

MACHINE LEARNING & DATA MINING

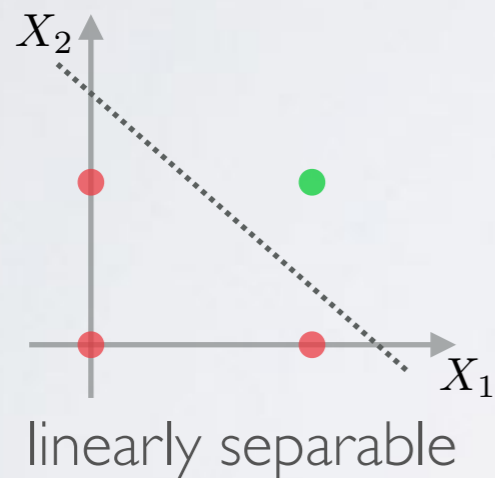
CS/CNS/EE 155

Deep Learning
Part II

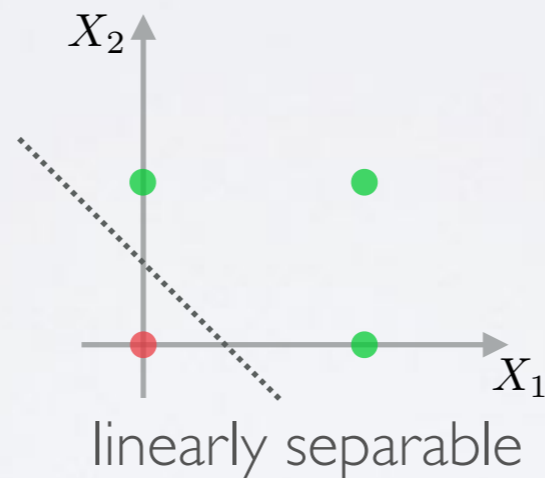
recap of last lecture

logistic regression can't handle non-linear data distributions

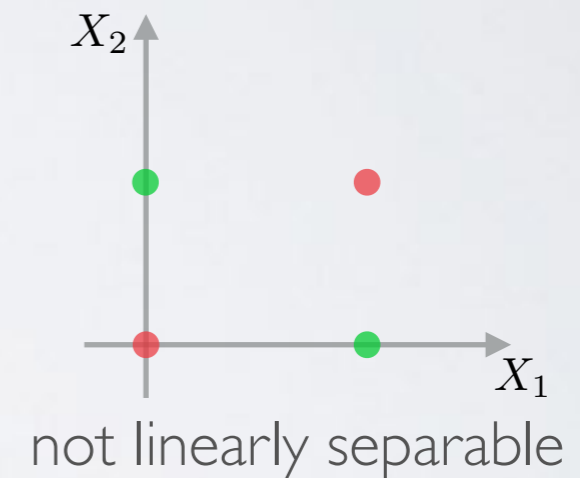
AND



OR



XOR



recap of last lecture

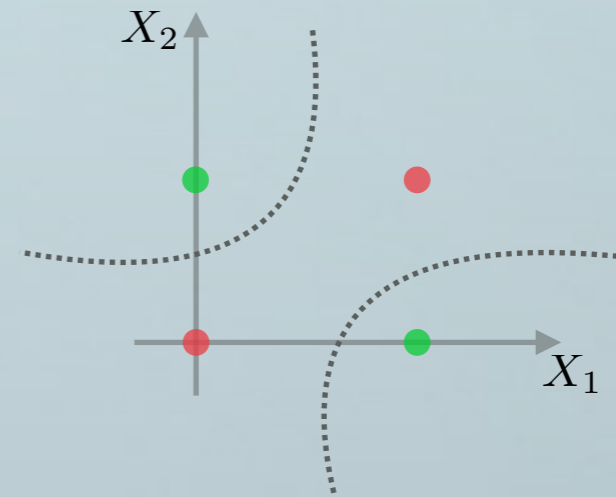
let's use non-linear features to linearize the problem!

