

# **Artificial Neural Networks: Lecture 8**

## **Convolutional Neural Networks**

Johanni Brea  
EPFL, Lausanne, Switzerland

### **Objectives for today:**

- Review: No free-lunch theorem and inductive biases
- Convolution Layers
- Max Pooling Layers
- Inductive bias of ConvNets
- ImageNet competition and modern ConvNets
- Training ConvNets with AutoDiff
- Applications beyond object recognition

Miniproject 1 is due on April 29 at 23:59  
End of next week: handout of miniprojects 2 & 3

## Reading for this lecture:

**Goodfellow et al., 2016 *Deep Learning***

- Ch 9

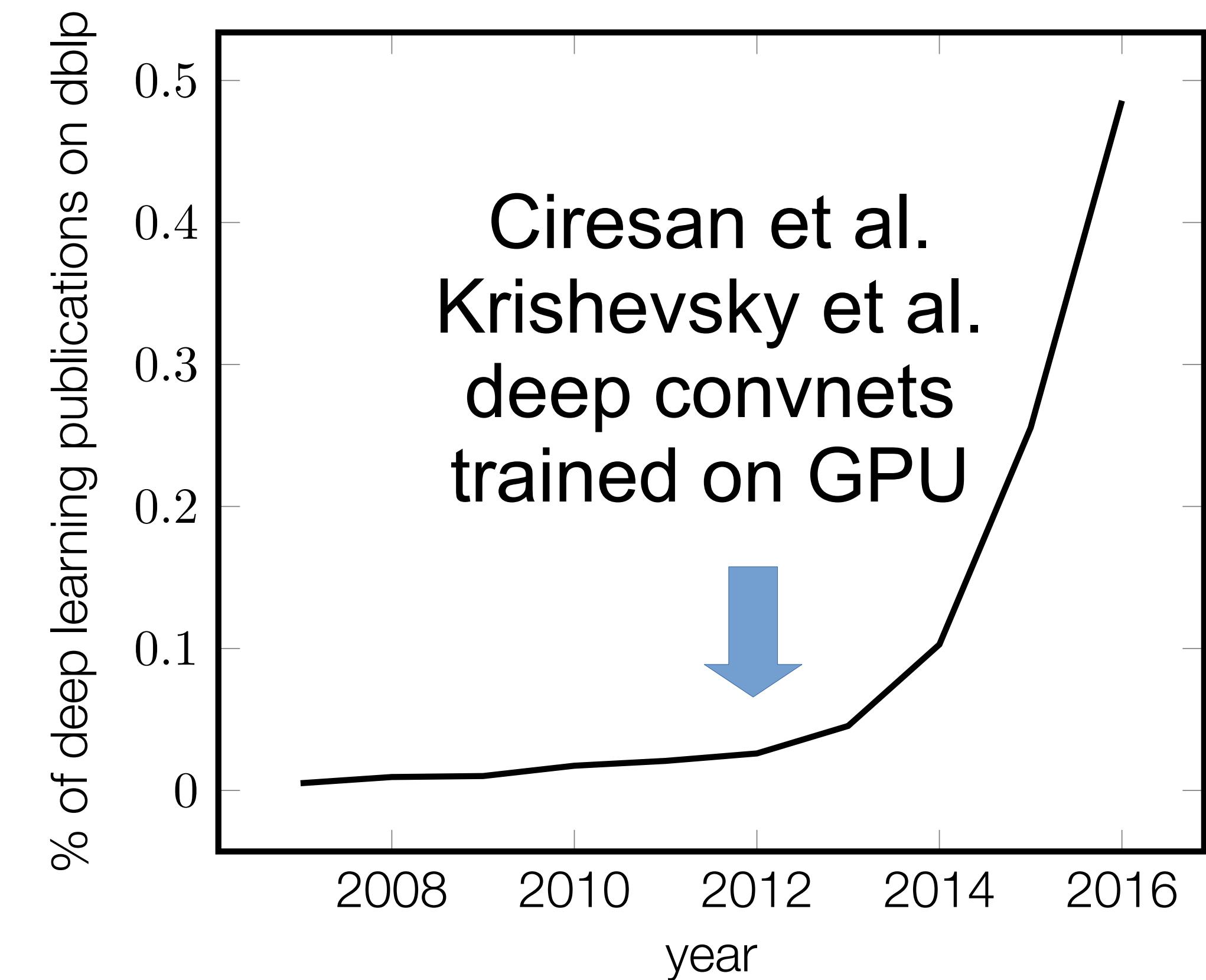
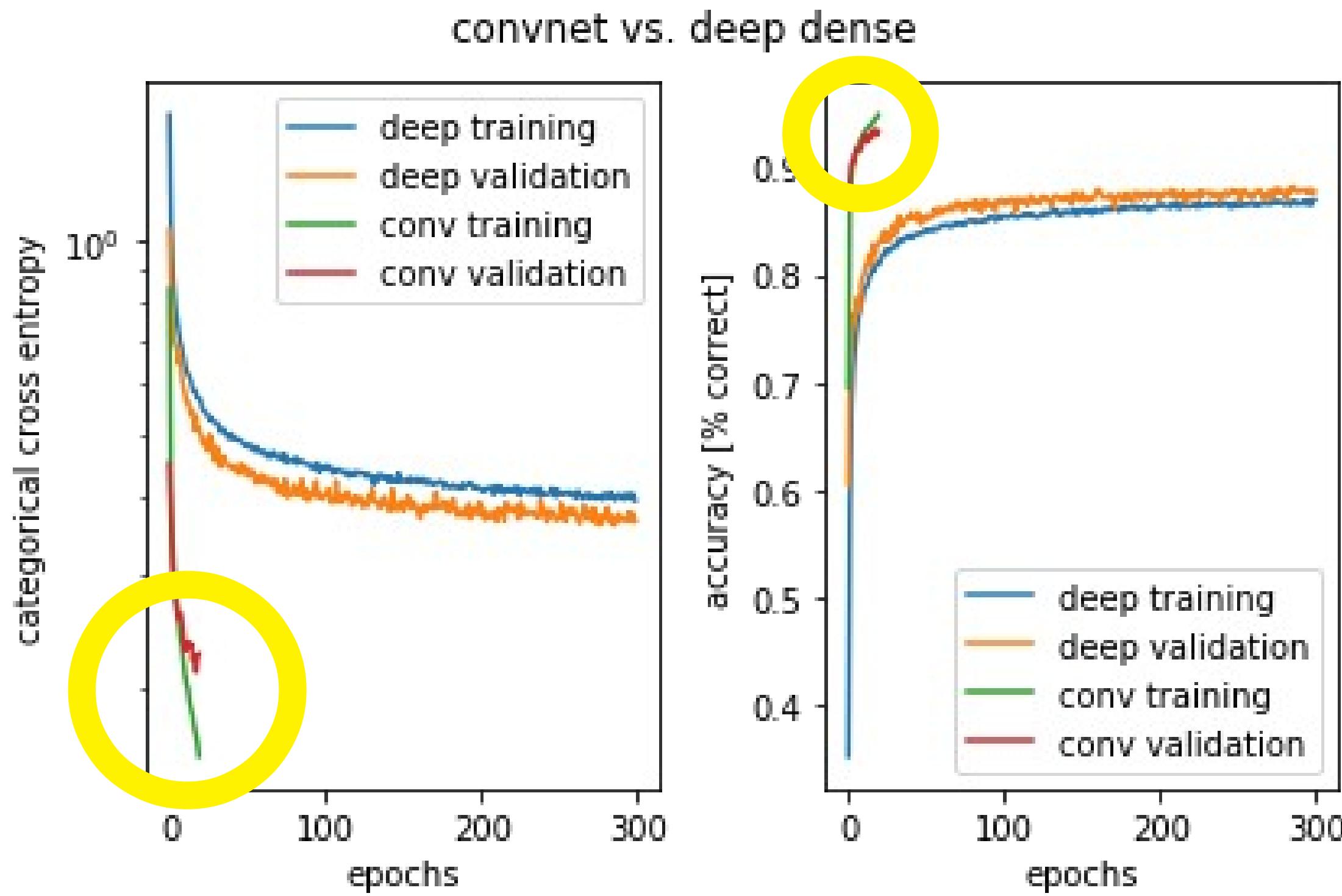
## Further (optional) reading/watching for this lecture:

See references in the slides.

Lectures by Andrew Ng

# Motivation

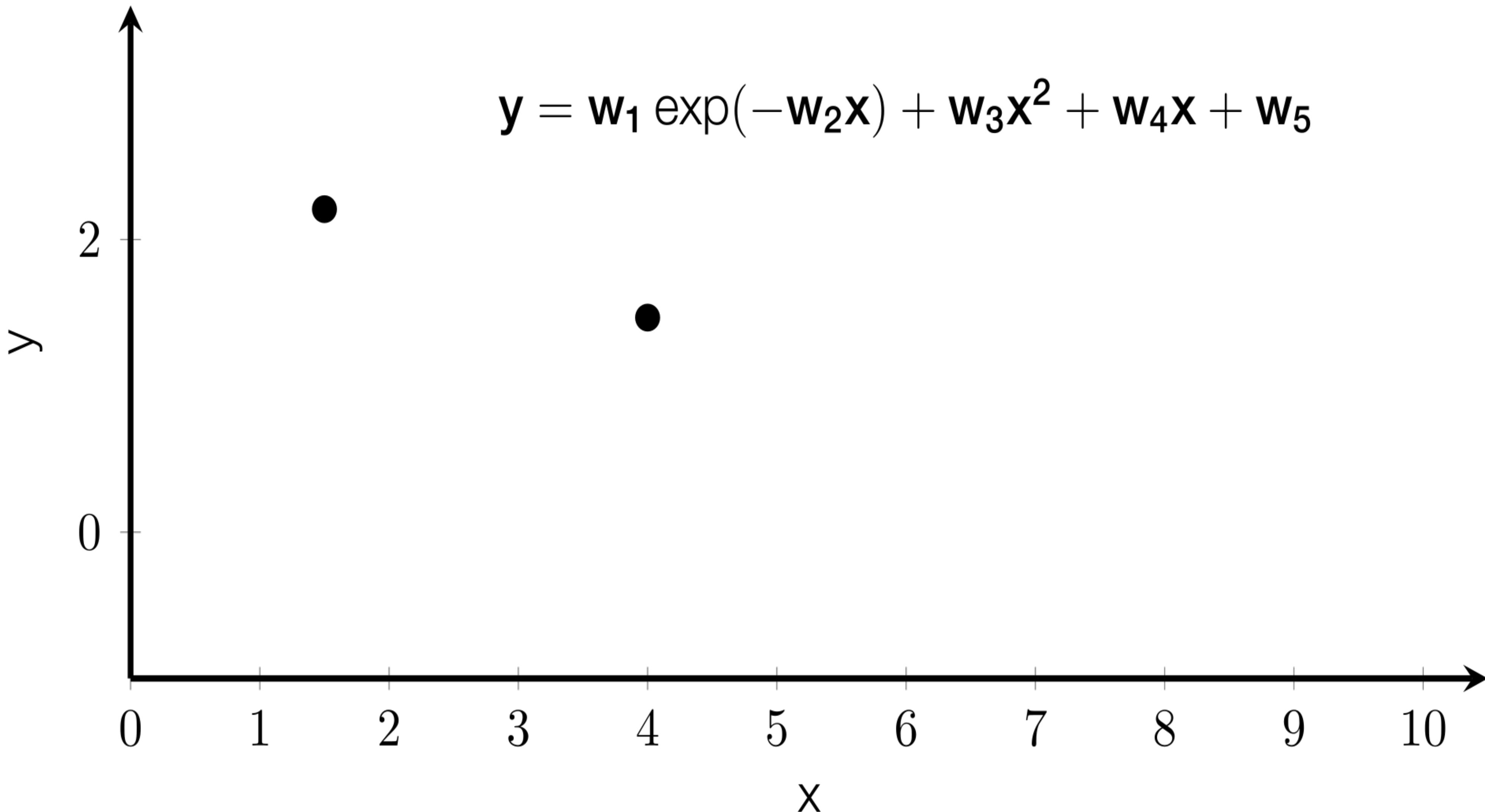
## Miniproject



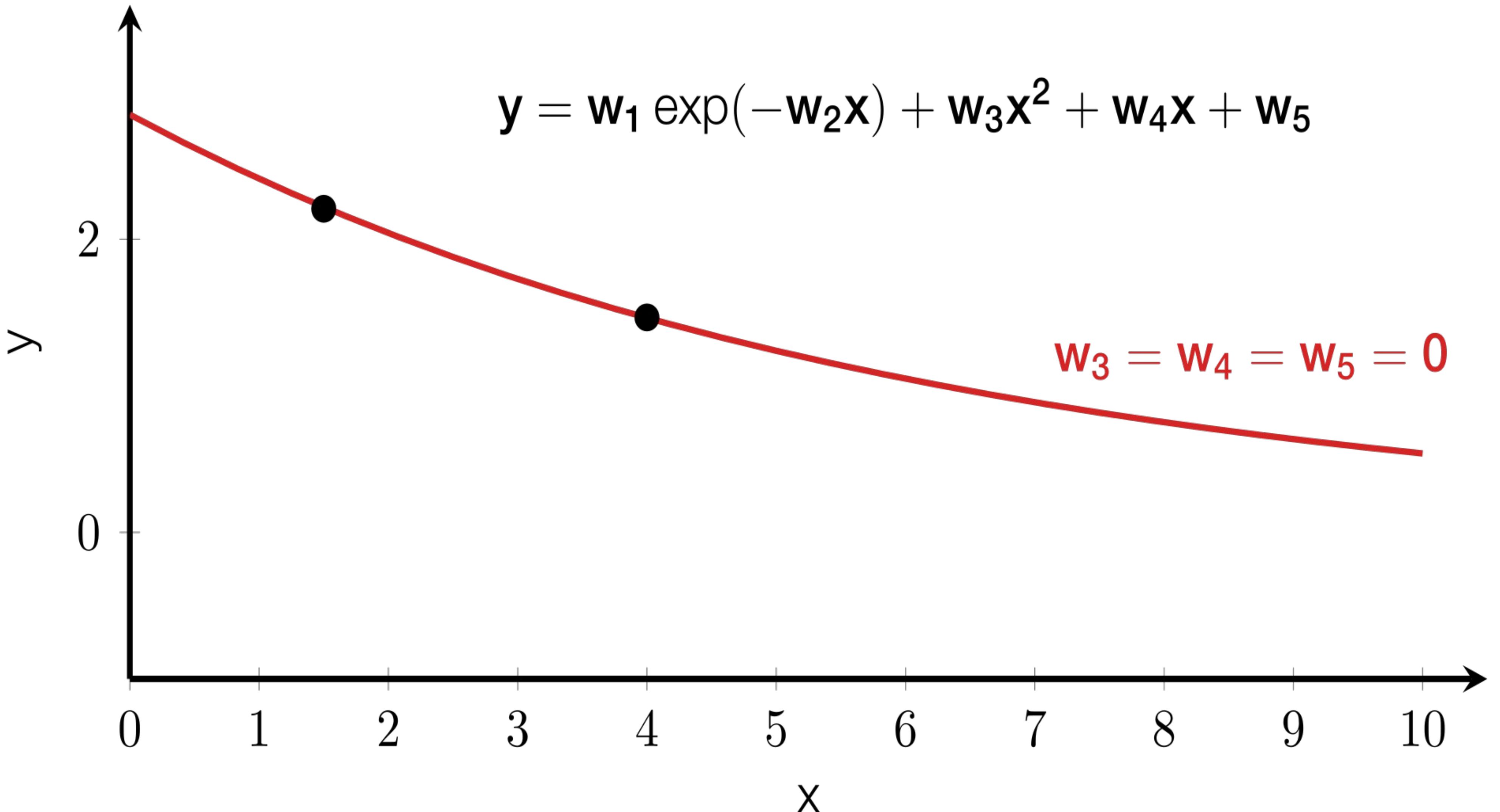
# **No free lunch theorem revisited**

## **Explicit inductive biases, data augmentation and transfer learning**

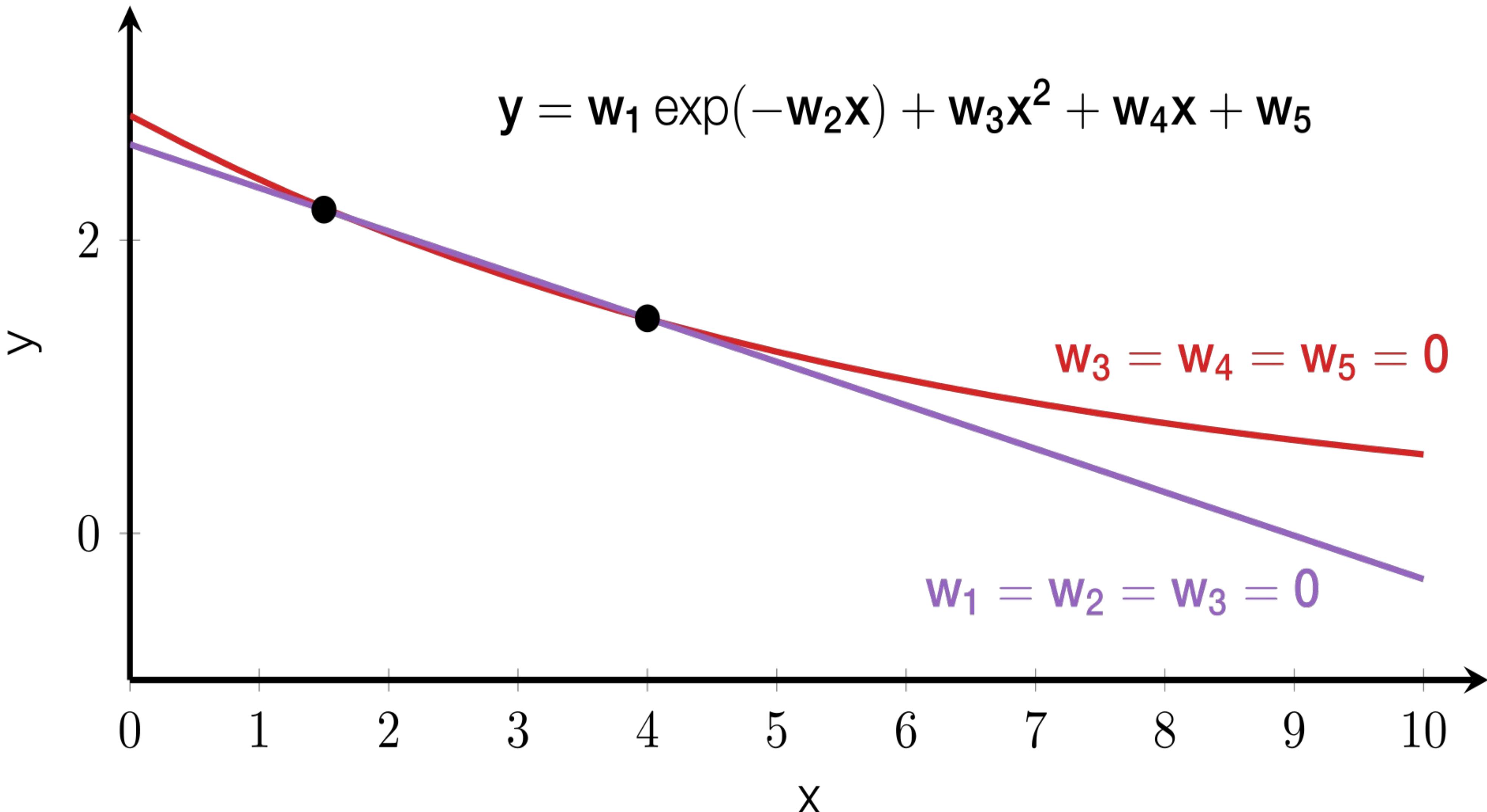
# review: No Free-lunch Theorem (weak inductive bias)



