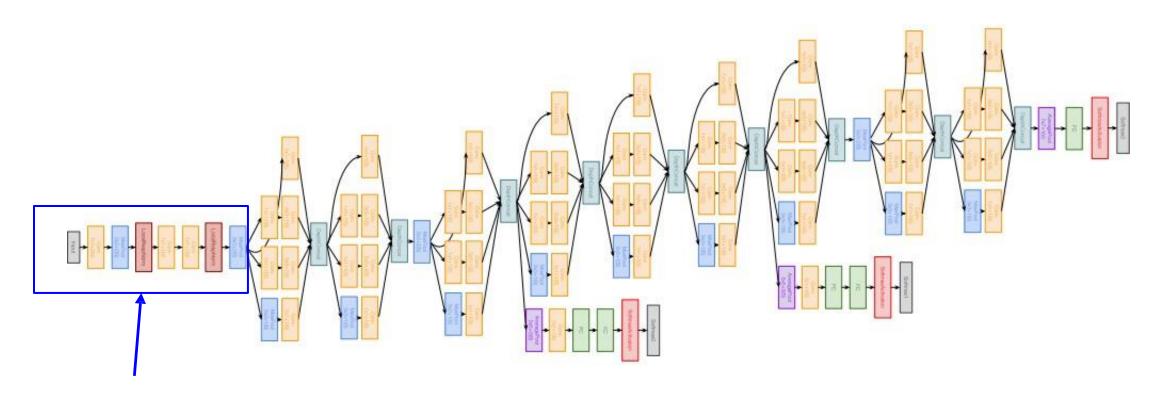
[3x3 conv, 192] 28x28x192x3x3x64

[5x5 conv, 96] 28x28x96x5x5x64

[1x1 conv, 64] 28x28x64x1x1x256

Total: 358M ops

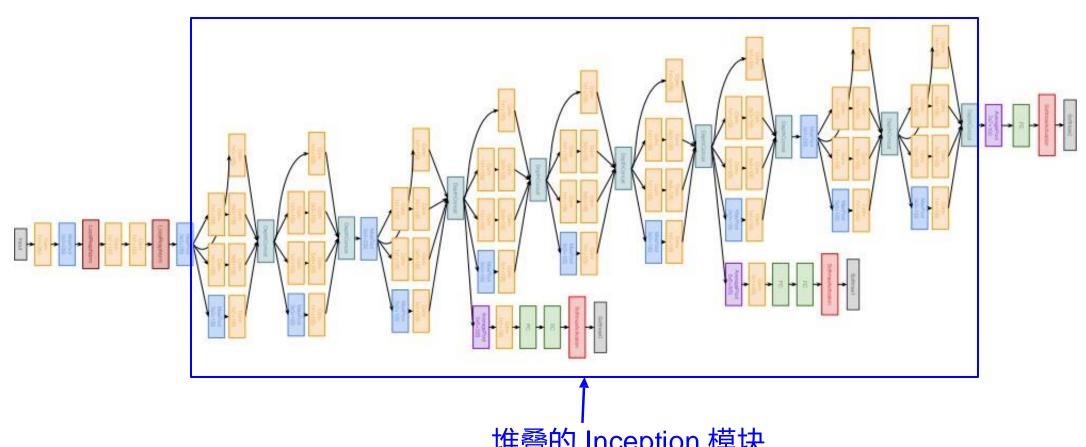




Stem Network: Conv-Pool-

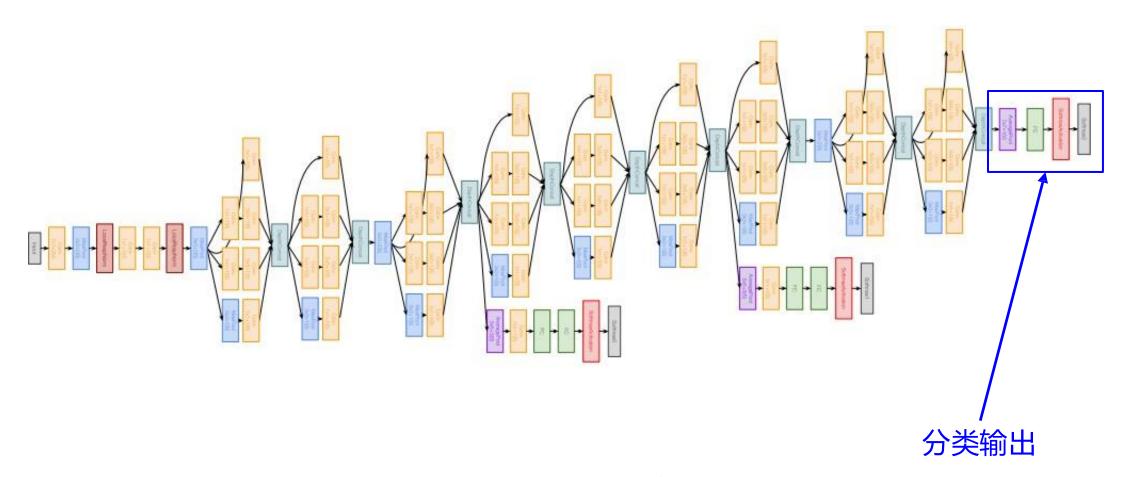