one approach: use a set of **hand-crafted** non-linear transformations

$$X_1, X_2 \rightarrow X_1, X_2, X_1 X_2$$

linear decision
boundary

hyperbolic decision
boundary

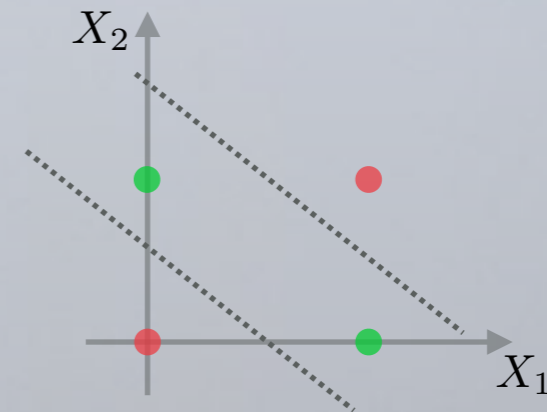


another approach: use a set of **learned** non-linear transformations

$$X_1, X_2 \rightarrow X_1 \wedge X_2, X_1 \vee X_2$$

linear decision
boundary

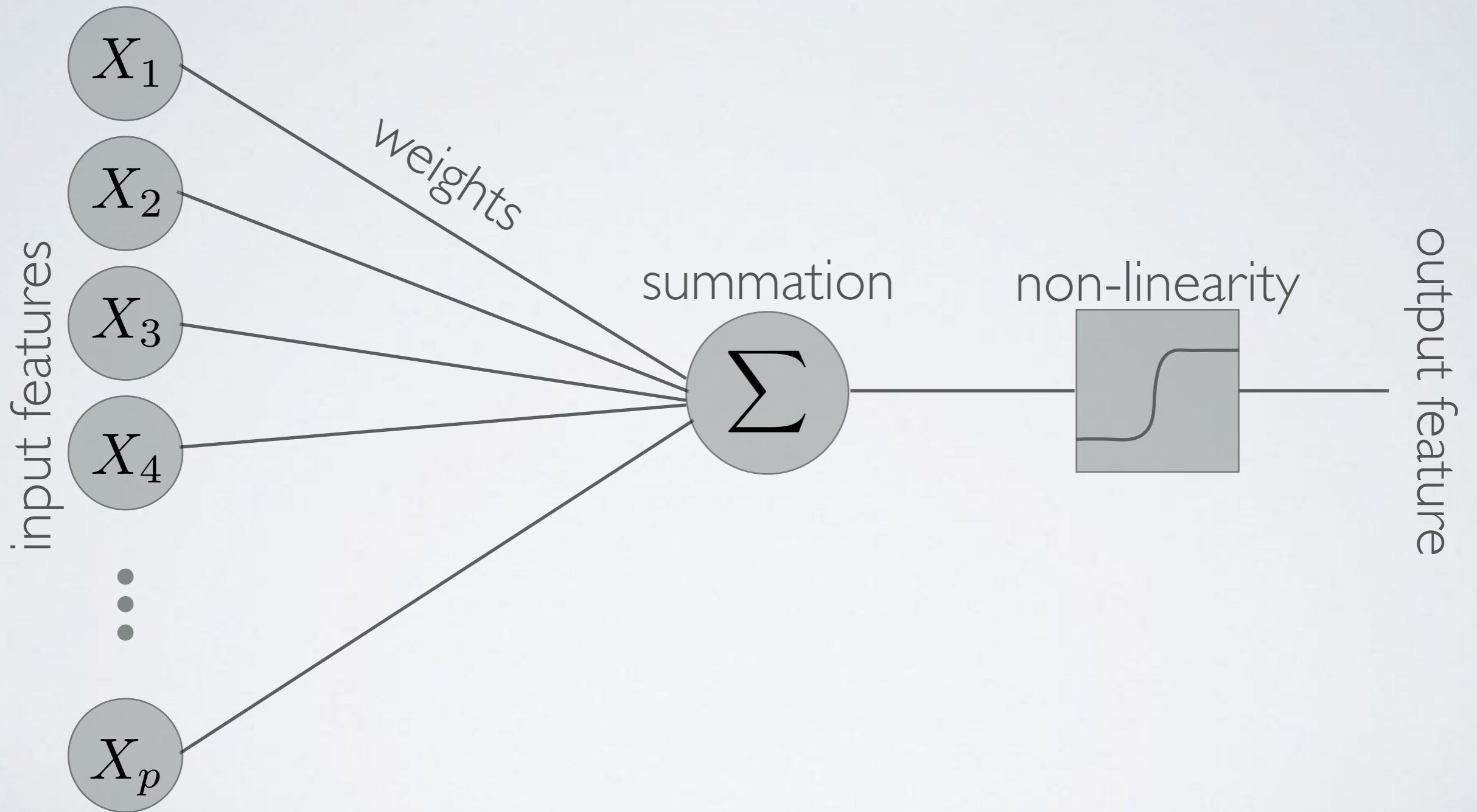
multiple linear decision
boundaries



$$\longrightarrow \neg(X_1 \wedge X_2) \wedge (X_1 \vee X_2)$$

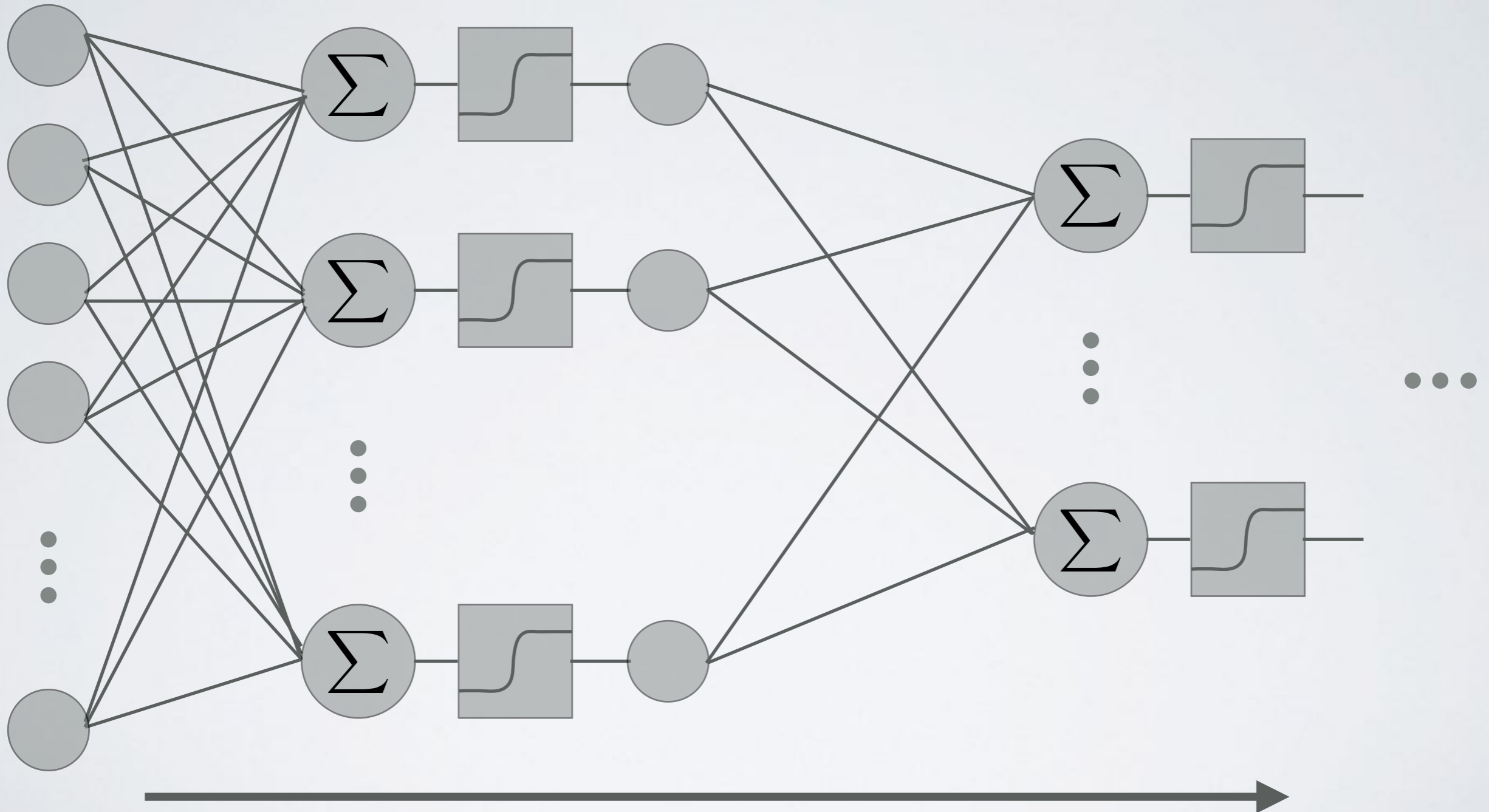
recap of last lecture

'neuron'



recap of last lecture

'neural network'



