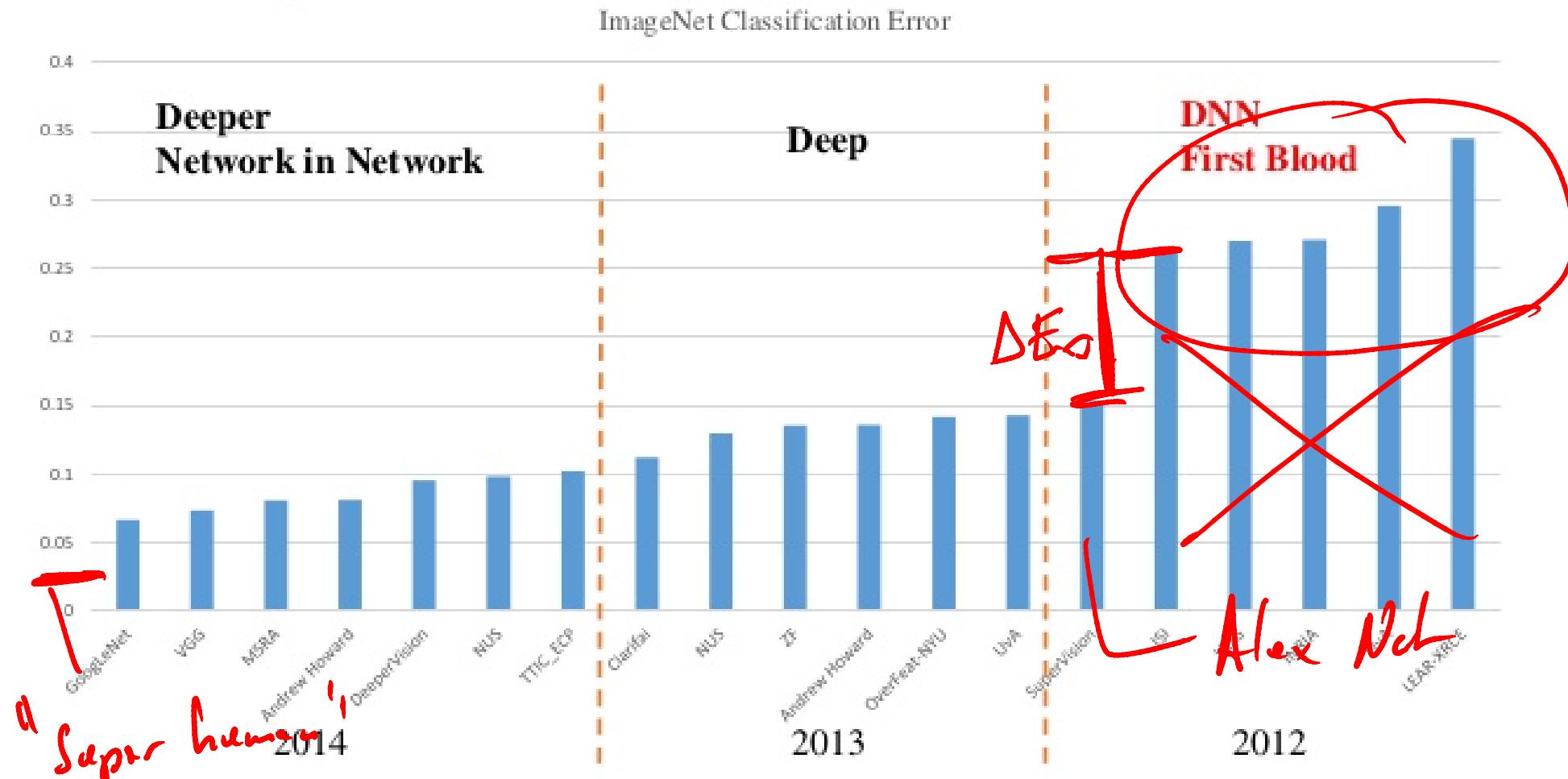


# Neural Networks

Jan Chorowski

# ImageNet Classification

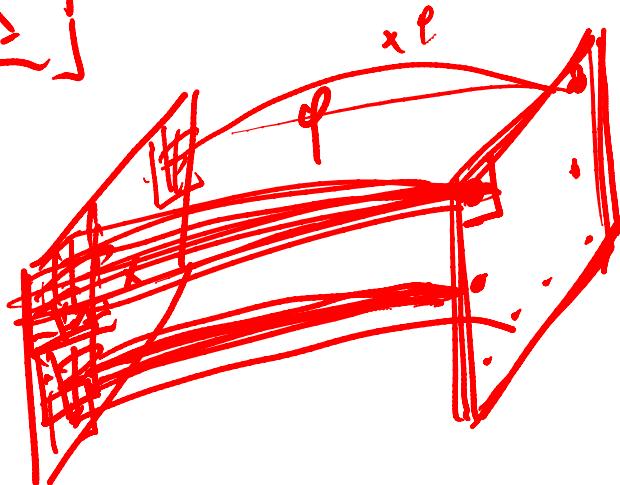
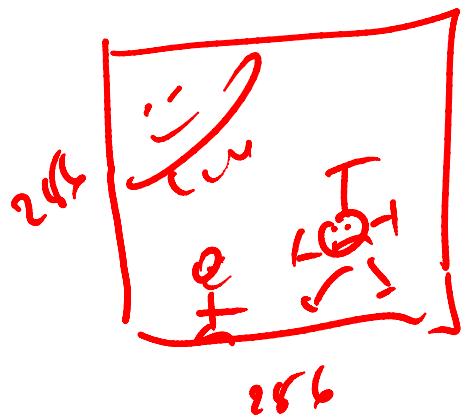
- 1000 categories and 1.2 million training images



Li Fei-Fei: ImageNet Large Scale Visual Recognition Challenge, 2014 <http://image-net.org/>



# Intuitions



①  $u_i$  to  $x^t$

$w \in \mathbb{R}$

$$256 \times 256$$

$\mathbb{R}$   
LOCAL  
CONNECTIVITY

$$256^2 \cdot q =$$

$$9 \cdot 2^{10} \cdot 2^6 = 6 \cdot 10^5 = 600 \text{ k word}$$

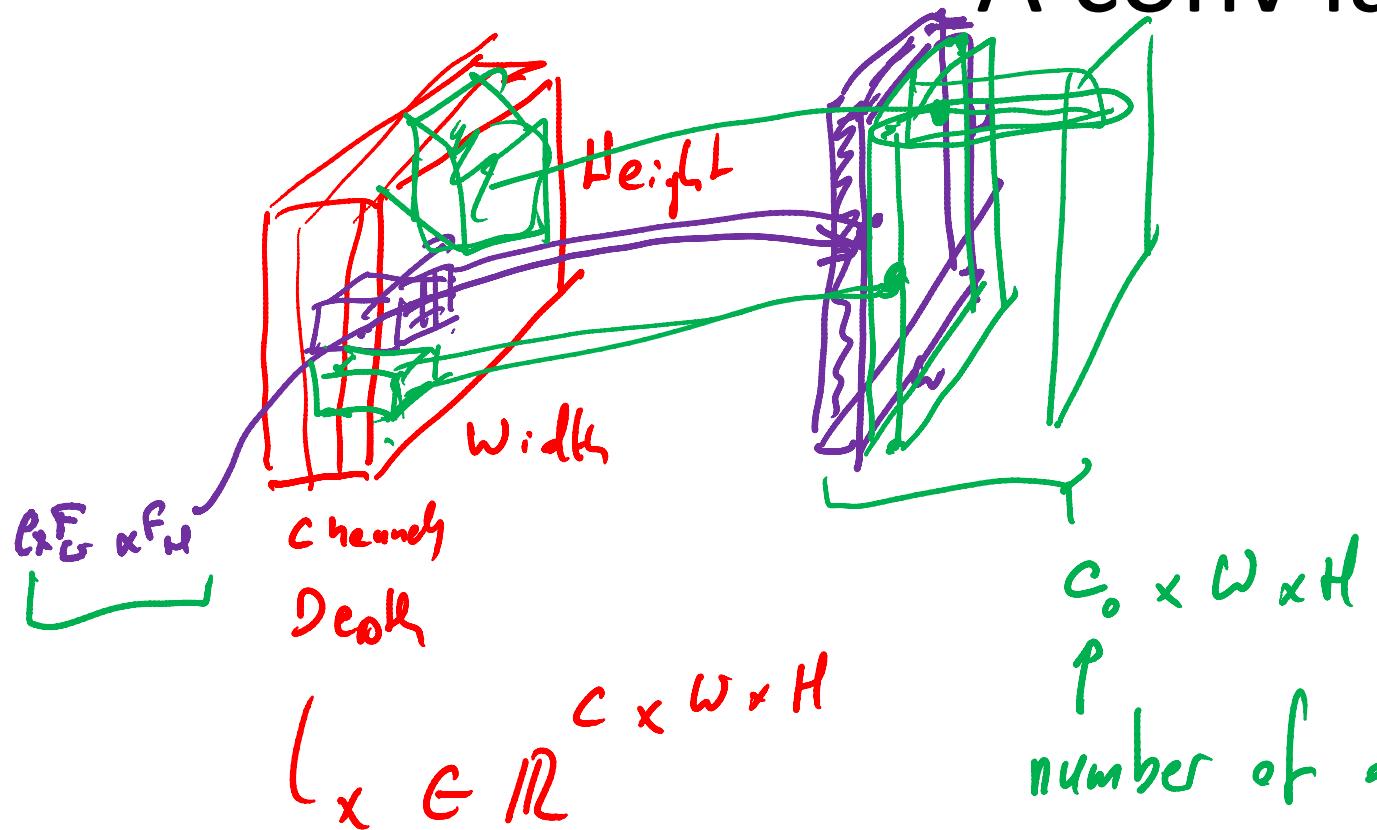
$$256 \times 256 = 256^2 = 2^{12}$$

$$(2^8)^4 = 2^{32} = 1 \cdot 10^9$$

LL Weight

③ Weigh Sharing + Local Connectivity  
q weight ???

# A conv layer



How many weights?

$C_o$  filters / sign process  
cores

each having  $C \times F_H \times F_W$

Total:

$$\underbrace{C_o \times C_m \times F_H \times F_W}_{\text{Total}}$$

How many multiplications?

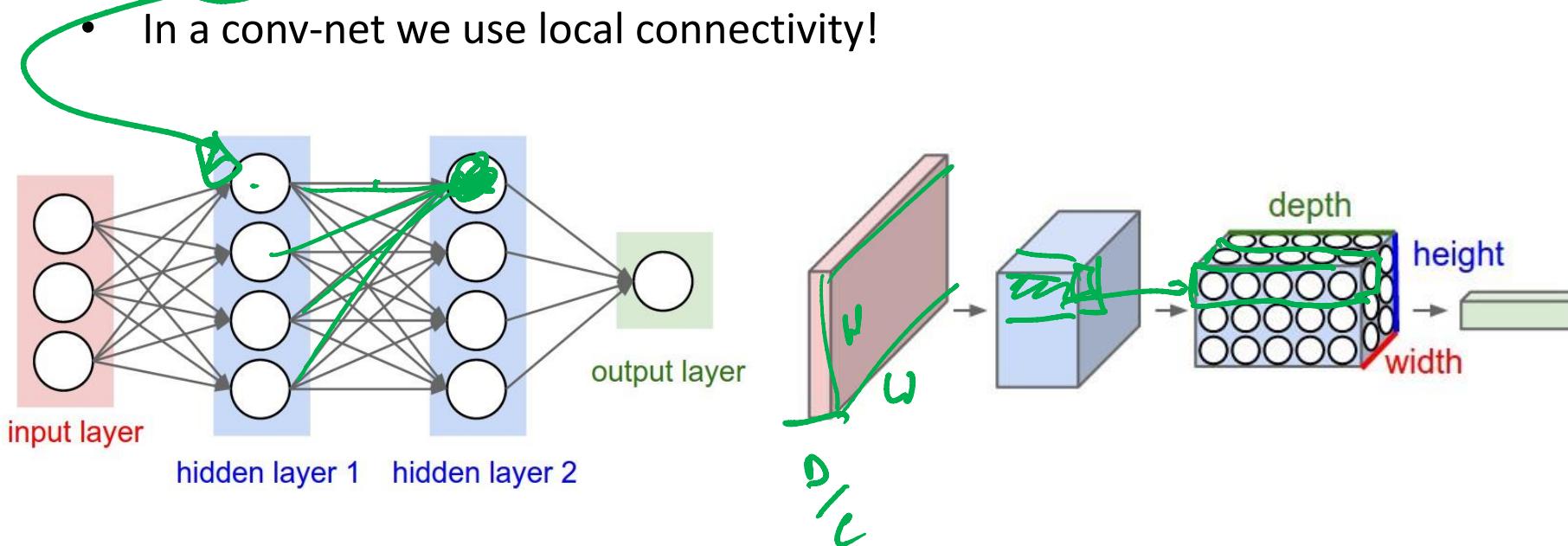
$$(W \times H) \times C_o \times C_m \times F_H \times F_W$$

# Sharing neurons - convolutions

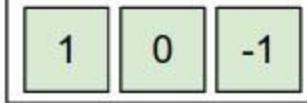
Note: material from <http://cs231n.github.io/convolutional-networks/>

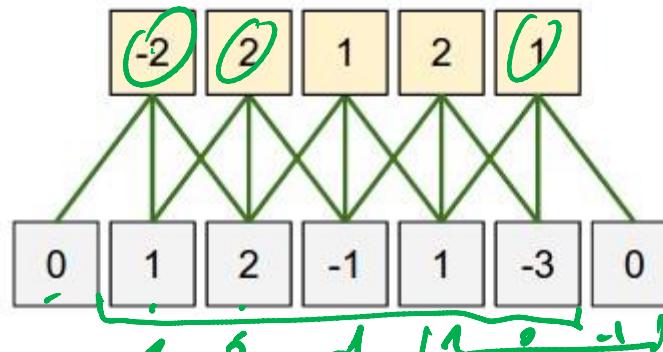
In a conv net we use a different connection pattern between layers:

- In a dense feed-forward layer used an all-to-all scheme
- In a conv-net we use local connectivity!

