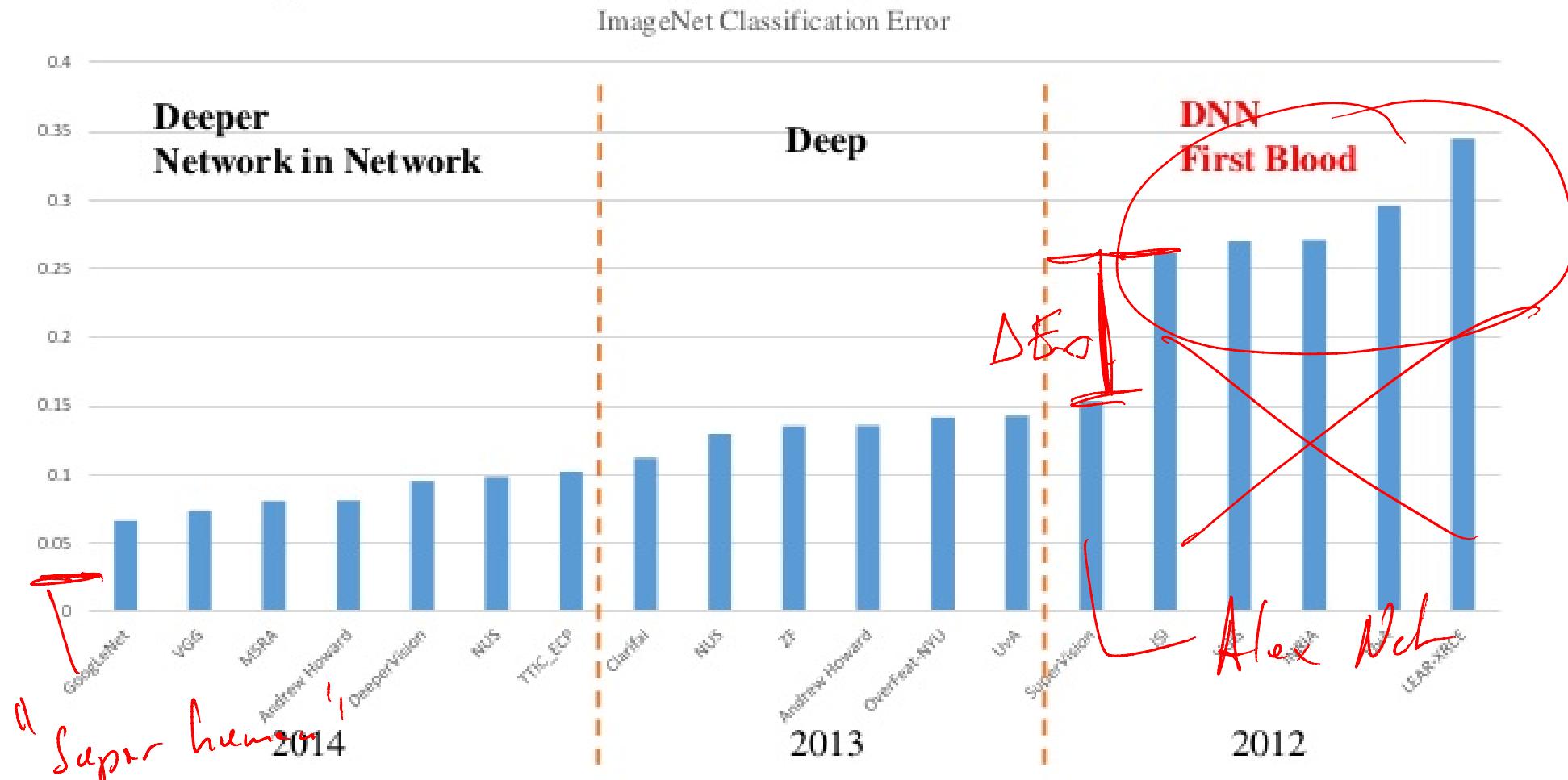


Neural Networks

Jan Chorowski

ImageNet Classification

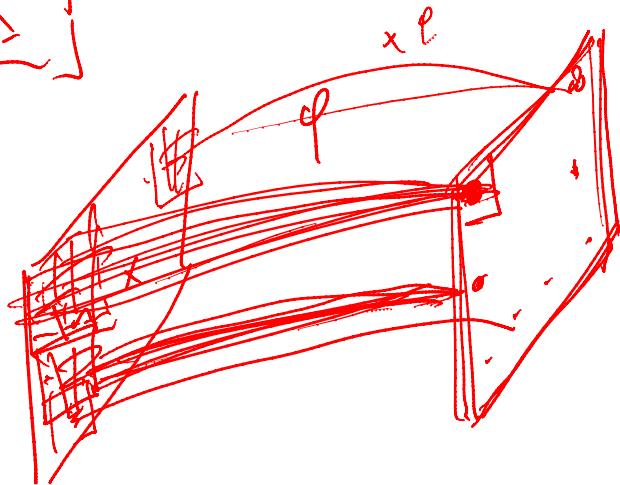
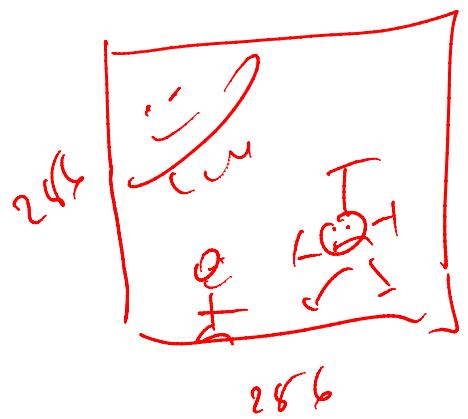
- 1000 categories and 1.2 million training images



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



Intuitions



① u_p to node i

$$w \in \mathbb{R}^{28 \times 28 \times 28 \times 28}$$

② LOCAL CONNECTIVITY

$$28^2 \cdot 9 =$$

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

③ Weight Sharing + Local Connectivity

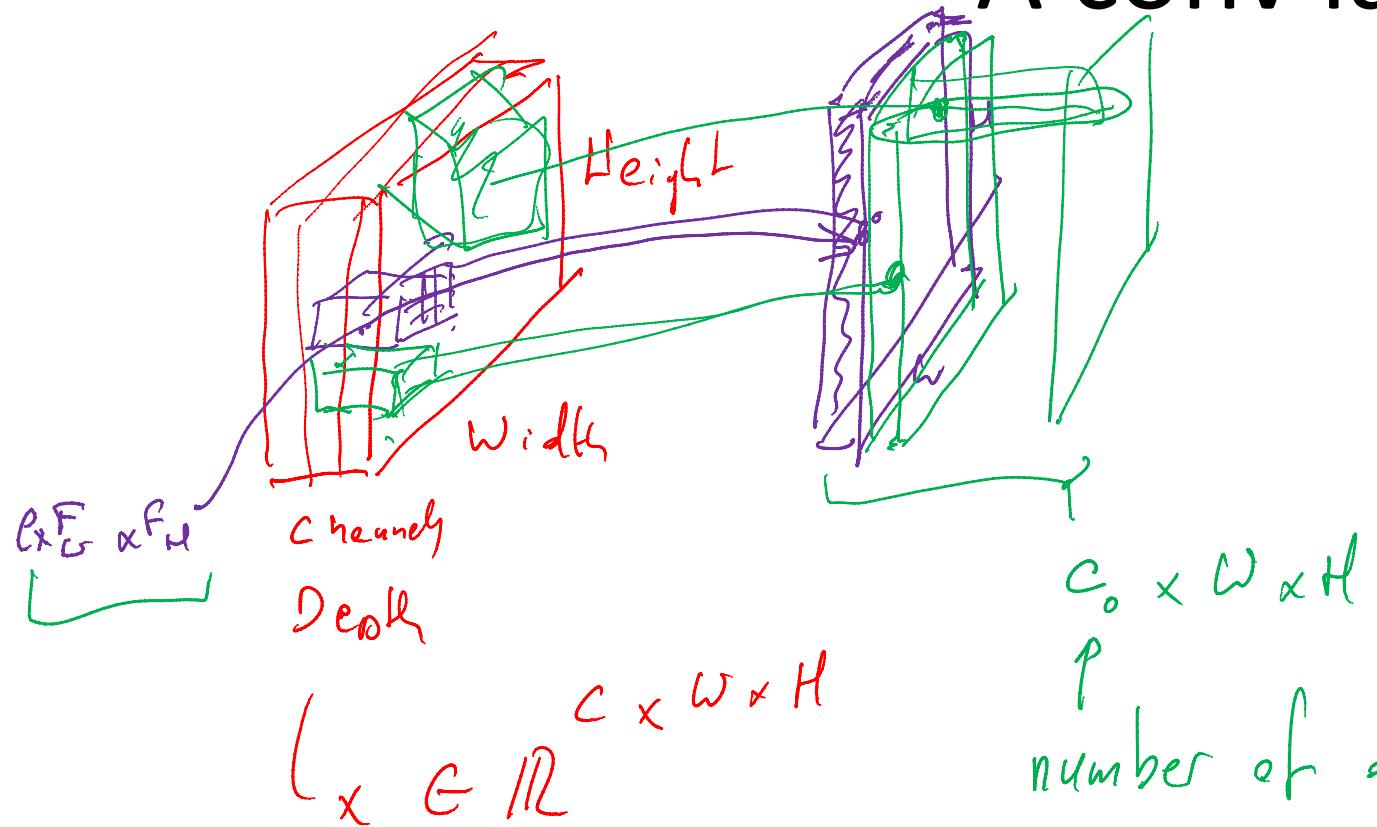
Weight ?? ??

$$w \in \mathbb{R}^{28 \times 28 \times 28 \times 28}$$

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

4G weight

A conv layer



How many weights?

C_o filters / signal process
cores

each having $C \times F_W \times F_H$

Total:

$$C_o \times C_i \times F_W \times F_H$$

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

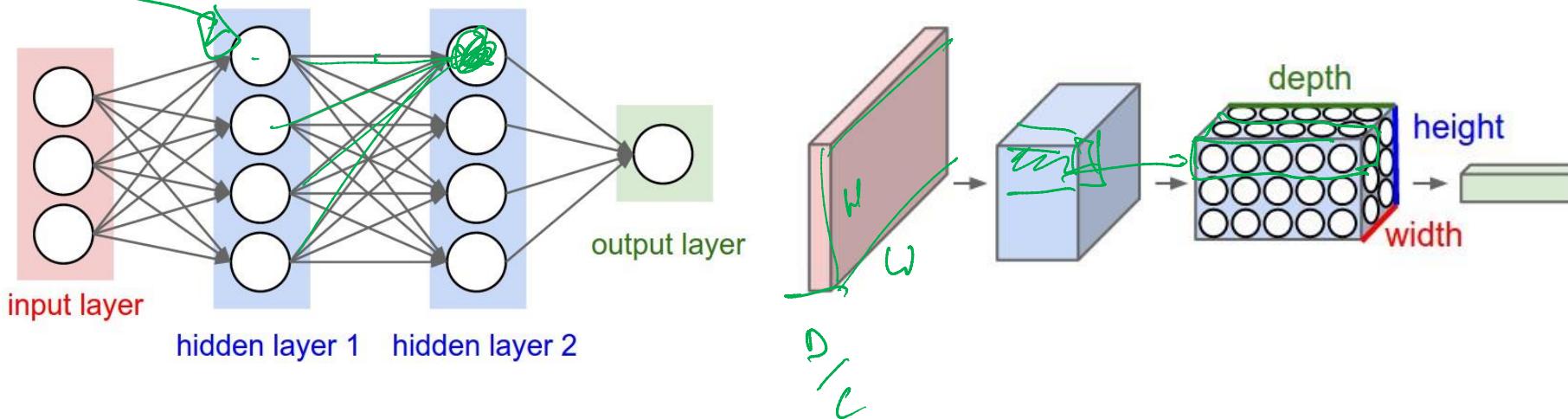
How many multiplications?

Sharing neurons - convolutions

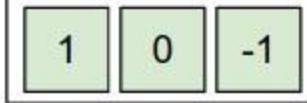
Note: material from <http://cs231n.github.io/convolutional-networks/>

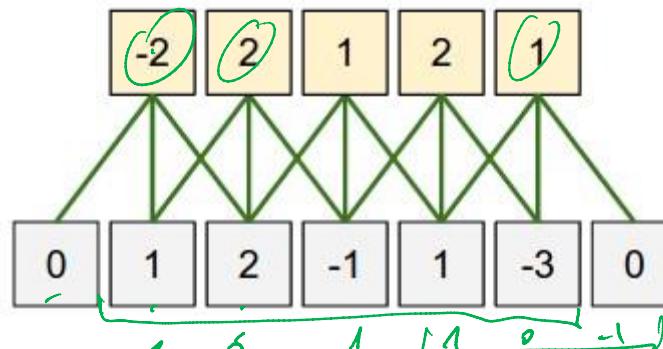
In a conv net we use a different connection pattern between layers:

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

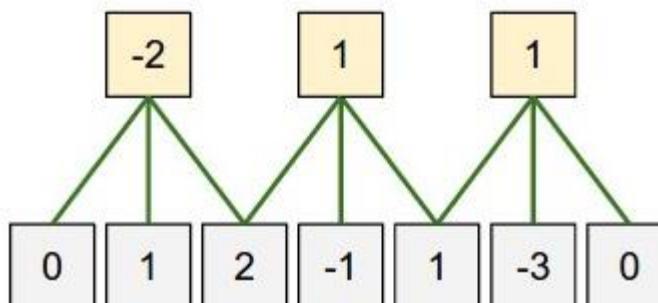


1D Convolution layer

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



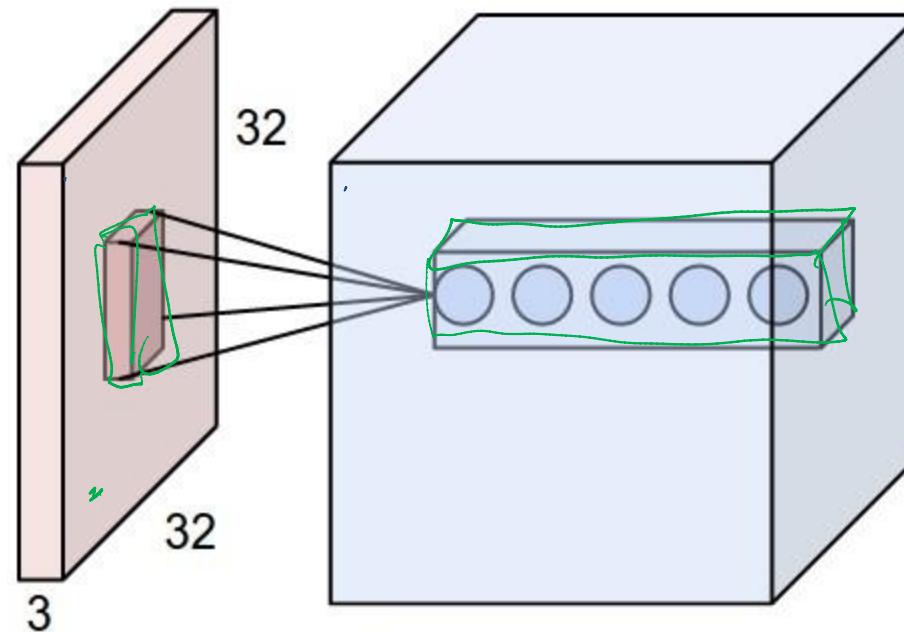
- Pool or stride (select only a few outputs)