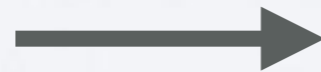
depth
5

big picture

neural networks are function approximators that can be trained to match the data's label distribution

$$**f(data) \sim P(label | data)**$$

more parameters,
depth



more expressive,
better approximation

(as long as you don't overfit)

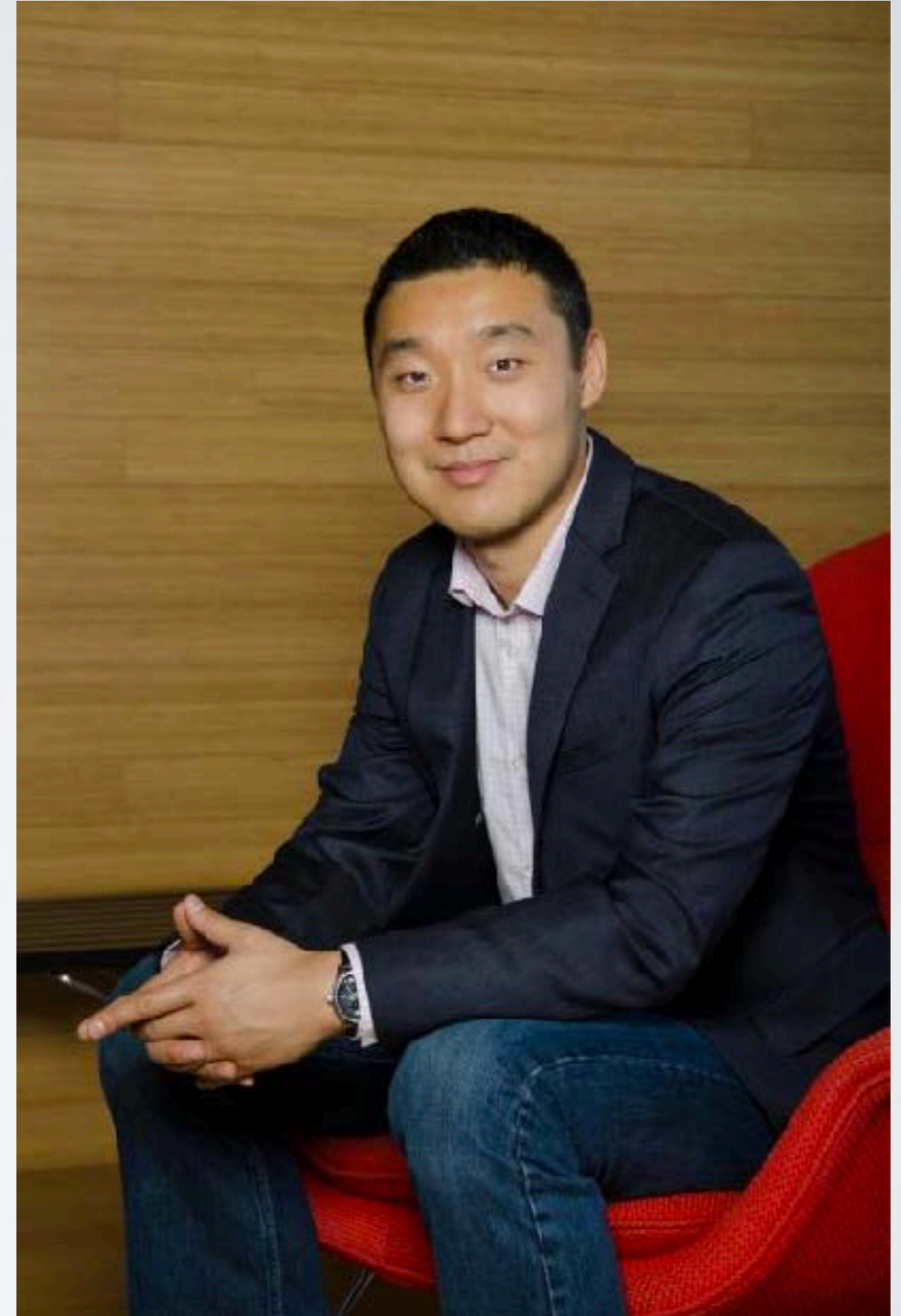
when is this useful?

Q: who is in this picture?

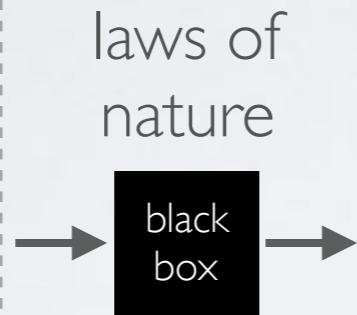
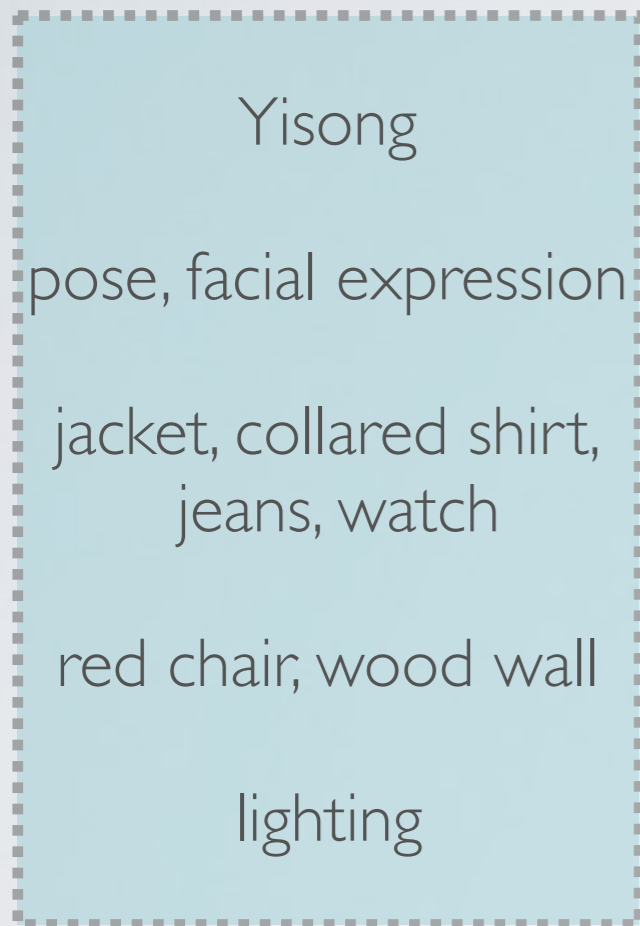
A: Yisong

Q: why? / how do you know?