2x Conv-Pool





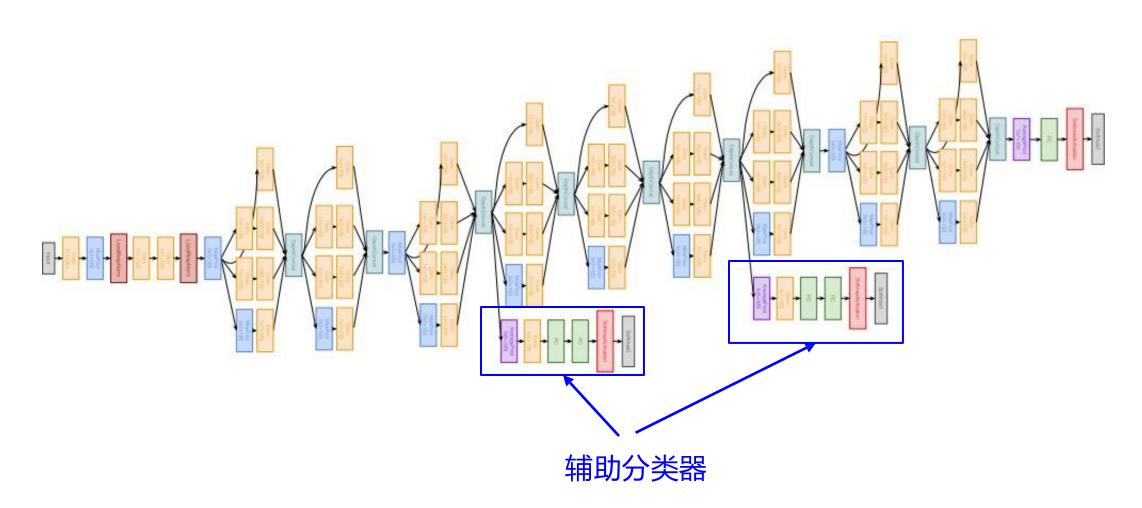
堆叠的 Inception 模块



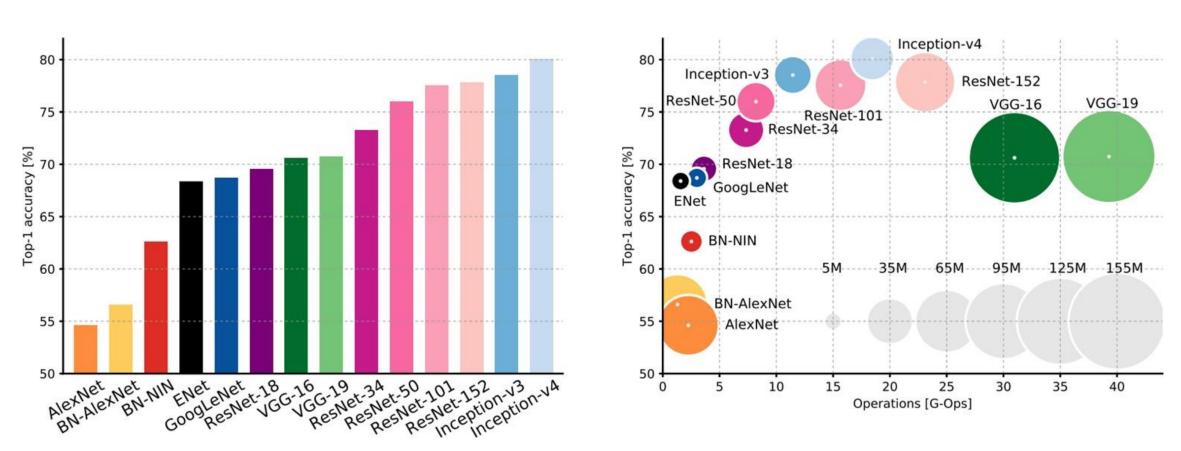


GoogLeNet 只在分类器中使用的 FC 层







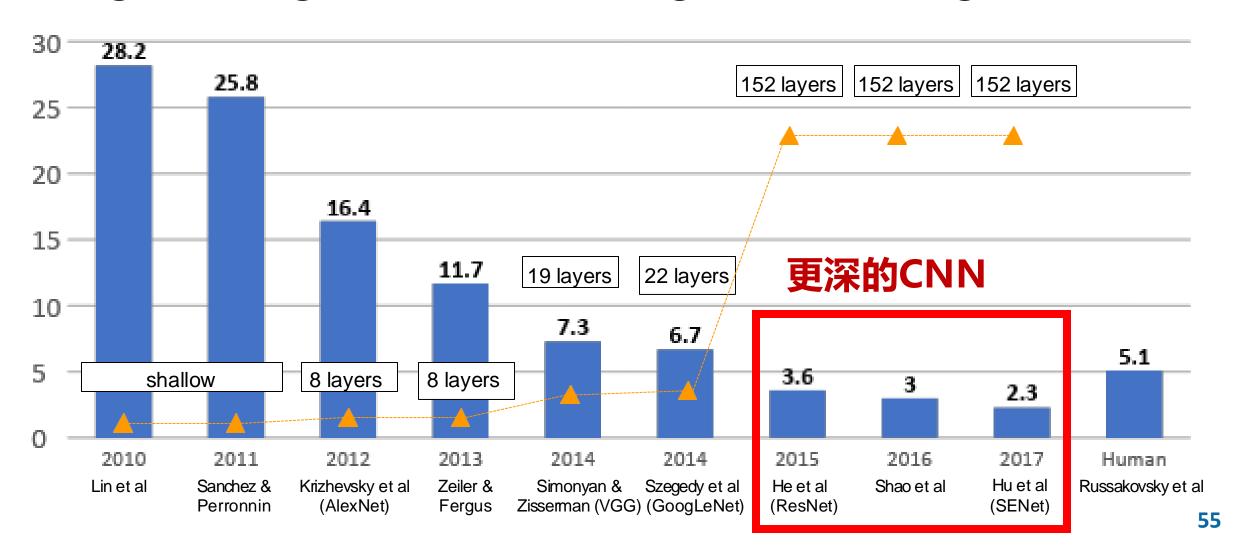


An Analysis of Deep Neural Network Models for Practical Applications, 2017.

## 残差卷积神经网络



### ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

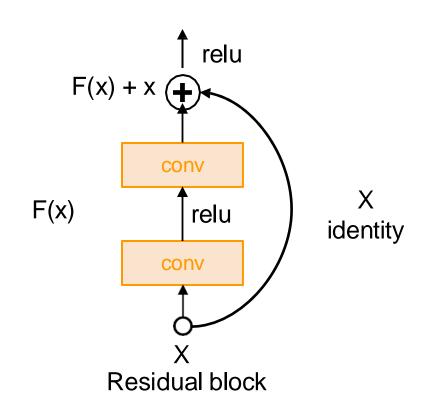


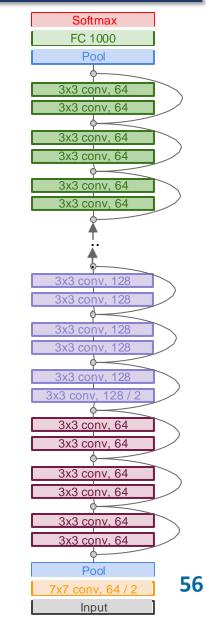


- 残差神经网络
- ILSVRC′15 冠军
- 残差连接
- 解决 CNN 训练问题
- 152层 (后续: 1001层)