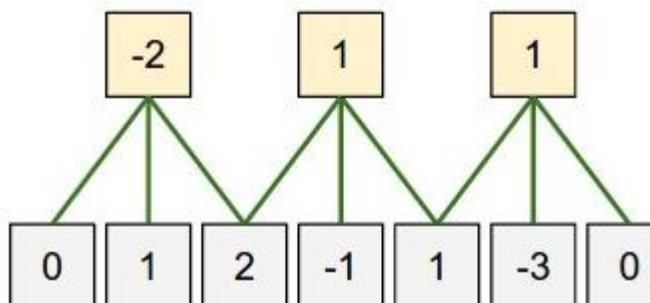


# 1D Convolution layer

- Small filter (neuron): 
- Swipe the filter over the sequence



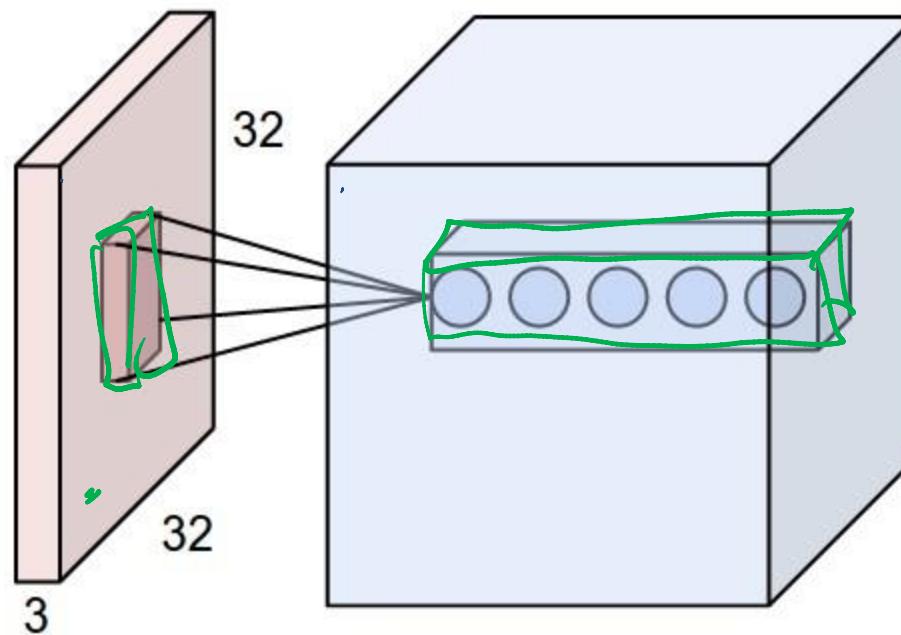
- Pool or stride (select only a few outputs)



Caveat:  
in math / signal  
processing we  
also flip the  
filter

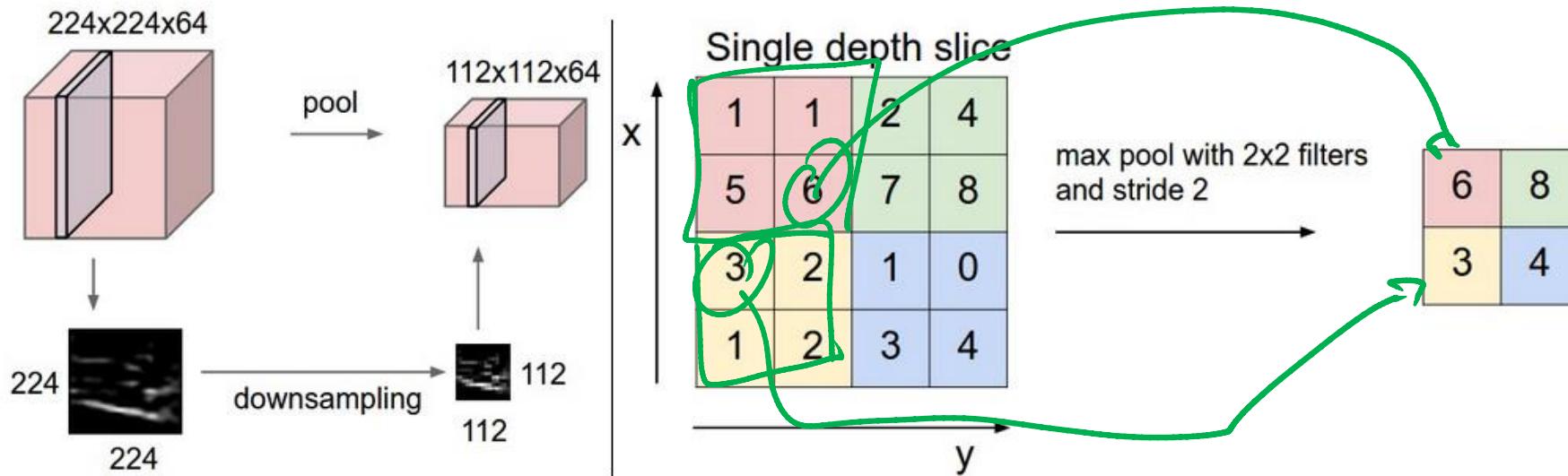
# 2D conv layer

- <http://cs231n.github.io/convolutional-networks/>



# Pooling

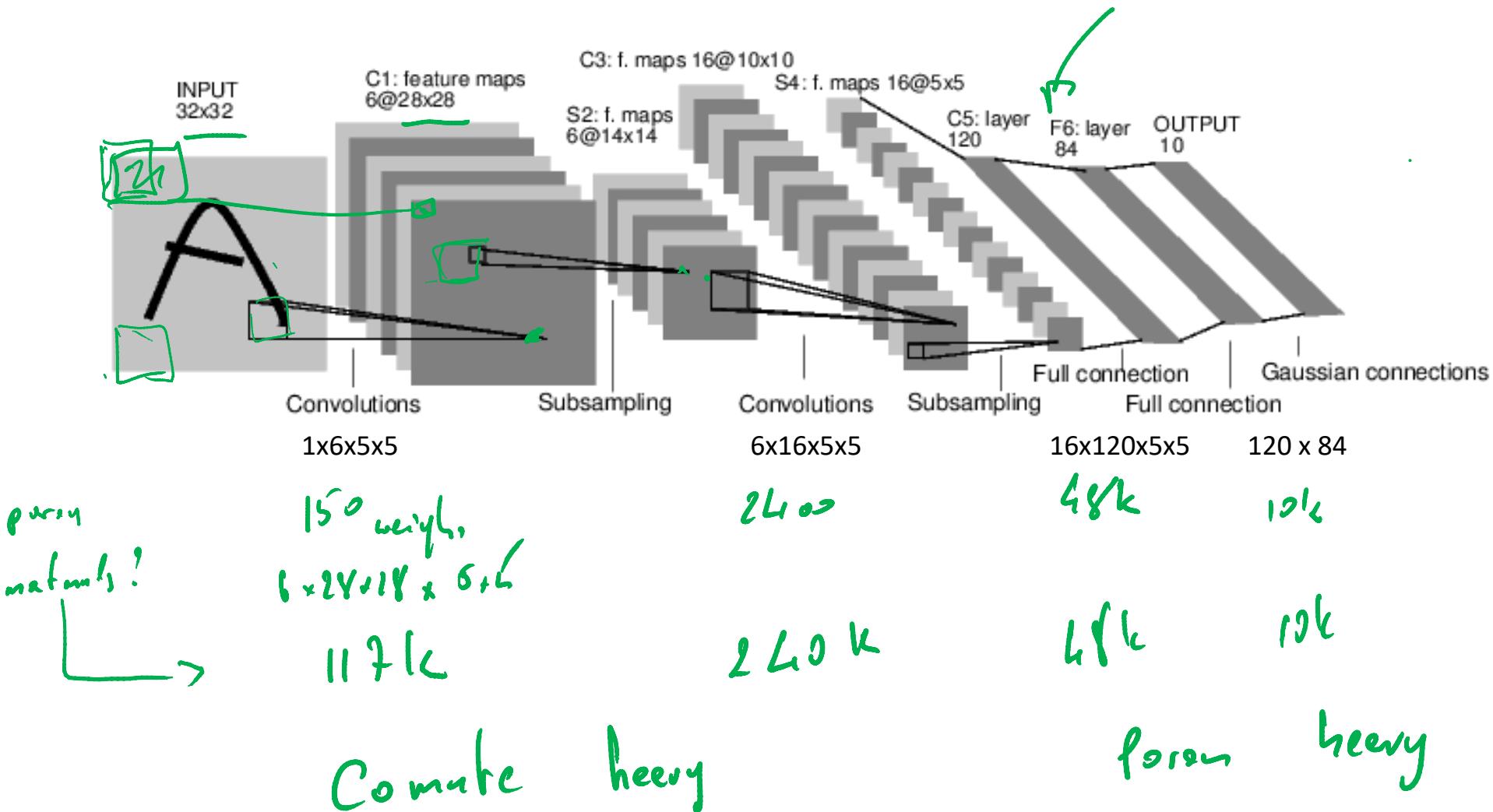
→ DD PALS !  
→ Reduces dimensions w & h