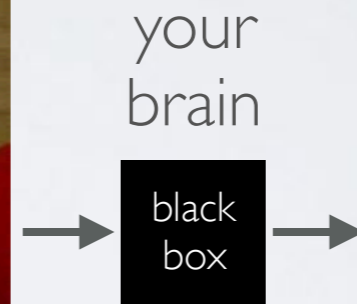
A: umm...



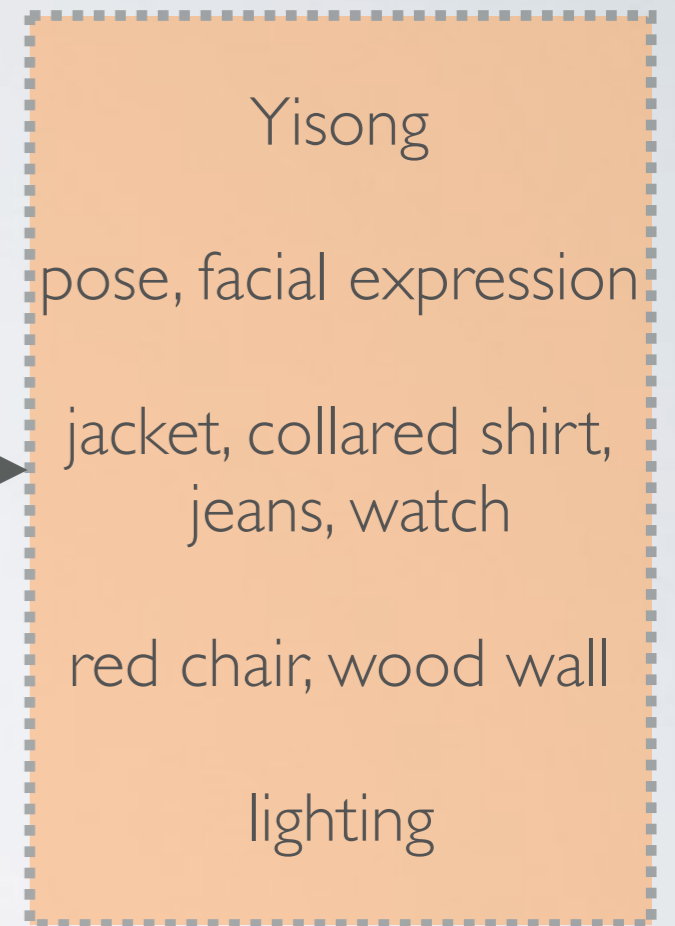
environment



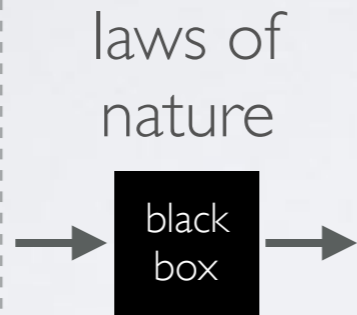
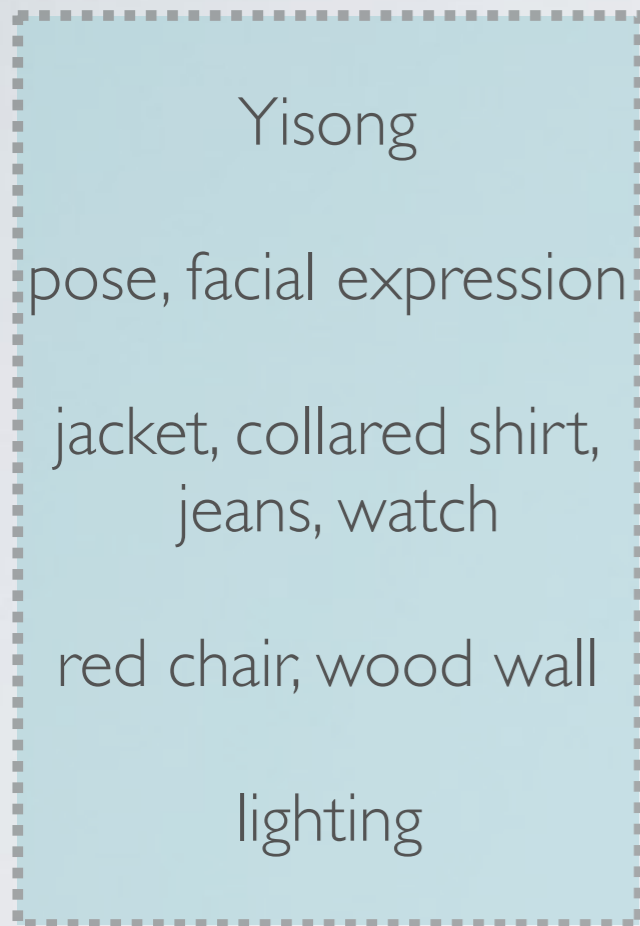
observation



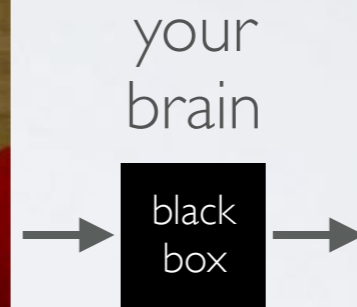
inference



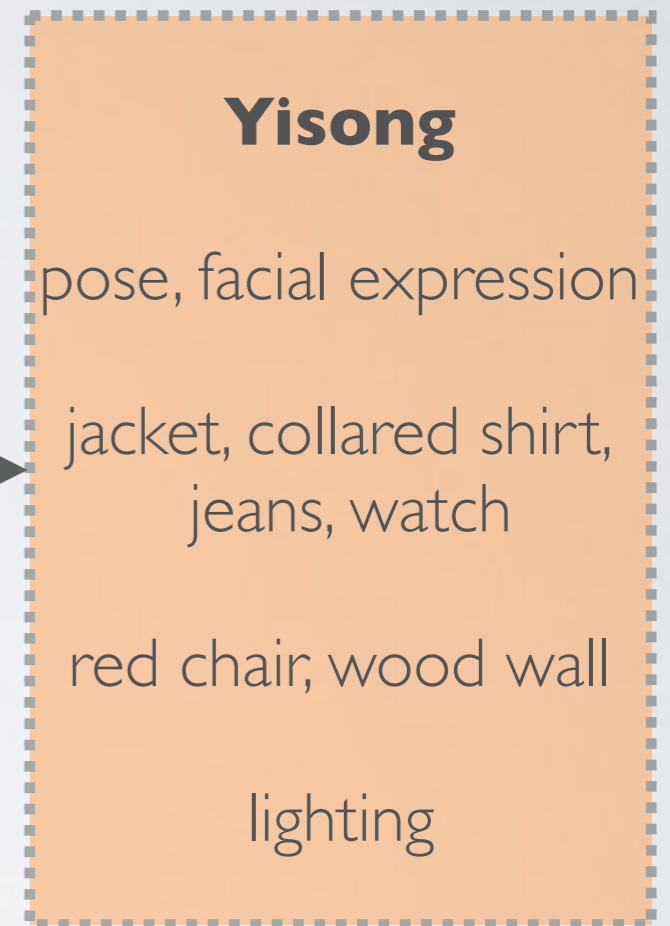
environment



observation



inference



'who is this?'