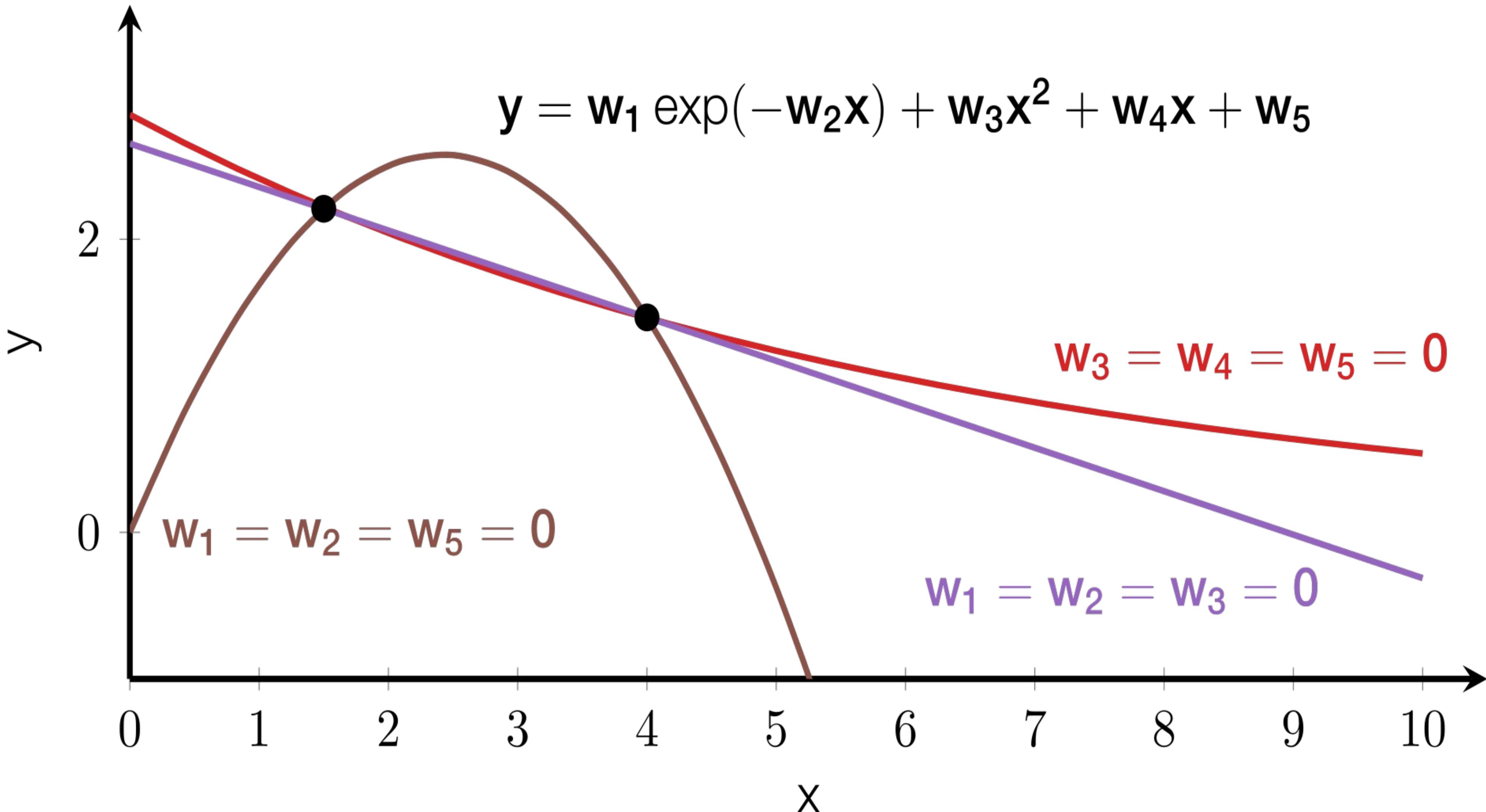
# review: No Free-lunch Theorem (strong inductive bias 1)



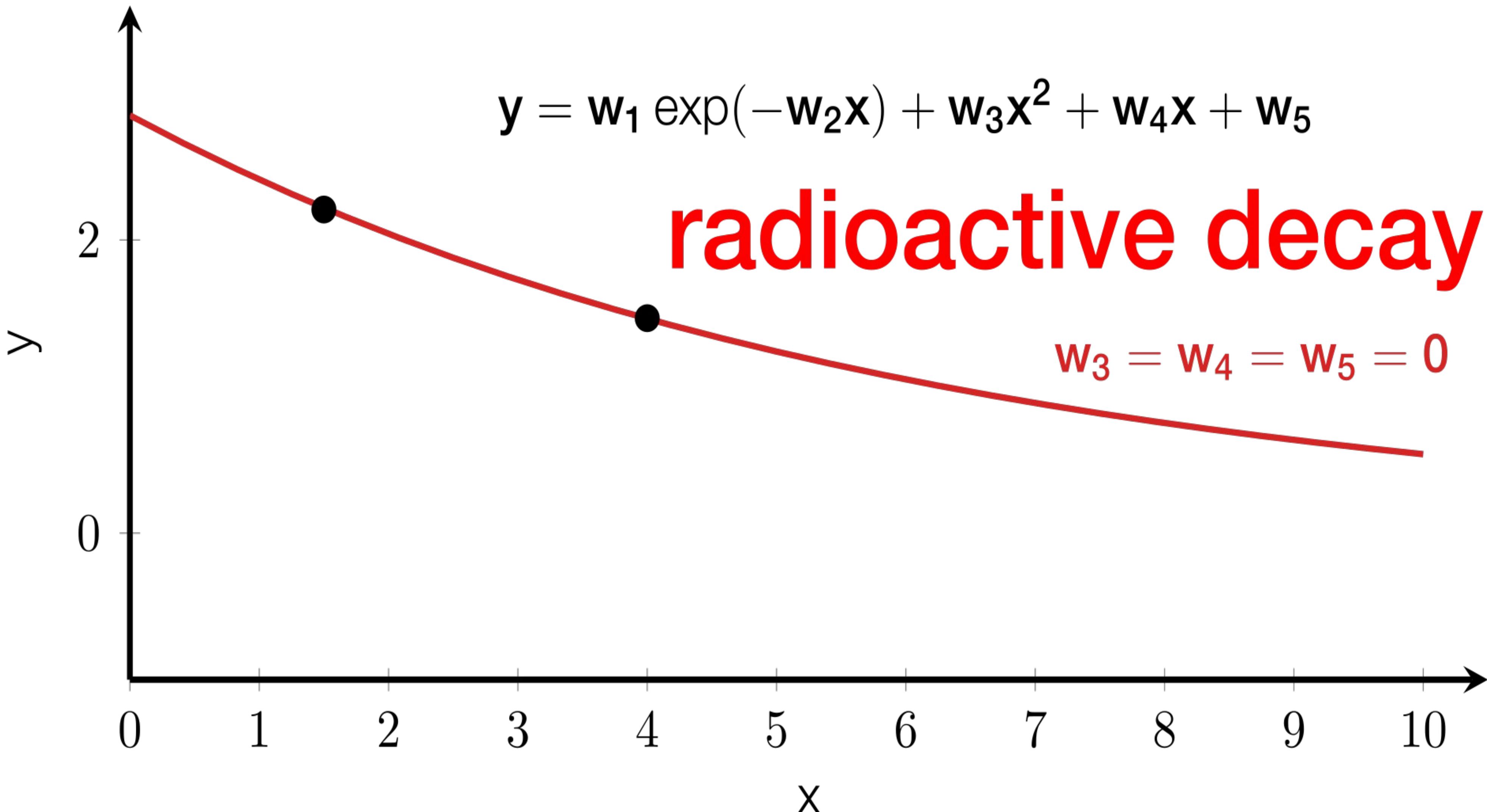
# review: No Free-lunch Theorem (strong inductive bias 2)



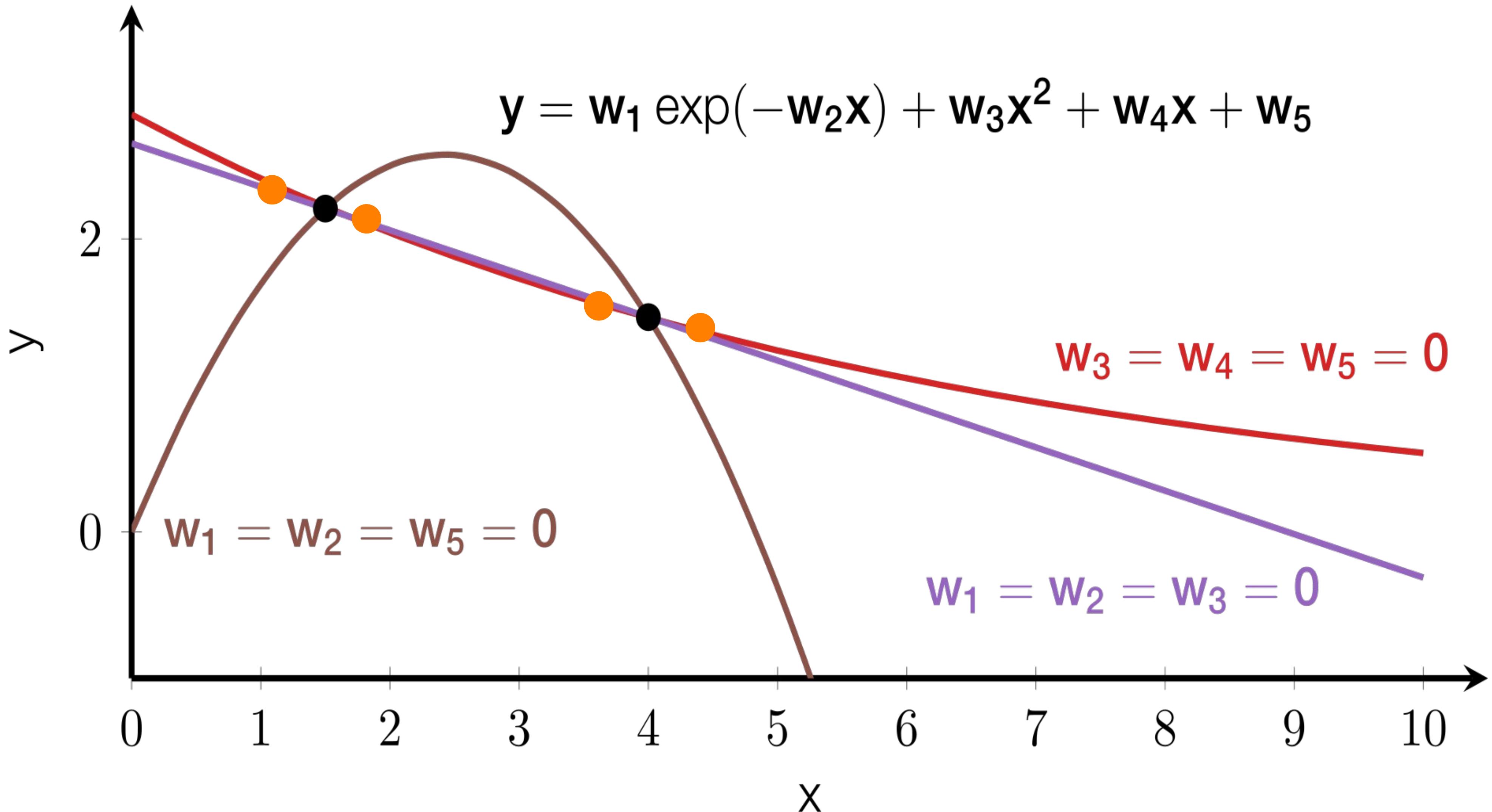
# review: No Free-lunch Theorem (strong inductive bias 3)



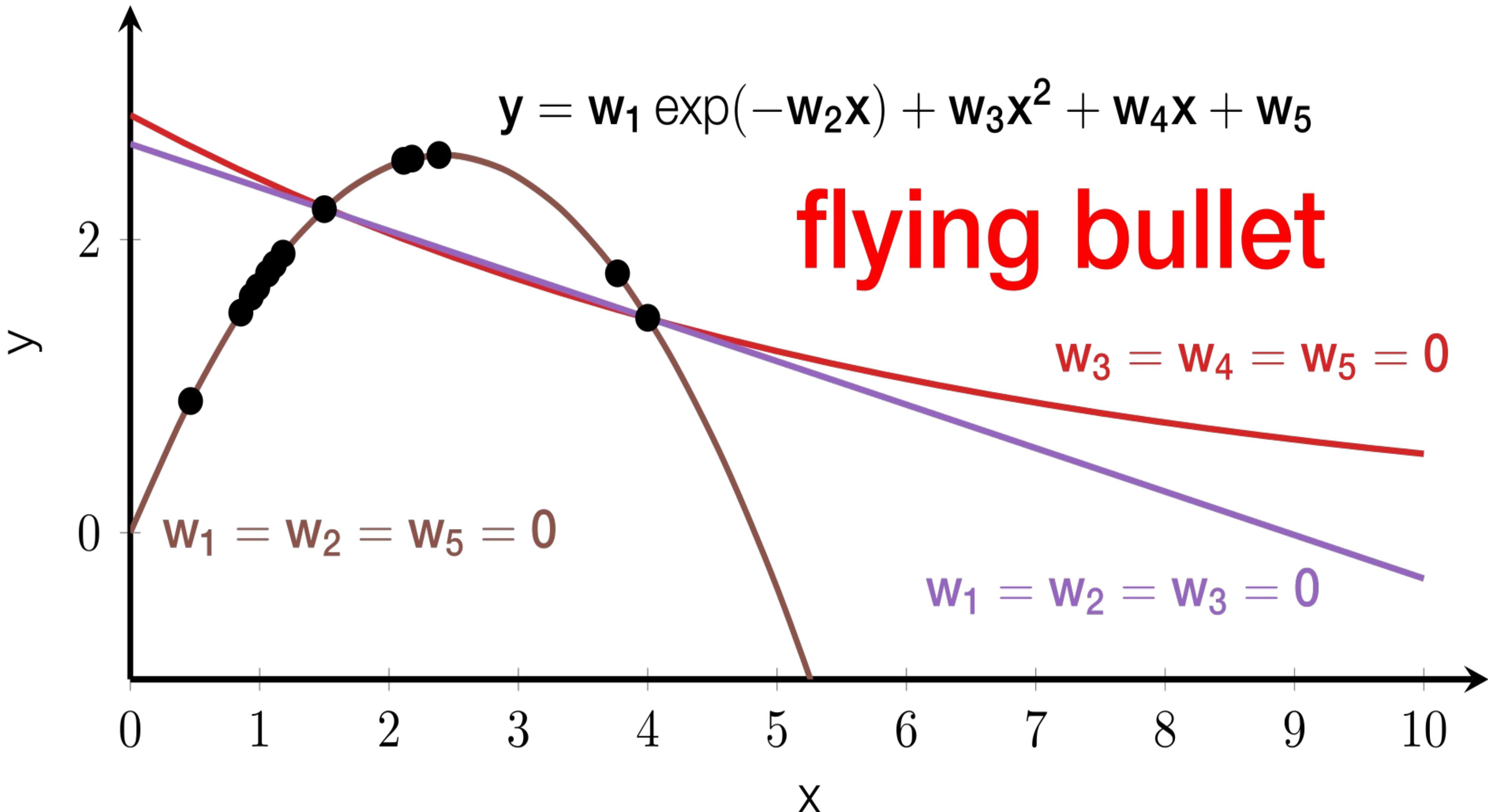
# review: No Free-lunch Theorem (transferred inductive bias)



# review: No Free-lunch Theorem (data augmentation)



# review: No Free-lunch Theorem (the wrong inductive bias)



# Inductive bias

Induction = finding a rule (function) from specific examples

Inductive bias = prior preference of rules (functions)

1) Explicit inductive bias (transfer reasoning)

“For radioactive decay I know that  $w_3 = w_4 = w_5 = 0$ ”

2) Inductive bias through transfer learning

“I train first different models on data from other radioactive elements and choose the best for my current case”

3) Inductive bias through data augmentation

“For radioactive decay neighboring points have similar values”

# Quiz: Inductive bias

- [ ] With a strong inductive bias one can reach a low test error with very little training data.
- [ ] With a strong inductive bias the test error will always be low.
- [ ] Data augmentation is a heuristic method to get more training data.
- [ ] In data augmentation there is an inductive bias in the form of our assumptions about reasonable transformations.
- [ ] Choosing a specific neural network architecture is choosing an explicit inductive bias.