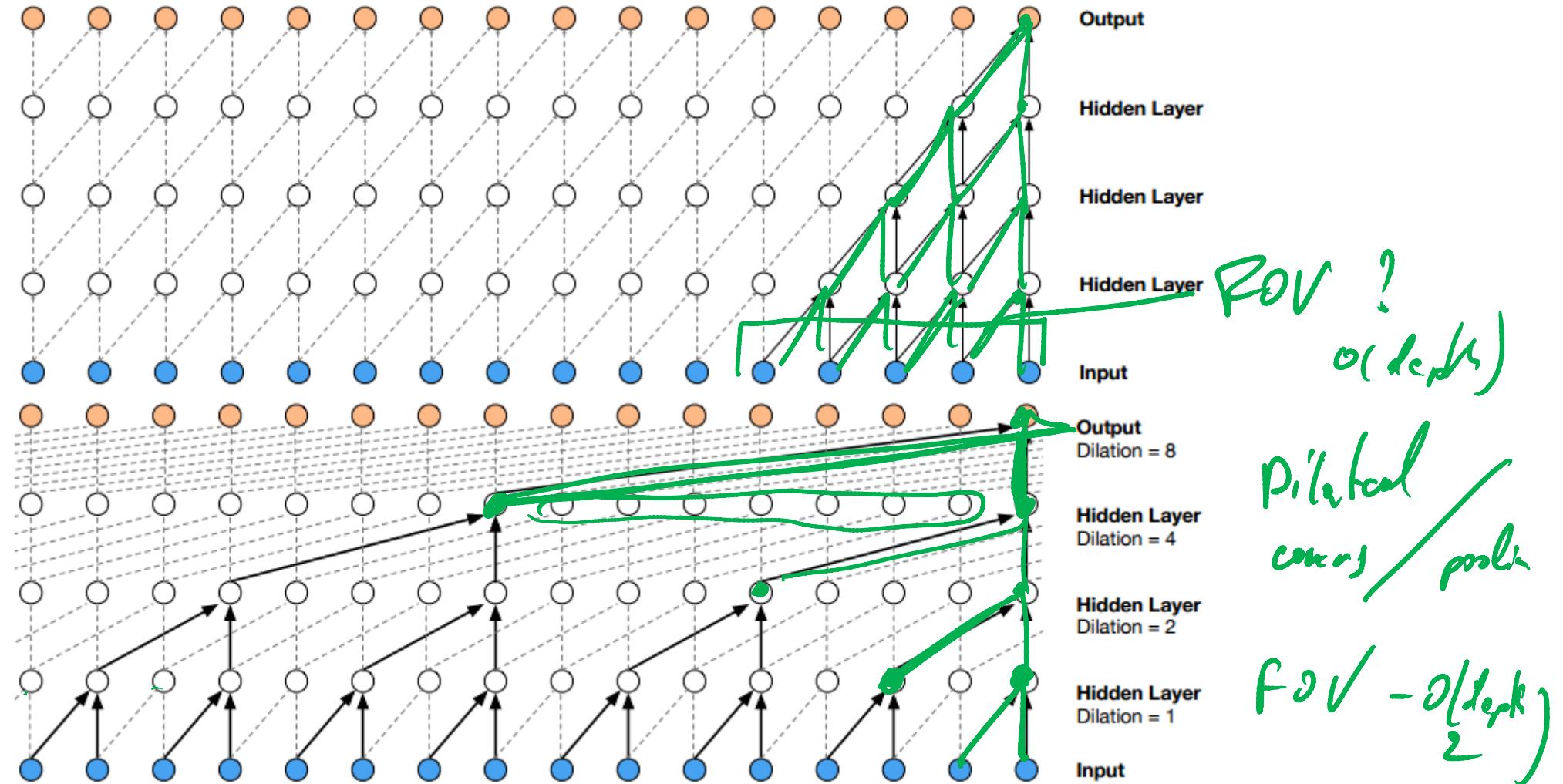
Pooling layer downsamples the volume spatially, independently in each depth slice of the input volume. **Left:** In this example, the input volume of size  $[224 \times 224 \times 64]$  is pooled with filter size 2, stride 2 into output volume of size  $[112 \times 112 \times 64]$ . Notice that the volume depth is preserved. **Right:** The most common downsampling operation is max, giving rise to **max pooling**, here shown with a stride of 2. That is, each max is taken over 4 numbers (little  $2 \times 2$  square).

# Full network - LeNet

Y. LeCun et al. „Gradient-Based Learning Applied to Document Recognition”  
<http://yann.lecun.com/exdb/publis/pdf/lecun-98.pdf>



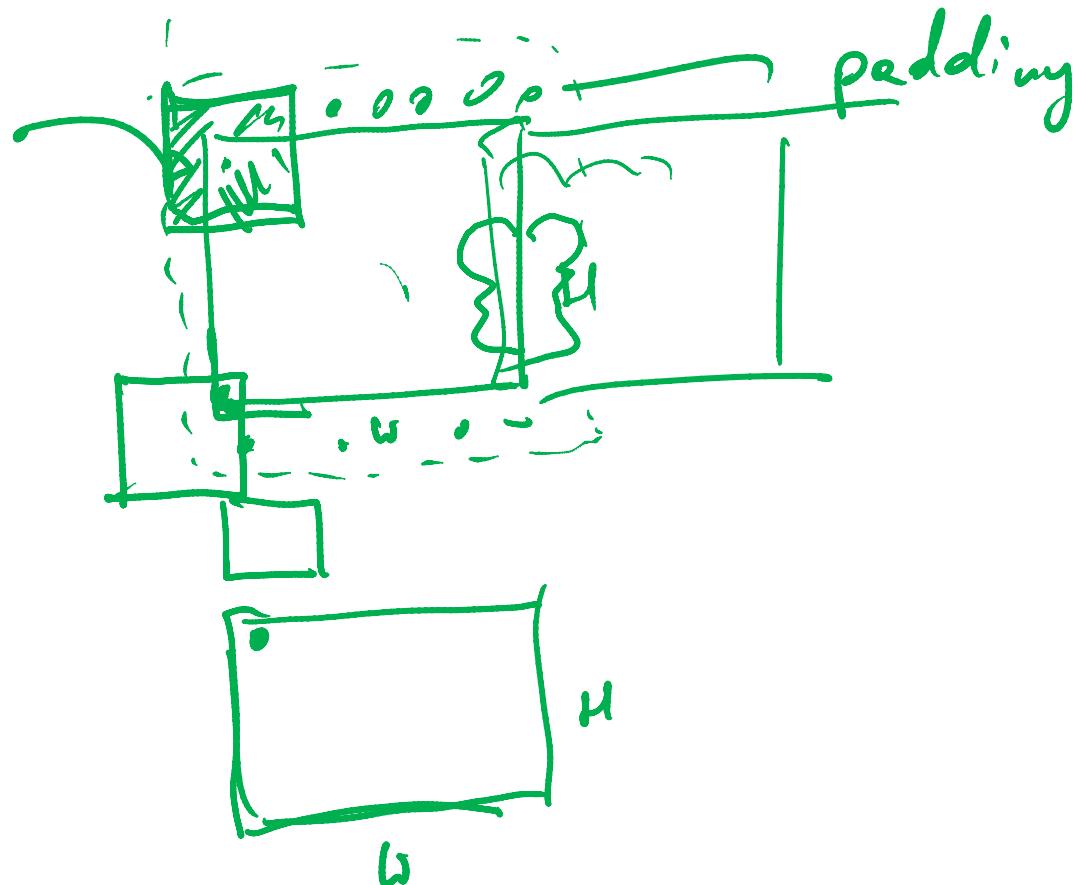
# Dilated convolutions: & troug wide receptive field at high resolutions



Wavenet: <https://arxiv.org/pdf/1609.03499.pdf>

<https://deepmind.com/blog/wavenet-generative-model-raw-audio/>

# Technicality: Edge padding



- VALID conv .

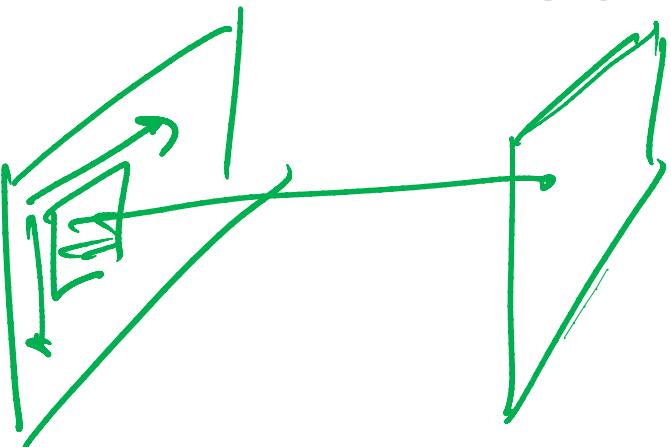
if input is  $W \times H$   $F_w \times F_h$

$$W - F_w + 1 \times H - F_h + 1$$

- FULL

$$(W + F_w - 1) \times (H + F_h - 1)$$

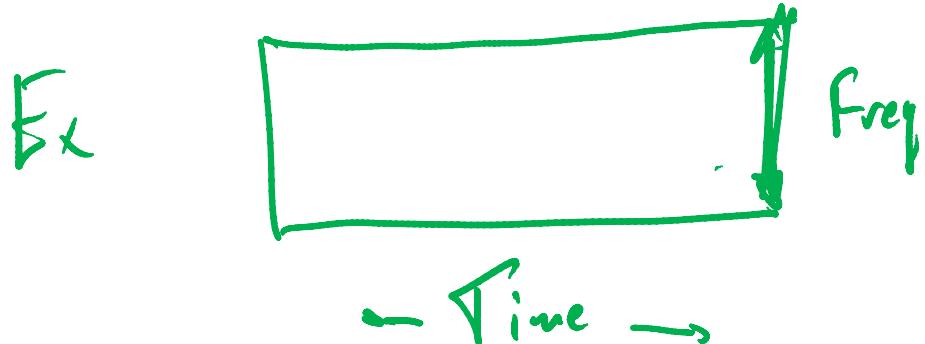
## Technicality: bias



$C \times W \times H$

$C \times 1 \times 1$

bias : a number per filter

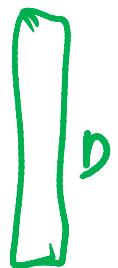
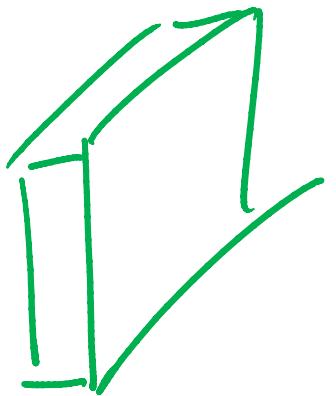


Weights : the same all over

bias :  $1 \times \# \text{Freqs}$

Time      Freq

# Convs can mimick FC layers



$$\text{Input: } C \times W \times H \rightarrow D \times 1 \times 1$$

Fully Connected or Dense

Filter size  $D \times C \times W \times H$

VALID over

$$W - P_W + 1 = W - w + t < 1$$

Filter:  $D \times C \times \underbrace{1 \times 1}_{\text{FC network}}$

Input:  $C \times W \times H$

# Conv Layers: All the Options

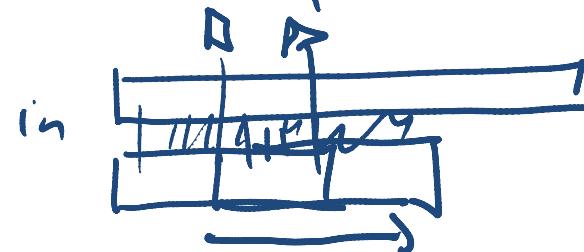
in 1 minibatch of multi-channel images - 4D tensor

`tf.nn.conv2d(input, filters, strides, padding, data_format='NHWC', dilations=None, name=None)`

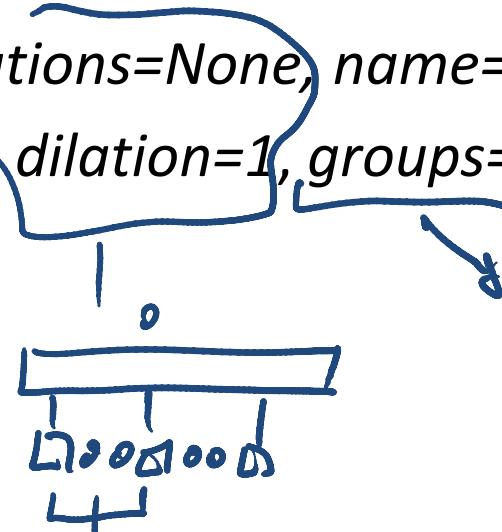
`torch.nn.functional.conv2d(input, weight, bias=None, stride=1, padding=0, dilation=1, groups=1)`