two sides of the same coin

generation

*there are latent properties
that result in specific patterns
in images*

Yisong



discrimination

*there are patterns in images
that allow us to infer latent
properties*



Yisong

the mappings between properties and images are
too complicated to define manually

deep learning to the rescue!

task:

train a deep neural network to discriminate whether or not an image contains Yisong

data



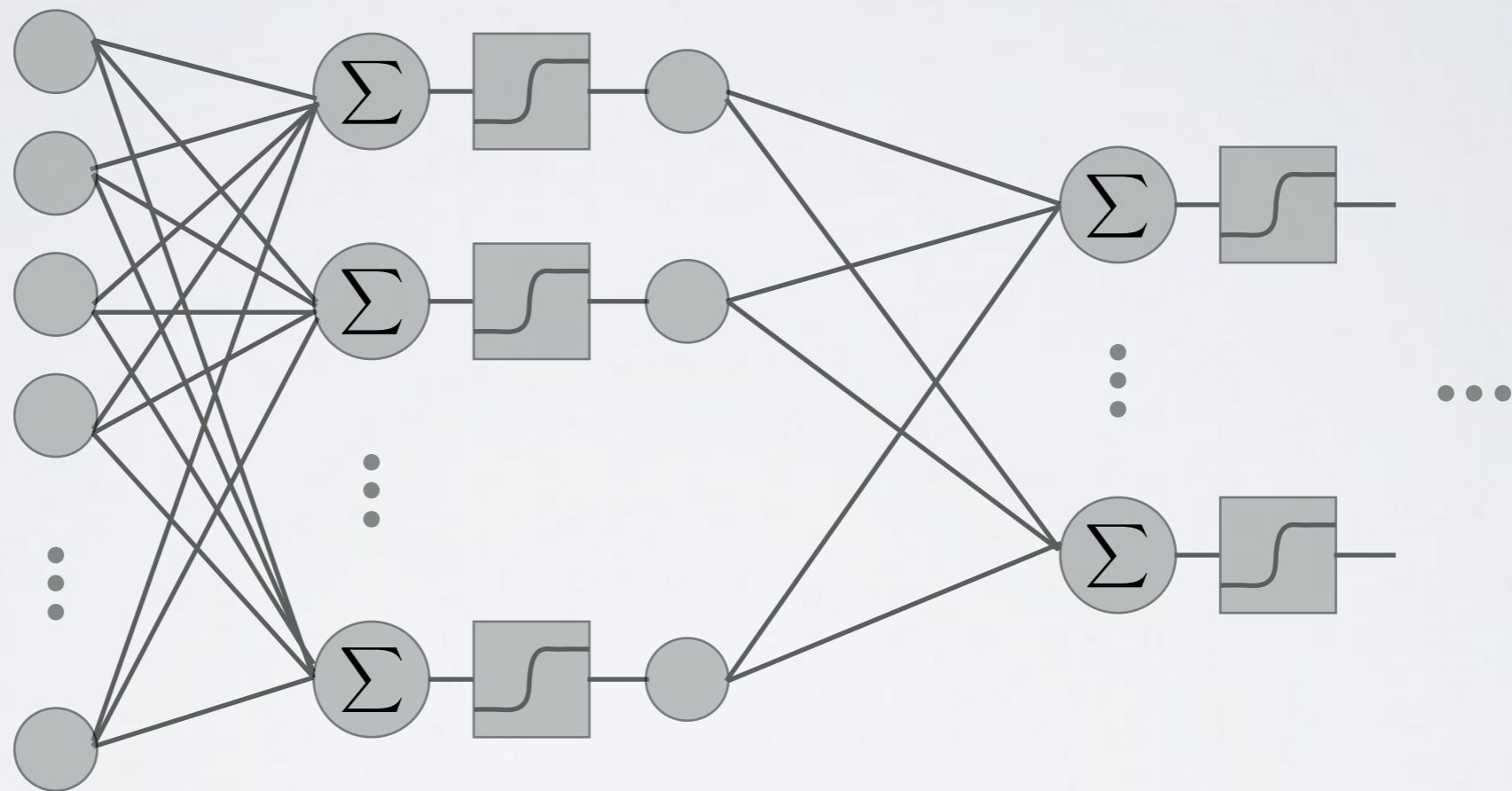
Yisong



not Yisong

labels

network architecture?