# **Convolutional Layers**

## **An explicit inductive bias for natural images**

# Convolution: one feature

1	0	1
0	1	0
1	0	1

Filter (or Feature)

1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature

$w_{xy}$

$$I_{xy} \quad a_{ij} = b + \sum_{x=1}^3 \sum_{y=1}^3 I_{i+x-1, j+y-1} w_{xy}$$

# Convolution: one feature

1	0	1
0	1	0
1	0	1

Filter (or Feature)

$w_{xy}$

1	1 <sub>x1</sub>	1 <sub>x0</sub>	0 <sub>x1</sub>	0
0	1 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	0
0	0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	1
0	0	1	1	0
0	1	1	0	0

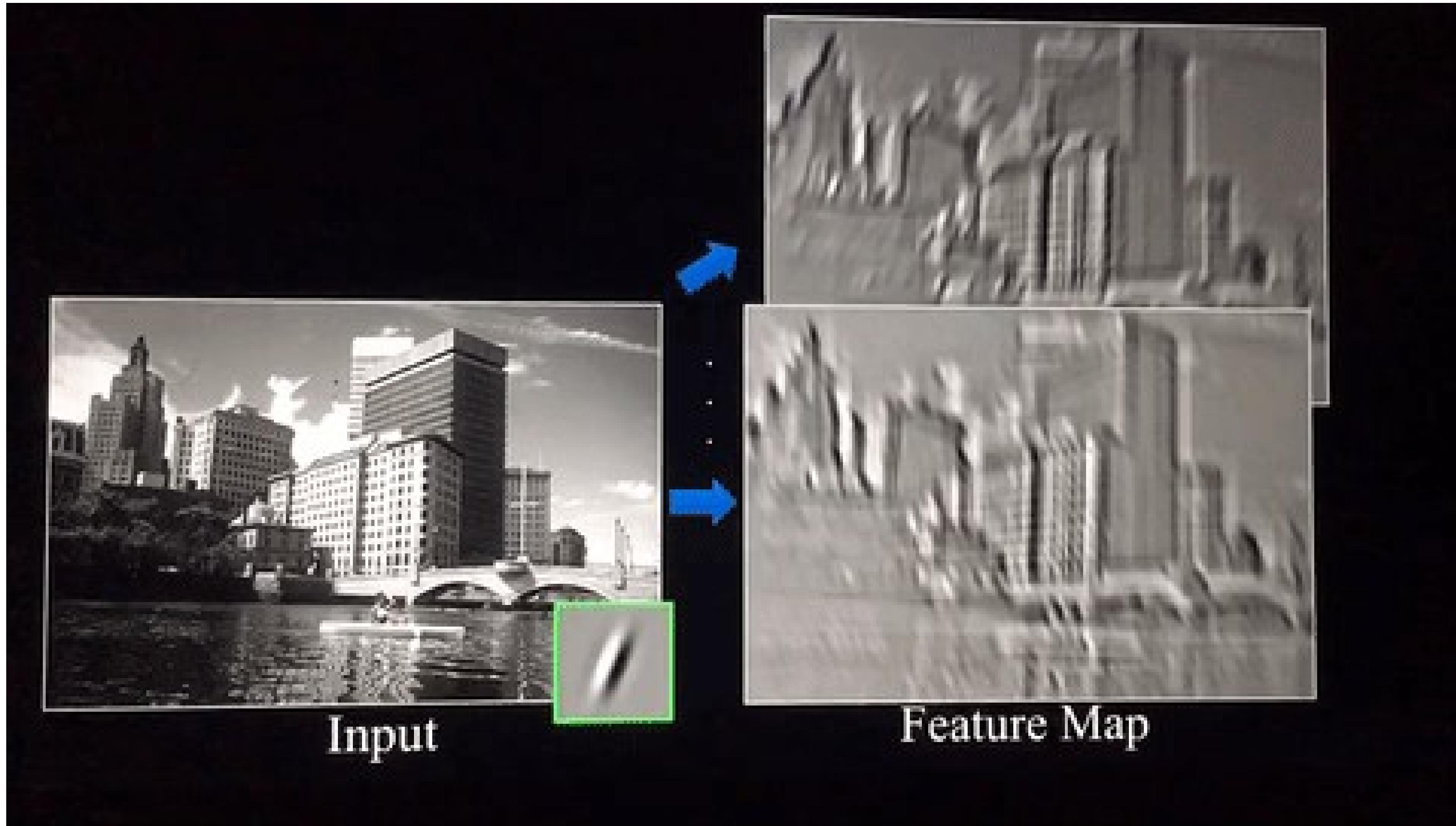
Image

4	3	

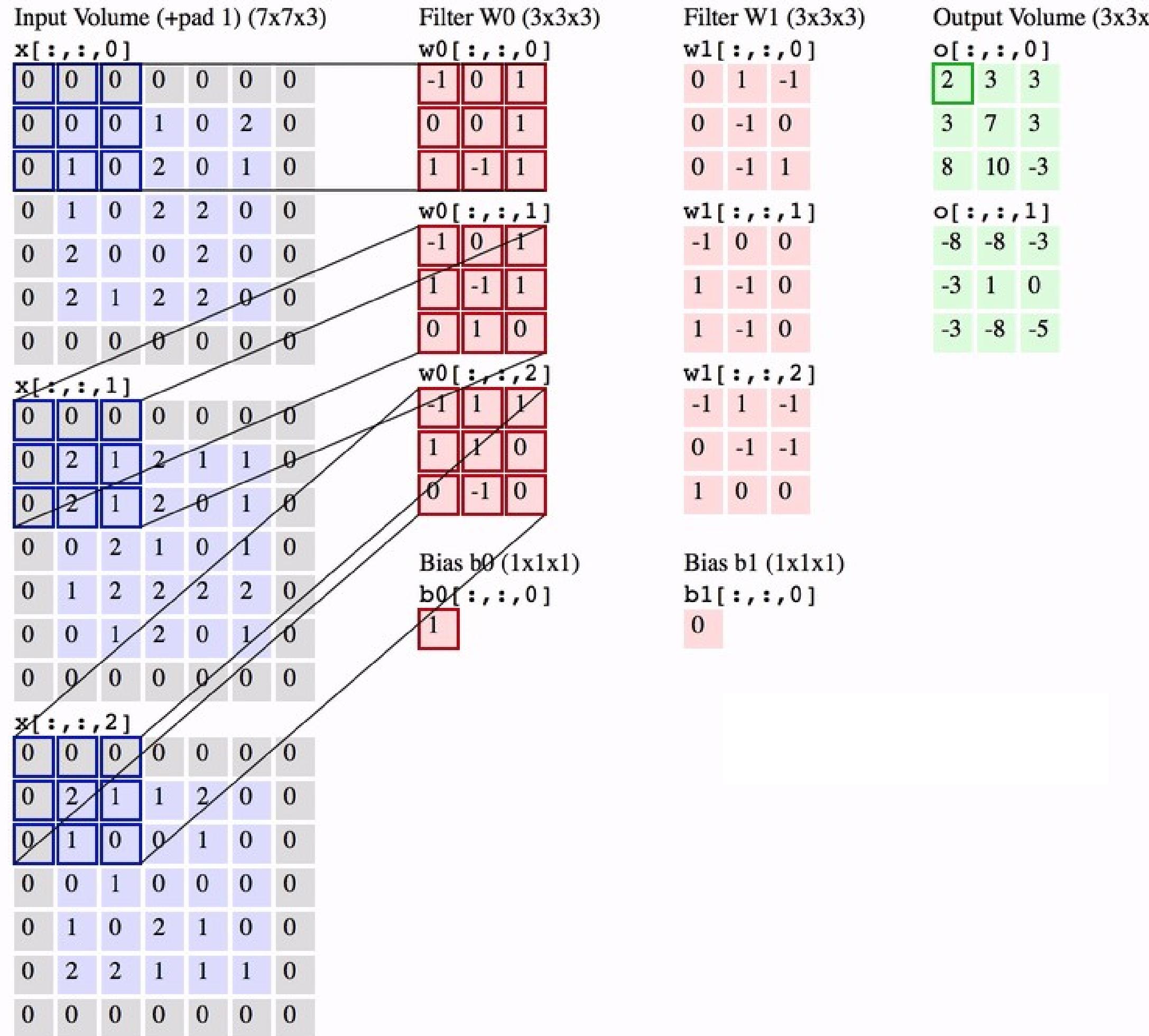
Convolved  
Feature

$$I_{xy} \quad a_{ij} = b + \sum_{x=1}^3 \sum_{y=1}^3 I_{i+x-1, j+y-1} w_{xy}$$

# Convolution: multiple features



# Convolution: stride and padding

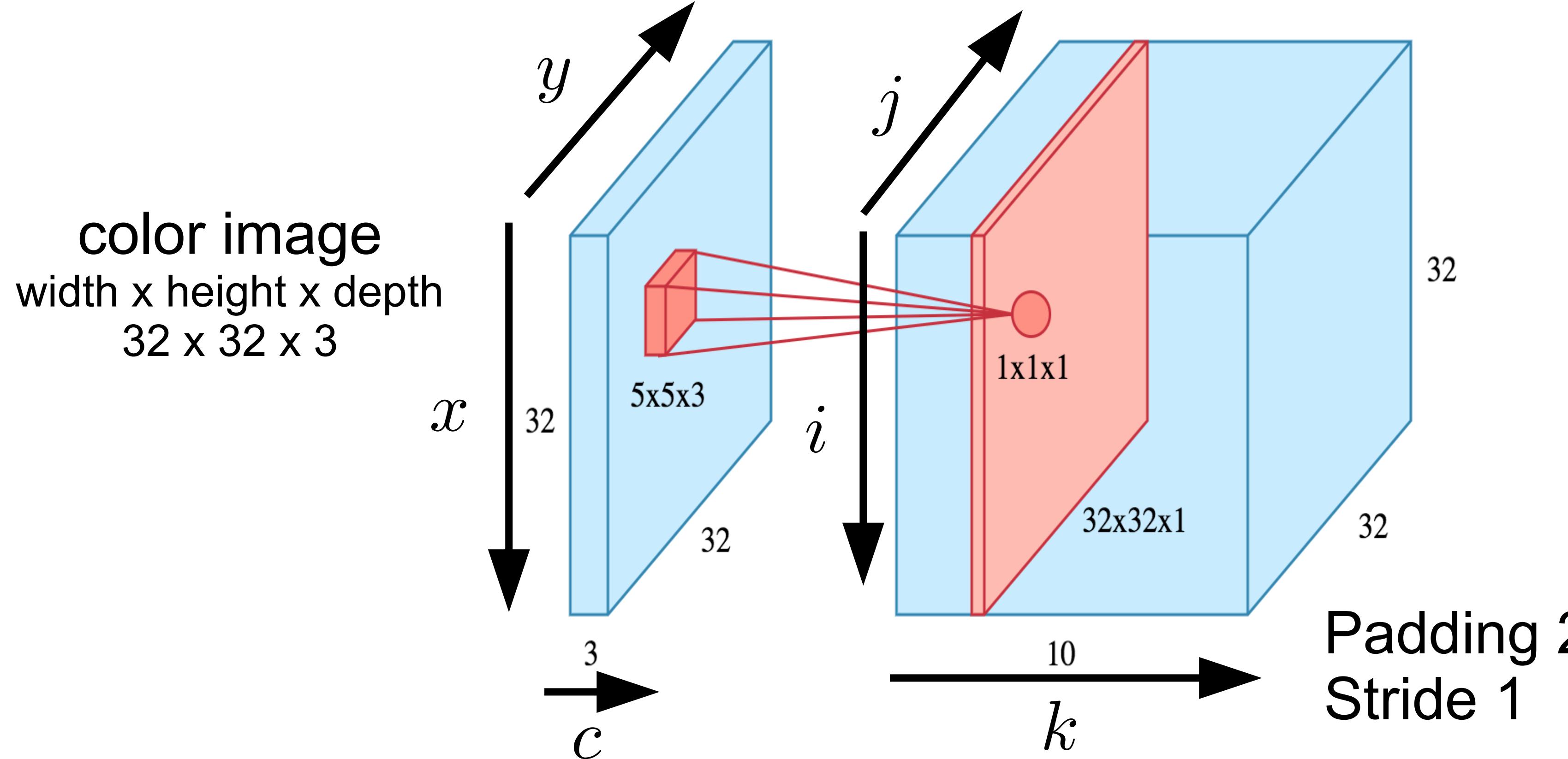


- Stride 2

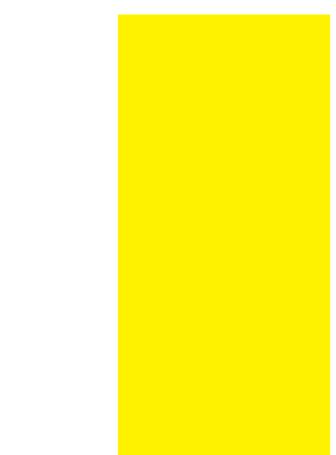
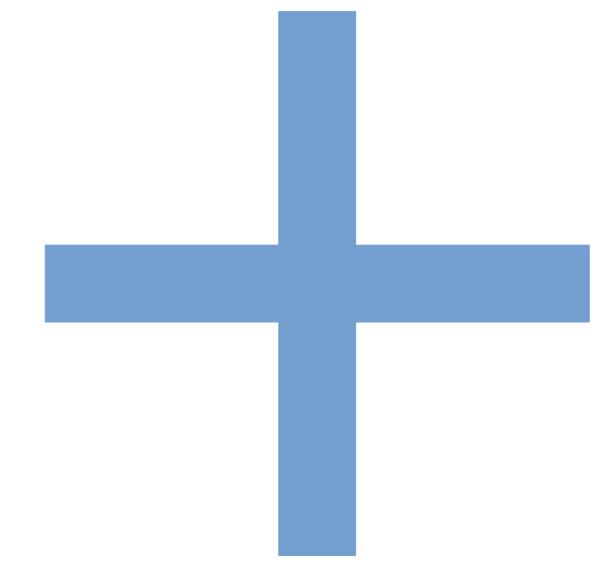
- filters jump 2 pixels at a time as we convolve

- Images are sometimes padded with 0 at the borders

# Convolution: one layer



filters, e.g



$$I_{xyc} \quad a_{ijk} = b_k + \sum_{x=1}^5 \sum_{y=1}^5 \sum_{c=1}^3 I_{i+x-1, j+y-1, c} w_{xyck}$$

# Convolution: computing the volumes

$n \times n \times c$  input volume (e.g. image)     $k$  filters of size  $f \times f \times c$

stride  $s$

padding  $p$

$$\left( \frac{n + 2p - f}{s} + 1 \right) \times \left( \frac{n + 2p - f}{s} + 1 \right) \times k$$

Proof: exercise

# Inductive bias of a convolutional layer

## 1) Equivariance to translation

$f(x)$  is equivariant to  $g$  if  $f(g(x)) = g(f(x))$

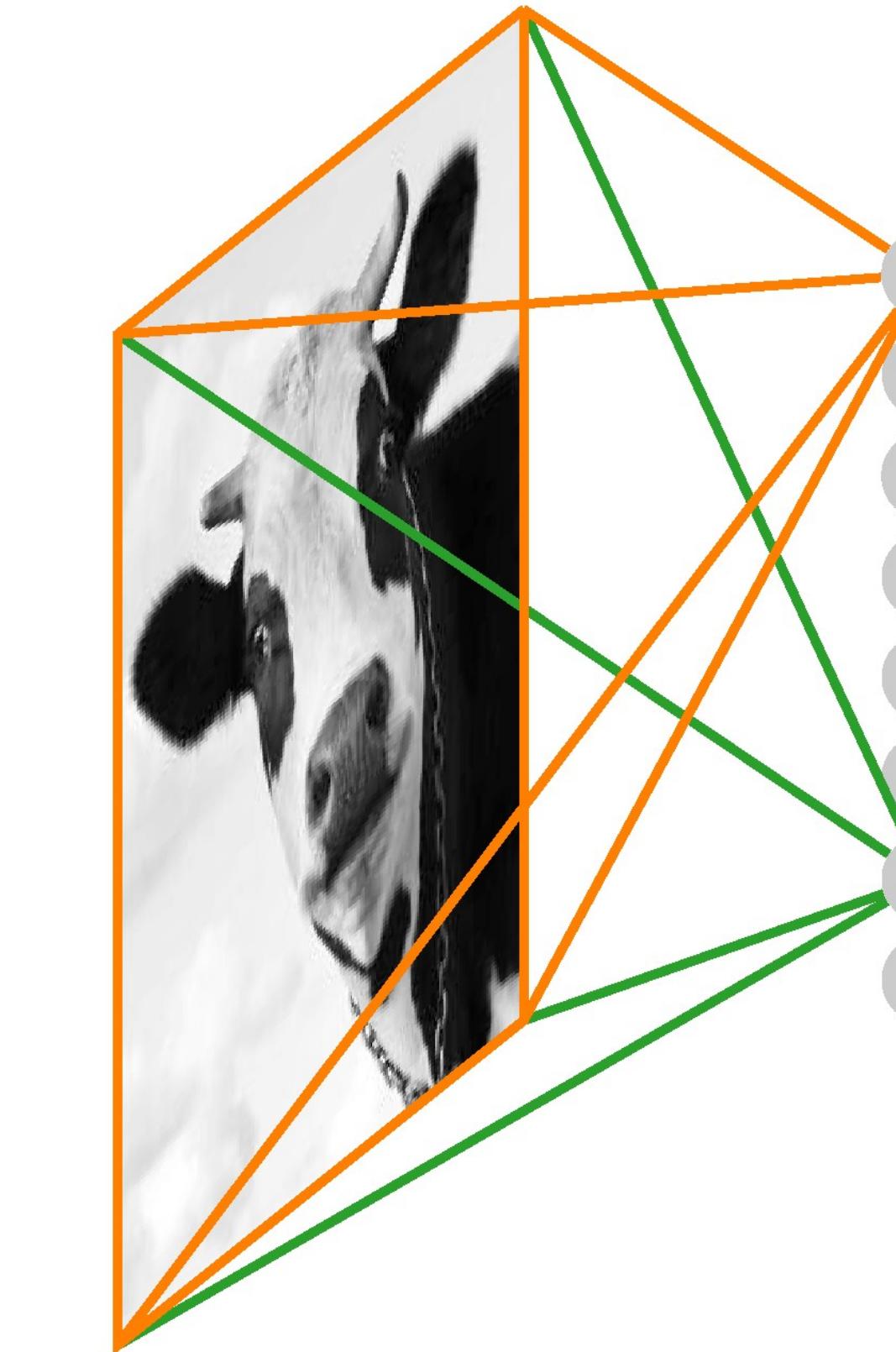
“A feature detector (such as a vertical edge detector) that’s useful in one part of the image is probably useful in another part of the image.”