MSRA @ ILSVRC & COCO 2015 Competitions

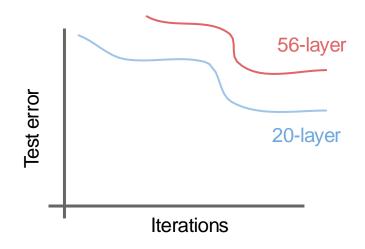
- 1st places in all five main tracks
  - ImageNet Classification: "Ultra-deep" (quote Yann) 152-layer nets
  - ImageNet Detection: 16% better than 2nd
  - ImageNet Localization: 27% better than 2nd
  - COCO Detection: 11% better than 2nd
  - COCO Segmentation: 12% better than 2nd

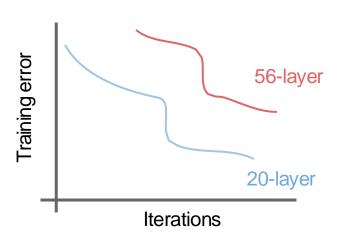






- 训练深层卷积网络存在问题
  - 深层网络的性能比浅层的更差
  - 但不是由过拟合造成的







- 训练深层卷积网络存在问题
  - 深层网络的性能比浅层的更差
  - 但不是由过拟合造成的
- 深层网络比浅层网络具有更强的表达能力 (更多参数)



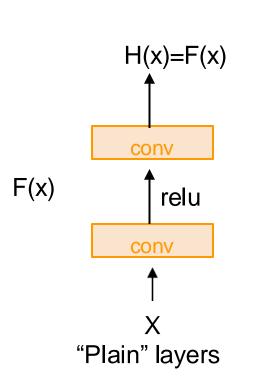
- 训练深层卷积网络存在问题
  - 深层网络的性能比浅层的更差
  - 但不是由过拟合造成的
- 深层网络比浅层网络具有更强的表达能力(更多参数)
- 猜想1: 优化问题,深度网络更难优化

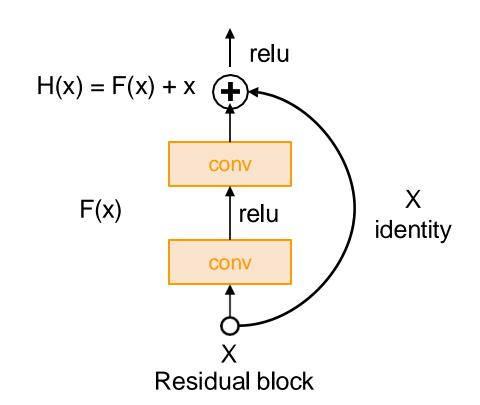


- 训练深层卷积网络存在问题
  - 深层网络的性能比浅层的更差
  - 但不是由过拟合造成的
- 深层网络比浅层网络具有更强的表达能力 (更多参数)
- 猜想1:优化问题,深度网络更难优化
- 关键问题: 让深层网络的能力至少与浅层网络一样好(恒等映射)



#### ■ 关键问题: 让深层网络的能力至少与浅层网络一样好

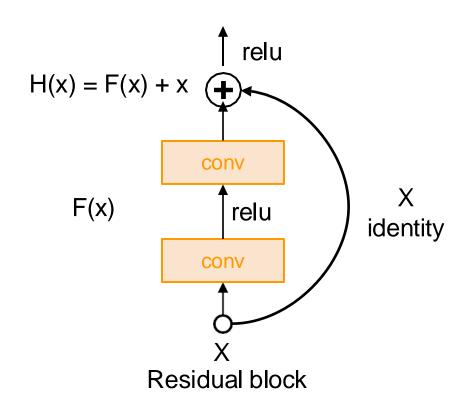


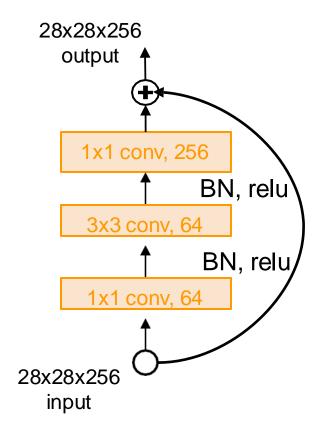


让网络拟合残差 F(x)=H(x)-x, 而不是直接拟合H(x)



#### ■ 使用 Bottleneck 提升计算效率





### ResNet 结构





arXiv https://arxiv.org > cs :

#### [1512.03385] Deep Residual Learning for Image Recognition

by K He · 2015 · Cited by 239504 — We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously.