4D tensor  $N \times C \times H \times W$

4D tensor  $Cout \times Cin \times f_H \times f_W$

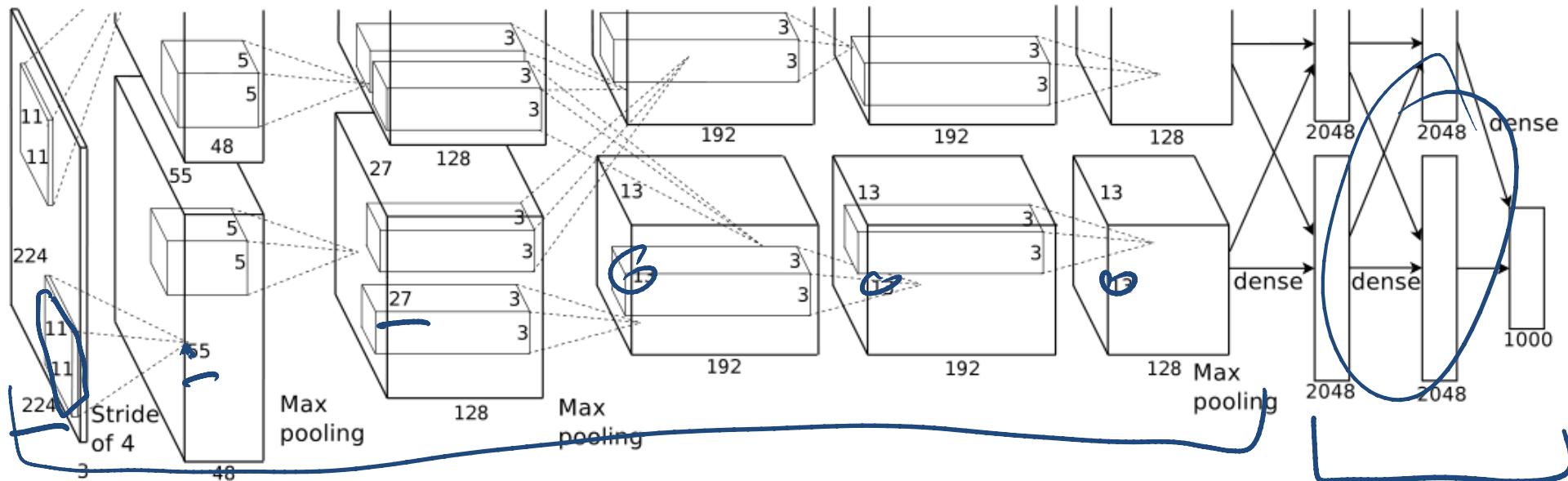


pooling, stride  
every n-th output



dilation-controlled  
by dilation

# AlexNet (2012)



60M parameters

The last dense layers takes:

$$4096 \times 4096 + 4096 \times 1000 = 20M \text{ params!}$$

# Separable Filters

- Sometimes a 2D conv == 2\* 1D conv:

$$\frac{1}{3} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} * \frac{1}{3} [1 \quad 1 \quad 1] = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$
$$\mathbf{G}_x = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} * A = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} * [+1 \quad 0 \quad -1] * A$$

- Inception and VGGnet use filter approximate separation, with more nonlinearity

$$\text{Conv}(x_1, 3 \times 3, f_i, l_i) = \text{Conv}(\text{Conv}(x_1, 3 \times 1), 1 \times 3)$$

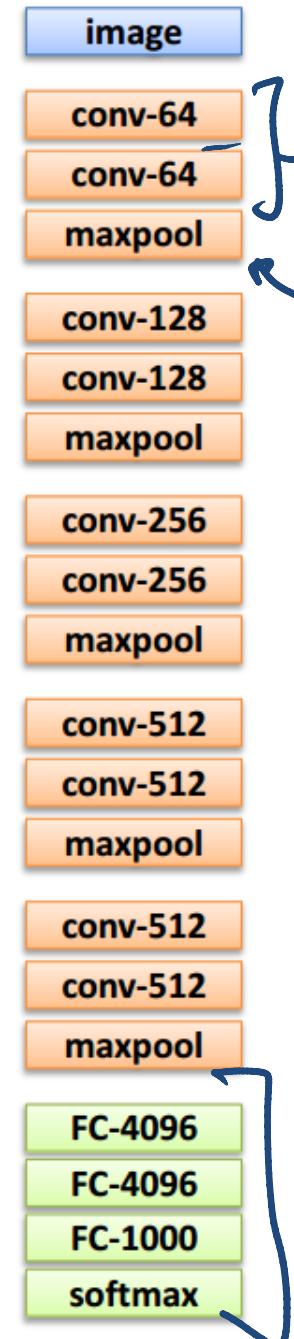
# VGGNet: Network Design

## Key design choices:

- 3x3 conv. kernels – very small
- conv. stride 1 – no loss of information

## Other details:

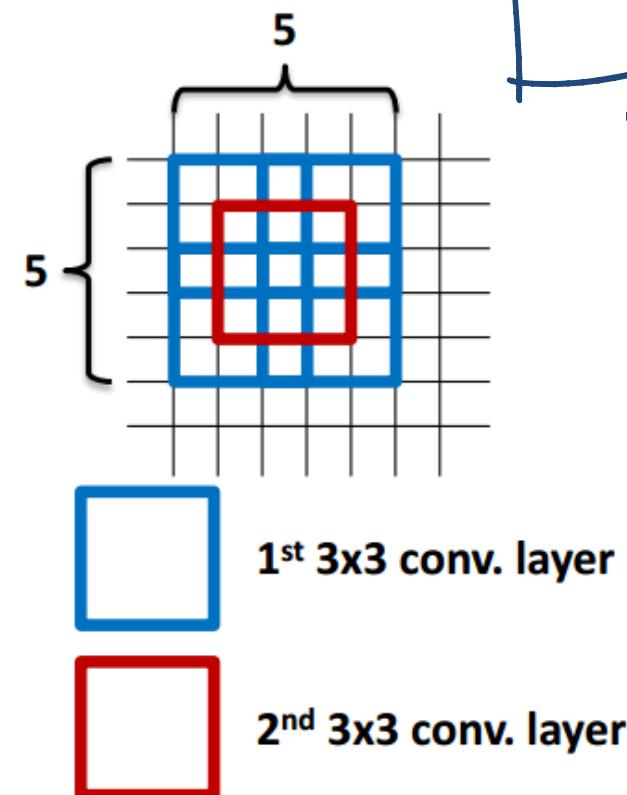
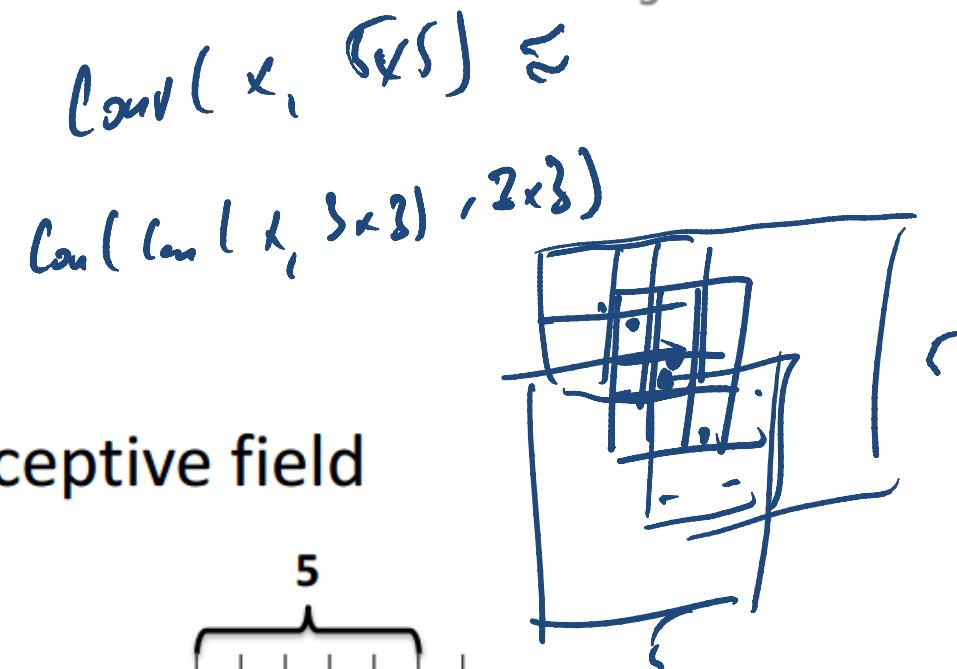
- Rectification (ReLU) non-linearity
- 5 max-pool layers (x2 reduction)
- no normalisation
- 3 fully-connected (FC) layers



# VGGNet: Discussion

Why 3x3 layers?

- Stacked conv. layers have a large receptive field
  - two 3x3 layers – 5x5 receptive field
  - three 3x3 layers – 7x7 receptive field
- More non-linearity
- Less parameters to learn
  - ~140M per net



# VGGNet (2014)

INPUT: [224x224x3] memory: 224\*224\*3=150K weights: 0

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M weights:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M weights:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: 112\*112\*64=800K weights: 0

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M weights:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M weights:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: 56\*56\*128=400K weights: 0

CONV3-256: [56x56x256] memory: 56\*56\*256=800K weights:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K weights:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K weights:  $(3*3*256)*256 = 589,824$

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

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

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

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

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

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

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

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

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

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

FC: [1x1x4096] memory: 4096 weights:  $4096*4096 = 16,777,216$

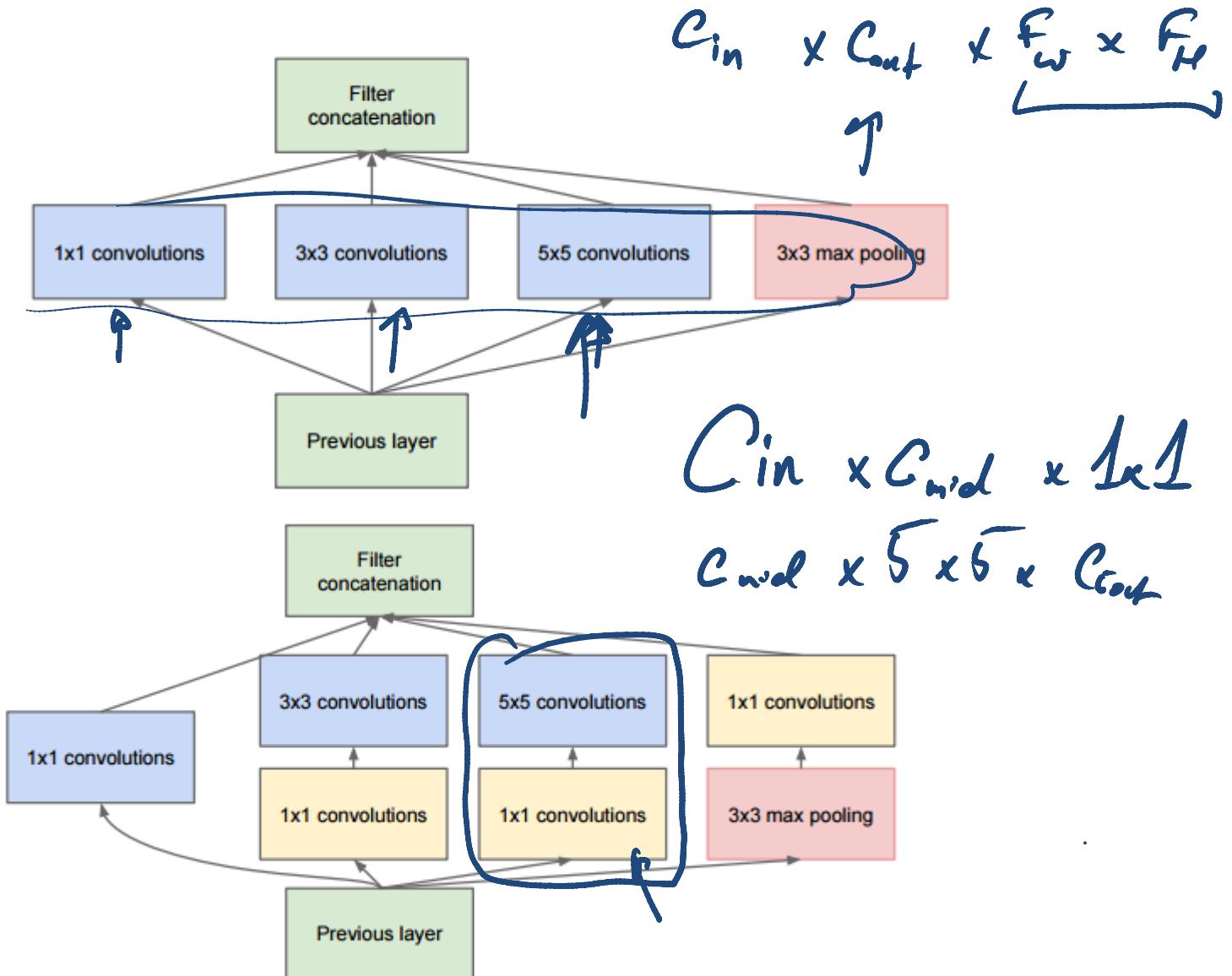
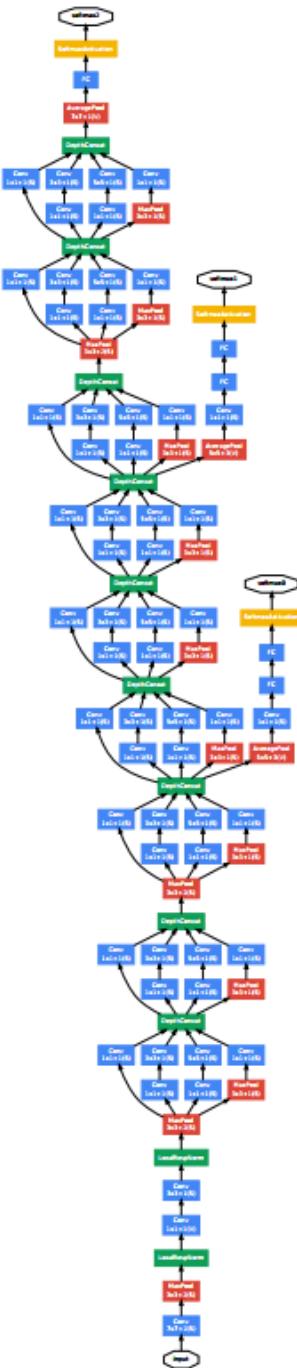
FC: [1x1x1000] memory: 1000 weights:  $4096*1000 = 4,096,000$