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

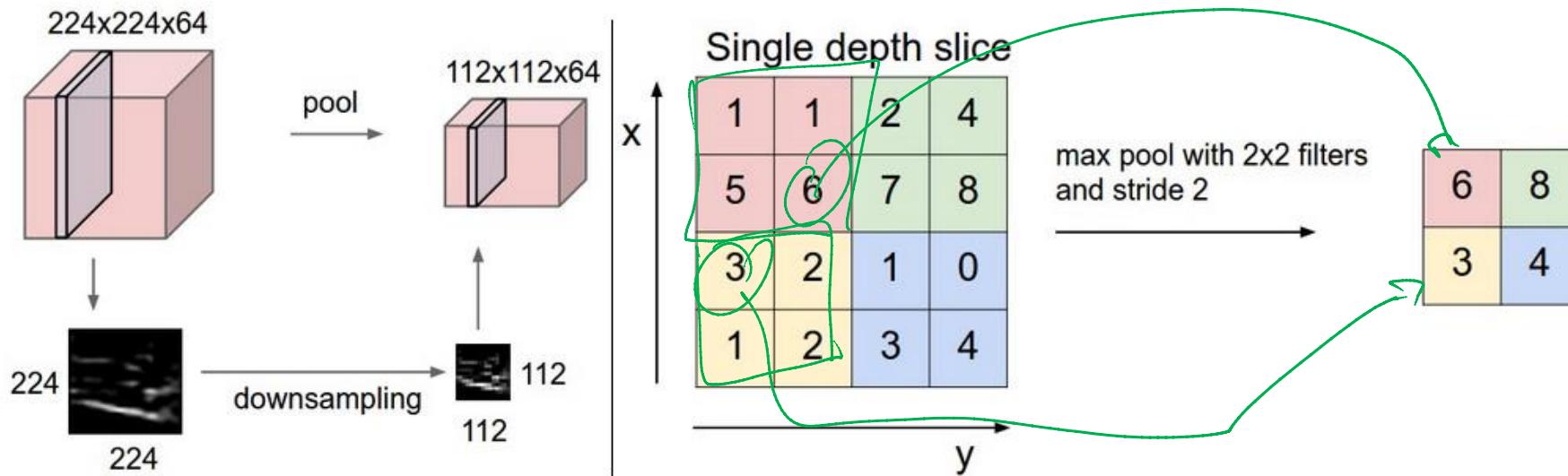
2D conv layer

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



Pooling

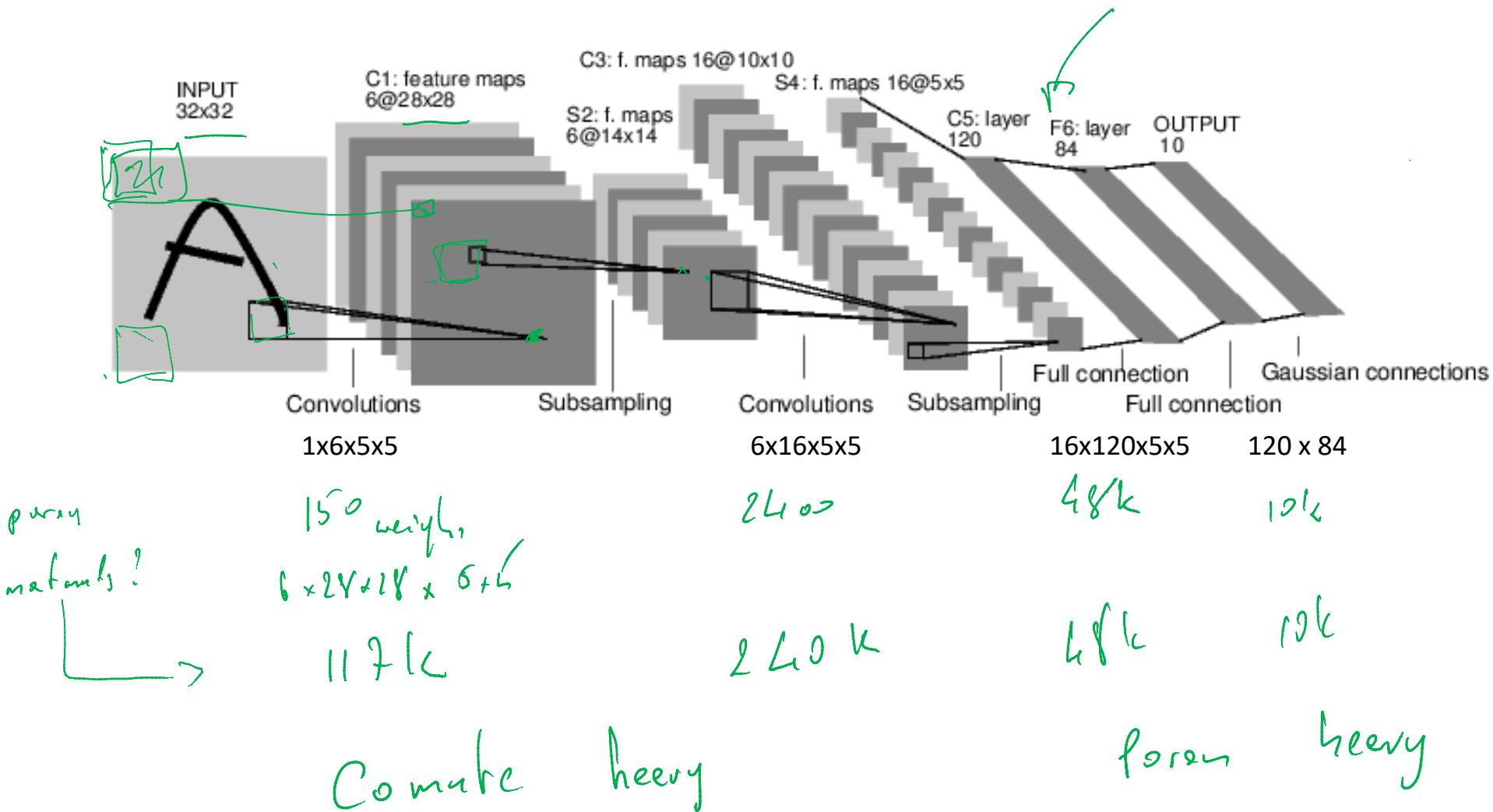
→ NO PARAMS ↗
→ Reduces dimensions w & h



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 2×2 square).

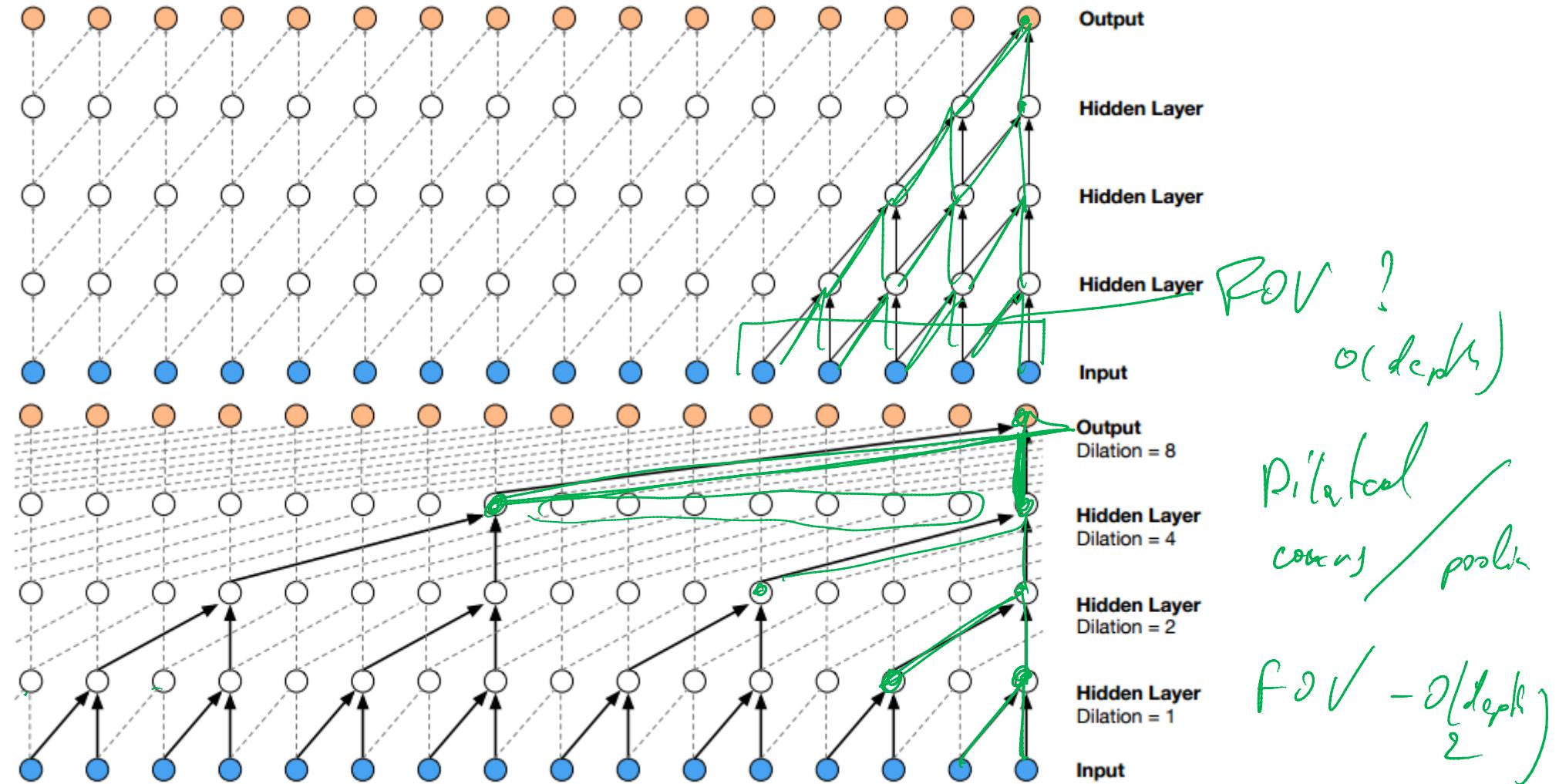
Full network - LeNet

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



Dilated convolutions: *& trouy*

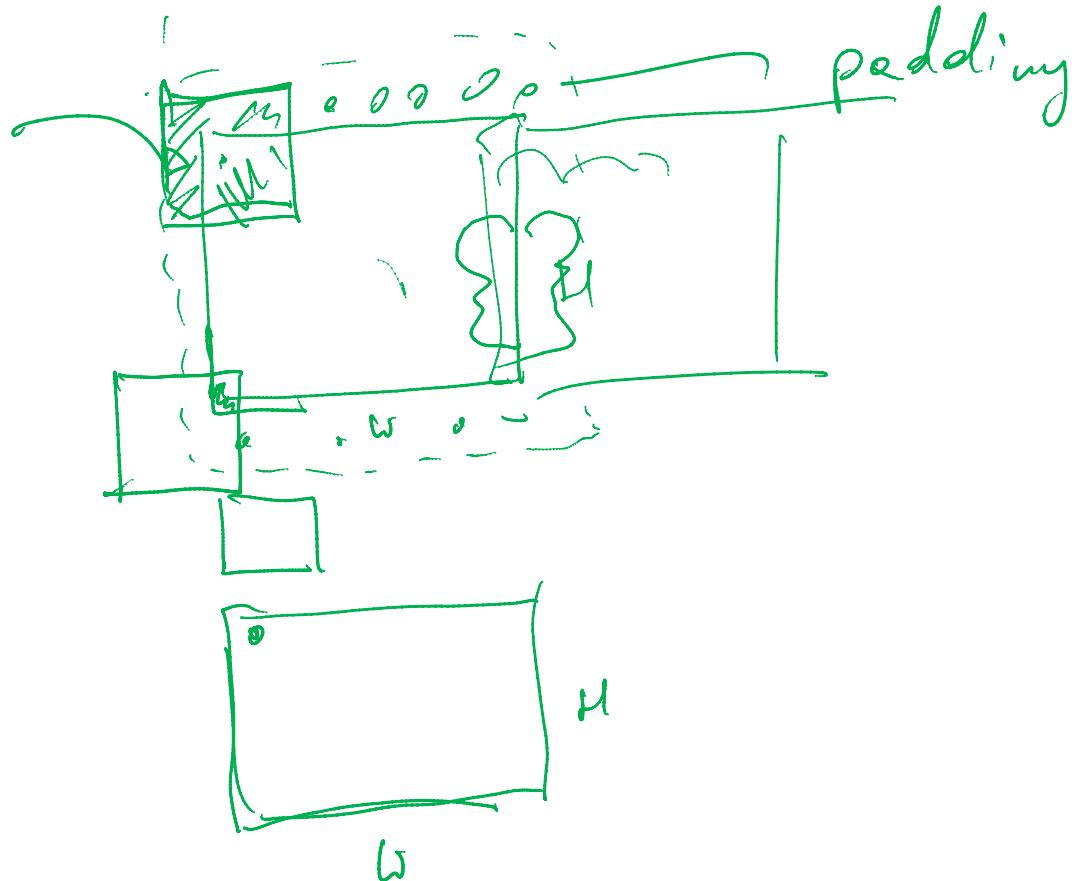
wide receptive field at high resolutions



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

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

Technicality: Edge padding



= VALID conv .

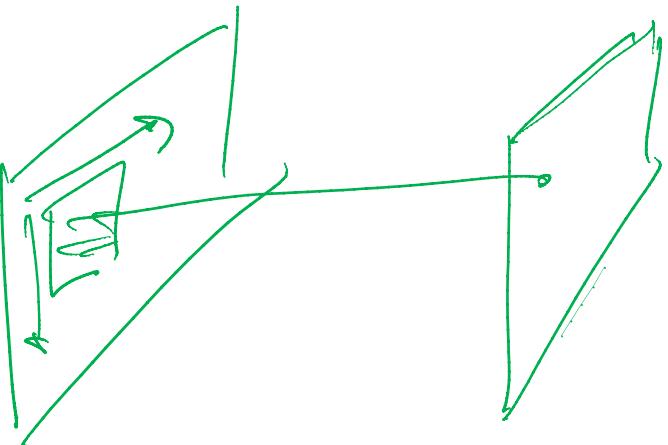
if input is $W \times H$ $F_w \times F_h$

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

- FULL

$$(W + P_w - 1) \times (H + P_h - 1)$$

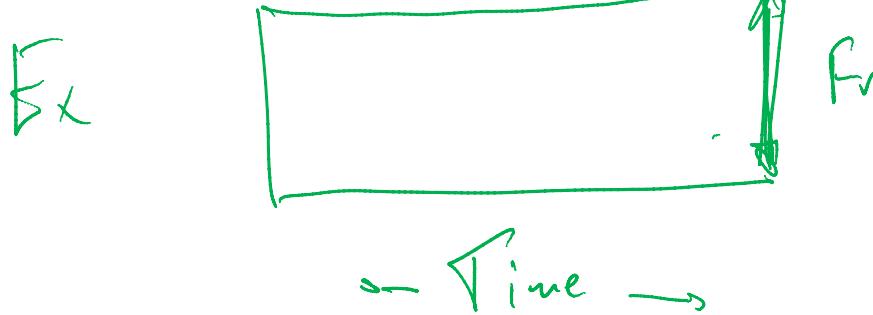
Technicality: bias



$$C \times W \times H$$

$$C \times 1 \times 1$$

bias : a number per filter

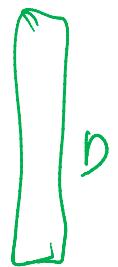
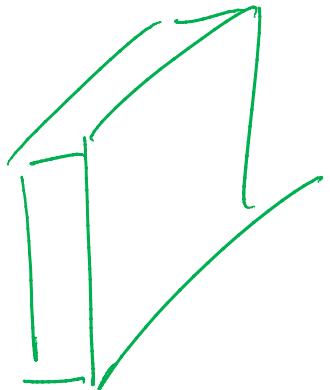


Weights : the same all over

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

Fine Prey

Convs can mimick FC layers



Fully Connected or Dense

$$C \times W \times H \rightarrow D \times 1 \times 1$$

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

VALID over

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

Input: $C \times W \times H$

Filters $D \times C \times \underbrace{1 \times 1}_{\text{Pixelwise}} \rightarrow$ pixelwise FC network

Conv Layers: All the Options

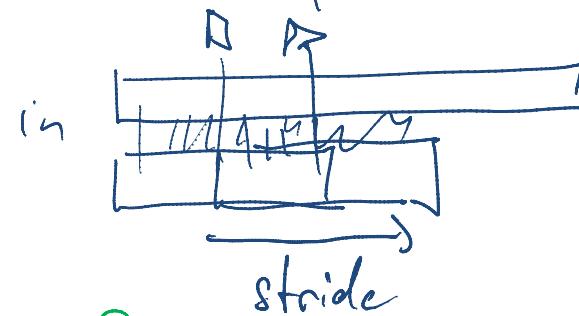
in 1 minibatch of multi-channel images - 4D tensor