decide on an input size

larger input:
more parameters
clearer patterns

15 25 35
× × ×
15 25 35
× × ×
3 3 3

$50 \times 50 \times 3$

$75 \times 75 \times 3$

$100 \times 100 \times 3$

$150 \times 150 \times 3$

$205 \times 205 \times 3$

$280 \times 280 \times 3$

675

1,875

3,675

7,500

16,875

30,000

67,500

126,075

235,200

smaller input:
fewer parameters
noisier patterns

decide on an input size

larger input:
more parameters
clearer patterns

15 25 35
× × ×
15 25 35
× × ×
| | |

50 × 50
× |

75 × 75 × |

100 × 100 × |

150 × 150 × |

205 × 205 × |

280 × 280 × |

225

625

1,225

2,500

5,625

10,000

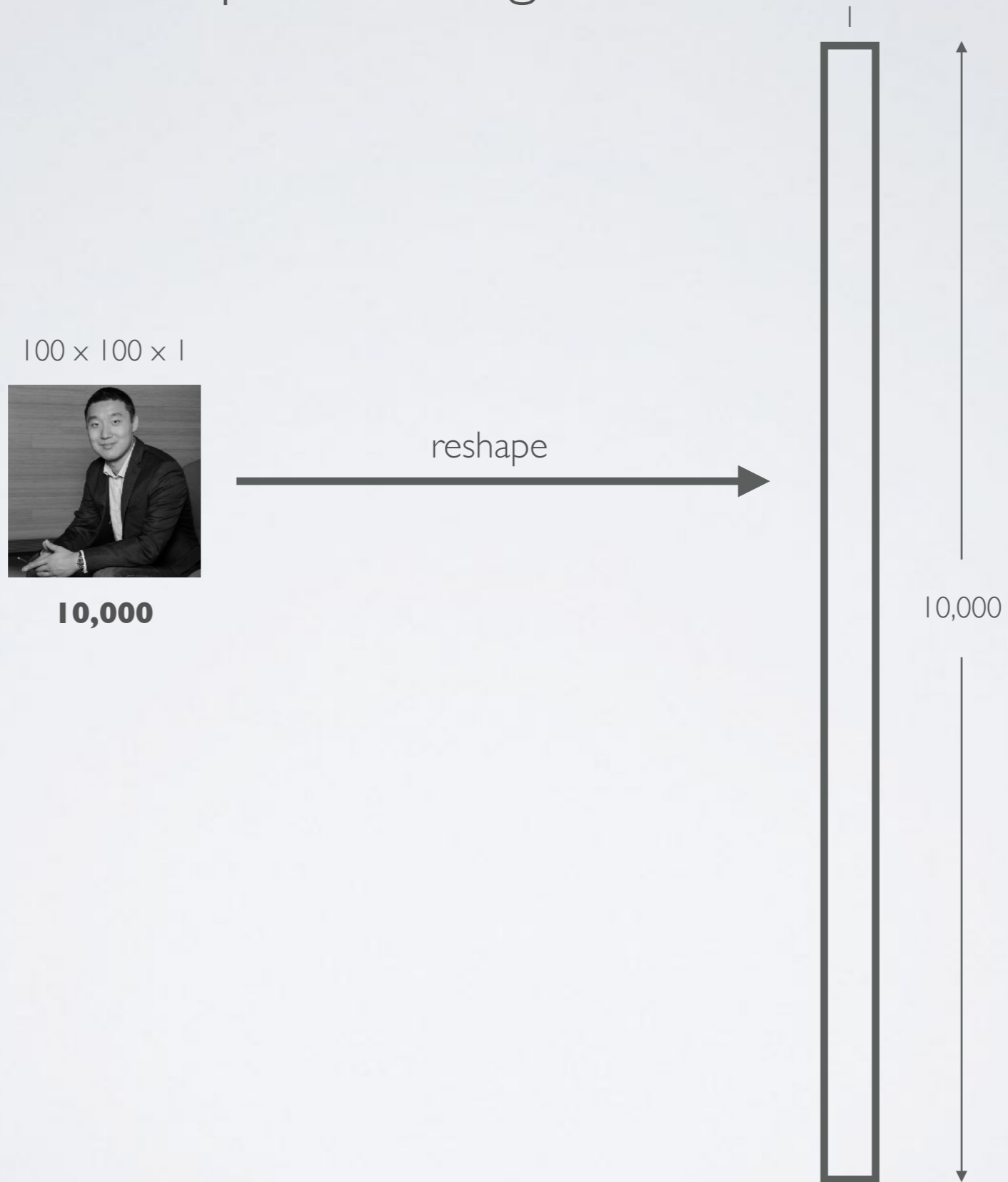
22,500

42,025

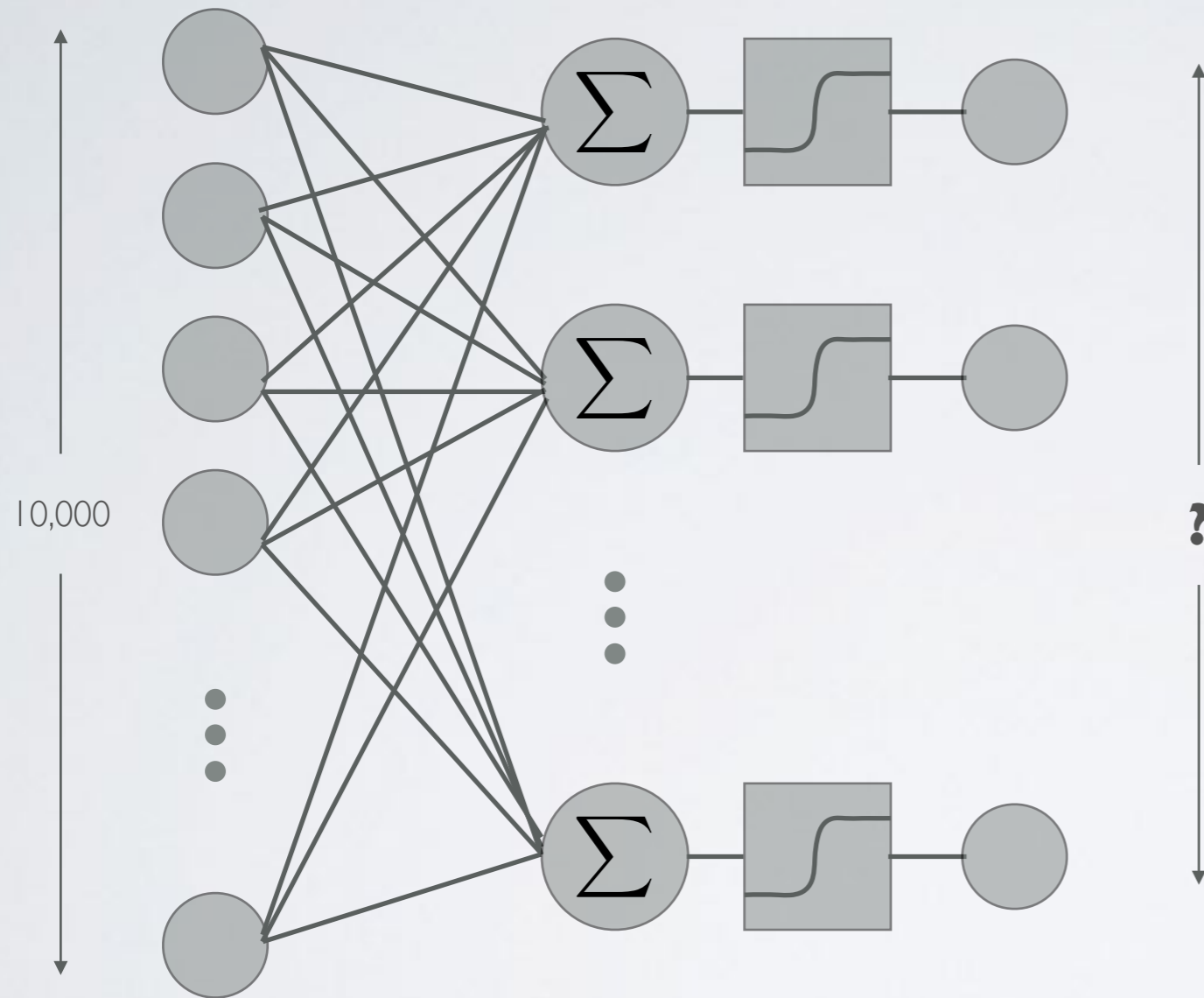
78,400

smaller input:
fewer parameters
noisier patterns

reshape the image into a vector



what about the rest of the architecture?



how many units can we afford?