## 2) Only local interactions matter

A conv layer is like a standard dense layer with

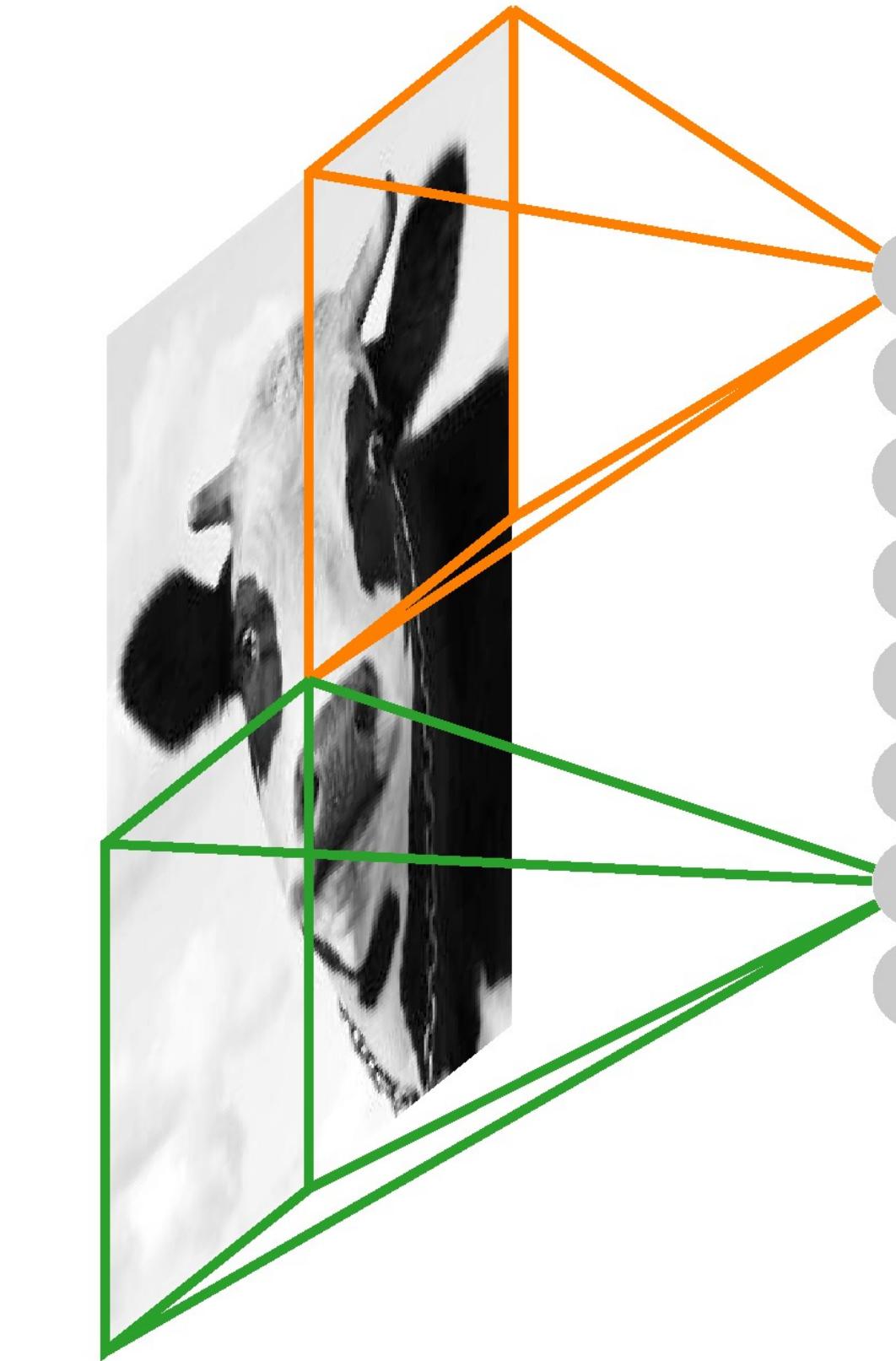
- 1) many neurons having the same weights
- 2) many weights zero (except in a small neighborhood)

# From dense to convolutional layers

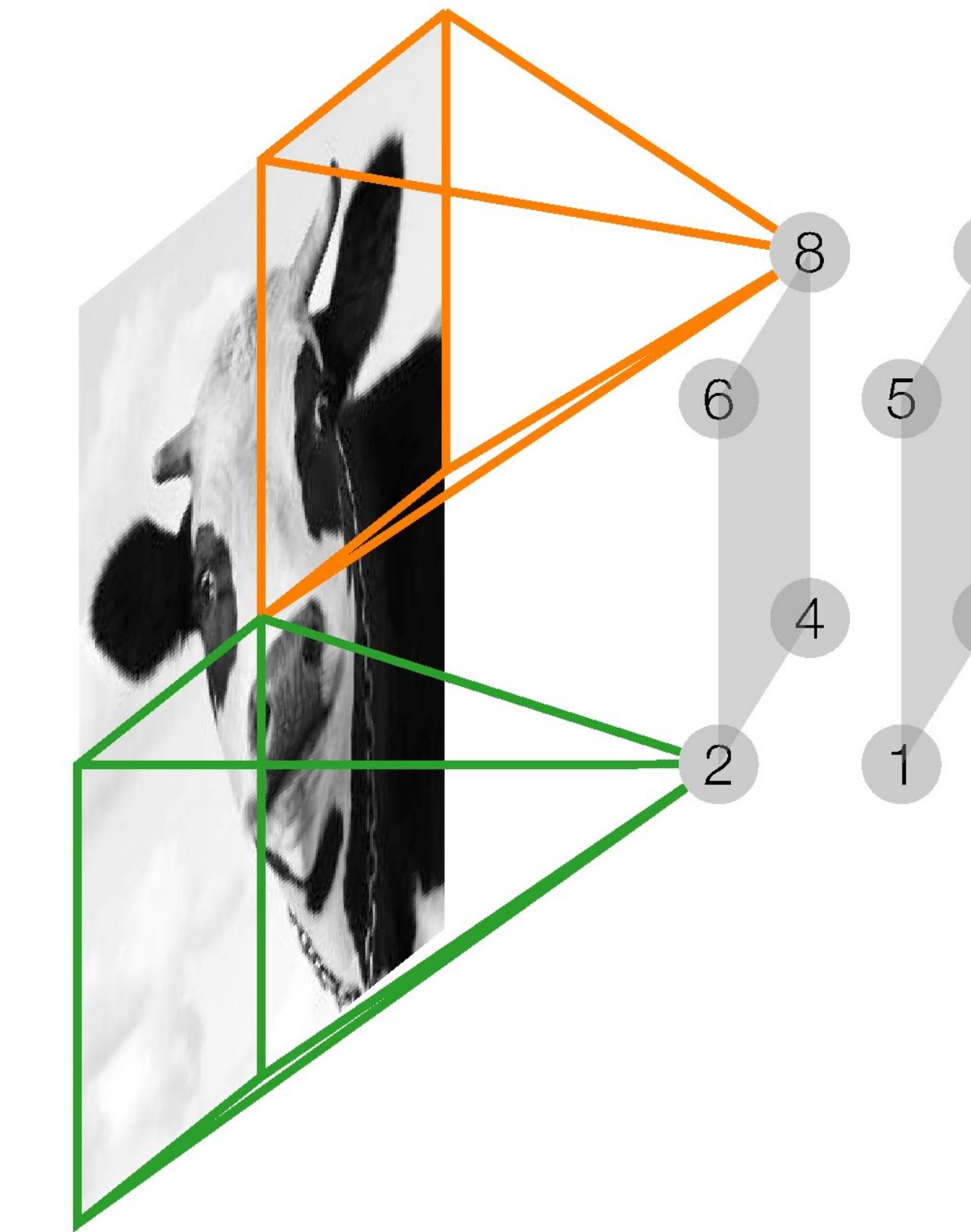


dense (all-to-all) connectivity

# From dense to convolutional layers

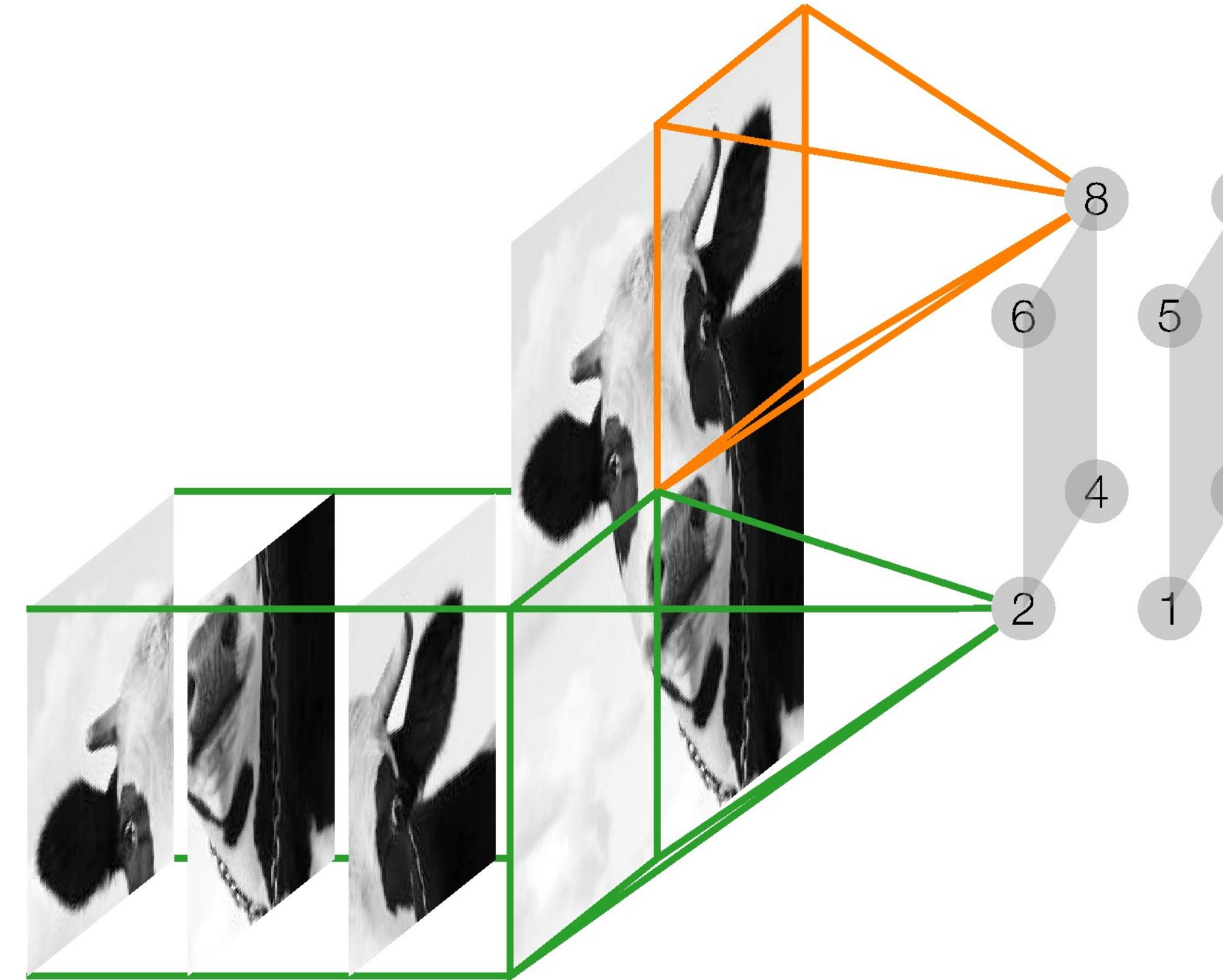


# From dense to convolutional layers



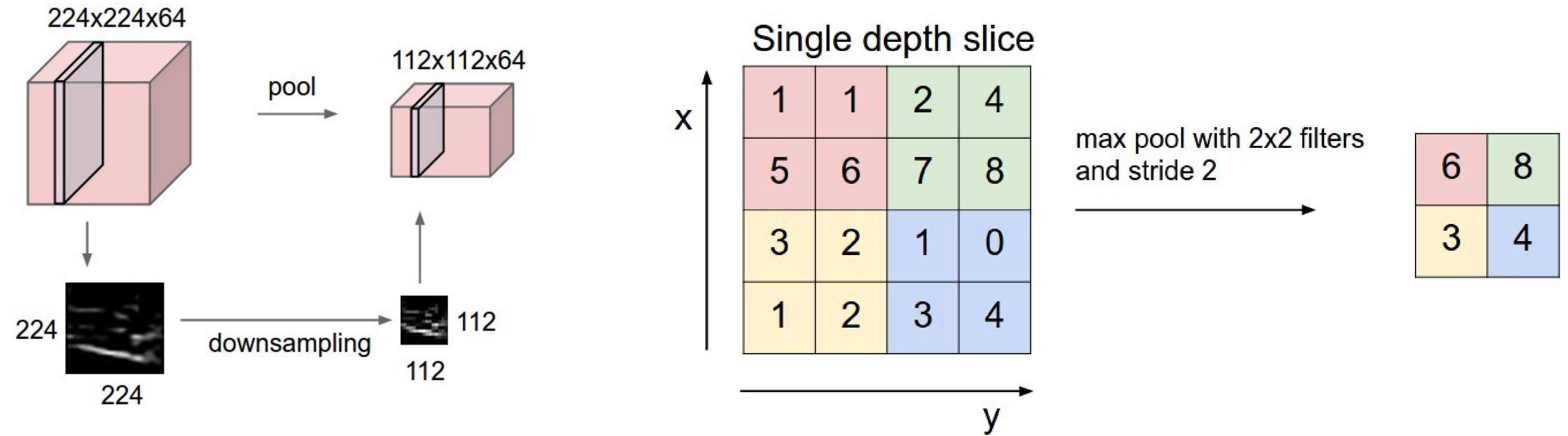
patchy connectivity

# From dense to convolutional layers



patchy connectivity & data augmentation

# Max Pooling



## Inductive bias of max pooling

Invariance to small translations with some probability

$f(x)$  is invariant to  $g$  if  $f(g(x)) = f(x)$

“quite a few activations remain the same when wiggling the image”

# Critic of max-pooling

“The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster.” Geoff Hinton

[https://www.reddit.com/r/MachineLearning/comments/2lmo0l/ama\\_geoffrey\\_hinton/clyj4jv/](https://www.reddit.com/r/MachineLearning/comments/2lmo0l/ama_geoffrey_hinton/clyj4jv/)

<https://openreview.net/pdf?id=HJWLfGWRb>

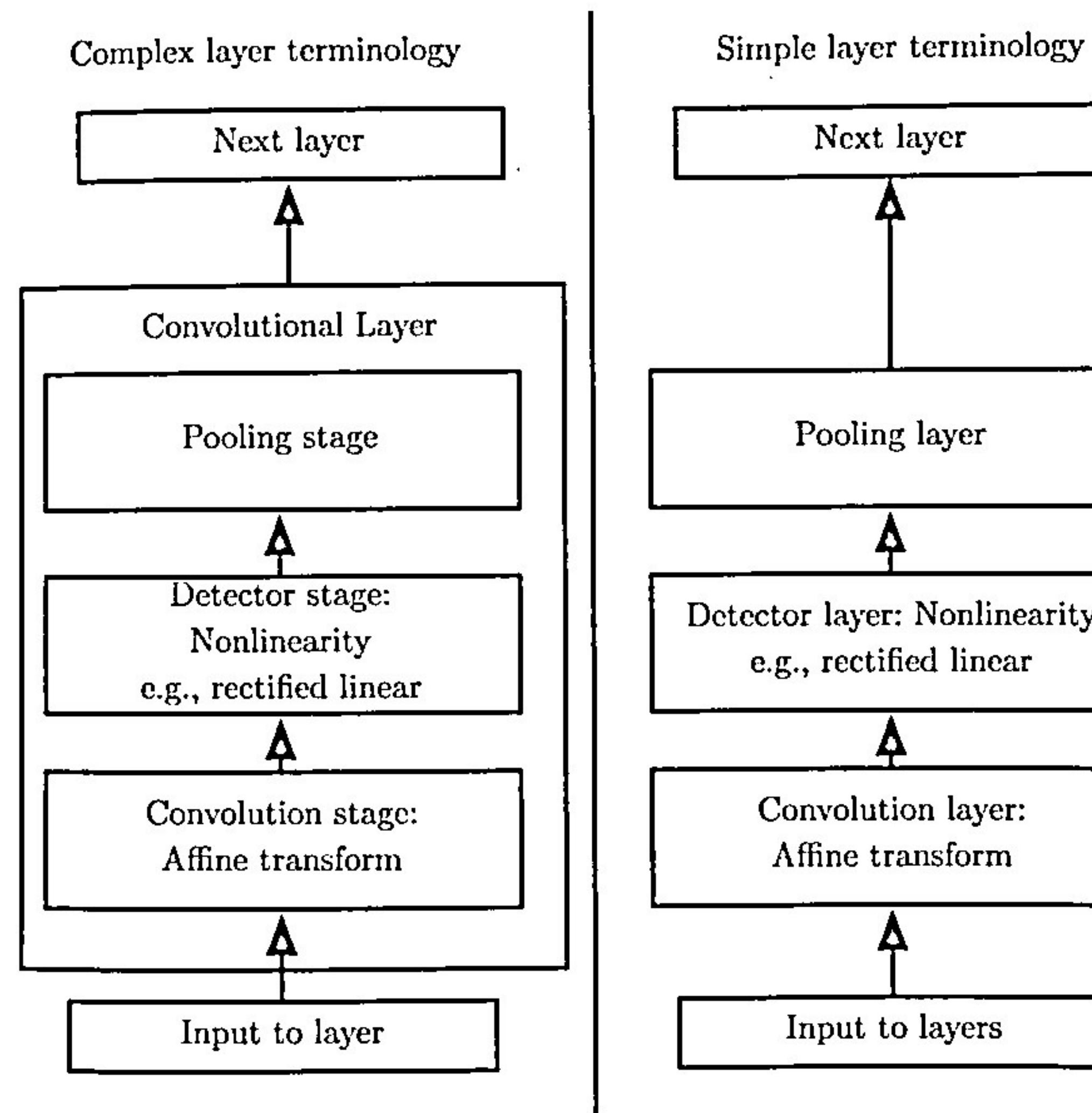
<https://www.youtube.com/watch?v=pPN8d0E3900>

“We find that max-pooling can simply be replaced by a convolutional layer with increased stride without loss in accuracy on several image recognition benchmarks.”