TOTAL memory: 24M \* 4 bytes ~ 93MB / image (only forward! ~\*2 for bwd) TOTAL params: 138M parameters

- up to 2M

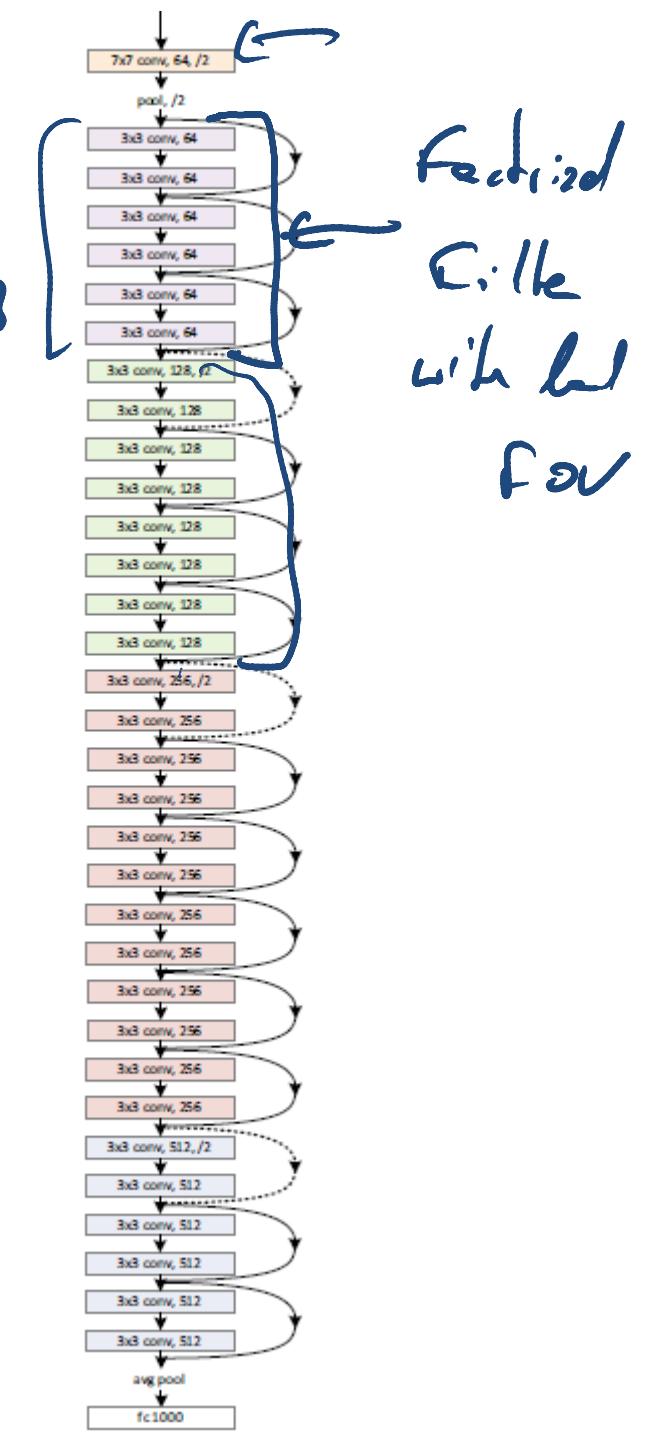
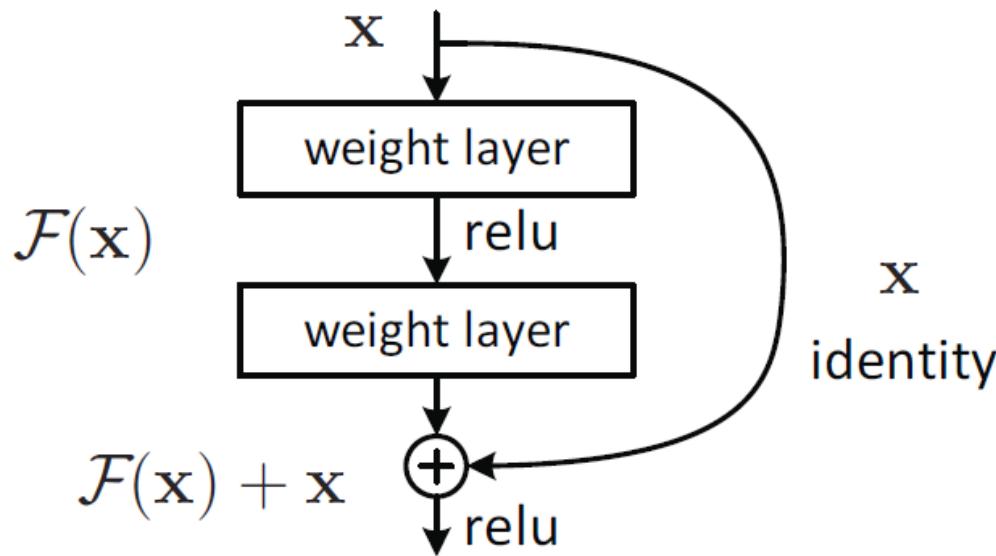
no weight

# GoogLeNet (Inception) 2014



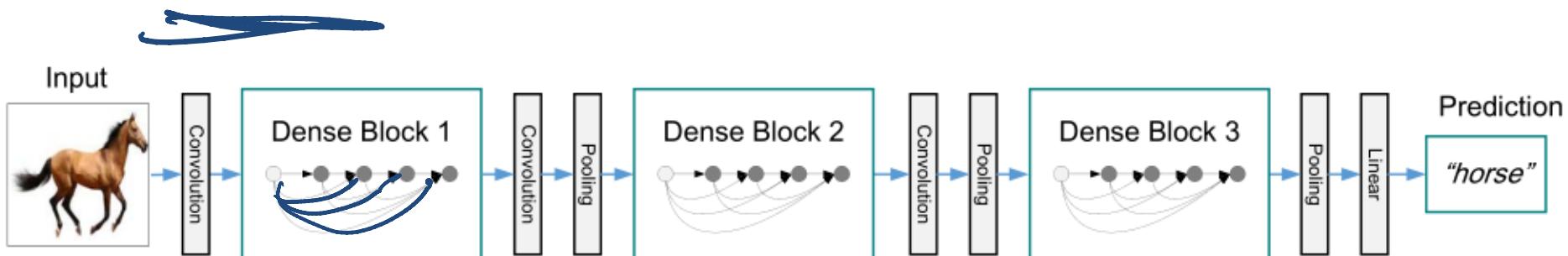
# ResNest (Residual Nets, 2015)

- Add bypasses to ease gradient flow
- Networks with more than 150 layers!
- 3x3 convs with number of feature maps doubling after every pooling

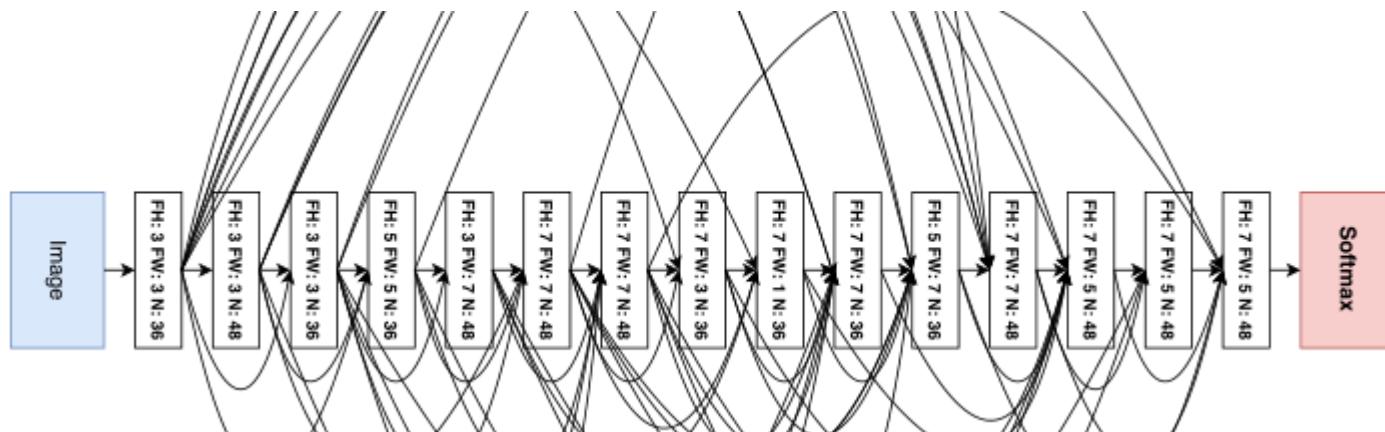


# Quest for models

- DenseNets (<https://arxiv.org/pdf/1608.06993>)



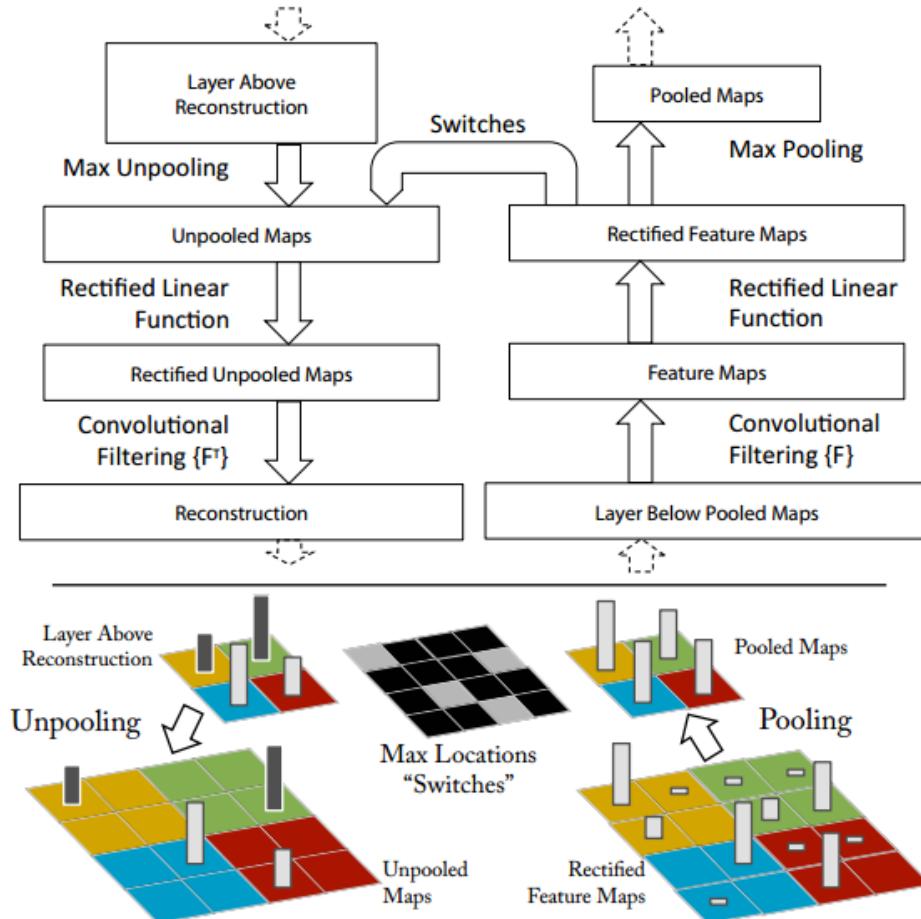
- Automatic architecture search  
(<https://arxiv.org/pdf/1611.01578>)



# Which ConvNet to use?

- Take current best on Imagenet
- Many networks are available for download  
(google tensorflow model and caffe model zoo)
- For your use – take a pretrained net and retrain only the last layers!
- On ImageNet (1M images) data augmentation and regularization is as important as the network!