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

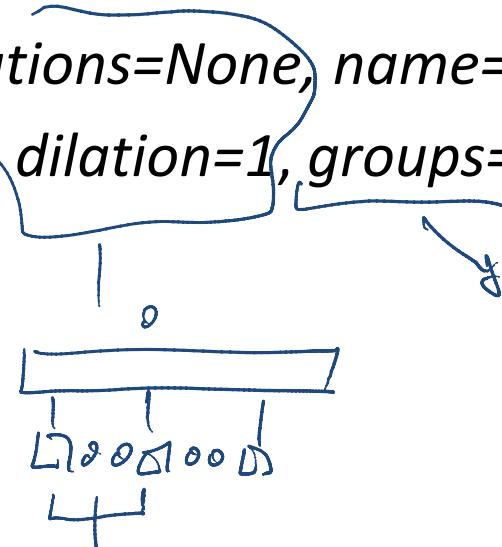
`torch.nn.functional.conv2d(input, weight, bias=None, stride=1, padding=0, dilation=1, groups=1)`

4D tensor $NCHW$

4D tensor $Cout \times Cin \times F_H \times F_W$

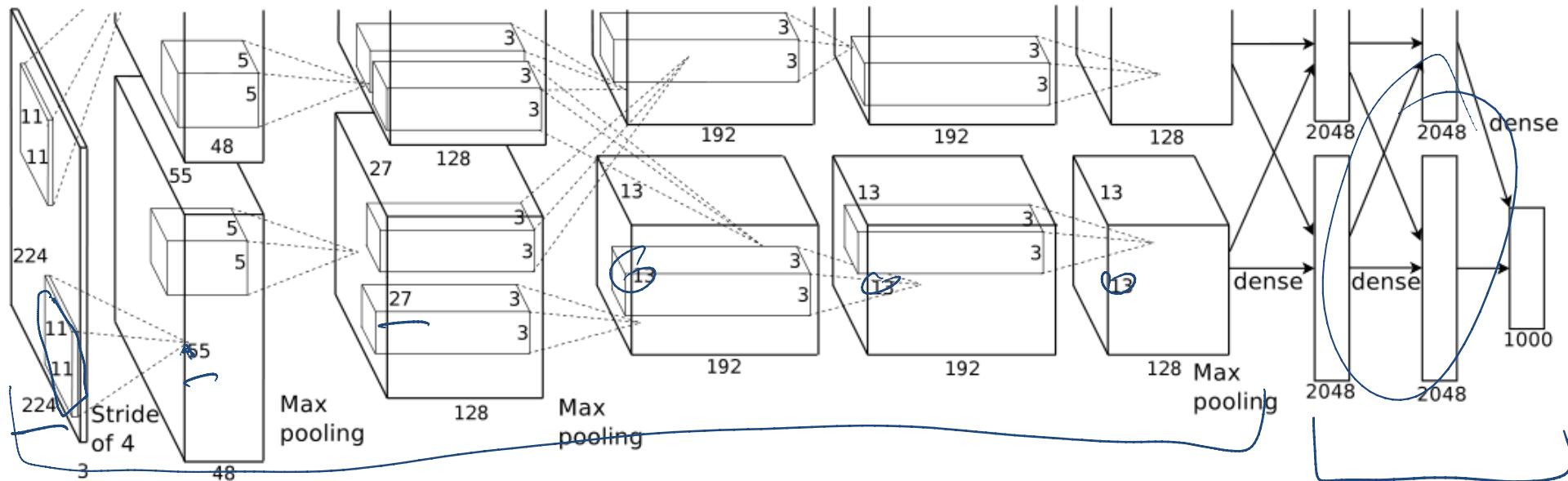


pooling, stride
every n-th output



2 way -
controlled
by dilation

AlexNet (2012)



60M parameters

The last dense layers takes:

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

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

$$\text{Conv}(x, q \times q \text{ filter}) = \text{Conv}(\text{Conv}(x, 3 \times 1), 1 \times 3)$$

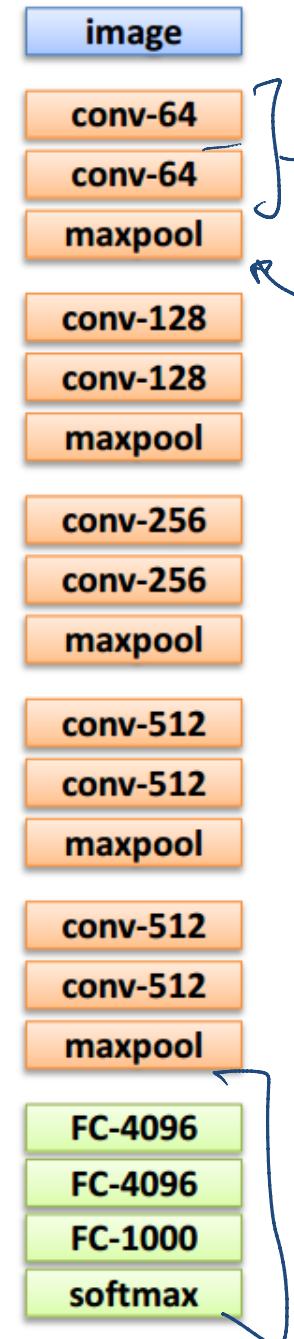
VGGNet: Network Design

Key design choices:

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

Other details:

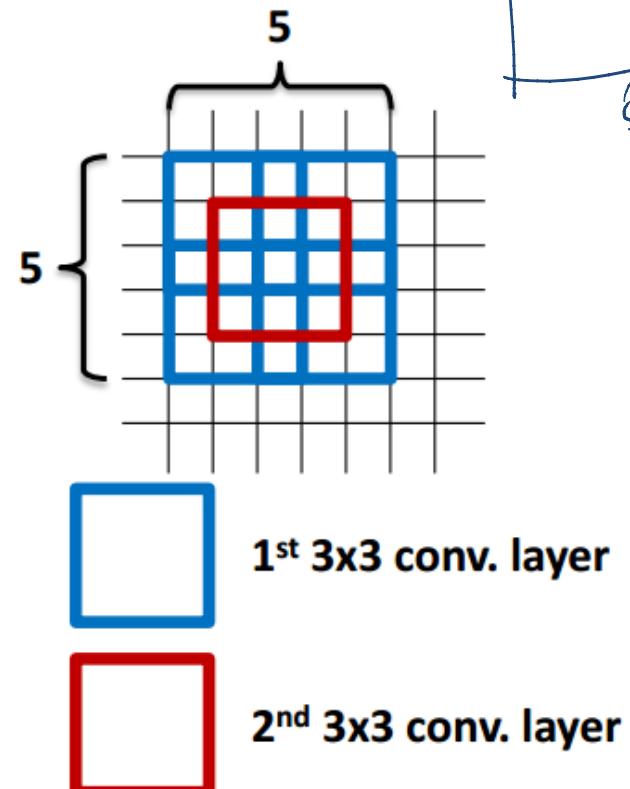
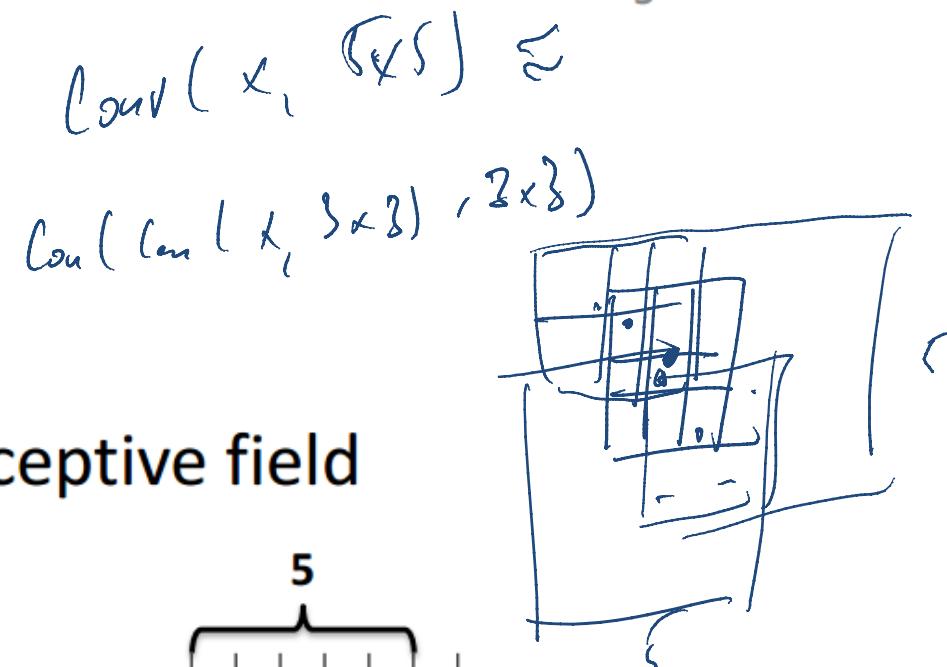
- Rectification (ReLU) non-linearity
- 5 max-pool layers (x2 reduction)
- no normalisation
- 3 fully-connected (FC) layers



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: $7*7*512*4096 = 102,760,448$

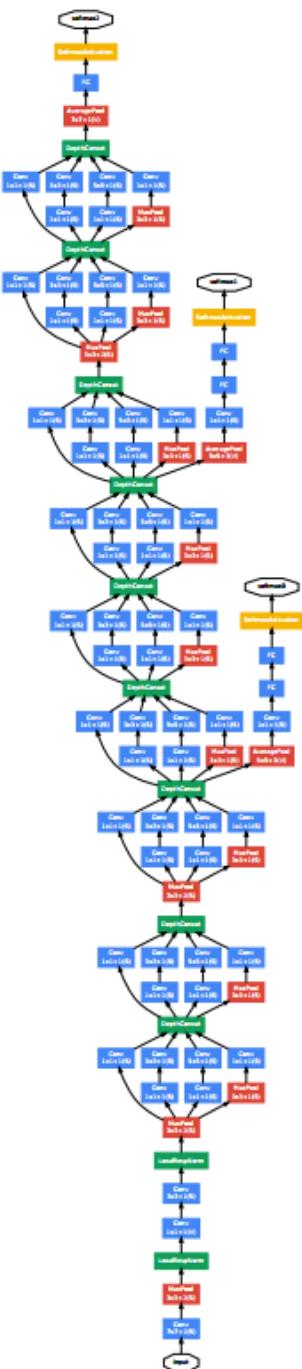
FC: [1x1x4096] memory: 4096 weights: $4096*4096 = 16,777,216$

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

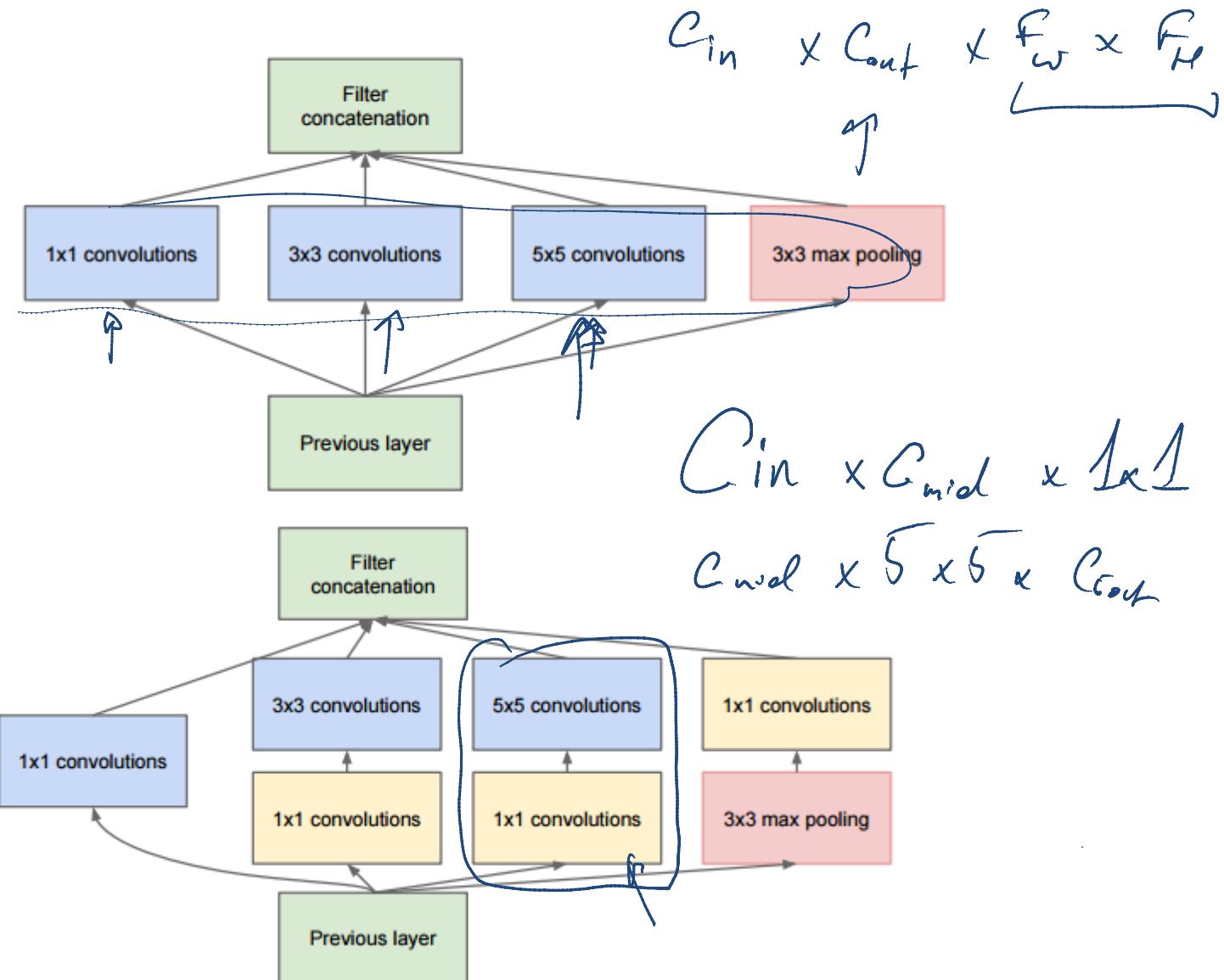
TOTAL memory: 24M * 4 bytes ~ 93MB / image (only forward! ~*2 for bwd) TOTAL params: 138M parameters