$10,000 \times 100 = 1$ million weights

$10,000 \times 1,000 = 10$ million weights

$10,000 \times 10,000 = 100$ million weights

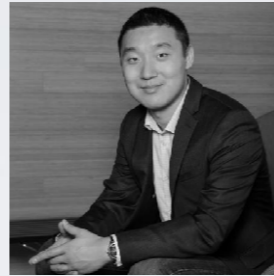
$10,000 \times 100,000 = 1$ billion weights

⋮

how many basic patterns do we expect to find in the image?

how many patterns can the image contain?

100 × 100 × 1



10,000

upper bound

if we consider all values (1 - 256) of all 10,000 pixels, there are **256^{10,000}** possible patterns

this is more than the number of atoms in the known, observable universe.

in reality, the actual number will be much less

lower bound

if we want to recognize *multiple* basic, low-level patterns (e.g. edges, gradients, etc.) *anywhere in the image*, I estimate there will be a total of *at least* **10,000** of these.

the actual number may be far more

—————→ at least 100 million weights in the first layer alone

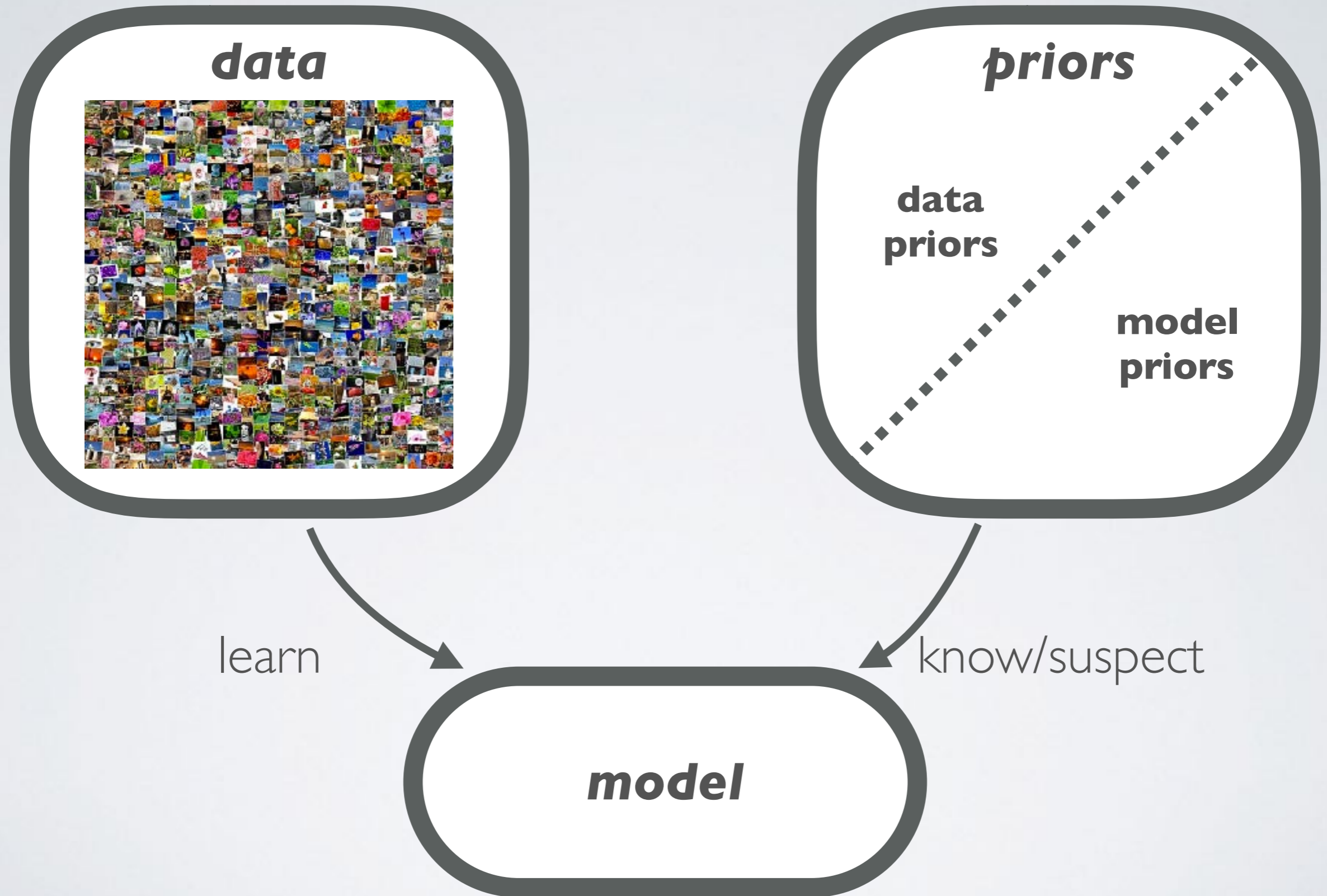
our approach requires a huge number of weights (parameters)