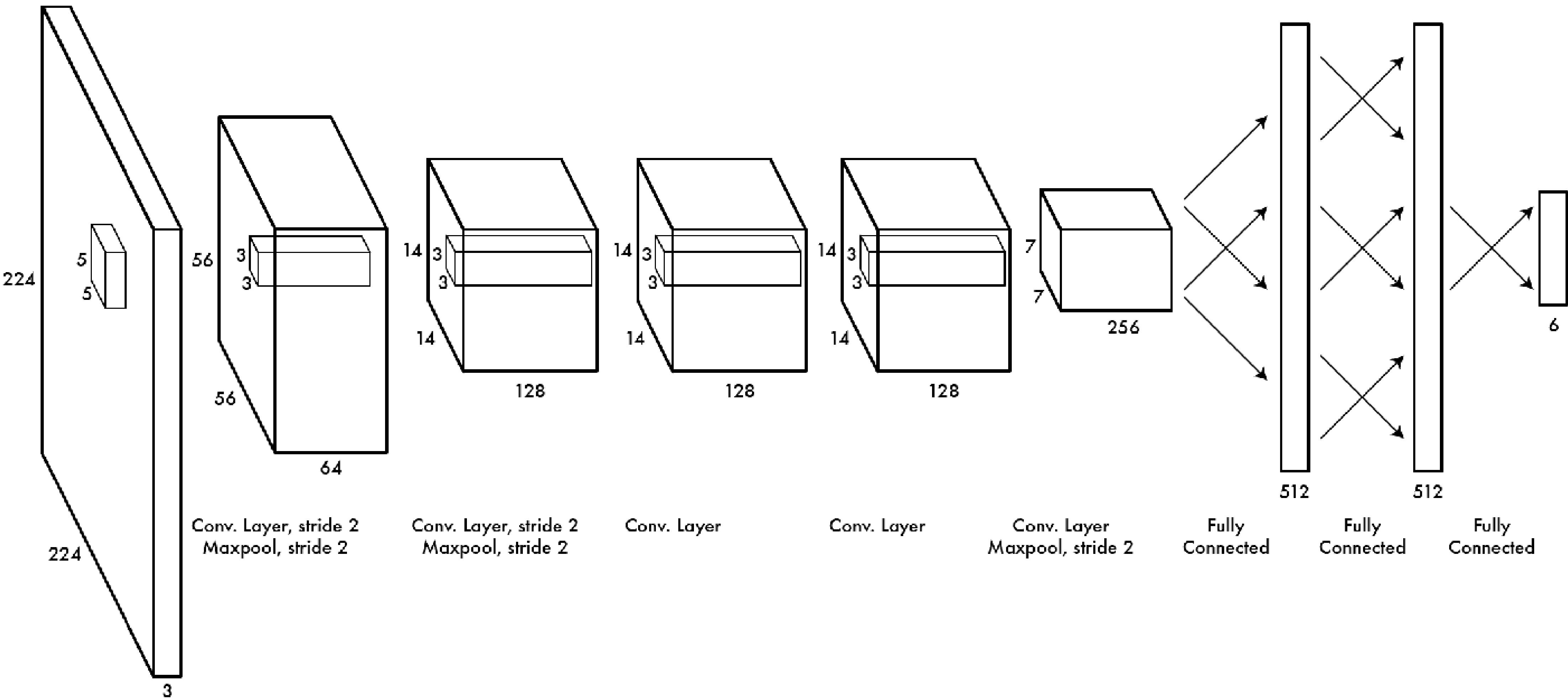
<https://arxiv.org/abs/1412.6806>

See also <https://openreview.net/forum?id=HJeuOiRqKQ>

# Components of a typical conv-net layer



# Convolutional network



# Quiz:

- [ ] The number of weights in a convolutional layer with 10 filters of size  $5 \times 5 \times 3$  is 750 (excluding biases).
- [ ] The number of weights in a convolutional layer does not depend on the size of the input layer.
- [ ] The number of outputs of a convolutional layer with 10 filters of size  $5 \times 5 \times 3$  is 250.
- [ ] Convolutional layers are also equivariant under rotations of the image.
- [ ] Max pooling layers are also likely to be invariant under rotations of the image.
- [ ] You have a dataset of centered portraits (passport photos) on white background. The equivariance property of ConvNets is a good inductive bias for this dataset.

# **Modern ConvNets and ImageNet Competitions**

# Imagenet challenge 2012

example images of australian terriers



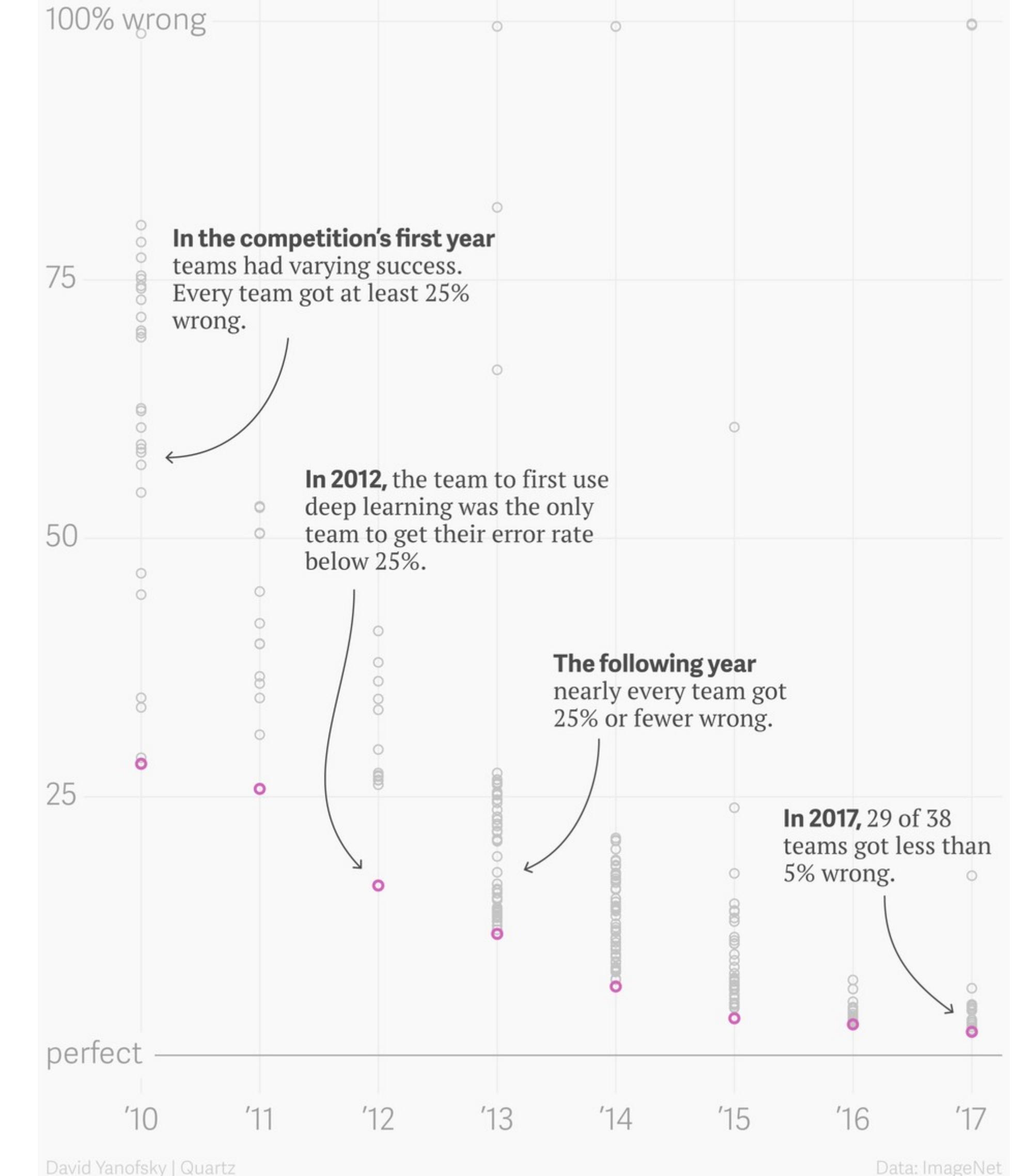
- 1.2 million images training set
- 1000 categories
- 50 thousand validation set
- 100 thousand test set

# Imagenet challenges



<http://www.image-net.org/>  
<https://www.kaggle.com/c/imagenet-object-detection-challenge>

ImageNet Large Scale Visual Recognition Challenge results

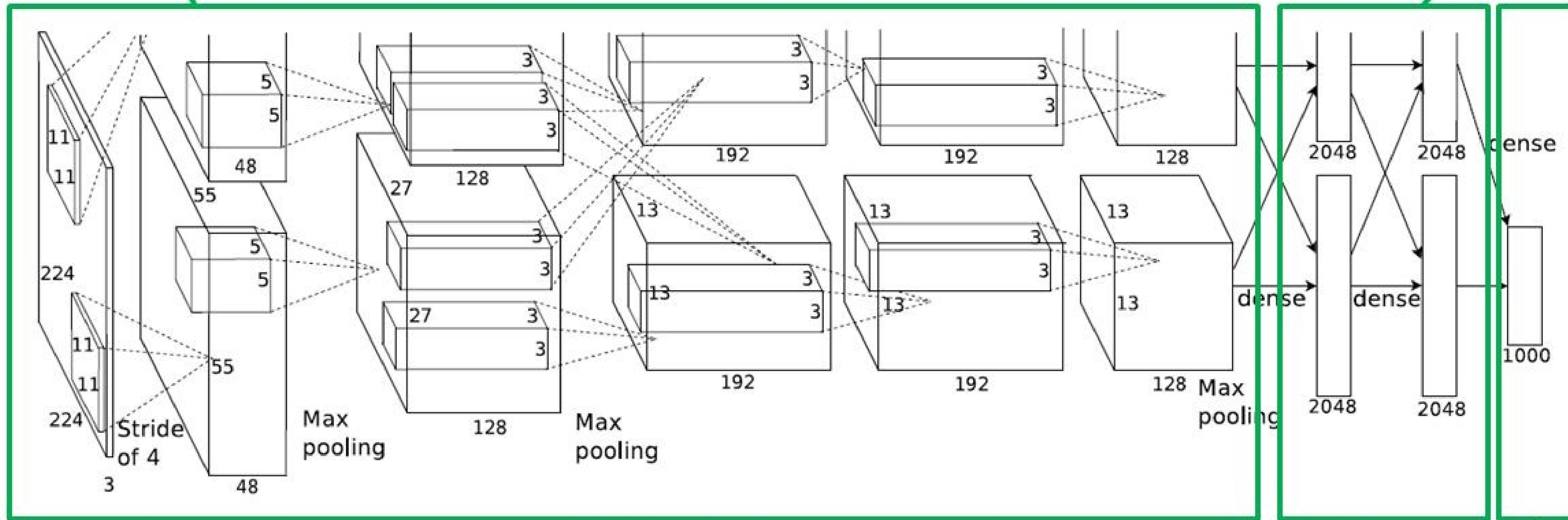


# AlexNet: winner of Imagenet challenge 2012

Convolutional Layers.  
Over 90% of  
computation time.

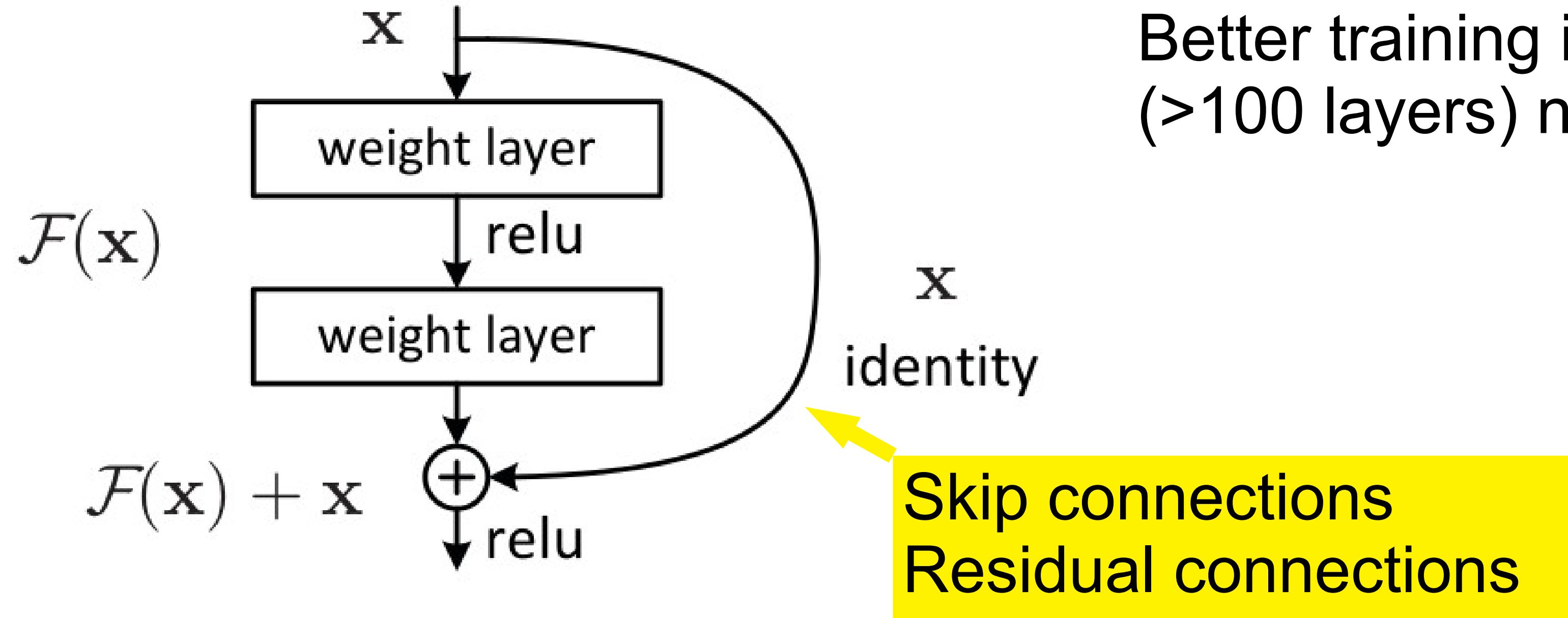
## AlexNet

Fully Connected Layers.  
They can extract local  
and global features.