# Visualizing CNN features



1

1

2

2

3

3

4

4

5

5

6

6

7

7

8

8

9

9

10

10

11

11

12

12

13

13

14

14

15

15

16

16

17

17

18

18

19

19

20

20

21

21

22

22

23

23

24

24

25

25

26

26

27

27

28

28

29

29

30

30

31

31

32

32

33

33

34

34

35

35

36

36

37

37

38

38

39

39

40

40

41

41

42

42

43

43

44

44

45

45

46

46

47

47

48

48

49

49

50

50

51

51

52

52

53

53

54

54

55

55

56

56

57

57

58

58

59

59

60

60

61

61

62

62

63

63

64

64

65

65

66

66

67

67

68

68

69

69

70

70

71

71

72

72

73

73

74

74

75

75

76

76

77

77

78

78

79

79

80

80

81

81

82

82

83

83

84

84

85

85

86

86

87

87

88

88

89

89

90

90

91

91

92

92

93

93

94

94

95

95

96

96

97

97

98

98

99

99

100

100

101

101

102

102

103

103

104

104

105

105

106

106

107

107

108

108

109

109

110

110

111

111

112

112

113

113

114

114

115

115

116

116

117

117

118

118

119

119

120

120

121

121

122

122

123

123

124

124

125

125

126

126

127

127

128

128

129

129

130

130

131

131

132

132

133

133

134

134

135

135

136

136

137

137

138

138

139

139

140

140

141

141

142

142

143

143

144

144

145

145

146

146

147

147

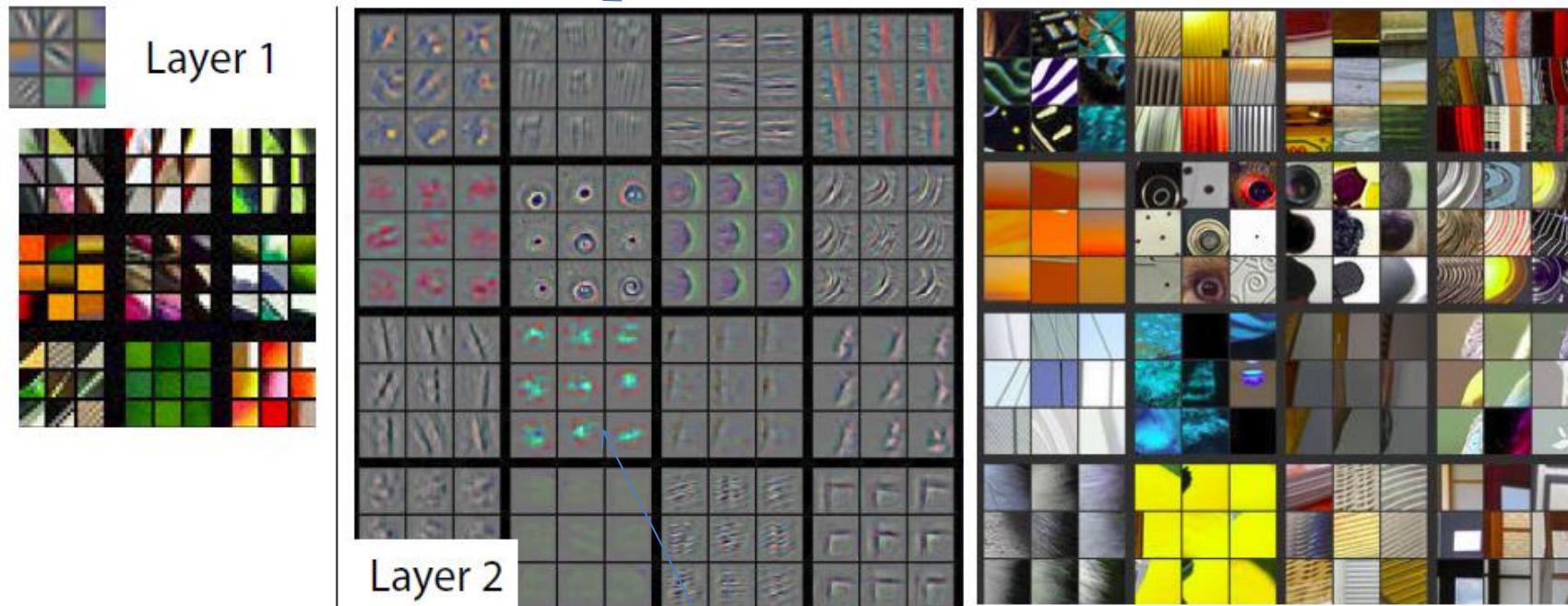
148

148

149

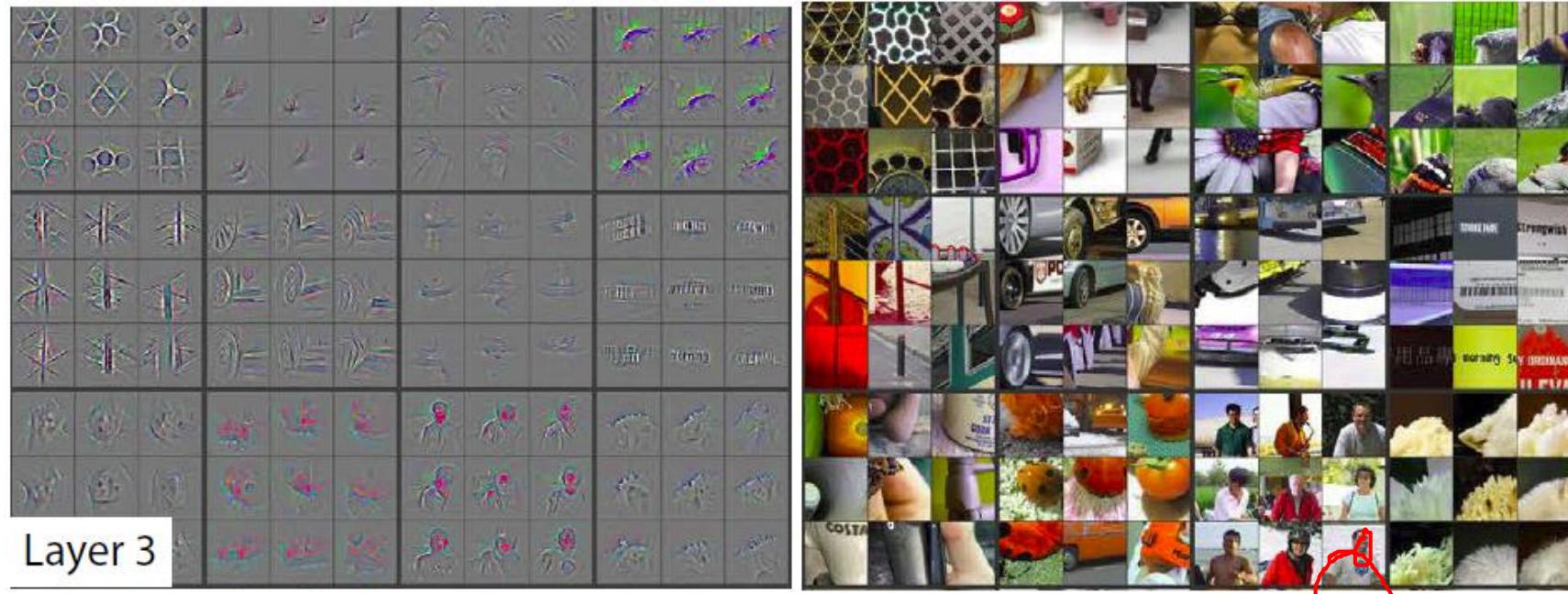
# Low-level features

What the neuron (feature-detector looks for)