up to 2M

no memory

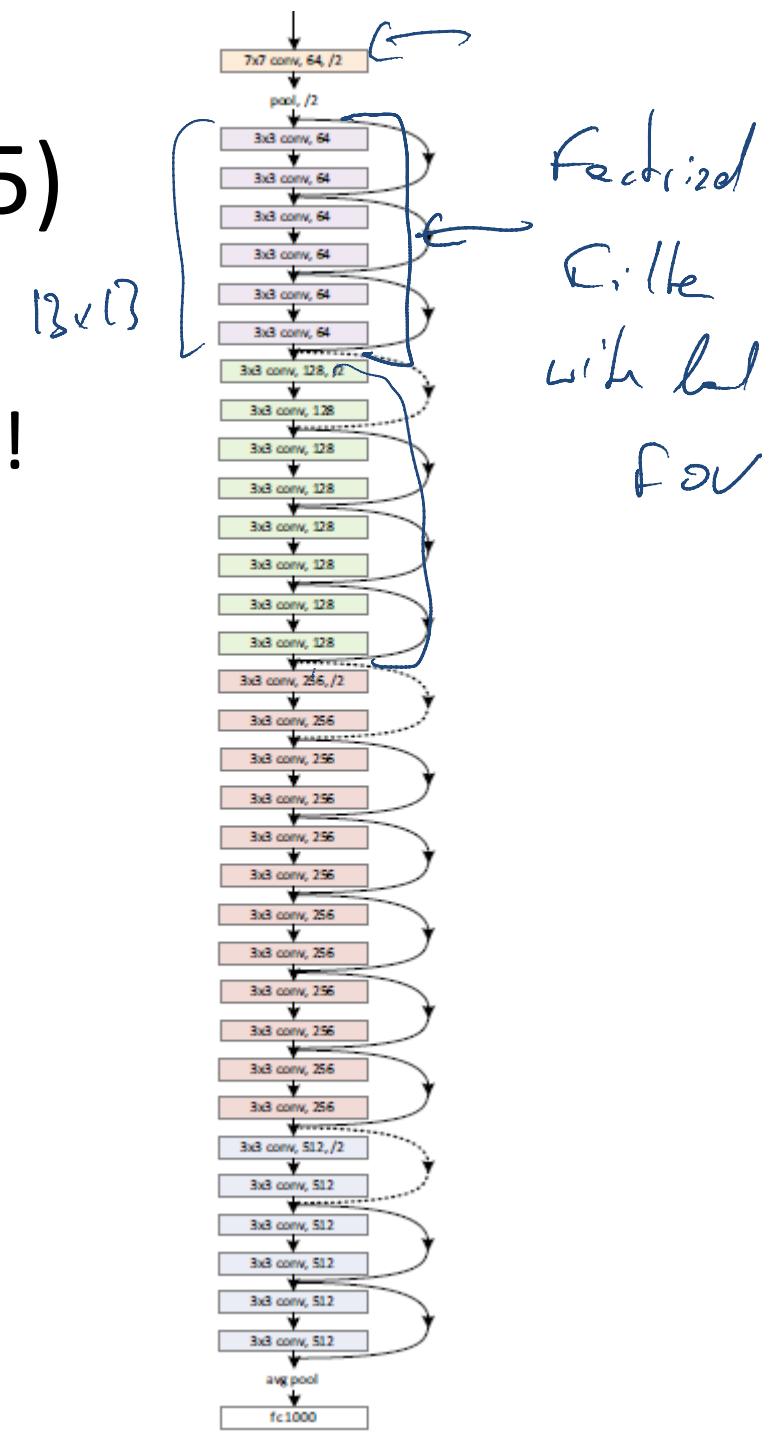
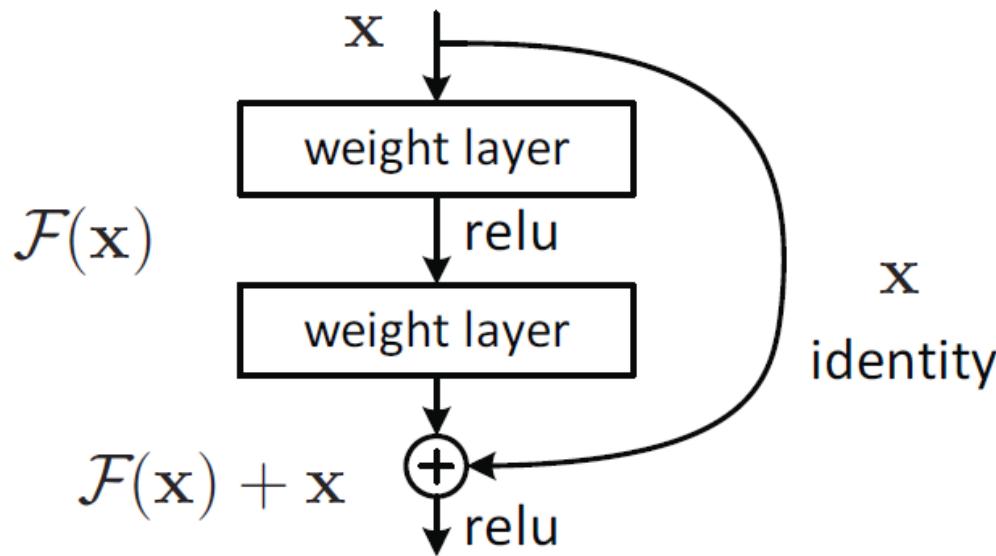


GoogLeNet (Inception) 2014



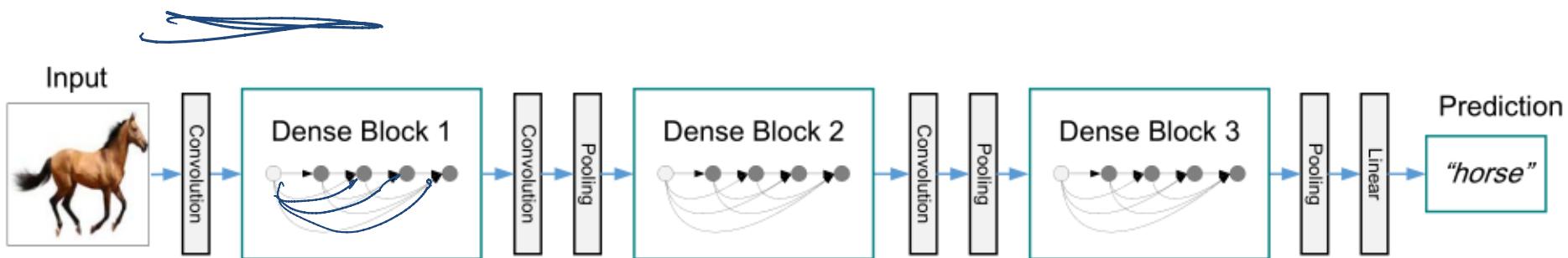
ResNest (Residual Nets, 2015)

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

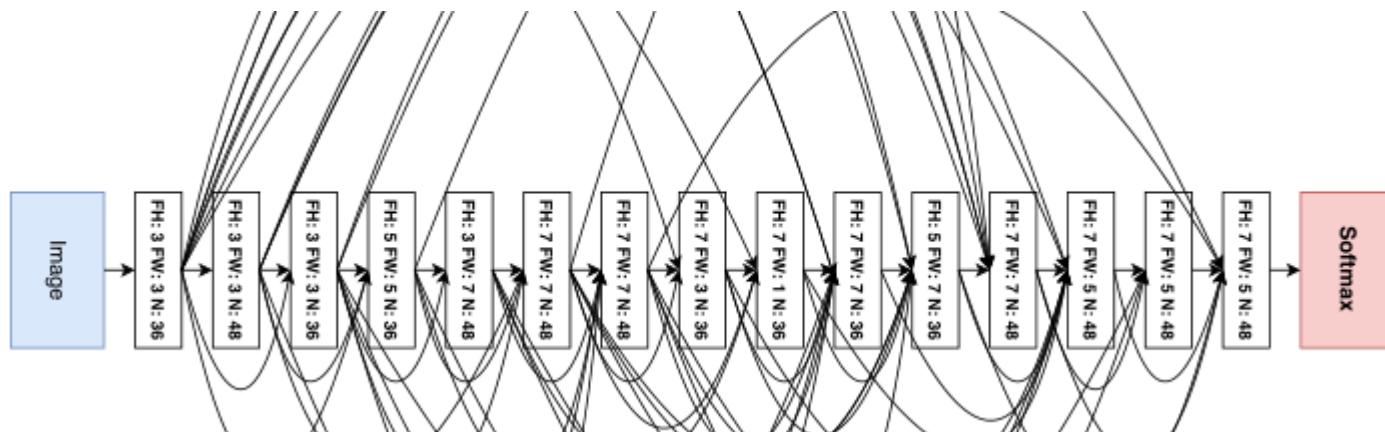


Quest for models

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



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

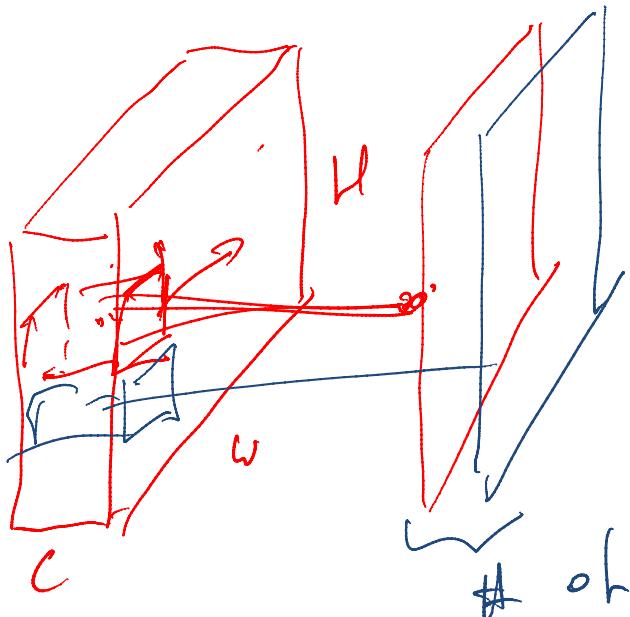


Which ConvNet to use?

- 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!

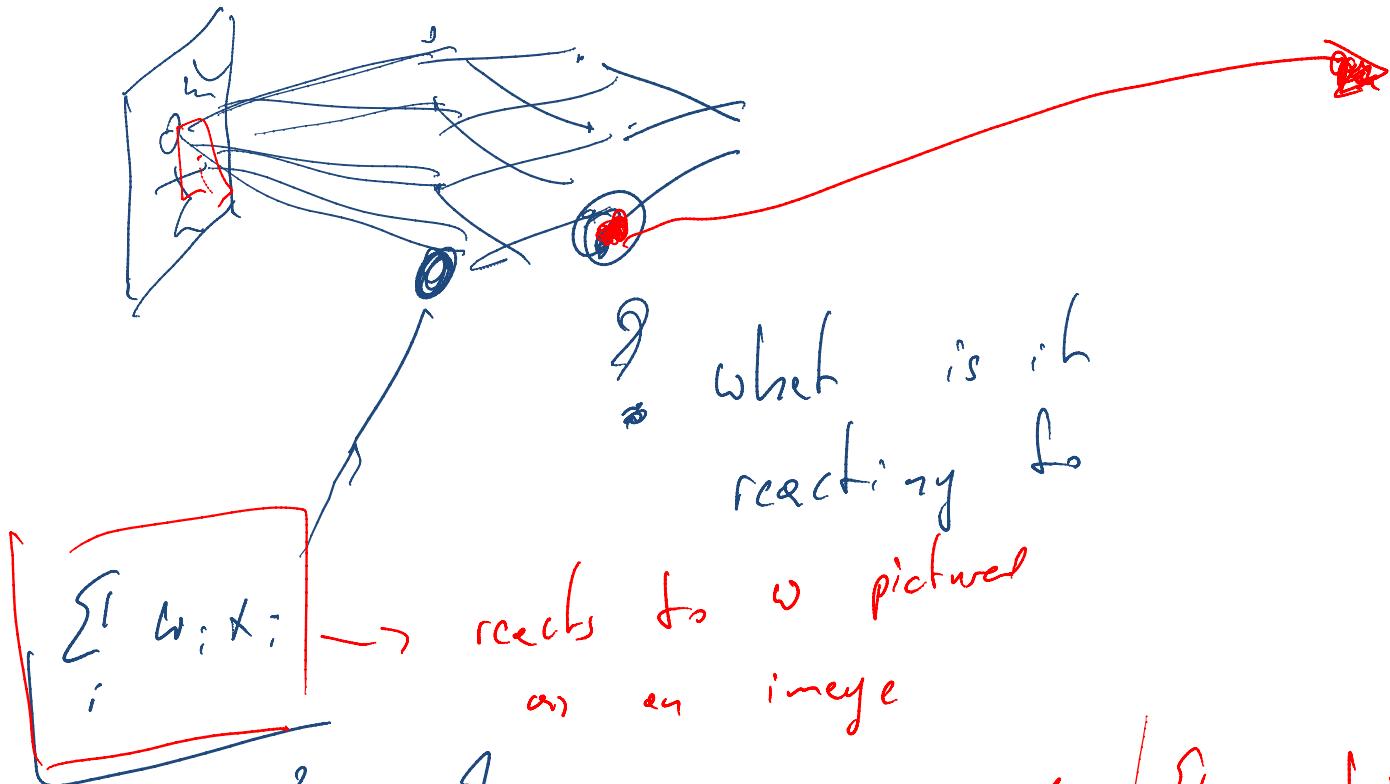
Conv recap

- Local connectivity \rightarrow small, localized features make sense
- Weight sharing \rightarrow same features are detected everywhere



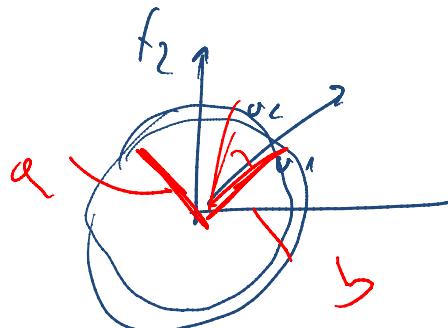
of feature maps = # number of output
desired = # filters

How to understand a model?



? What is it reacting to
as an image

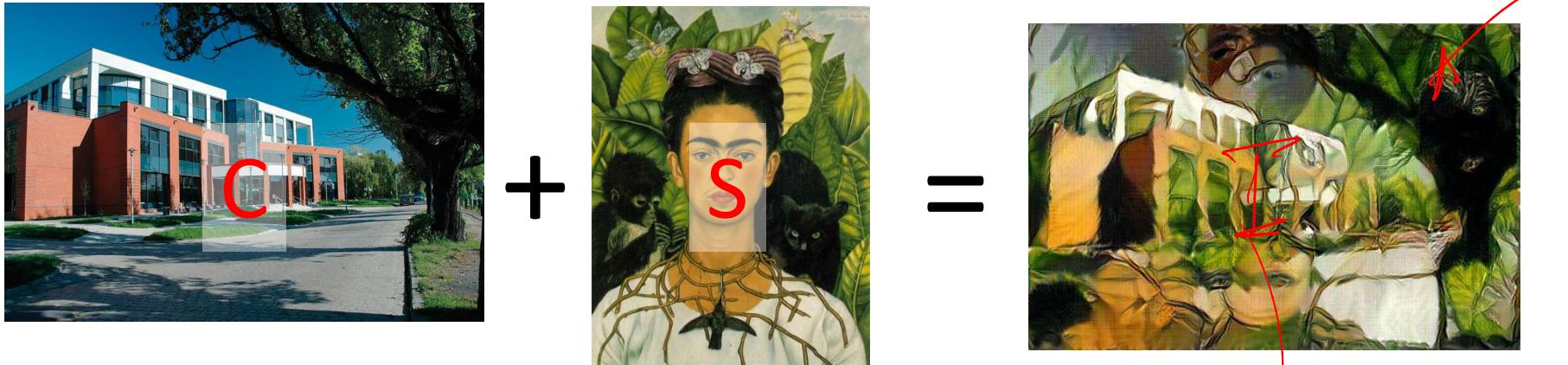
$$\|x\|_2^2 \leq 1$$



$$\sum_i q_i w_i < \left(\begin{cases} a, b \end{cases} \right)$$

- 1 Build a linear approx to this unit $o(x+\Delta x) =$
$$o(x) + \frac{\partial o}{\partial x} \Delta x + \epsilon_m$$
- 2 Solve for an input that maximally excites this neuron

Sidenote Style Transfer



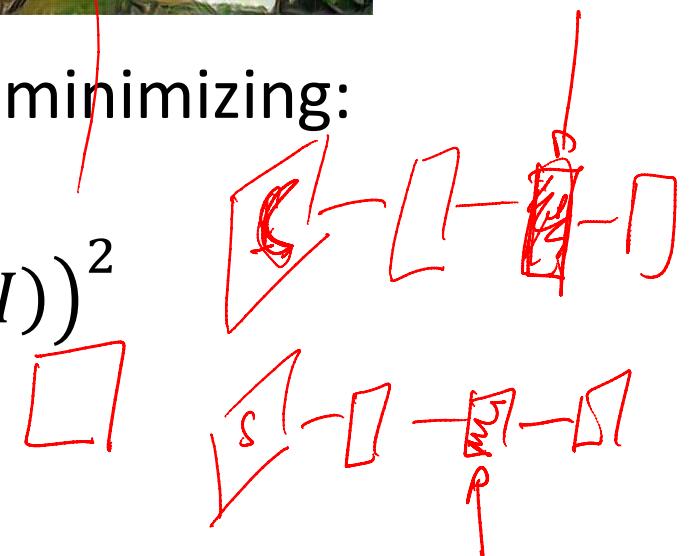
Find image (backpropagation toward pixels) minimizing:

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

Where:

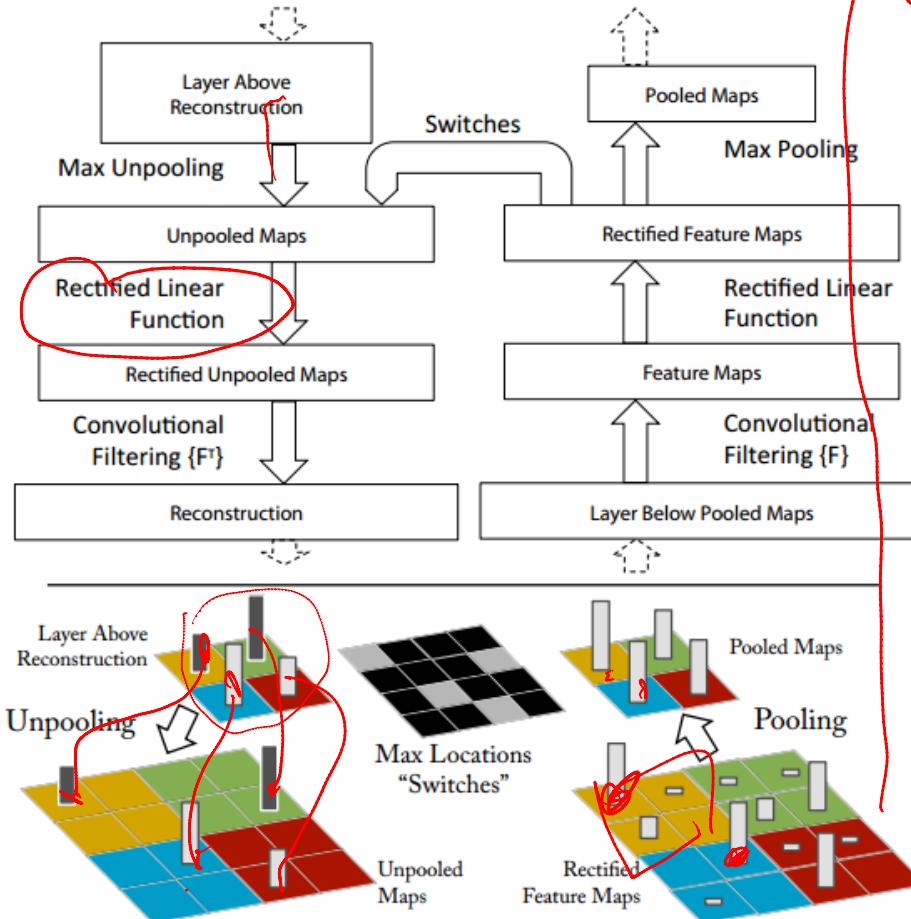
$F(x)$ features of a CNN layer on image x

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



Visualizing CNN features

Vanilla PWD rev



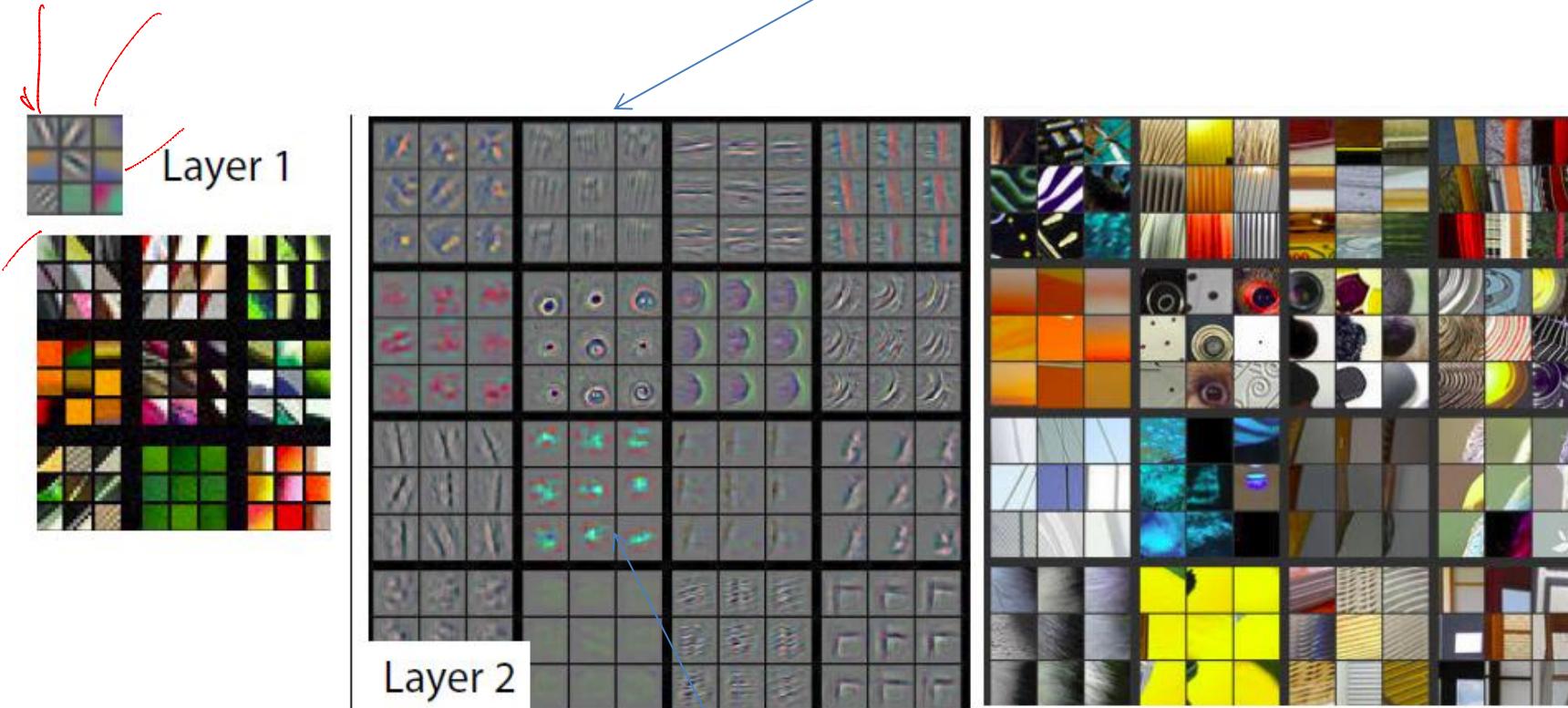
Main idea:

1. Do a forward propagation,
2. Save state of max-poolings,
3. „Deconvolve“: trace back the activations of chosen units to inputs

Note: similar to gradient computation, but applies ReLU during deconvolution
(gradient would use the saved jacobian of relu)

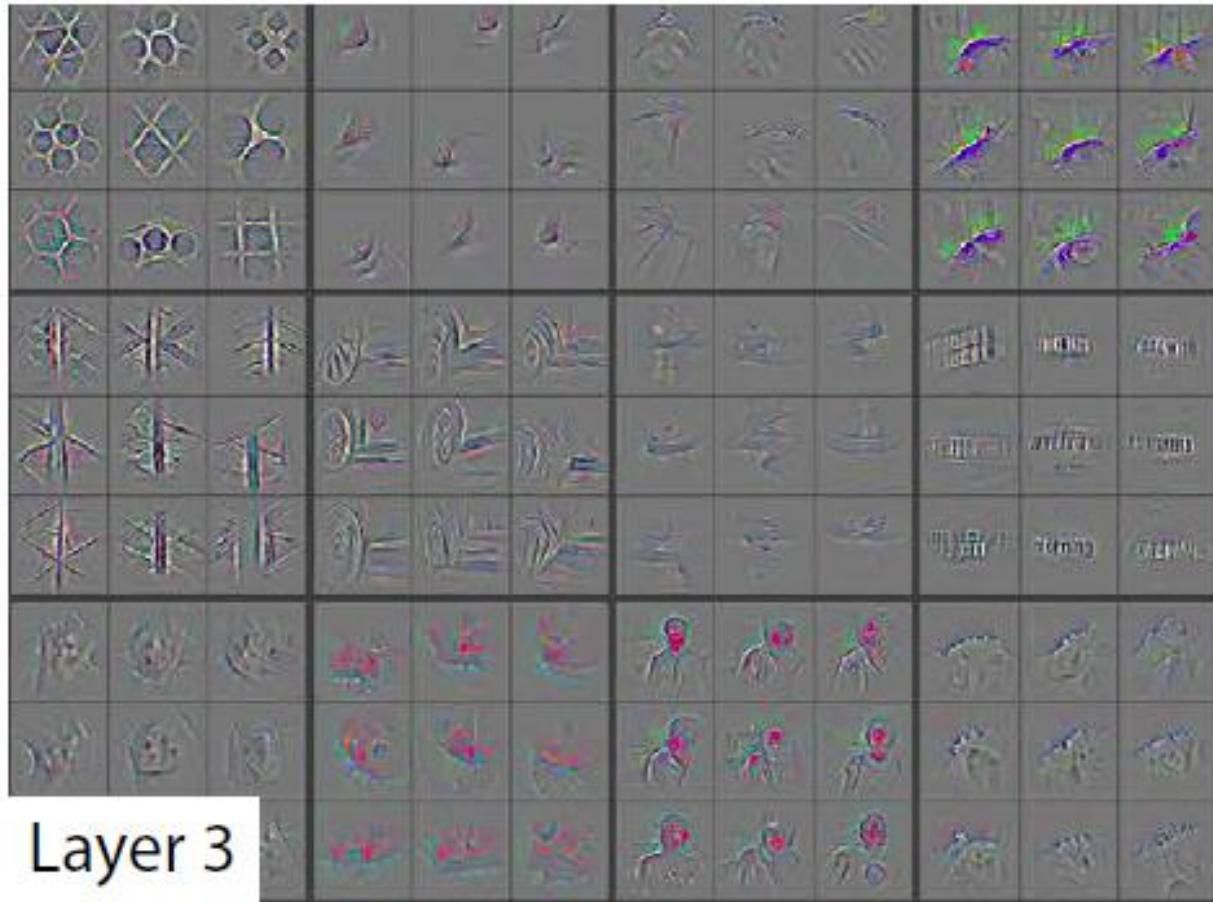
Low-level features

What the neuron (feature-detector looks for)