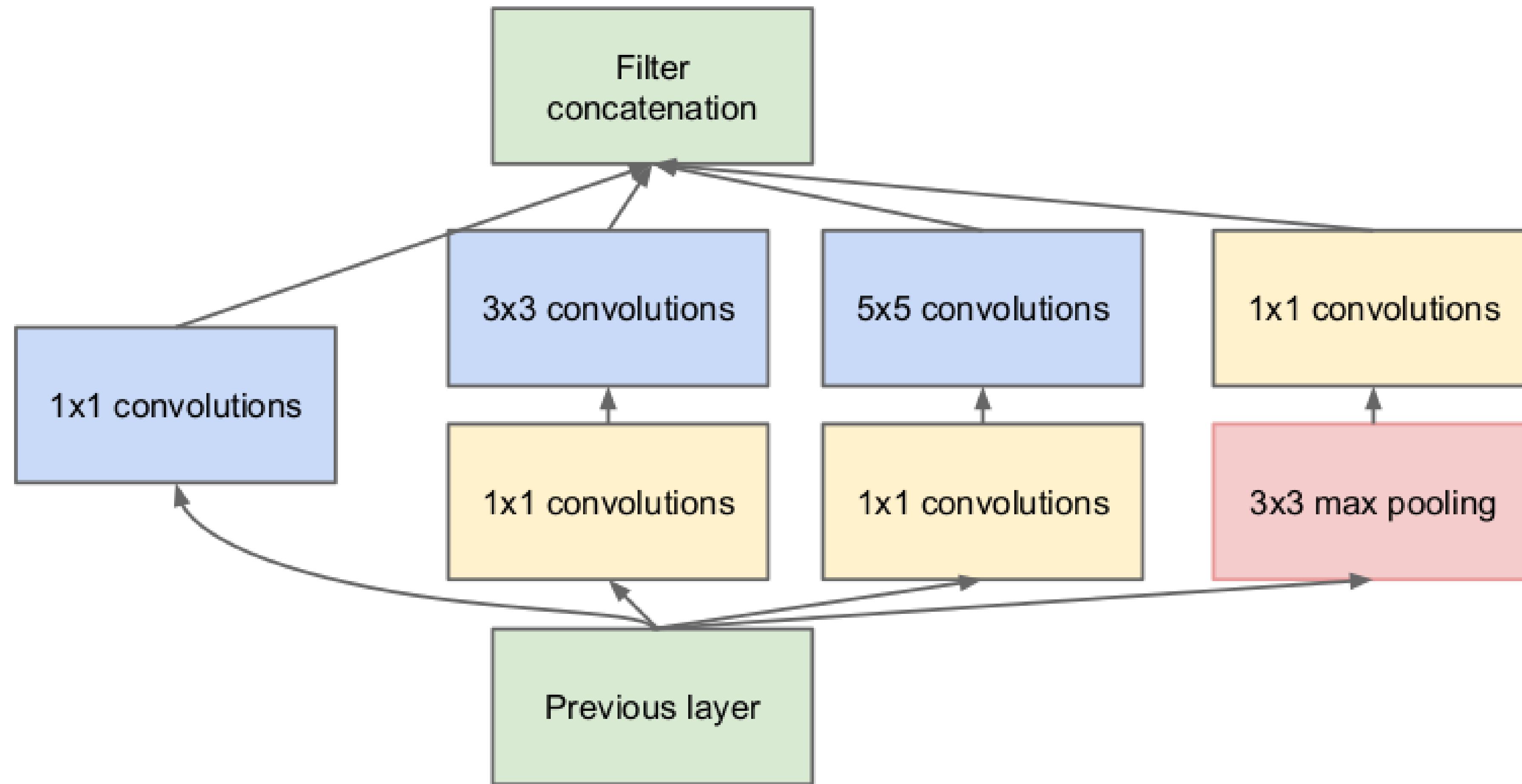
[Krizhevsky 2012]

# ResNet

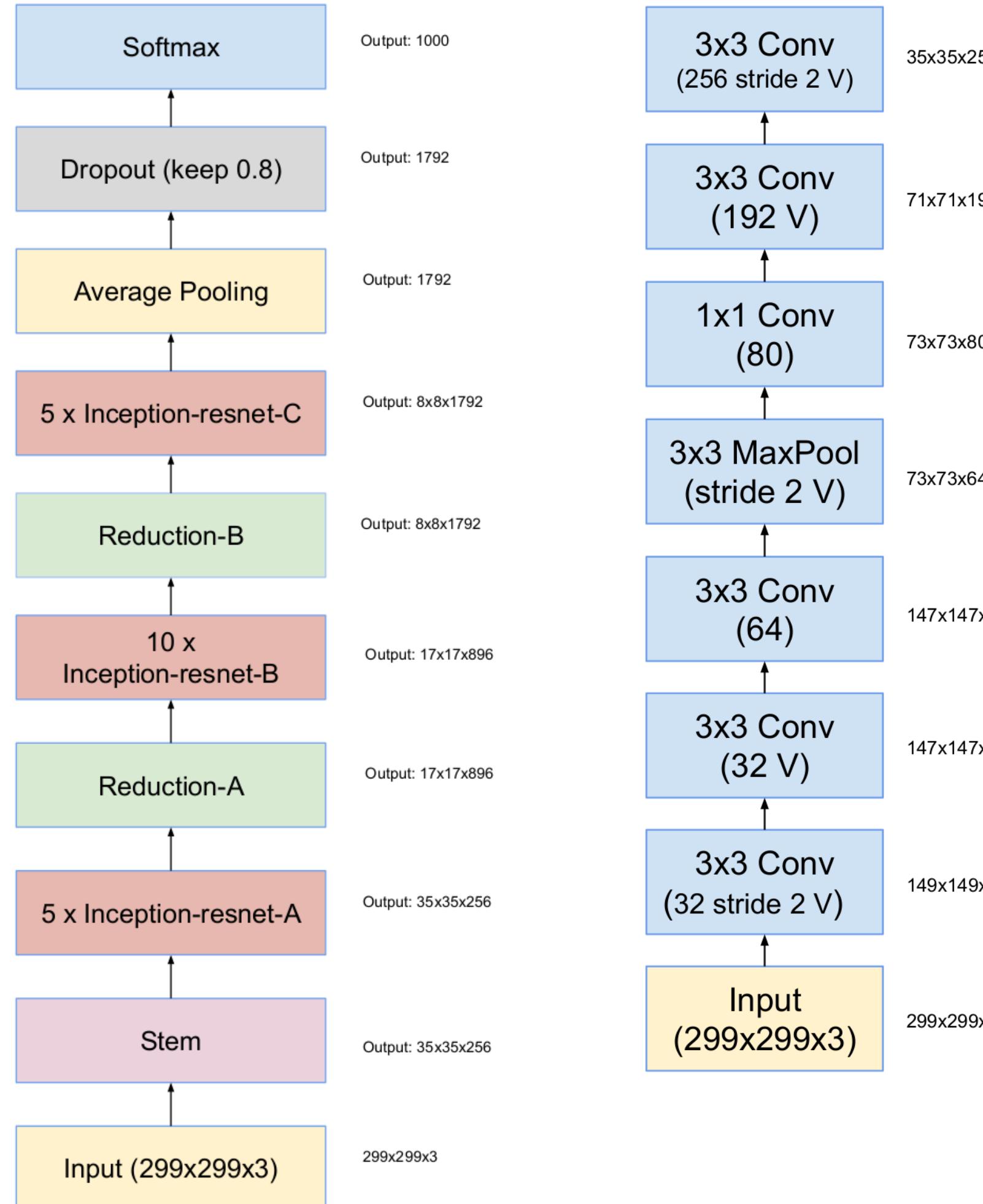


Better training in very deep  
(>100 layers) networks.

# Inception module



# Inception-ResNet-v2



Full model

Stem

Special Layers

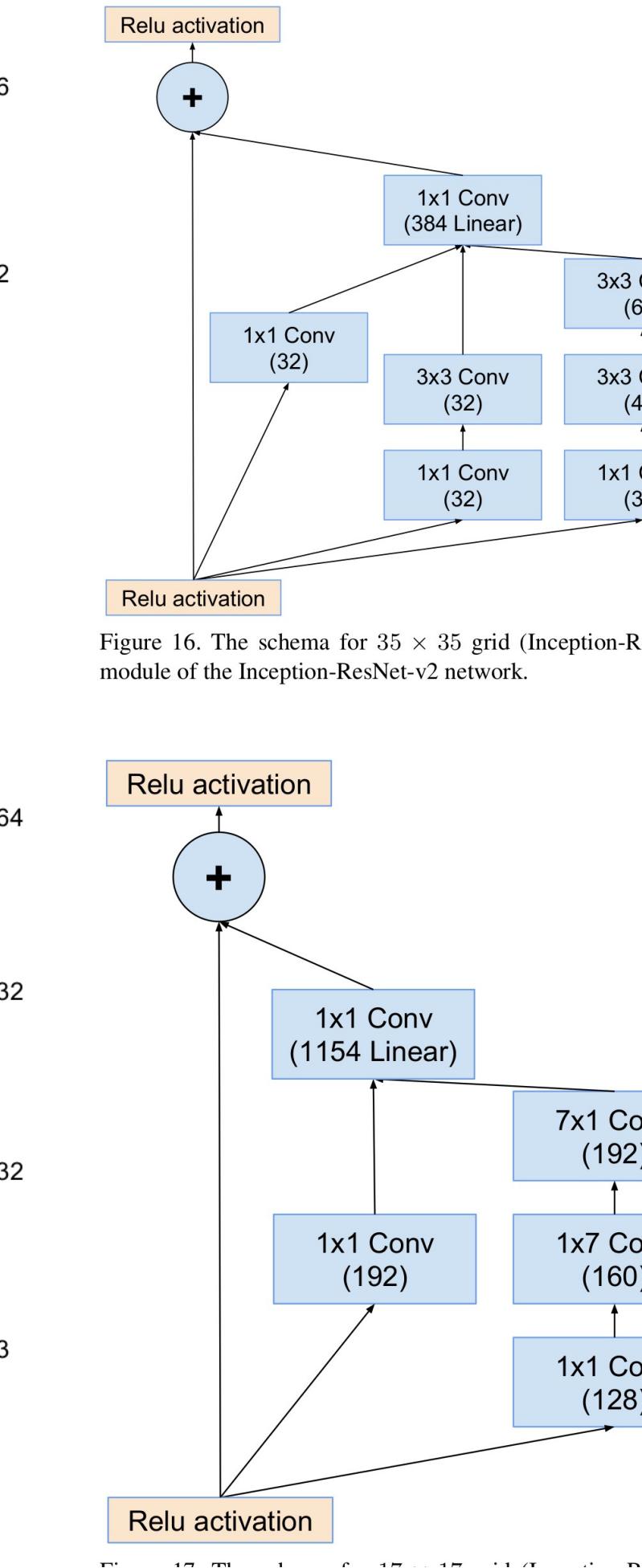


Figure 16. The schema for  $35 \times 35$  grid (Inception-ResNet-A) module of the Inception-ResNet-v2 network.

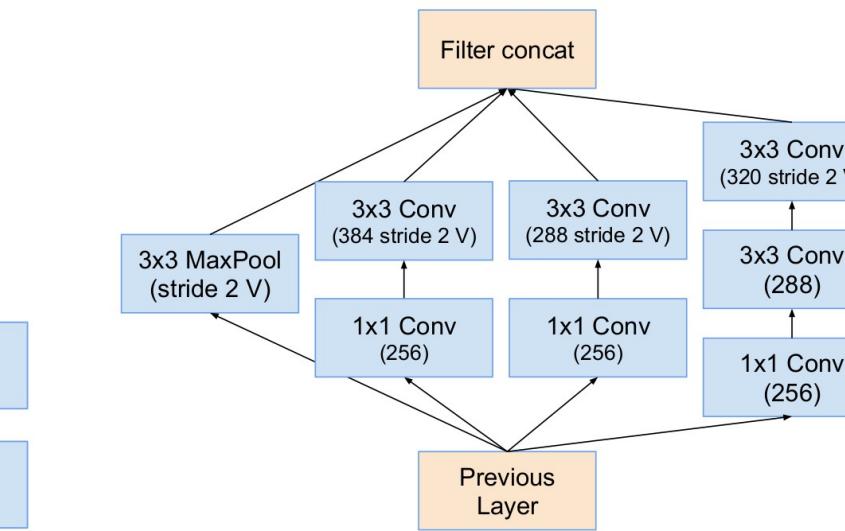


Figure 18. The schema for  $17 \times 17$  to  $8 \times 8$  grid-reduction module. Reduction-B module used by the wider Inception-ResNet-v1 network in Figure 15

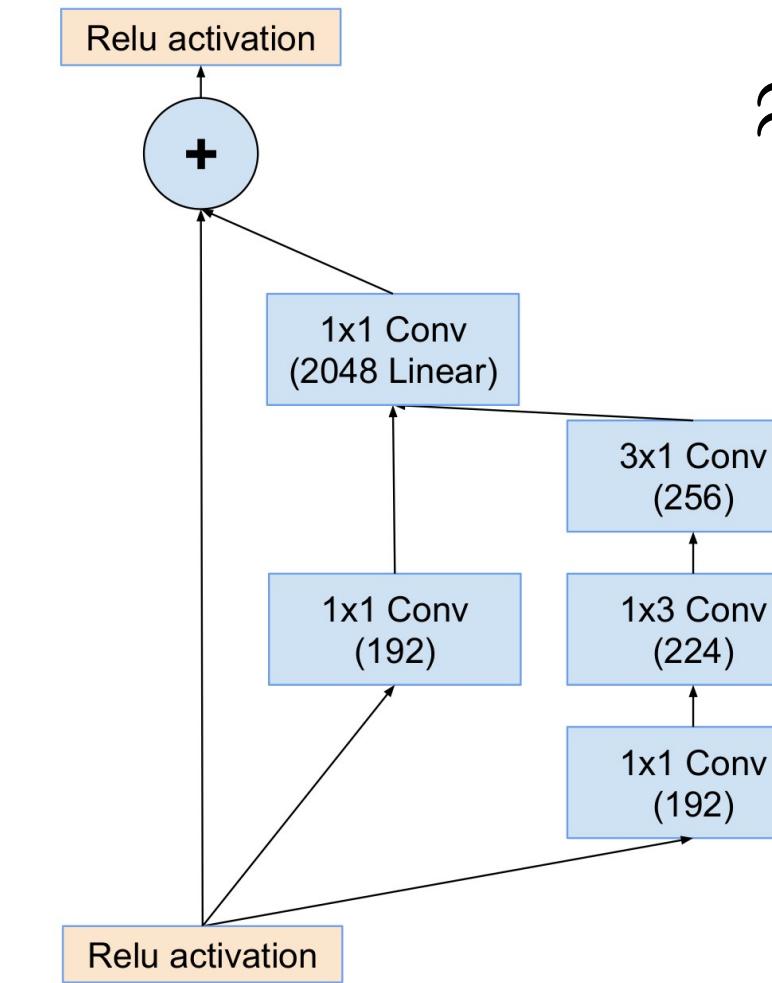


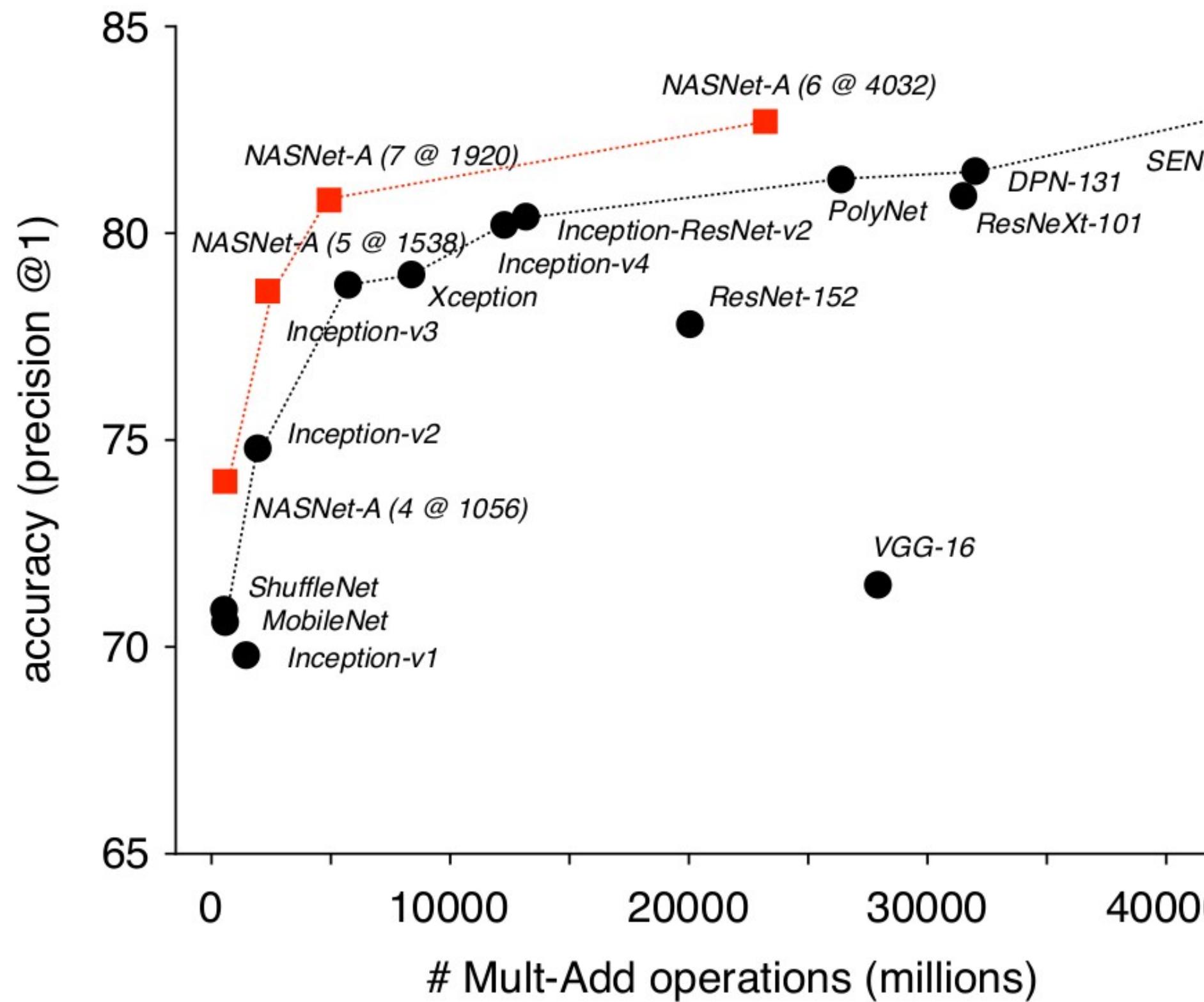
Figure 19. The schema for  $8 \times 8$  grid (Inception-ResNet-C) module of the Inception-ResNet-v2 network.