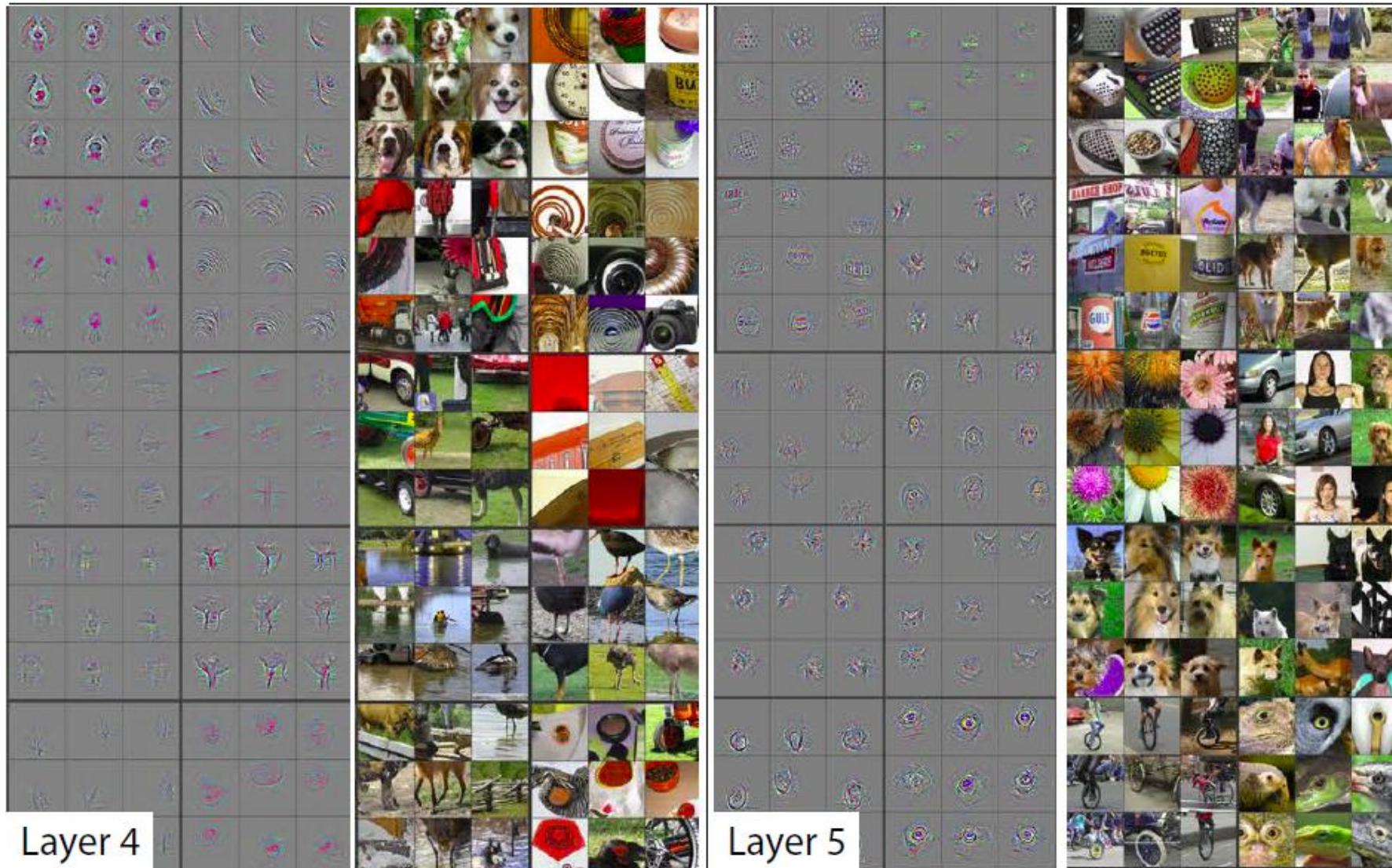


What images are selected by the neuron

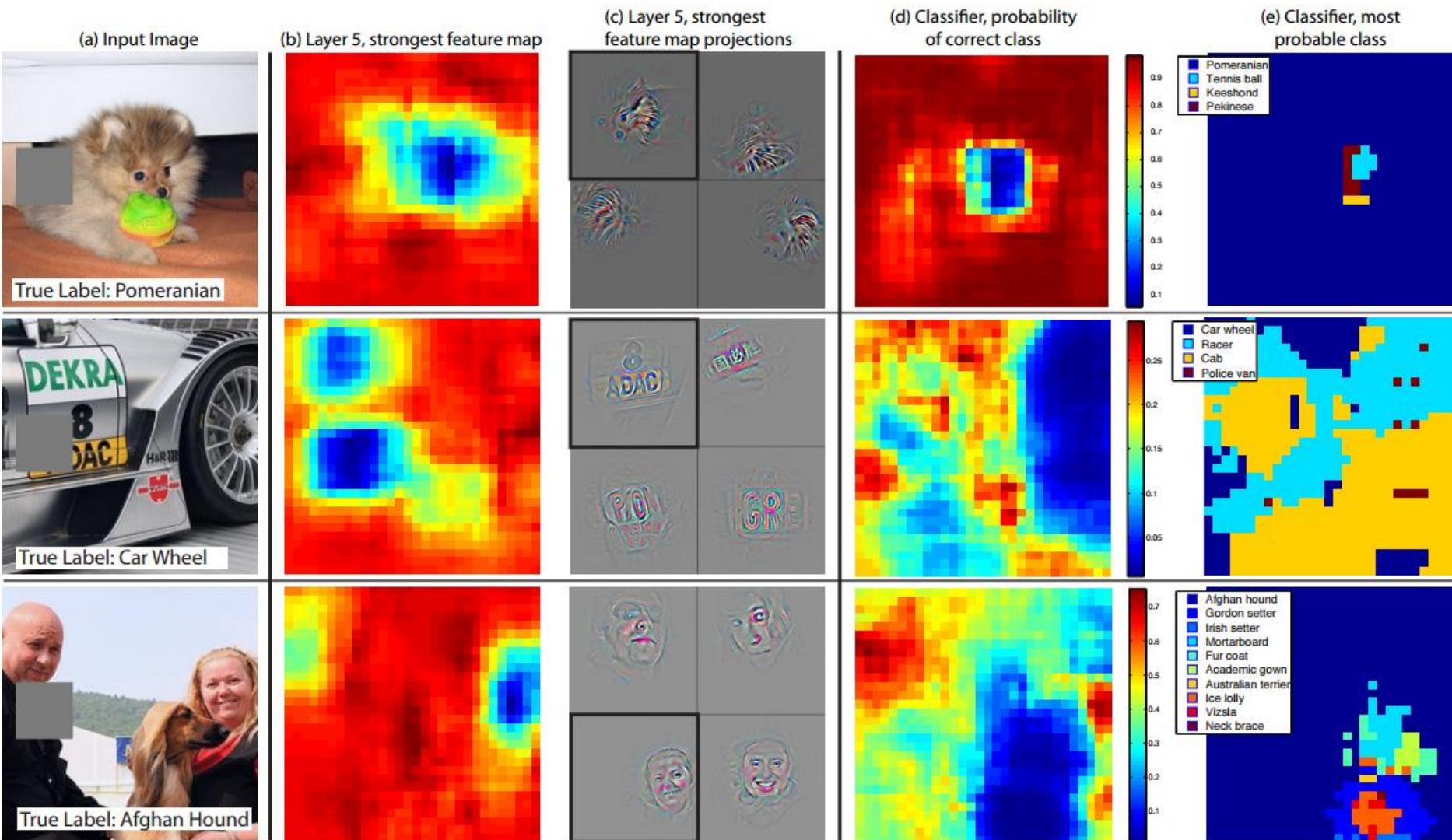
# Mid-level features



# High-level features



# Where Convnets look?



# Style Transfer



Find image (backpropagation toward pixels) minimizing:

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_c + \mathcal{L}_S = \\ &= (F(C) - F(I))^2 + (G(C) - G(I))^2\end{aligned}$$

Where:

$F(x)$  features of a CNN layer on image  $x$

$G(x)$  matrix of correlations between a layer's features over pixels of image  $x$

# **SPECIALIZED CONVNETS**



# MobileNets

## 1. Std convolution

Filter size:  $D \times C \times W \times H$

Matmuls per pixel:  $D \times C \times W \times H$

## 2. Depthwise convolution (each filter sees only one input channel)

Filter size:  $C \times W \times H$  (nb: no  $D$ )

Matmuls per pixel:  $C \times W \times H$

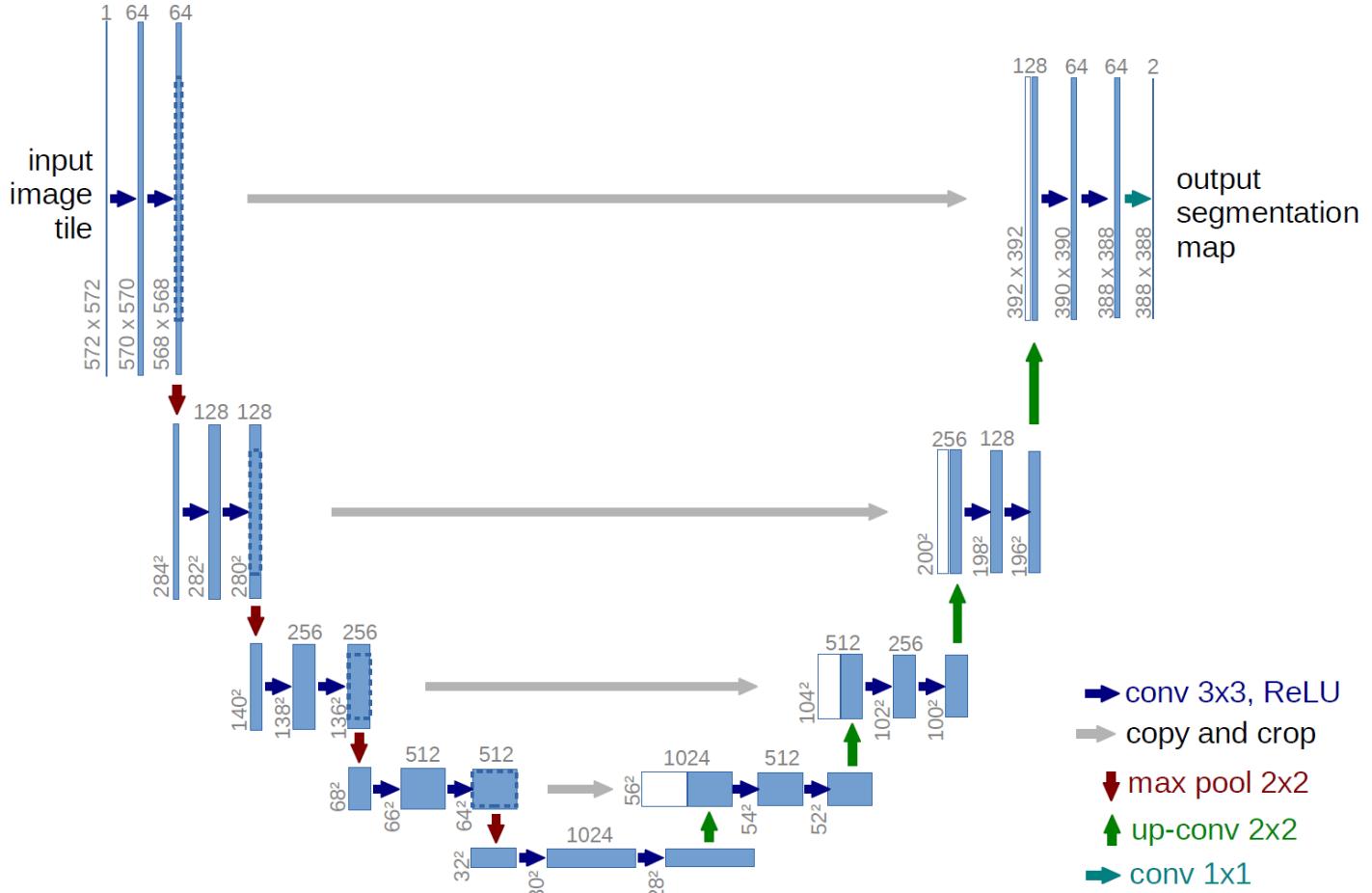
## 3. MobileNet: Depthwise + 1x1 convs

Why? Depthwise: fast, doesn't mix channels

1x1: mixes channels, still small

# Unet: Pixel labeling intro

# U-Net: Pixel labeling

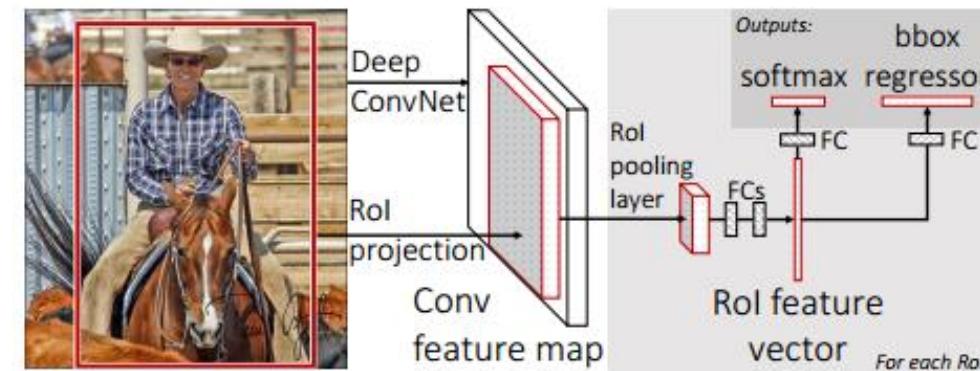
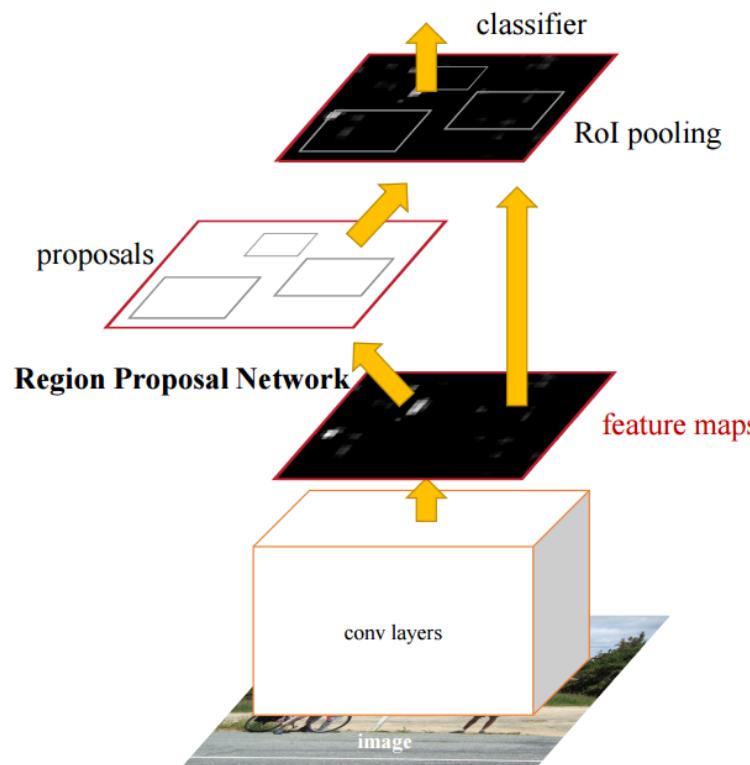
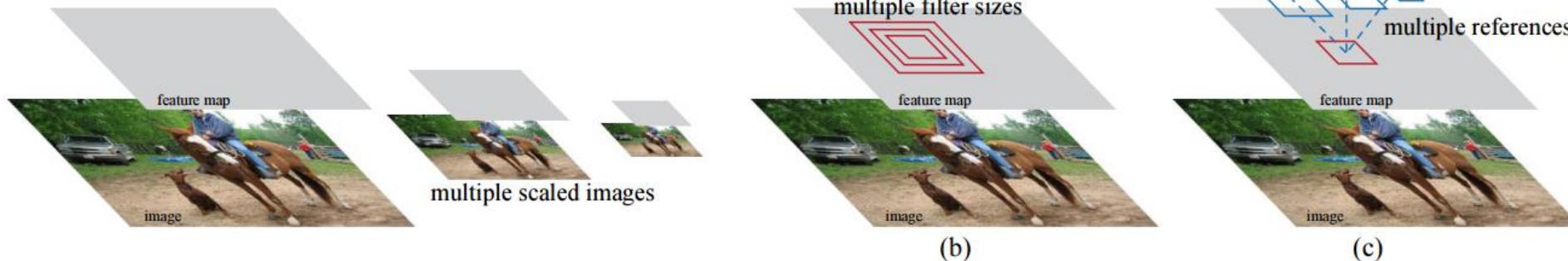