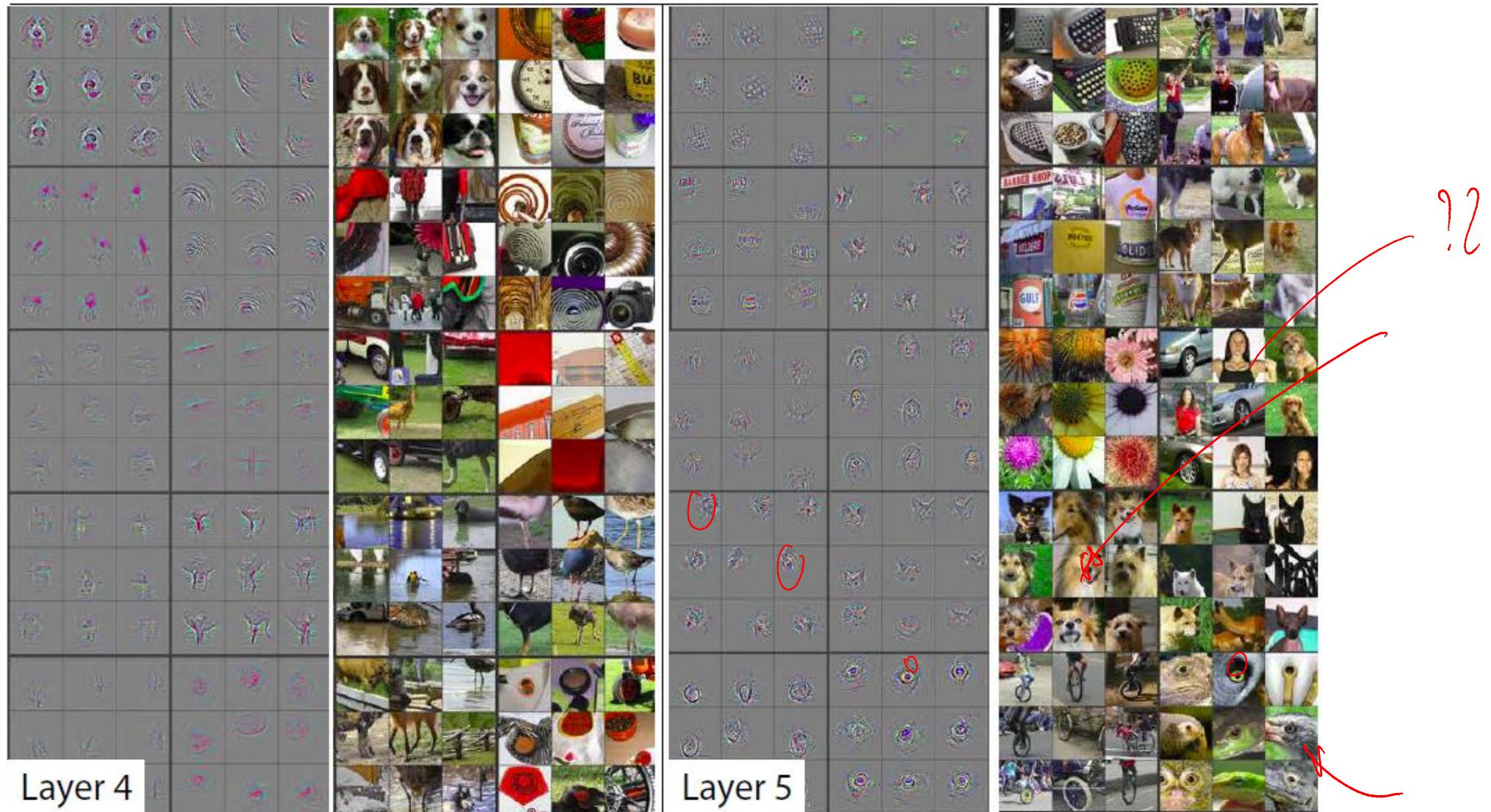


What images are selected by the neuron

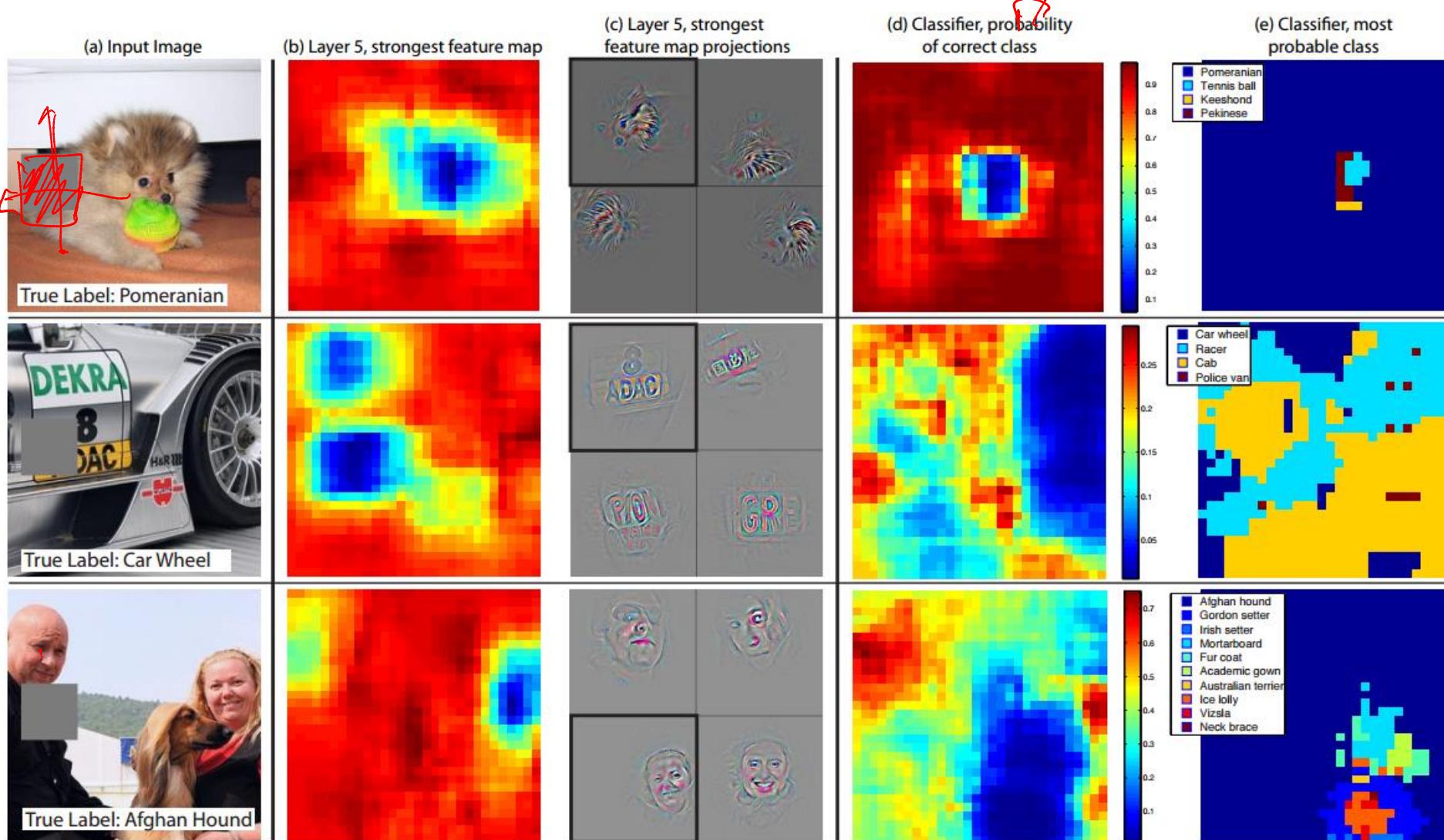
Mid-level features



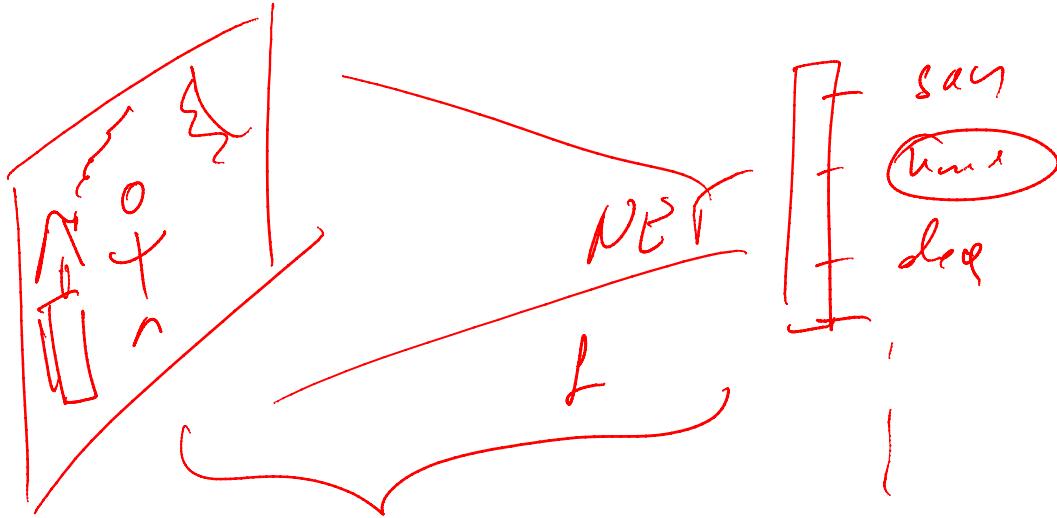
High-level features



Where Convnets look?

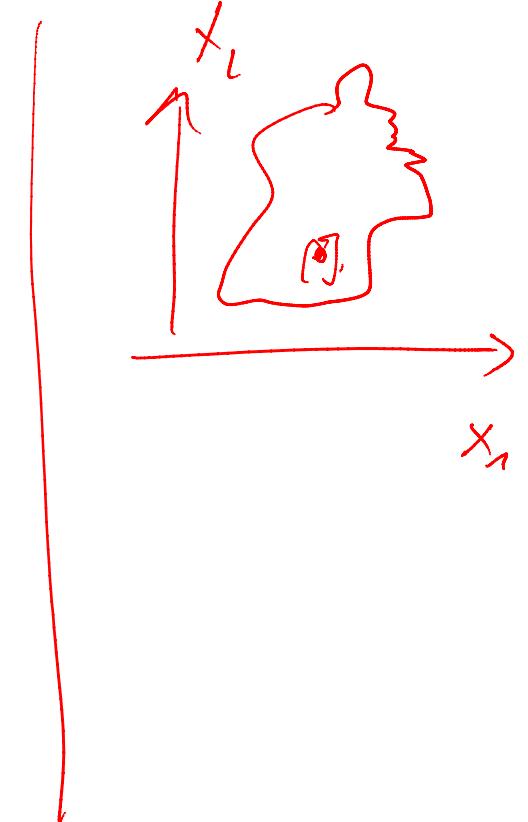


Gradient saliency map



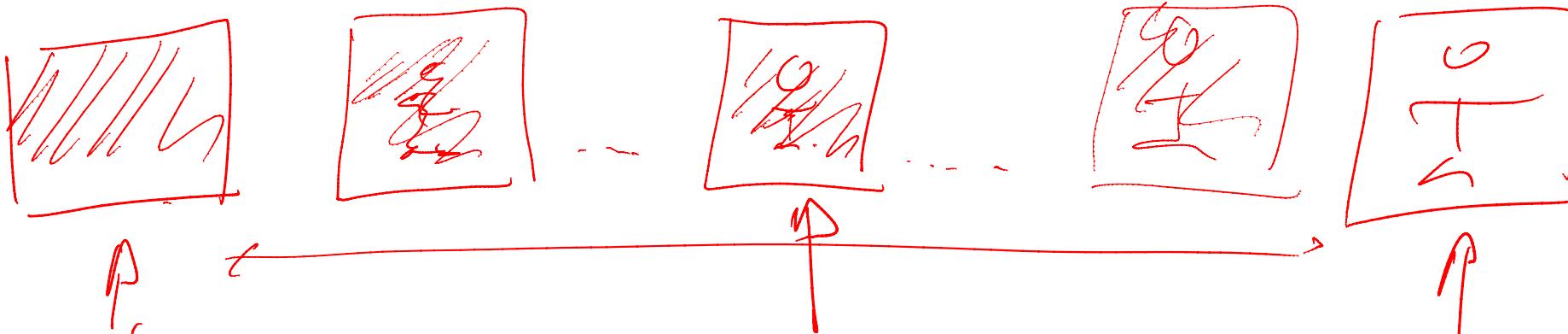
$$f(x + \Delta x) \approx f(x) + \Delta x \frac{\partial f}{\partial x}$$

Saturation or stability → we hope flat



$$f(x - \Delta x) \approx f(x)$$

Saliency maps with integrated gradients



uniform grey

uniform bin-

ary, $\frac{\partial L}{\partial x}$ should be
small

subject barely
visible

$\frac{\partial L}{\partial x}$ is large

Final img
pattern clearly
seen

to stable out-

$\frac{\partial L}{\partial x}$ is big

Idea: integrate (add) gradients over
linear interpolation from uniform grey to image

LIME: local surrogate models

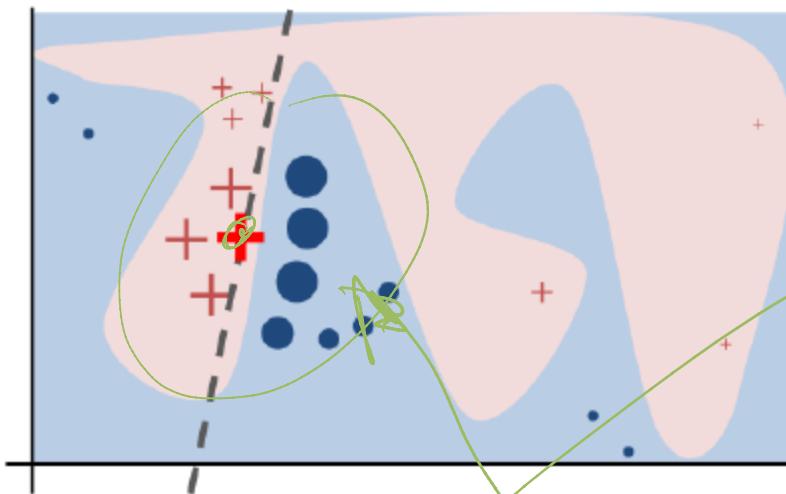


Several points sampled
close to x
Boil the net on them

Idea: build a surrogate (model of a model)
that explains locally the decision boundary
 $f \approx \text{small neighborhood}$

p_1	$f(p_1)$	great
p_2	$f(p_2)$	bad & train,
p_3	$f(p_3)$	train
\vdots	\vdots	in regression

LIME details

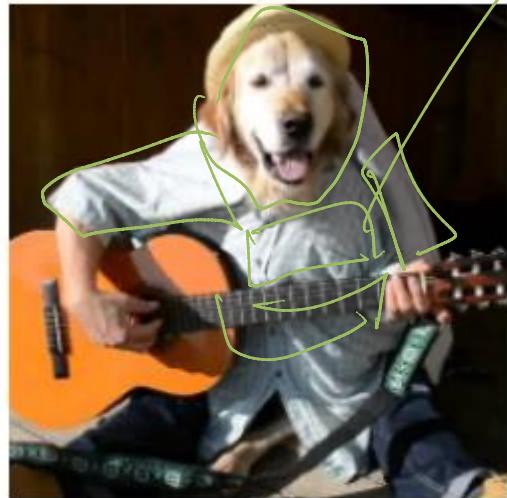


Focay on vicinity
of x
 π_x - weighs higher
points closer to x
 g is simple

$$\xi(x) = \underset{g \in G}{\operatorname{argmin}} \quad \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

↑
sample
class for g

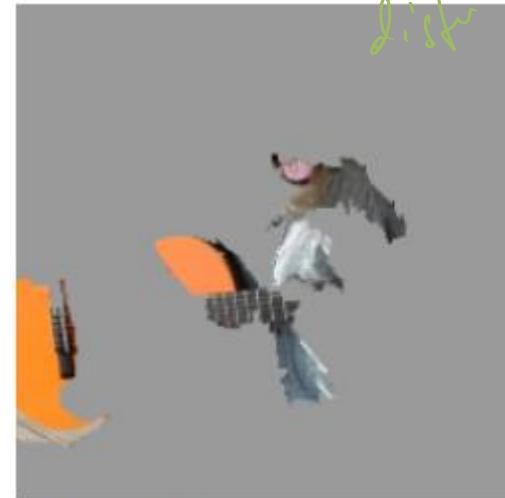
LIME on images



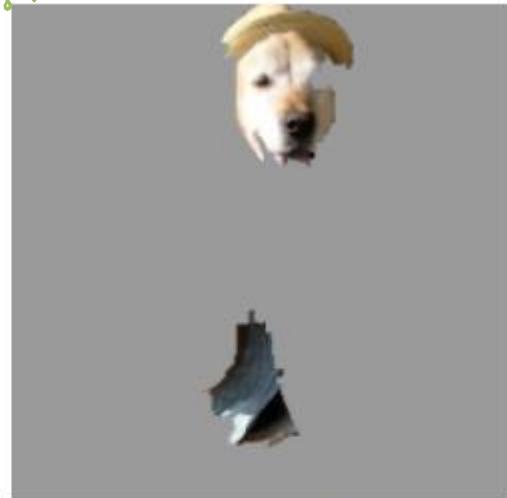
(a) Original Image



(b) Explaining *Electric guitar*



(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*

super pixels
standard CV task
lime tree set
disk
so hard on

sp1	sp2	sp3	...	f(x)
1	0	1		.
0	1	1		.
1	1	1		.
0	0	0		.

Figure 4: Explaining an image classification prediction made by Google’s Inception neural network. The top 3 classes predicted are “Electric Guitar” ($p = 0.32$), “Acoustic guitar” ($p = 0.24$) and “Labrador” ($p = 0.21$)

Beware: some methods lie

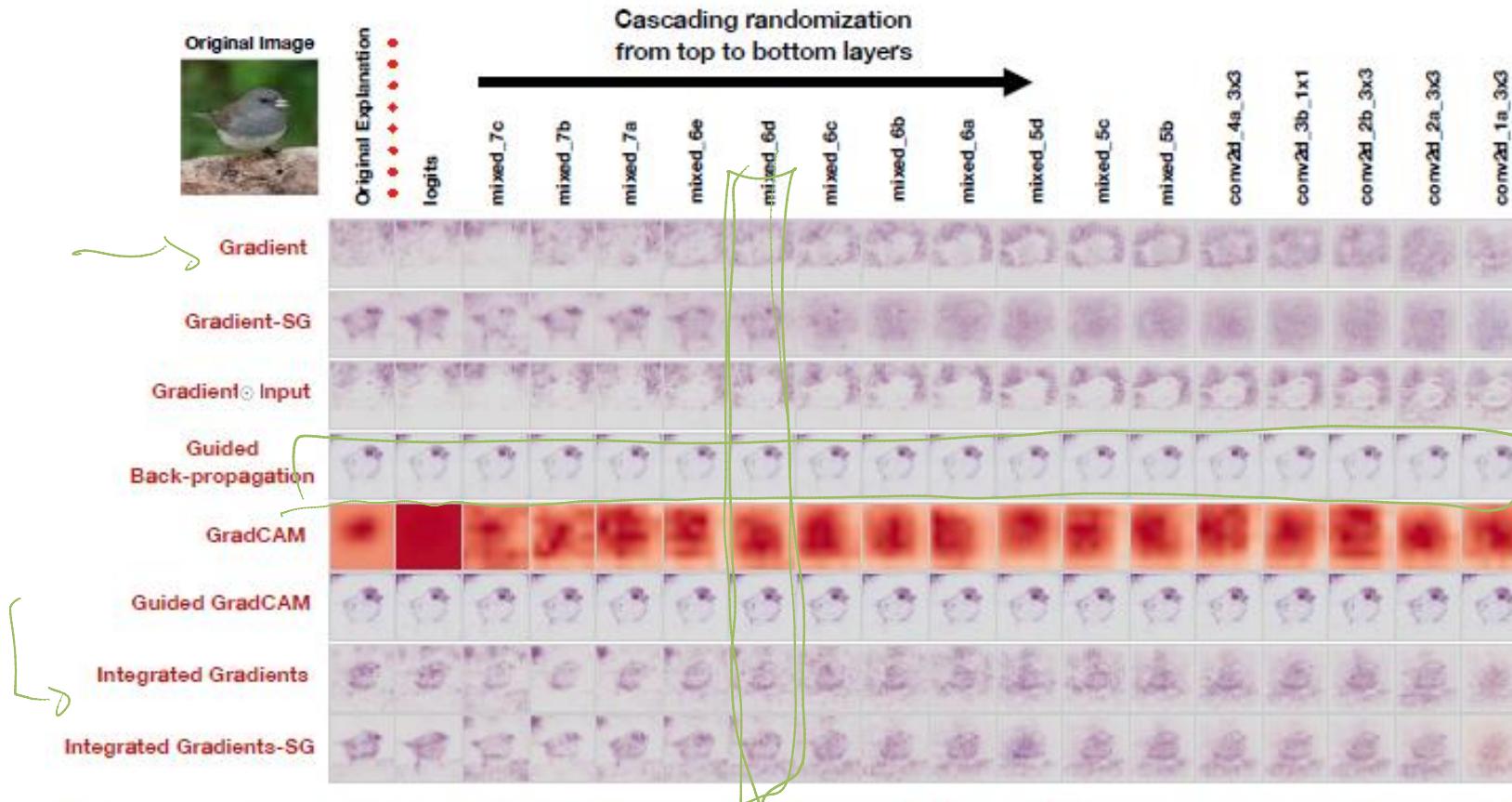


Figure 2: **Cascading randomization on Inception v3 (ImageNet).** Figure shows the original explanations (first column) for the Junco bird. Progression from left to right indicates complete randomization of network weights (and other trainable variables) up to that ‘block’ inclusive. We show images for 17 blocks of randomization. Coordinate (Gradient, mixed_7b) shows the gradient explanation for the network in which the top layers starting from Logits up to mixed_7b have been reinitialized. The last column corresponds to a network with completely reinitialized weights.