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer	
conv1	112×112	7×7, 64, stride 2					
7	6.95.76	72	76.25.75	3×3 max pool, stric	de 2	26536	
conv2_x	56×56	56×56	$\begin{bmatrix} 3 \times 3, 64 \end{bmatrix}_{\times 2}$	$\left[\begin{array}{c} 3\times3, 64\\ 3\times3, 64 \end{array}\right]\times3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$
		$\begin{bmatrix} 3 \times 3, 64 \end{bmatrix}^{2}$	$\begin{bmatrix} 3 \times 3, 64 \end{bmatrix}^{3}$	[ 1×1, 256 ]	1×1, 256	1×1, 256	
conv3_x	28×28	$\left[\begin{array}{c} 3\times3, 128\\ 3\times3, 128 \end{array}\right] \times 2$	$\left[\begin{array}{c} 3\times3, 128\\ 3\times3, 128 \end{array}\right] \times 4$	$ \left[\begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array}\right] \times 4 $	$ \left[\begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array}\right] \times 4 $	$ \left[\begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array}\right] \times 8 $	
conv4_x	14×14	$\left[\begin{array}{c} 3\times3,256\\ 3\times3,256 \end{array}\right]\times2$	$ \begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6 $	$   \begin{bmatrix}     1 \times 1, 256 \\     3 \times 3, 256 \\     1 \times 1, 1024   \end{bmatrix} \times 6 $	$ \left[\begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array}\right] \times 23 $	$ \begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36 $	
conv5_x	7×7	$\left[\begin{array}{c} 3\times3,512\\ 3\times3,512 \end{array}\right]\times2$	$\left[\begin{array}{c}3\times3,512\\3\times3,512\end{array}\right]\times3$	$ \left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array}\right] \times 3 $	$ \left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array}\right] \times 3 $	$ \left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array}\right] \times 3 $	
3/20/	1×1	average pool, 1000-d fc, softmax					
FLO	OPs	$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$	

## 经典卷积网络总结



■ AlexNet: CNN可以用于计算机视觉任务

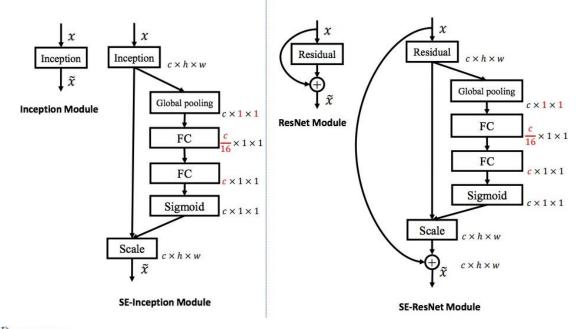
■ VGG: 网络深度的重要性

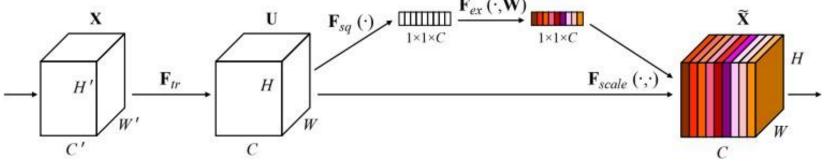
■ GoogLeNet: 网络宽度的重要性

■ ResNet:训练极深卷积网络,CNN 的表现优于人类!