Network	$k$	$l$	$m$	$n$
Inception-v4	192	224	256	384
Inception-ResNet-v1	192	192	256	384
Inception-ResNet-v2	256	256	384	384

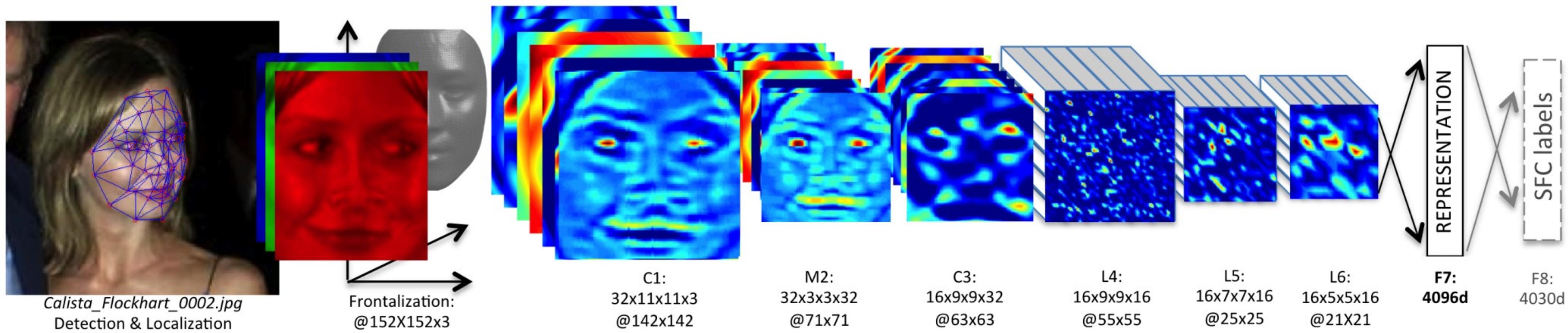
$\approx$  60 million parameters

# Transfer Learning for Image Recognition

- Find network architecture with reinforcement learning or evolutionary programming on a small dataset
- Use the same architectural elements on Imagenet



# DeepFace



## Closing the Gap to Human-Level Performance in Face Verification (2014)

- 4 million user - labeled faces on FaceBook images
- 4000 individuals
- Retrain fully-connected layers at the top on Labeled Faces in the Wild (LFW) dataset reaching (human level) accuracy of 97.35%

# **Automatic Differentiation**

# Shall we manually compute the gradients? (blackboard)

No. Use automatic (reverse mode) differentiation.

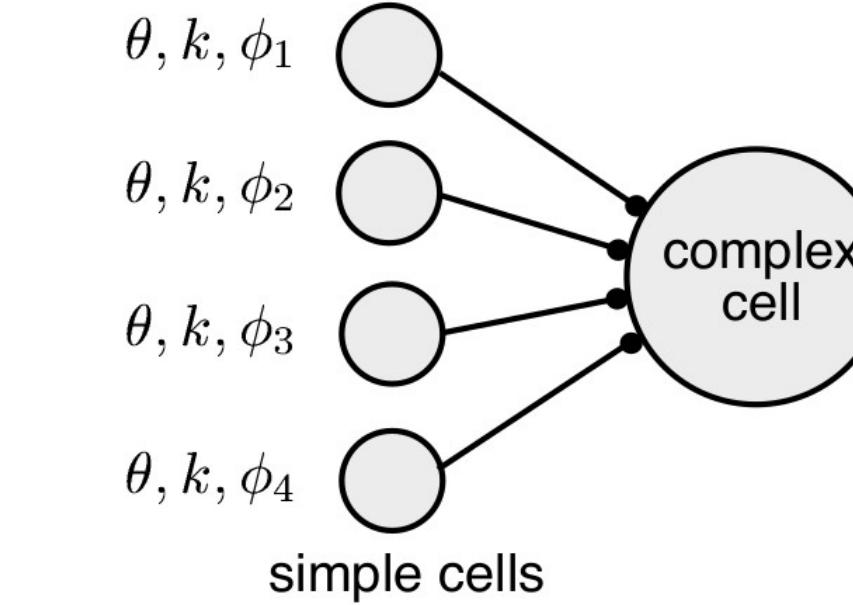
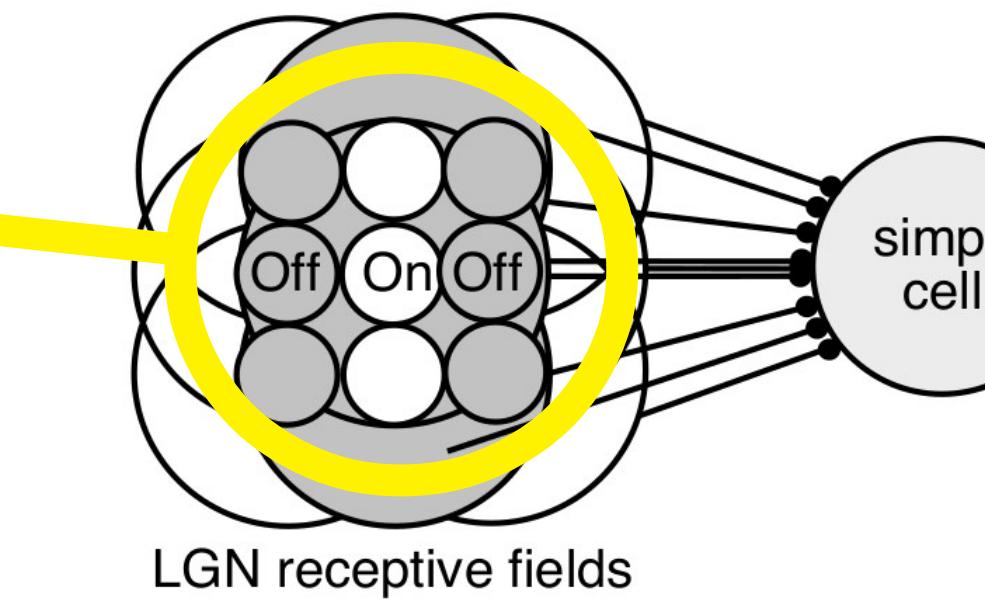
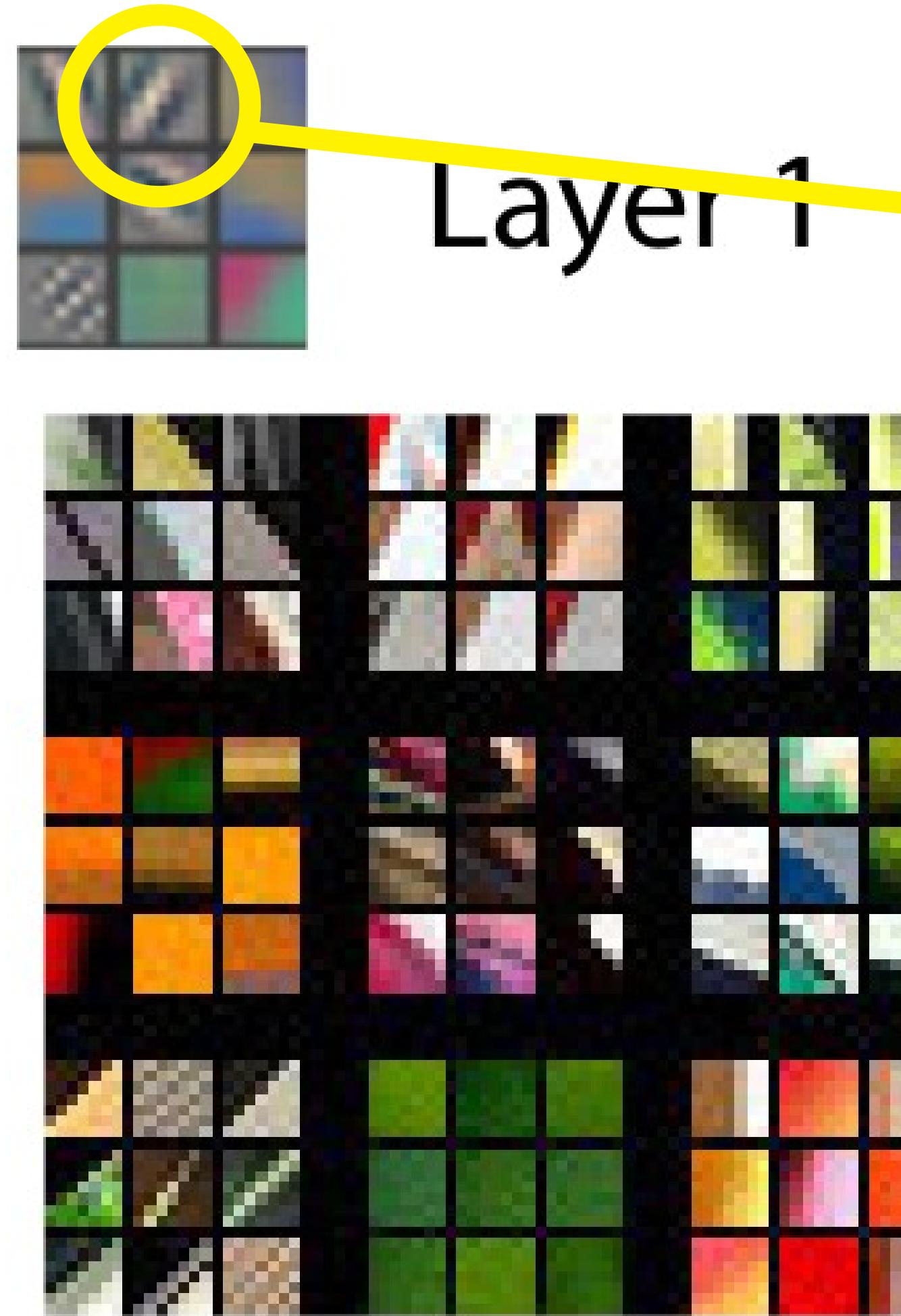
- Record computation tree on forward path to
- apply chain rule in reverse order.  
(Same idea as backprop in simple feedforward net)

An implementation requires just a Tracker (to record the computation tree) and analytical expressions of primitive operations

$$\frac{\partial f}{\partial x} \text{ e.g. } \frac{\partial \sin(x)}{\partial x} = \cos(x)$$

# **Applications beyond object recognition**

# Neuroscience

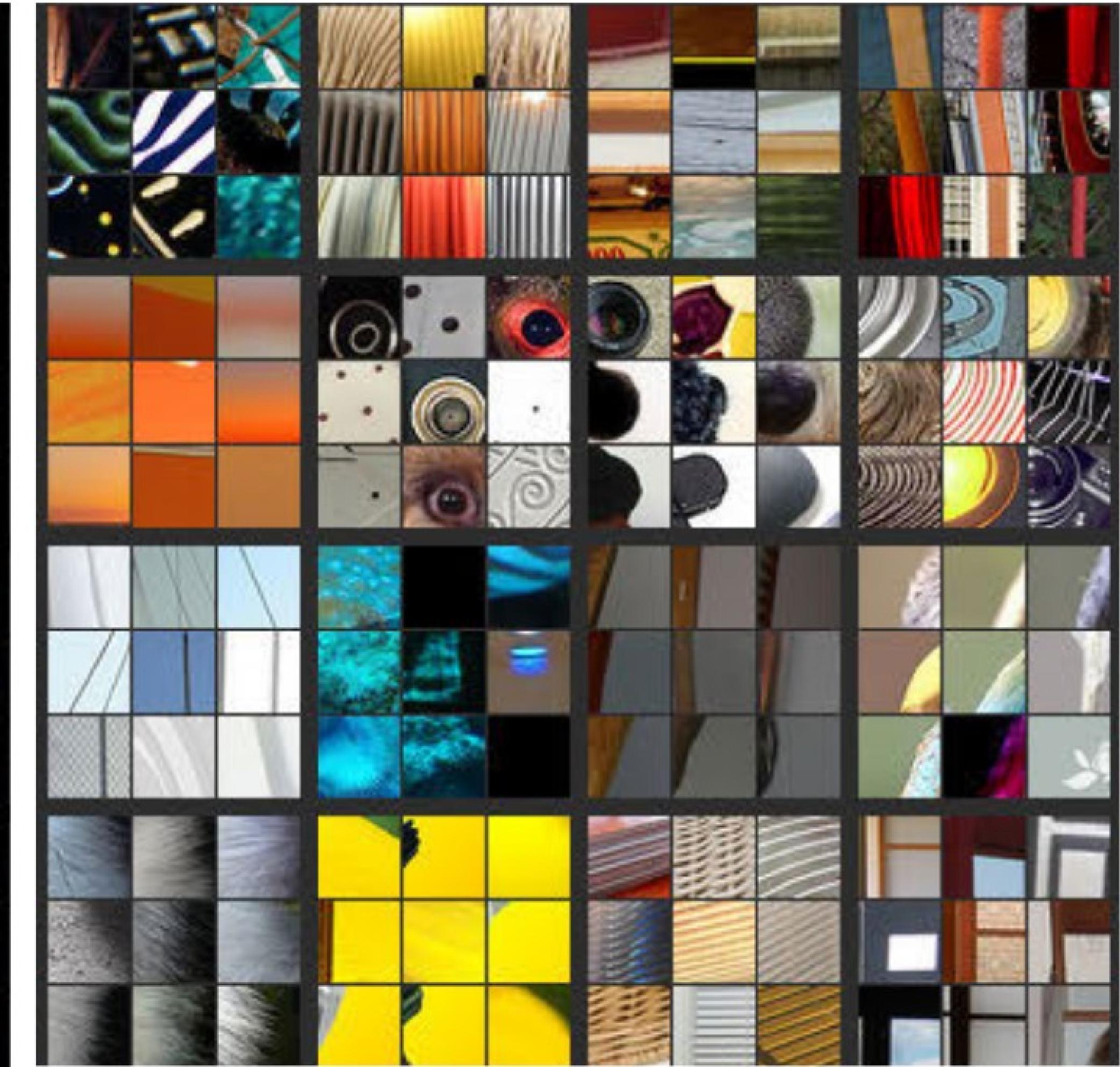
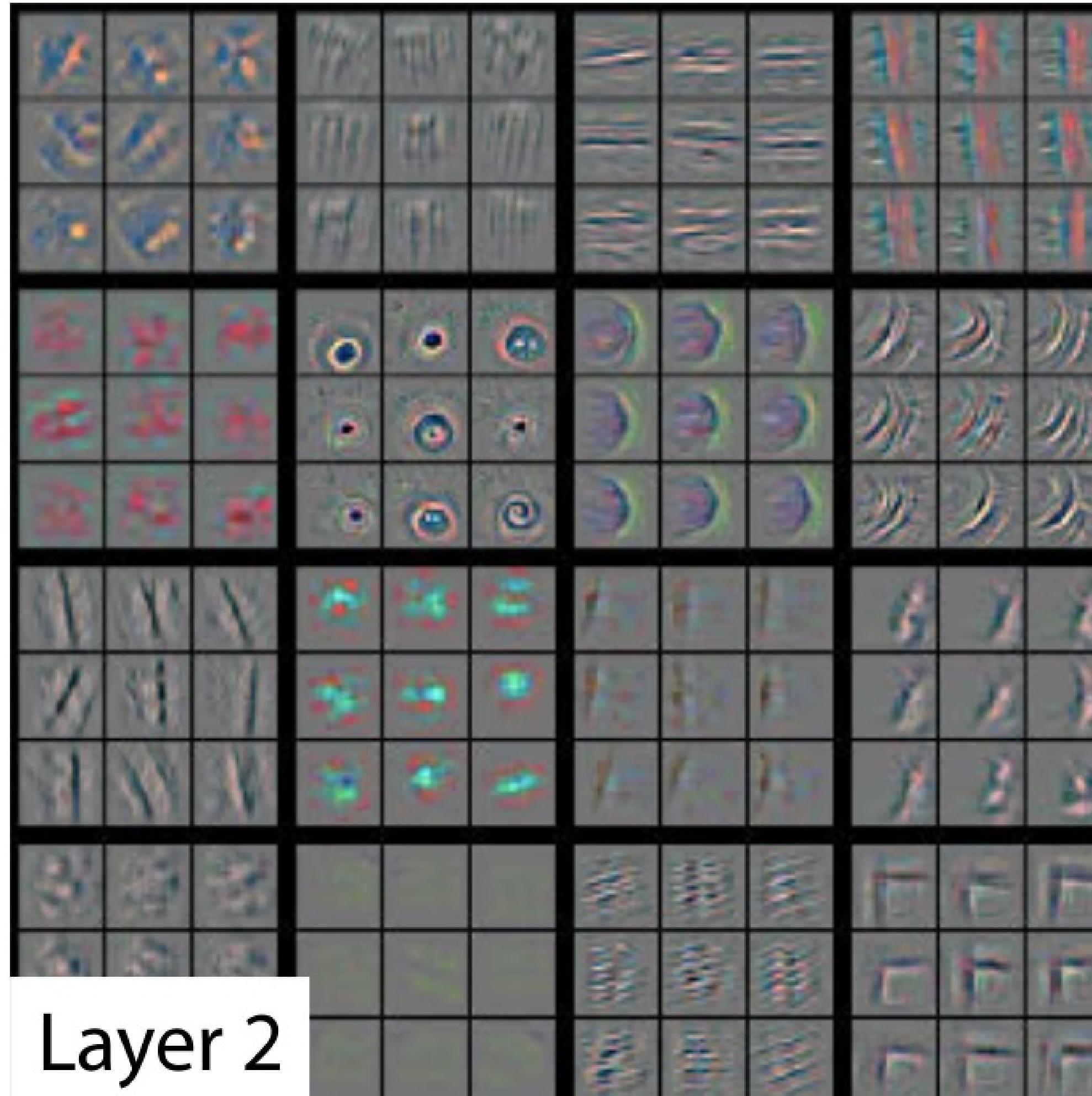


ConvNets have similar architecture and similar features as the visual system of animals and humans