Model explanations, further reading

- Integrated gradients: <https://arxiv.org/abs/1703.01365v2>
- Important analysis of saliency maps
<http://papers.nips.cc/paper/8160-sanity-checks-for-saliency-maps.pdf>
- LIME: <https://arxiv.org/abs/1602.04938v3>
- SHAP (similar to LIME, more stable and justified):
<http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>

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)

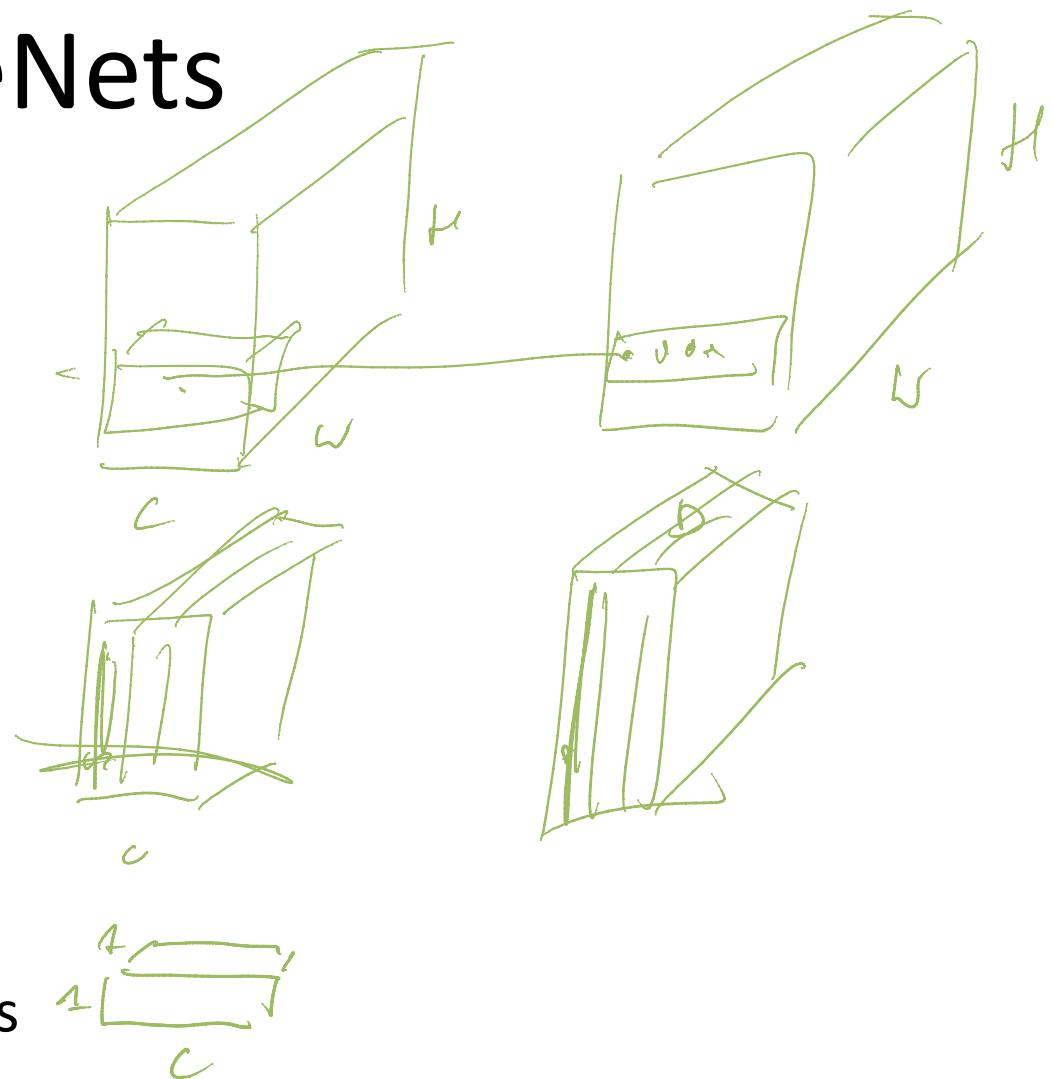
Filter size: $C \times W \times H$ (nb: no D)

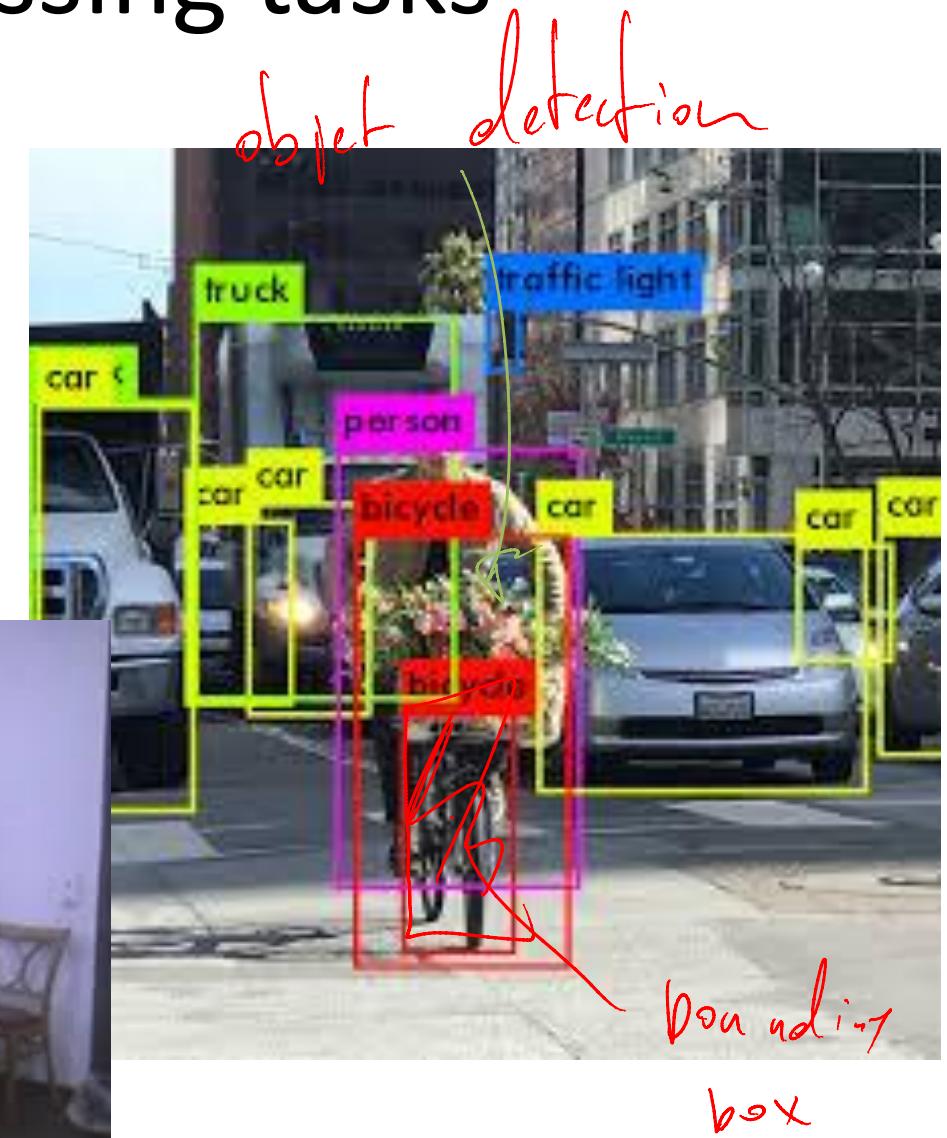
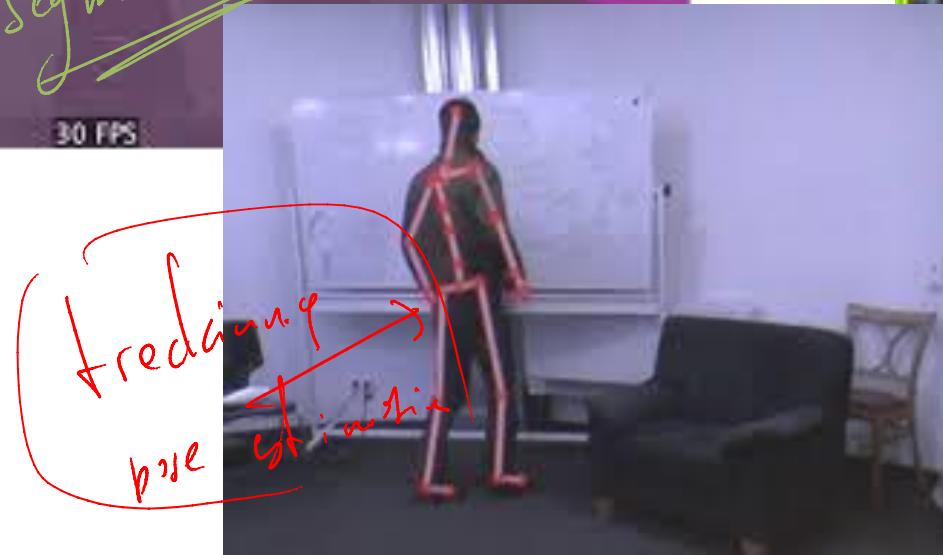
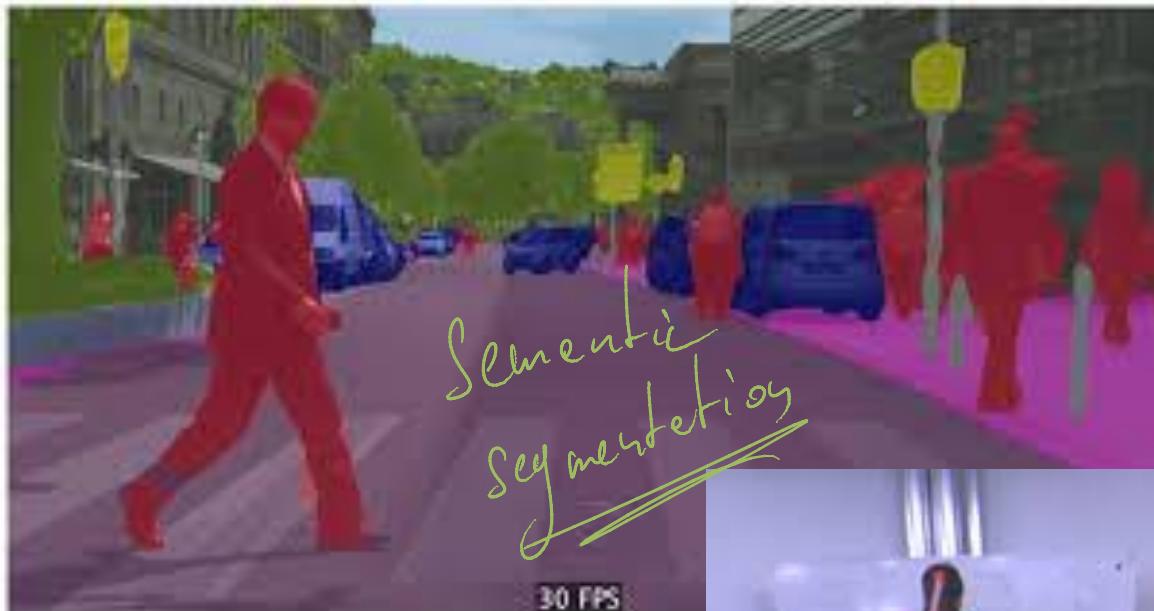
Matmuls per pixel: $C \times W \times H$

3. MobileNet: Depthwise + 1x1 convs

Why? Depthwise: fast, doesn't mix channels

1x1: mixes channels, still small

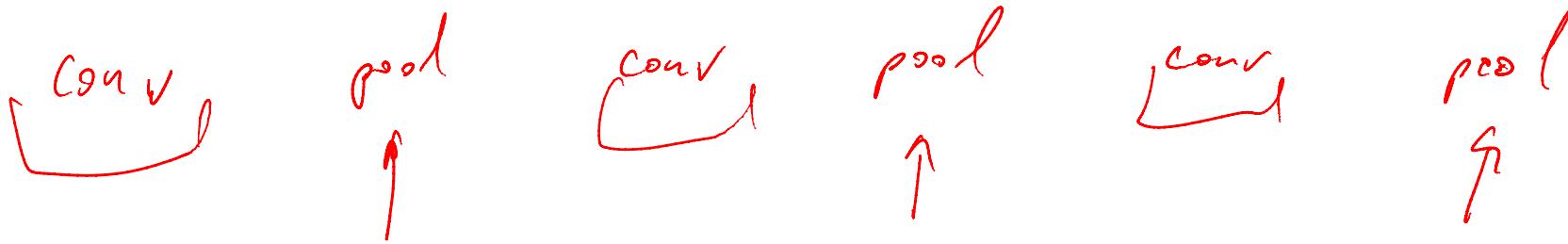




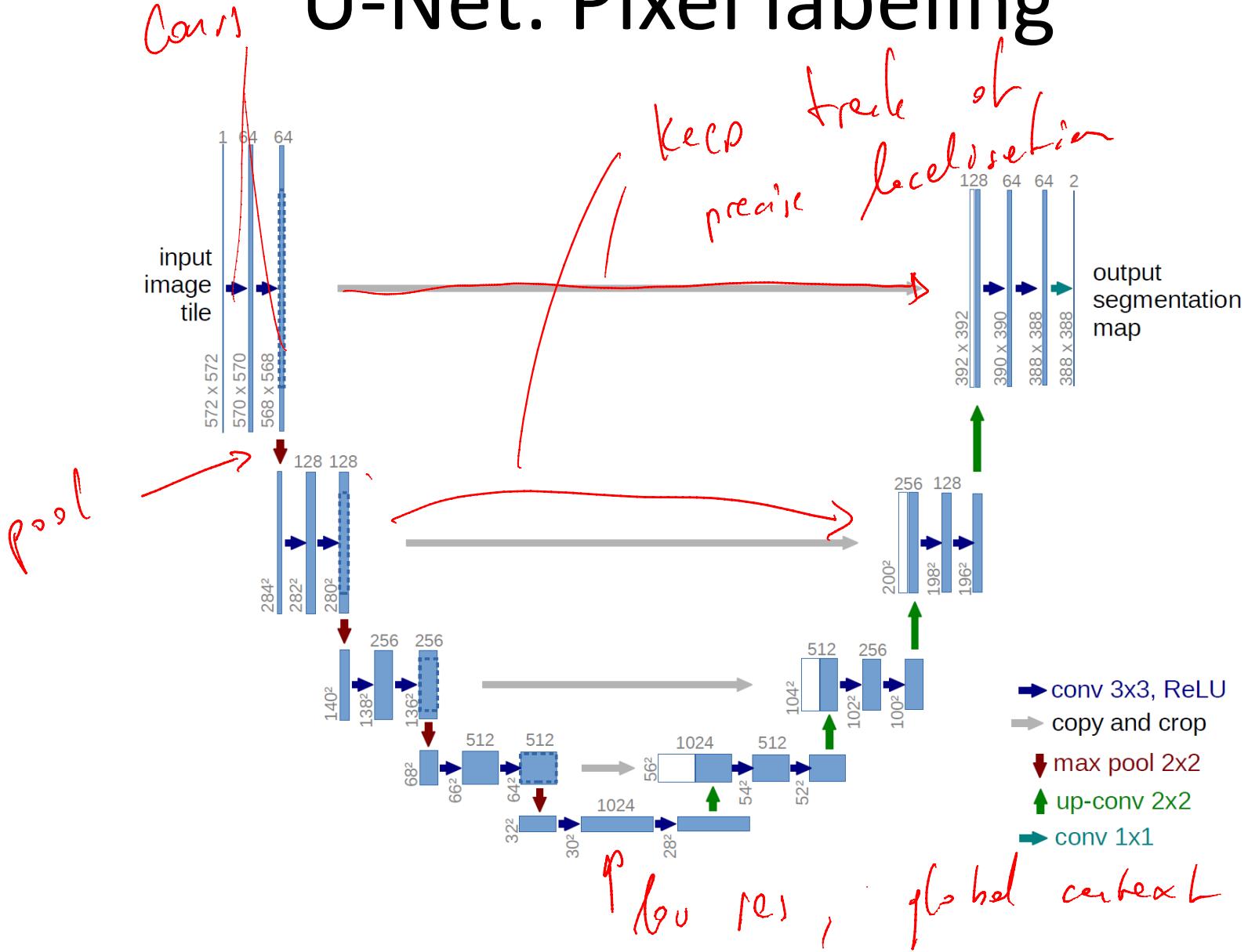
Unet: Pixel labeling intro

Schematic SEP.