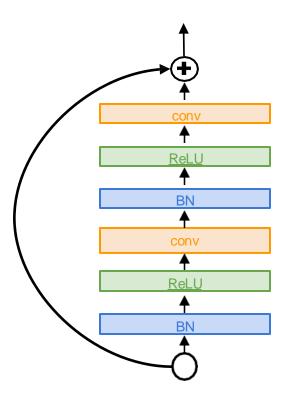
- Squeeze-and-Excitation Networks (SENet, 2017)
- 为 ResNet 加上注意力
- ILSVRC′ 17 分类任务冠军
- 引领后续各种注意力
- Transformer 也是基于注意力





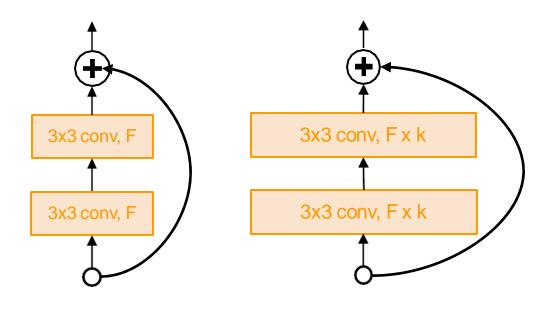


- Identity Mappings in Deep Residual Networks (2016)
- 改进 ResNet 模块设计
- 让网络中的信息更容易传播
- 深度达到1001层





- Wide Residual Networks (2016)
- 增加宽度,减少深度,计算效率更高(可并行化)
- Wide-ResNet-50 性能优于 ResNet-152





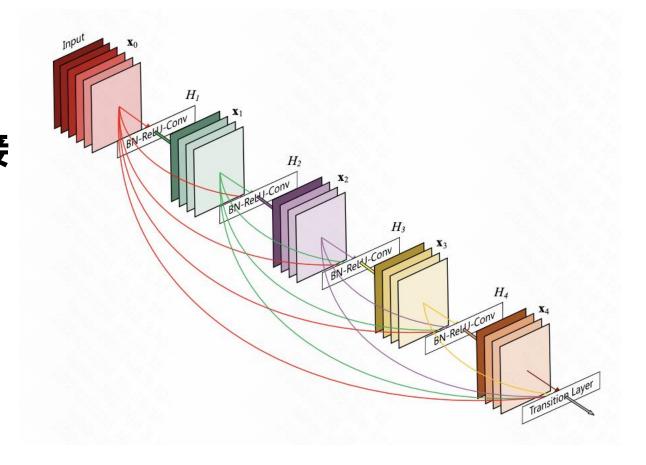
Aggregated Residual Transformations for Deep Neural Networks (ResNeXt, 2016)

增加宽度,保持深度,不增加参数量 256-d out 256-d out 1x1 conv, 256 1x1 conv, 1x1 conv, 1x1 conv, 32 paths 3x3 conv. 64 3x3 conv. 4 1x1 conv, 64 1x1 conv. 4 1x1 conv. 4 256-d in



Densely Connected Convolutional Networks (DenseNet, CVPR 2017 Best Paper)

- 密集使用残差模块
- 每一层都与前面所有层连接





- 将在分类任务上训练好的网络迁移到其他网络
  - 分类任务数据量大,标注难度低,标注成本低
  - 其他任务数据量小,标注难度高,标注成本高
  - 直接在小数据上训练网络,性能很差