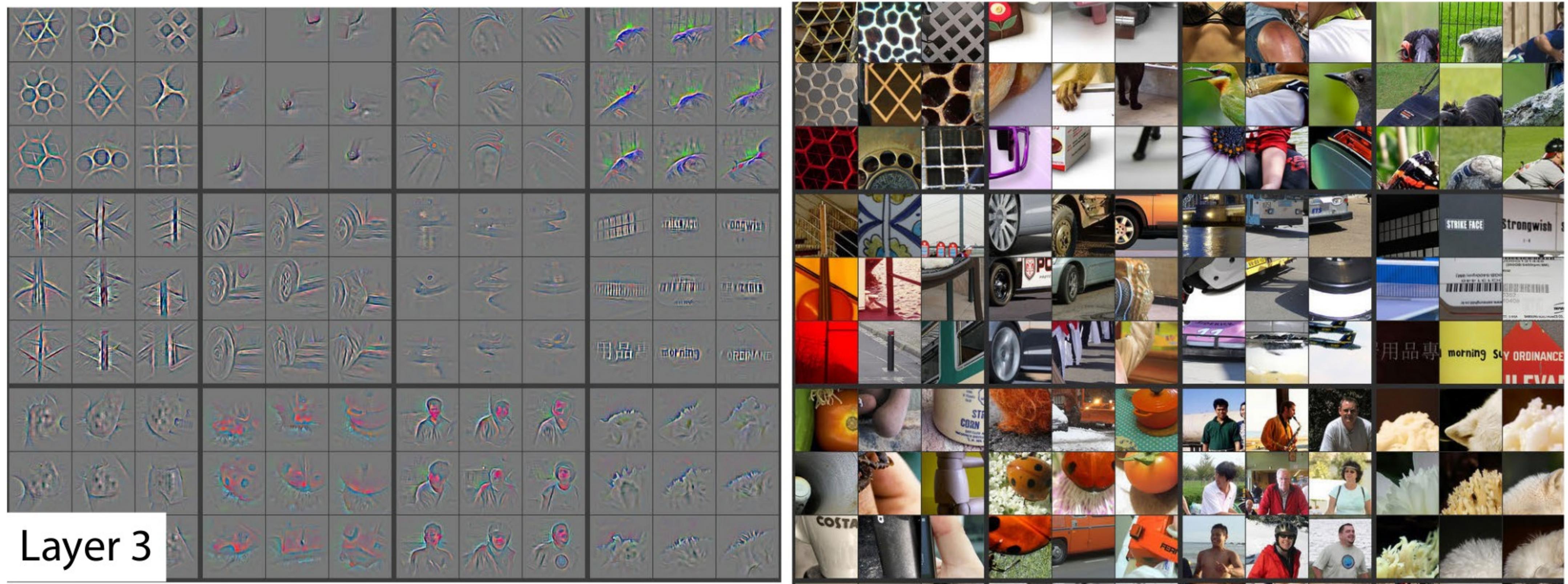
Dayan and Abbott, Theoretical Neuroscience, 2001

<https://arxiv.org/pdf/1311.2901.pdf>

# Visualizing convolutional networks



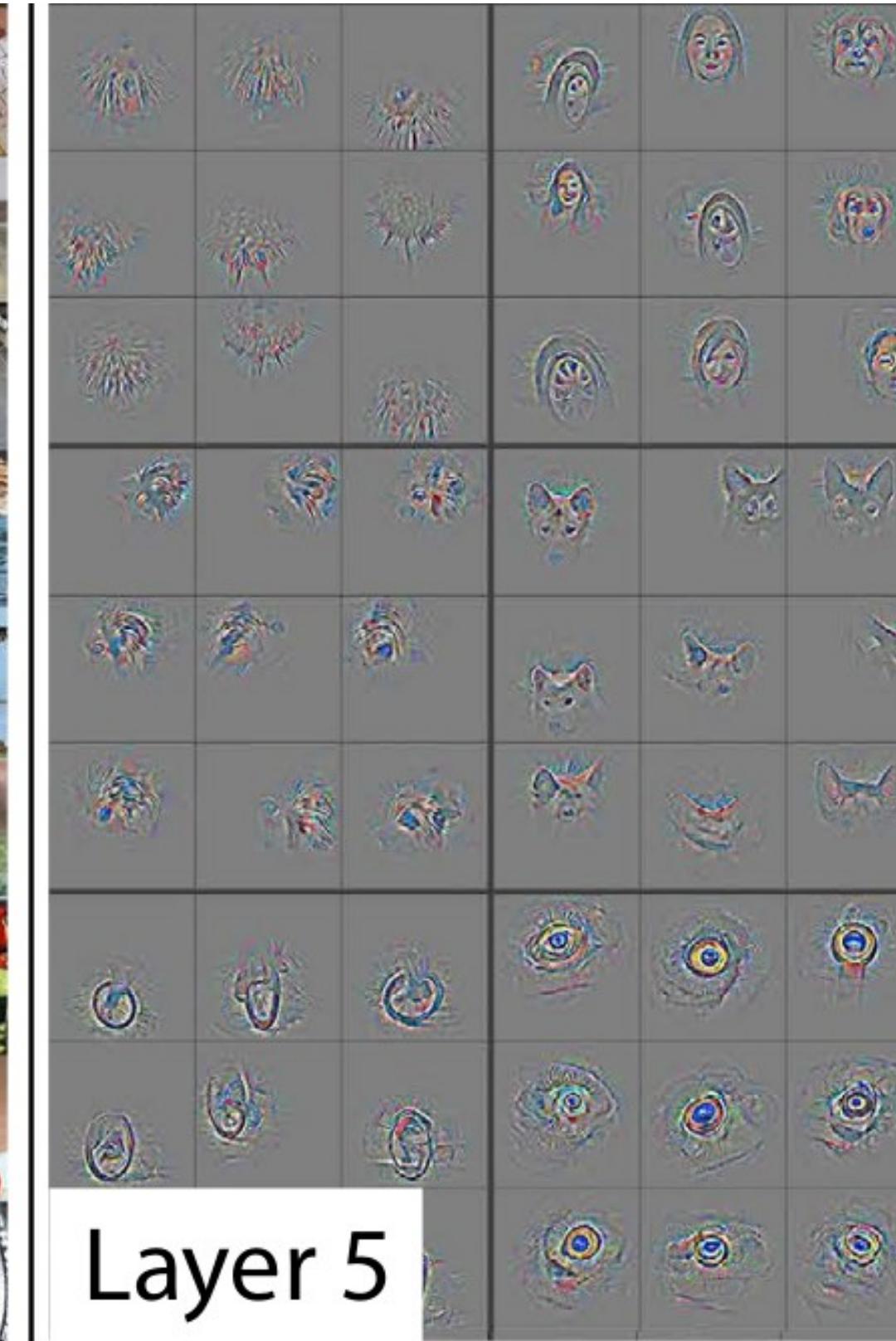
# Visualizing convolutional networks



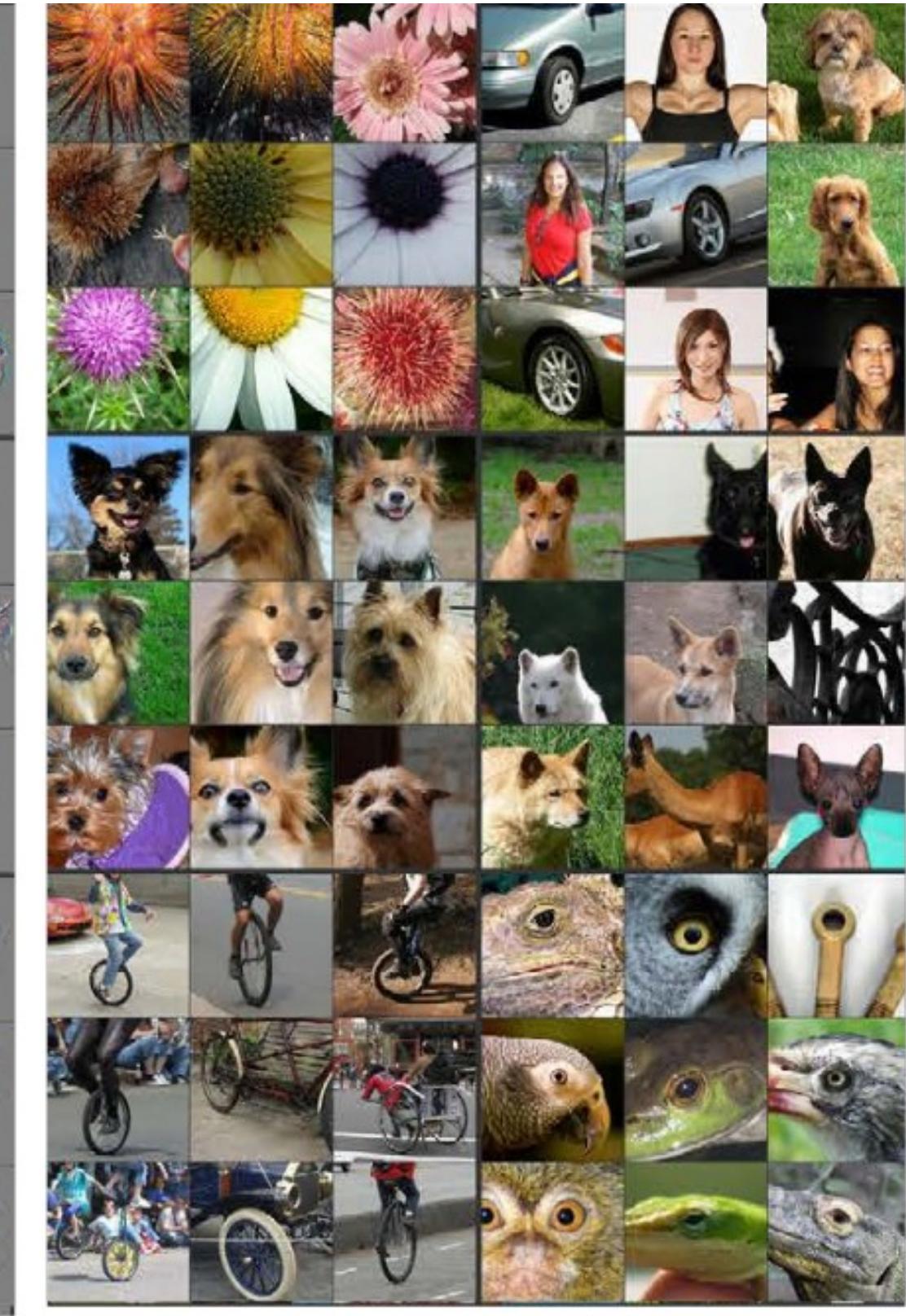
# Visualizing convolutional networks



Layer 4



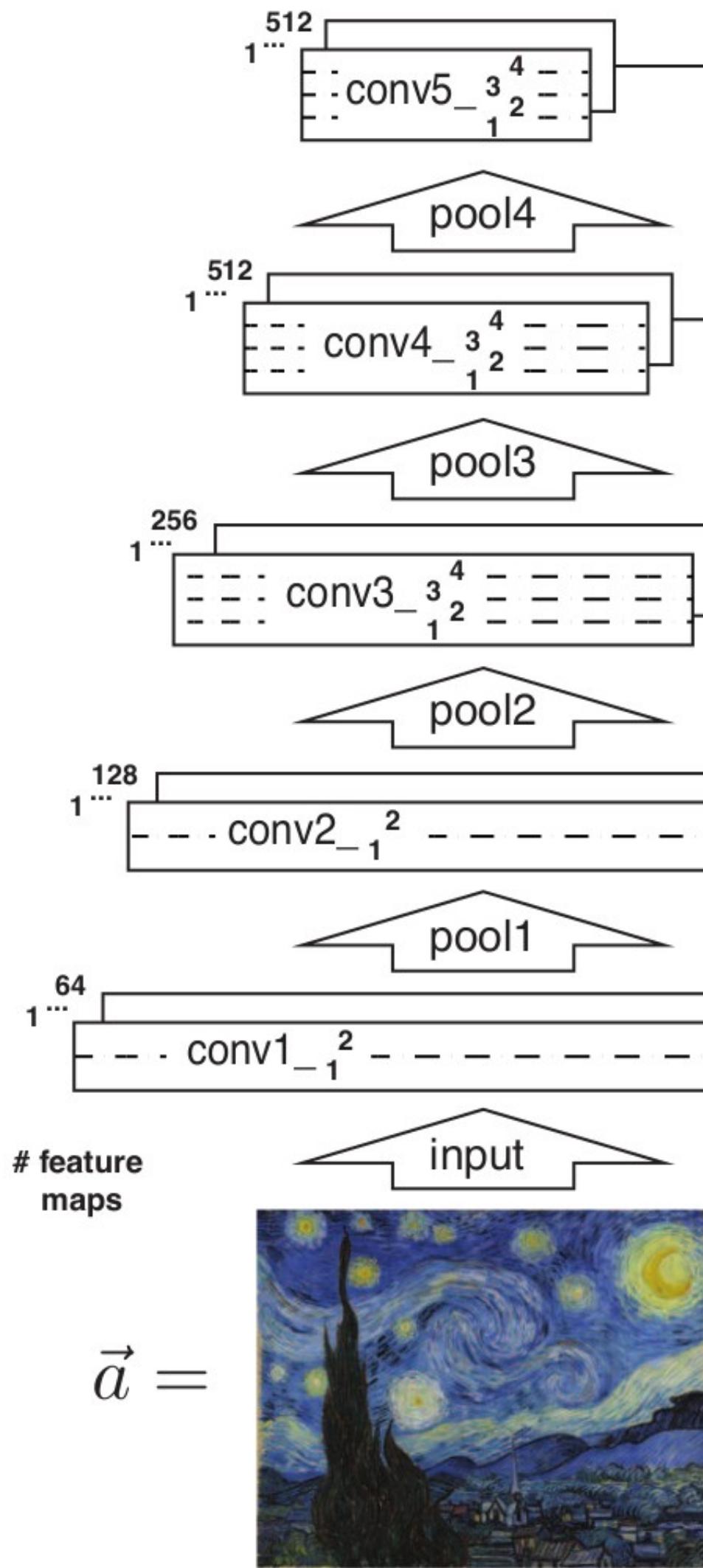
Layer 5



# Neural style



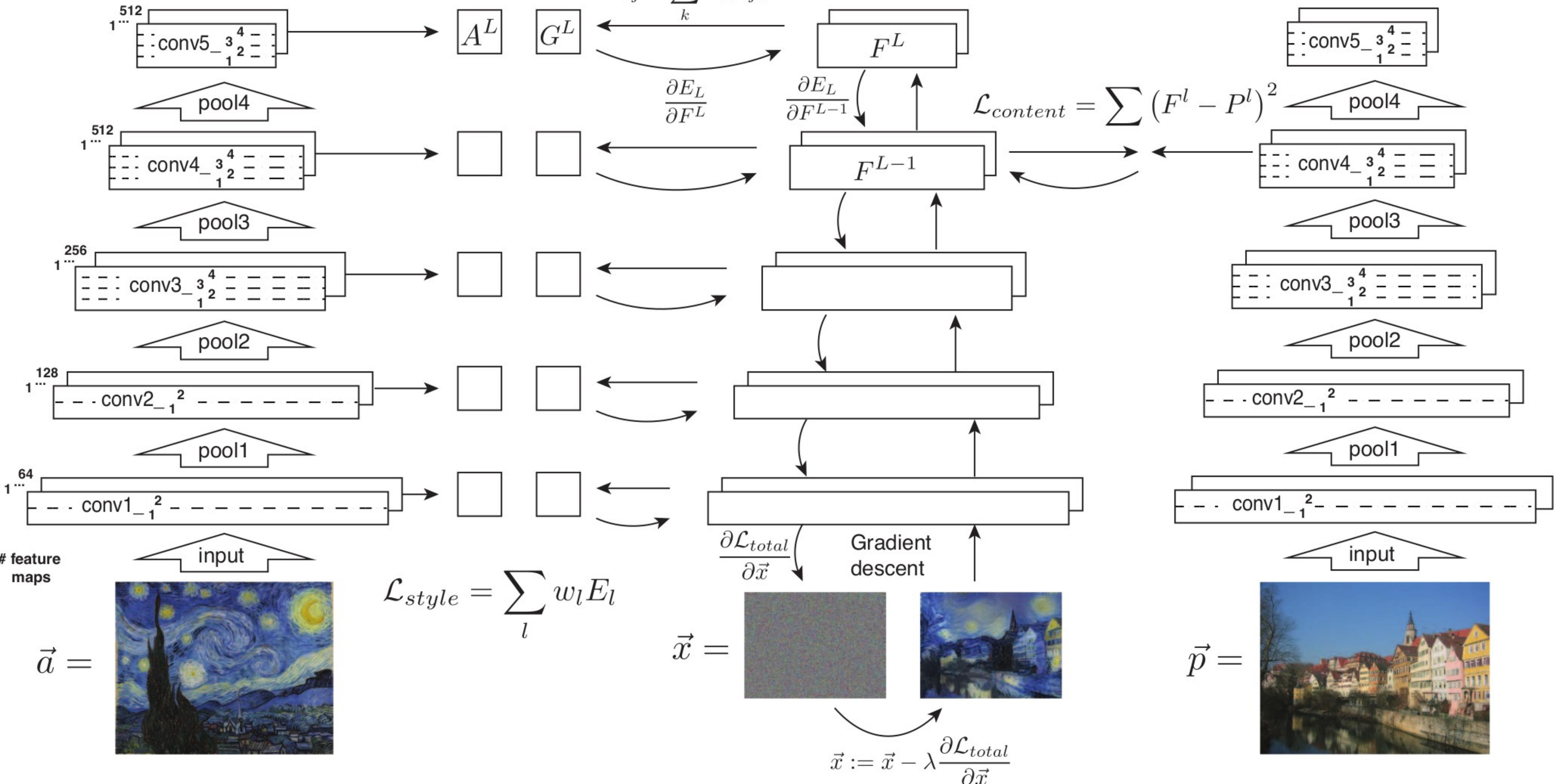
# Neural style



$$E_L = \sum (G^L - A^L)^2$$

$$G_{ij}^L = \sum_k F_{ik}^L F_{jk}^L.$$

$$\mathcal{L}_{total} = \alpha \mathcal{L}_{content} + \beta \mathcal{L}_{style}$$



# Caption generation and generative models of images



a cow is standing in the middle of a street

<https://cs.stanford.edu/people/karpathy/deepimagesent/>

ConvNet + bidirectional RNN



<https://thispersondoesnotexist.com/>

ConvNets in generative  
adversarial networks (GANs)

# Summary

- 1) Use good explicit inductive biases together with transfer learning and data augmentation.
- 2) Convolutions (and max-pooling) are reasonable inductive biases for natural images.
- 3) Residual connections help to learn in very deep networks.
- 4) Modern architectures like Inception-Resnet-v2, DeepFace and NASNet reach human level performance on recognition tasks.
- 5) AutoDiff is a flexible tool to train different architectures.
- 6) ConvNets and neurons in the visual systems have similar receptive fields.