# From Classification to Detection

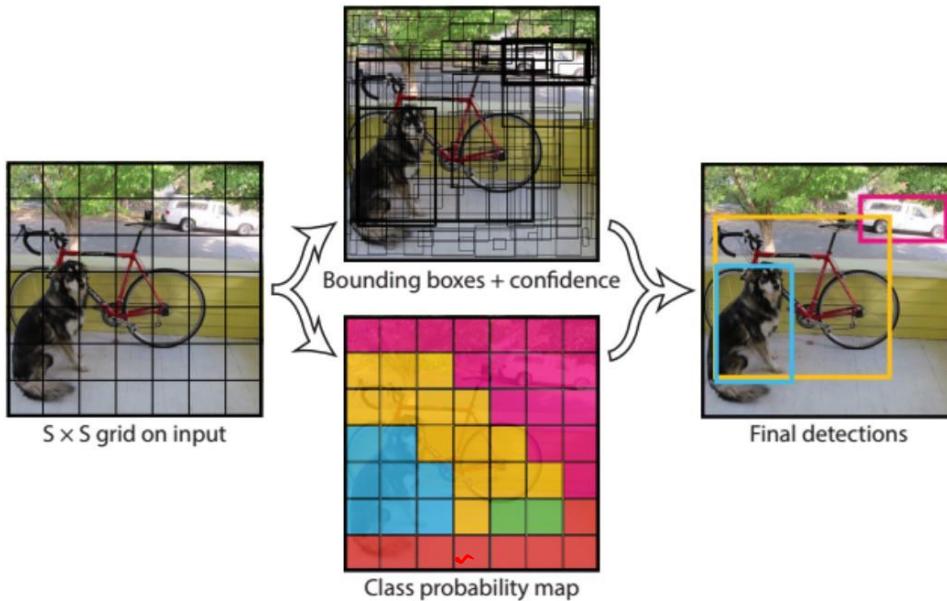
# Object detection and segmentation

Fast R-CNN (<https://arxiv.org/pdf/1504.08083.pdf>)

Faster R-CNN (<https://arxiv.org/pdf/1506.01497.pdf>)



# YOLO – fast object detection

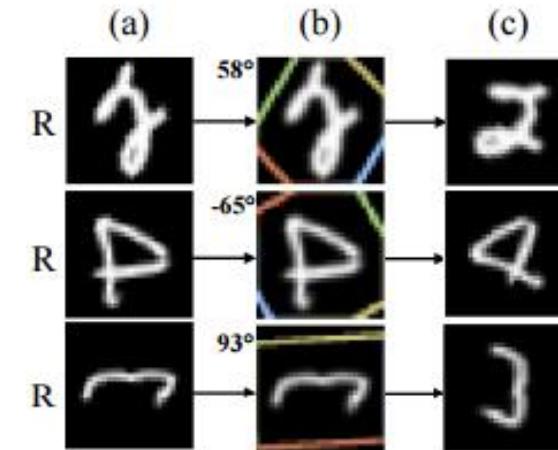
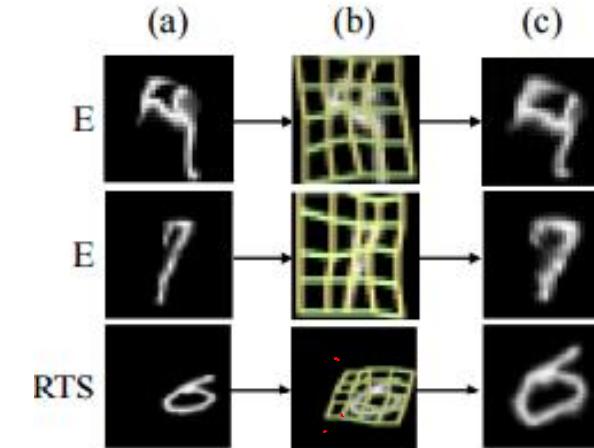
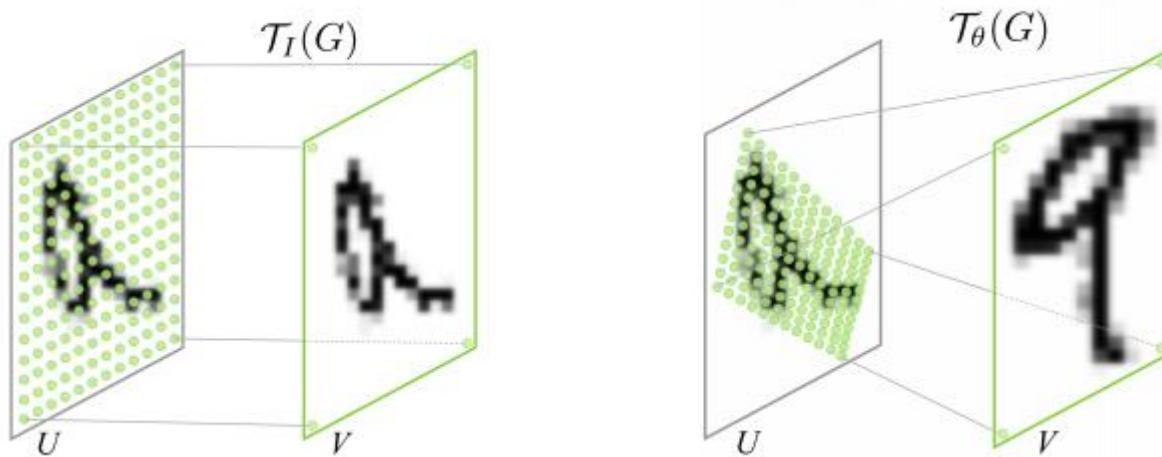
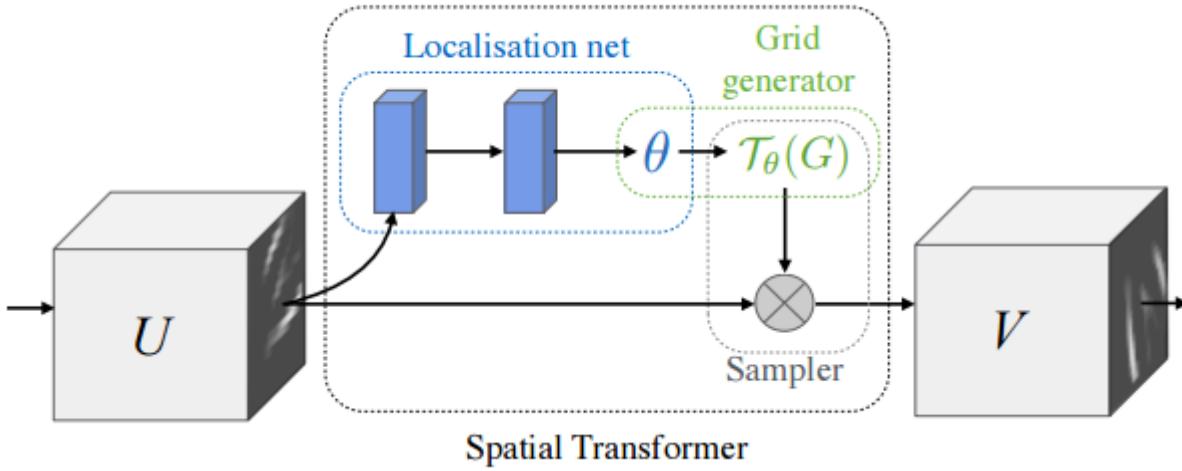


**Figure 2: The Model.** Our system models detection as a regression problem. It divides the image into an  $S \times S$  grid and for each grid cell predicts  $B$  bounding boxes, confidence for those boxes, and  $C$  class probabilities. These predictions are encoded as an  $S \times S \times (B * 5 + C)$  tensor.

- No region proposal and looping over ROIs!
- In a fully convolutional net predict for each pixel of a feature map:
  - Objectness – is something there
  - Bounding Box – how large it is

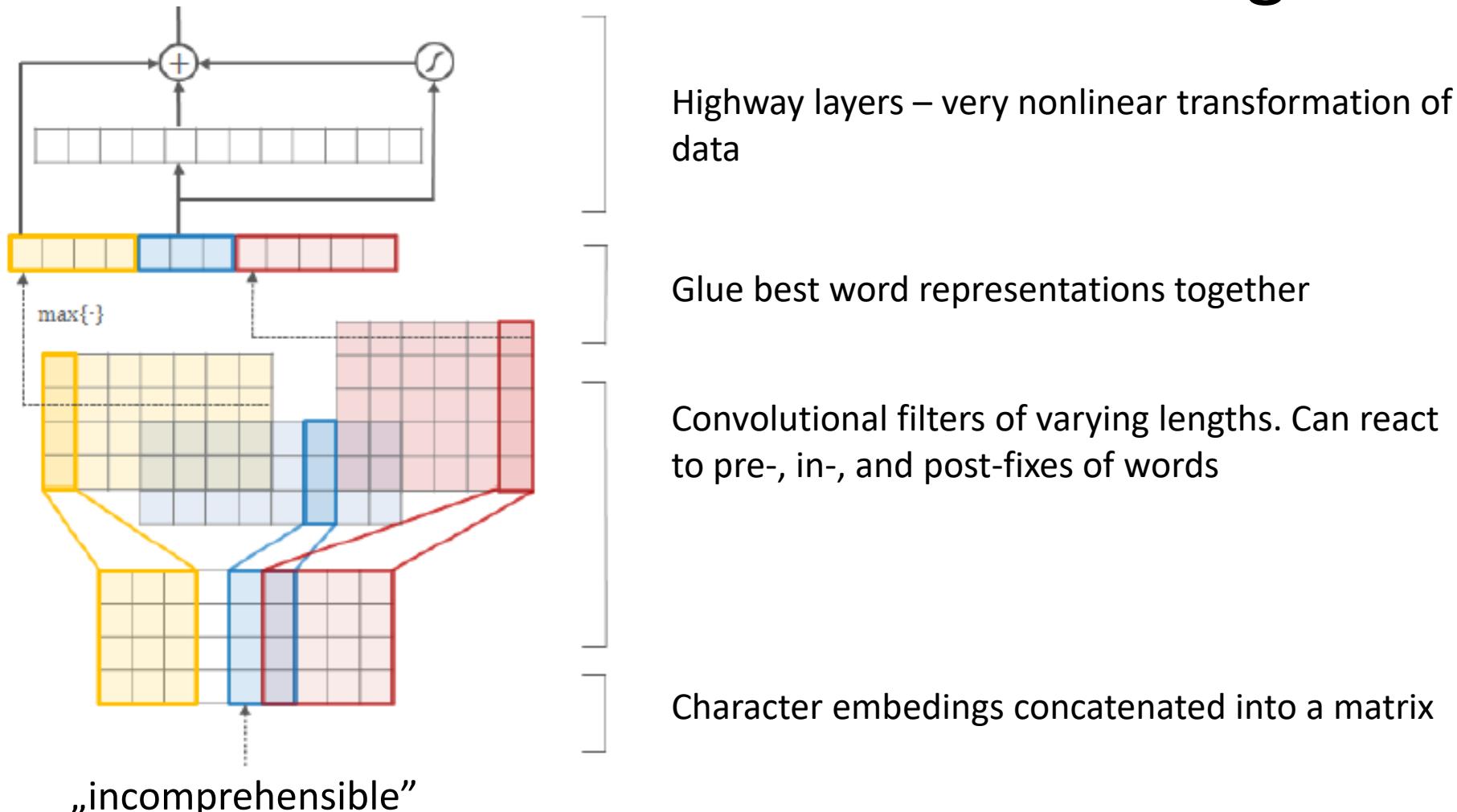
# Spatial Transformer Networks

<https://arxiv.org/pdf/1506.02025.pdf>



# CNN beyond images

## Or character-to-word embedding



Y. Kim, Y. Jernite, D. Sontag, and A. M. Rush, “Character-Aware Neural Language Models,”  
*arXiv:1508.06615 [cs, stat]*, Aug. 2015.

# Additional references

- Goodfellow Chp. 9
- <http://cs231n.github.io>
- [How transferable are features in deep neural networks?](#) studies the transfer learning performance in detail, including some unintuitive findings about layer co-adaptations.
- [CNN Features off-the-shelf: an Astounding Baseline for Recognition](#) trains SVMs on features from ImageNet-pretrained ConvNet and reports several state of the art results.