#### 1.在分类任务上训练

FC-1000

FC-4096

FC-4096

MaxPool

Conv-512

Conv-512

MaxPool

Conv-512

Conv-512

MaxPool

Conv-256

Conv-256

MaxPool

**Conv-128** 

Conv-128

**MaxPool** 

Conv-64

Conv-64

**Image** 



#### 1.在分类任务上训练

FC-1000 FC-4096 FC-4096

MaxPool Conv-512 Conv-512

MaxPool Conv-512 Conv-512

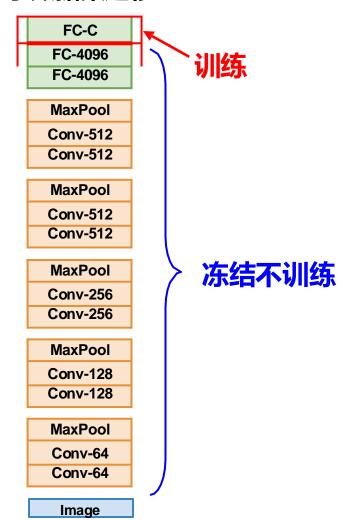
MaxPool Conv-256 Conv-256

MaxPool
Conv-128
Conv-128

MaxPool Conv-64 Conv-64

**Image** 

#### 2.小数据集迁移





#### 1.在分类任务上训练

FC-1000 FC-4096 FC-4096

**MaxPool Conv-512** Conv-512

MaxPool **Conv-512** Conv-512

MaxPool Conv-256 Conv-256

MaxPool Conv-128 Conv-128

**MaxPool** Conv-64 Conv-64

**Image** 

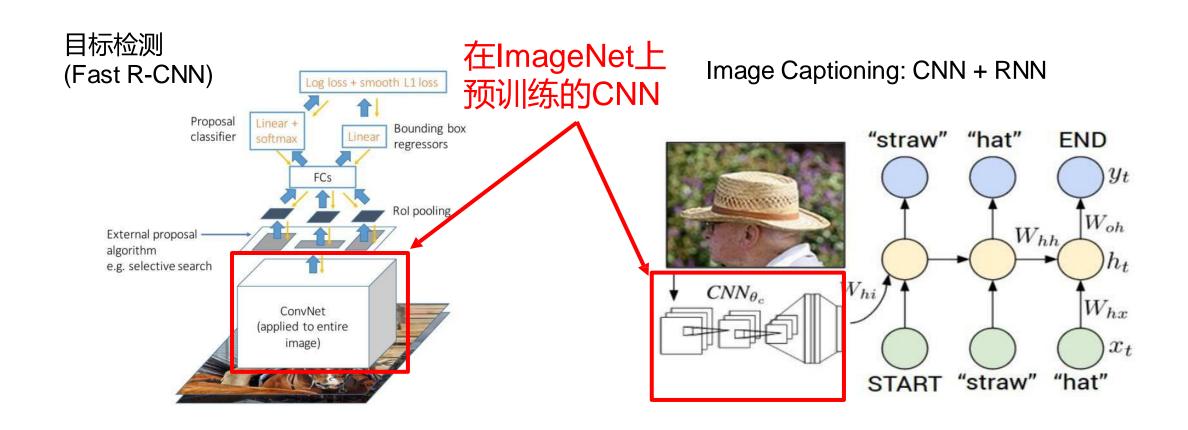
#### 2.小数据集迁移



#### 3.更大数据集迁移



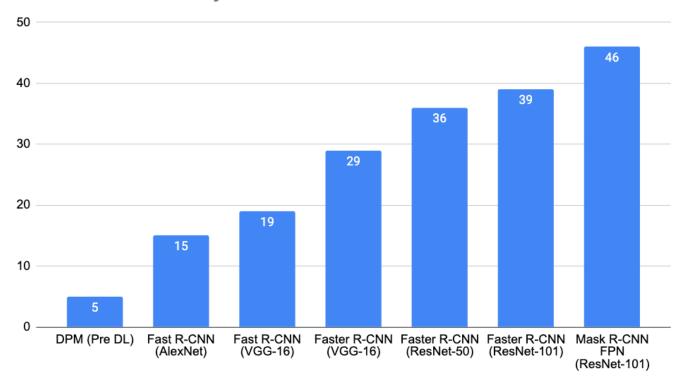




# 迁移学习-设计更好的模型

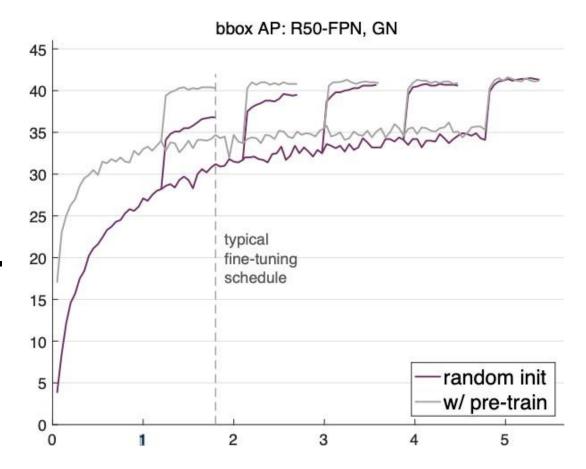


#### Object detection on MSCOCO





- 后续的研究发现预训练并不是必 须的
- 更长的训练时长+tricks可以达 到相同的性能
- 但是,使用合适的预训练模型一 定不会得到更差的结果
- 合适的预训练!



# 下节课



■ 训练卷积神经网络