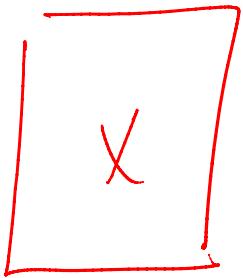
Usually



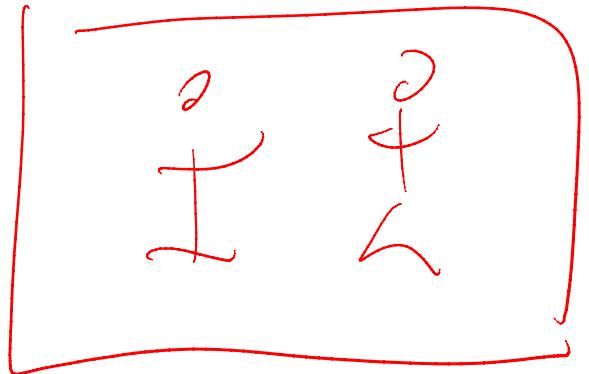
U-Net: Pixel labeling



From Classification to Detection



$x \rightarrow y \rightarrow$ what's in the img



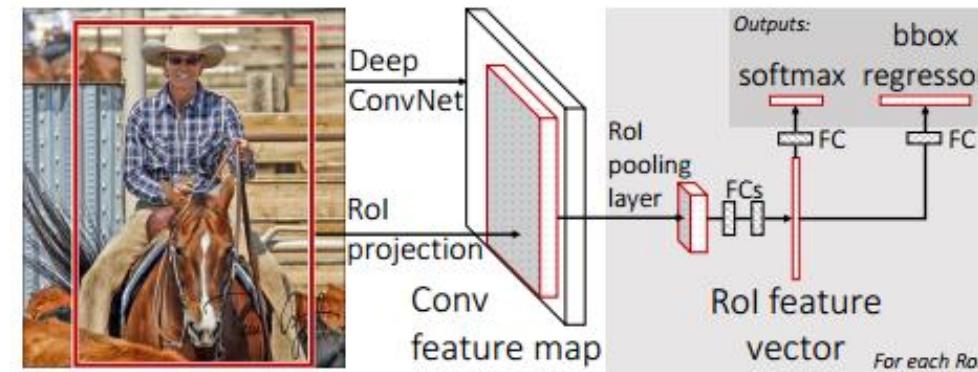
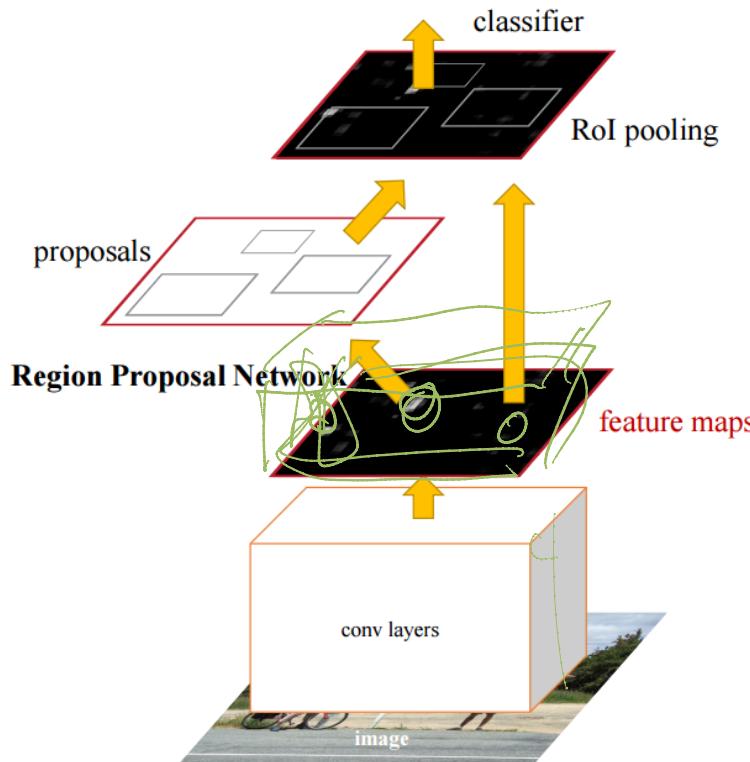
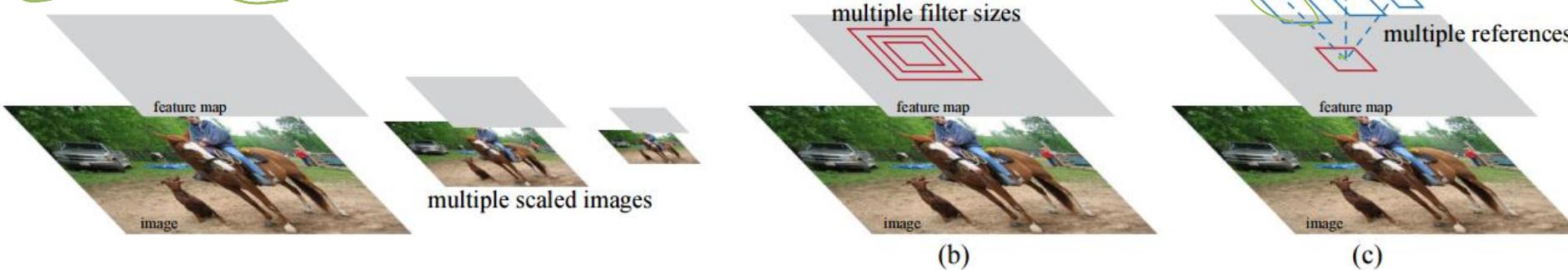
- For a correct detection
can be equivalent,
1. Take crops at all locations
run all through the network
indep. \Rightarrow slow

2. Run the network over large
image.

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

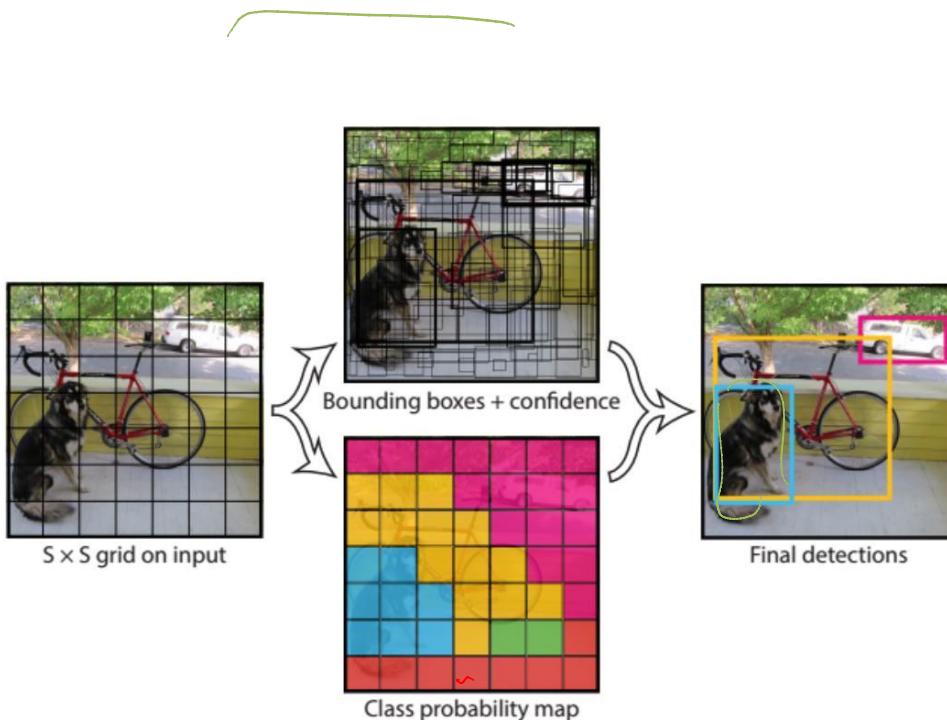
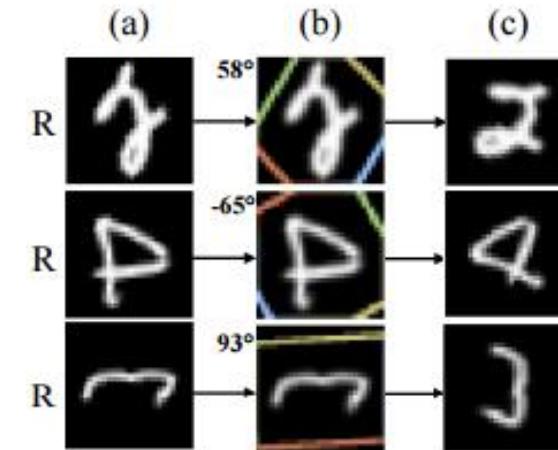
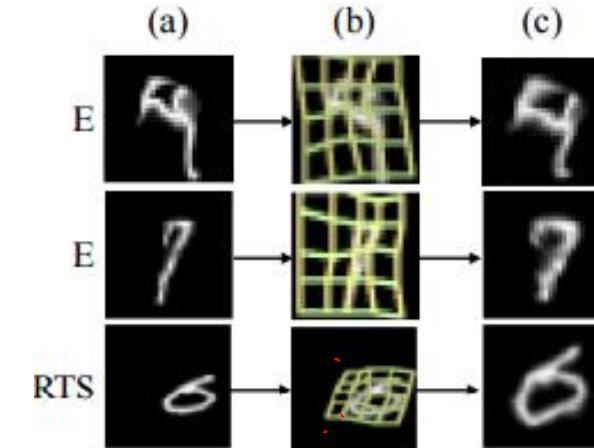
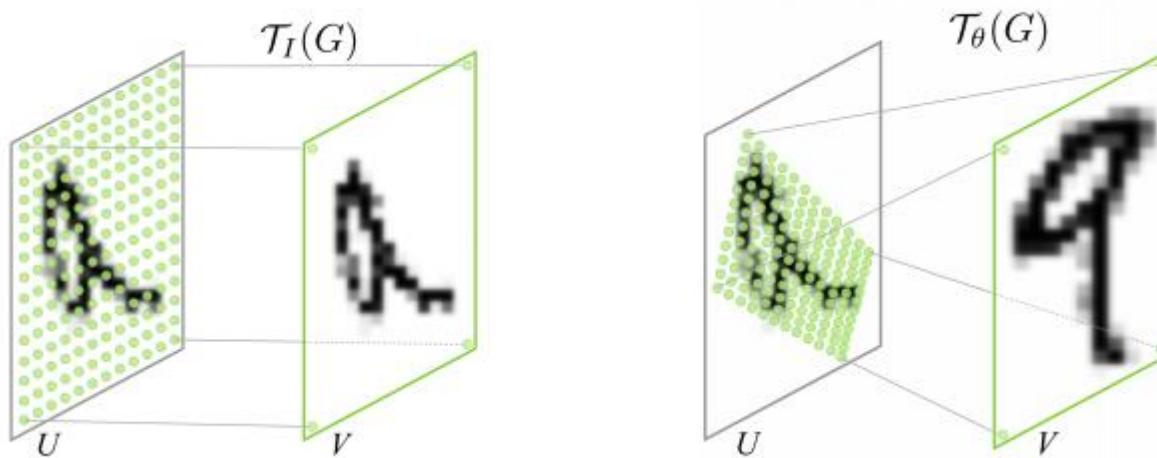
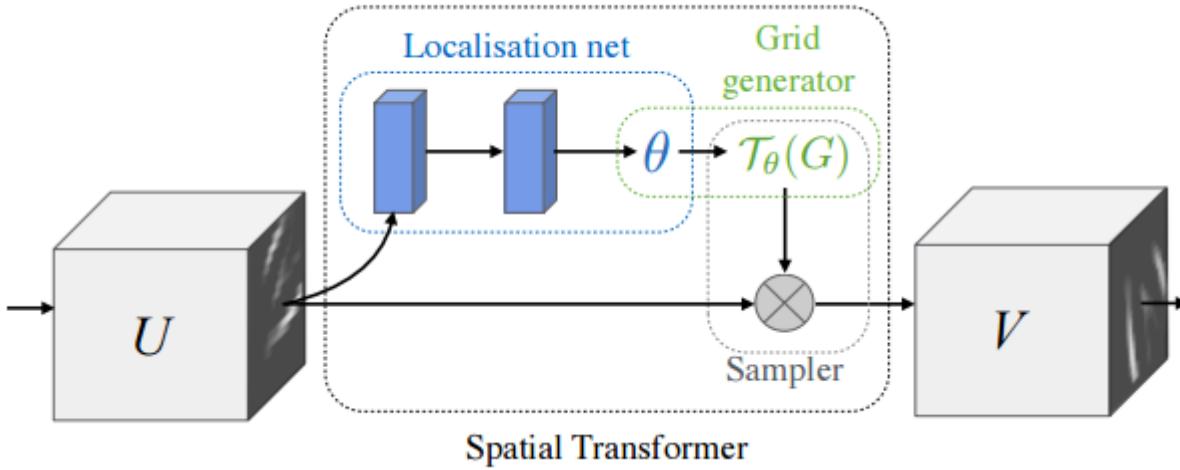


Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

- 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

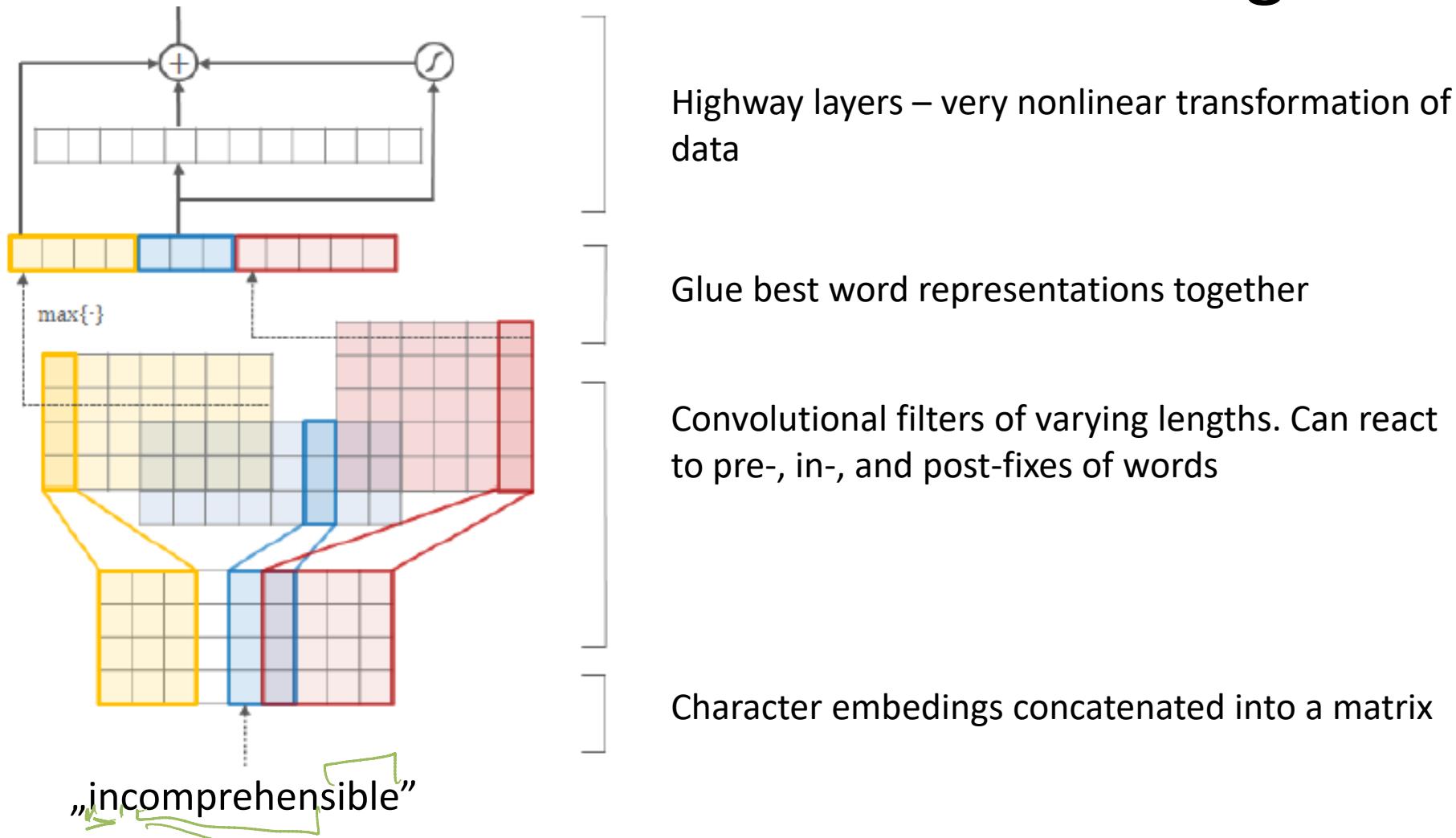
Spatial Transformer Networks

<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.