this dramatically increases the amount of data/labels we need to collect

as well as the amount of computation required for training

we need to re-evaluate our approach

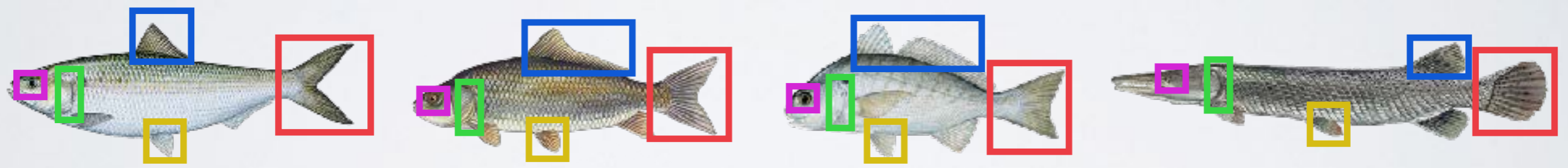
two sources for improving a model



two sources for improving a model



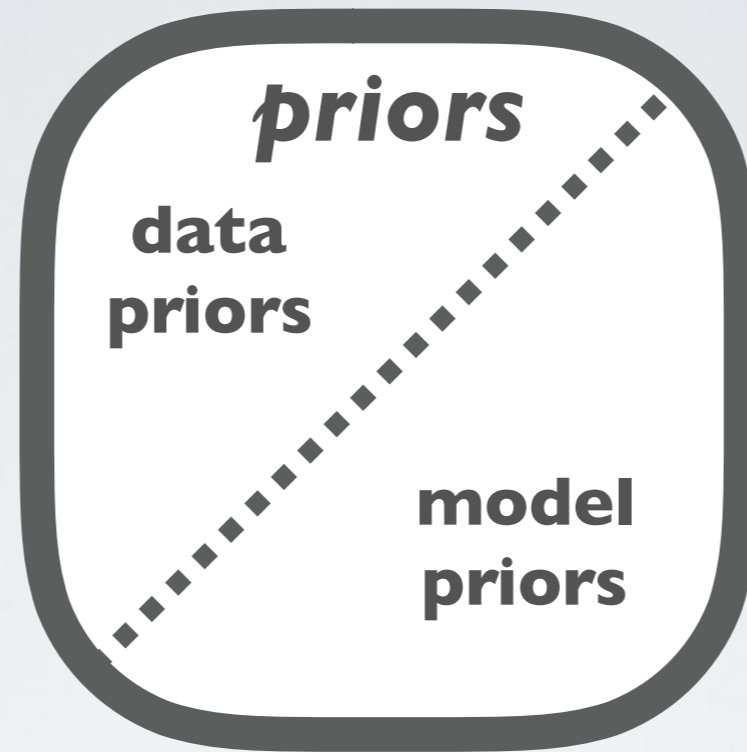
by observing data samples, we can learn about the data distribution



e.g. learn features common to fish and differences between them

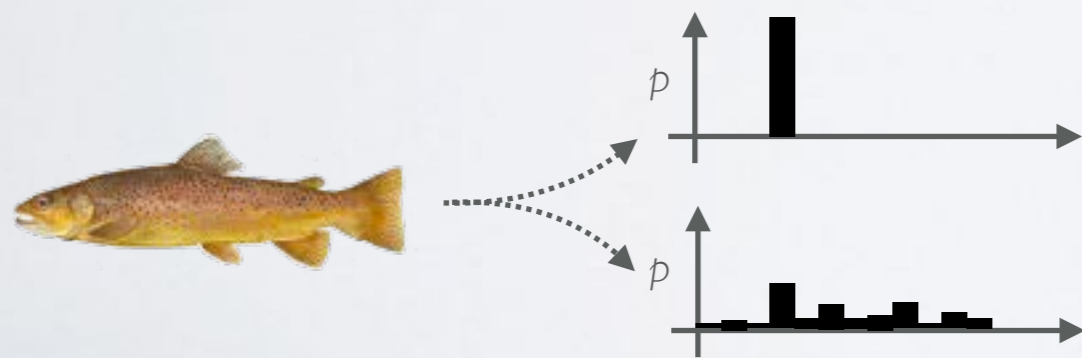
with this knowledge, we can more easily learn mappings to/from the data to latent quantities of interest (transfer learning from unsupervised features)

two sources for improving a model

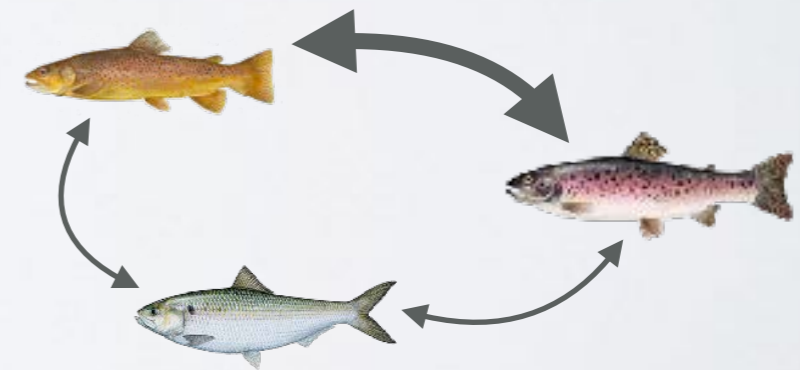


priors correspond to knowledge (or suspicions) that we already have about the task

a **data prior** is additional relevant information about a data example

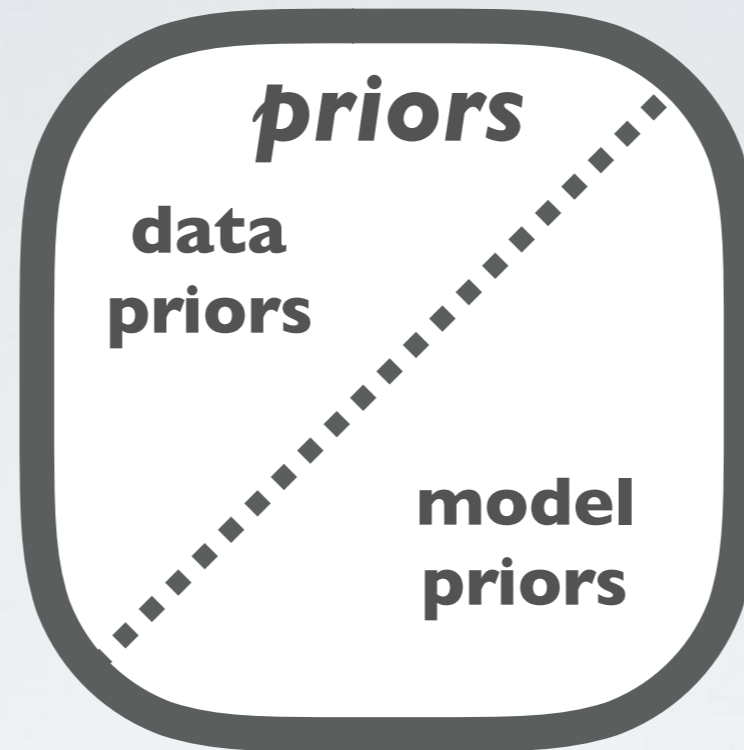


label or label distribution



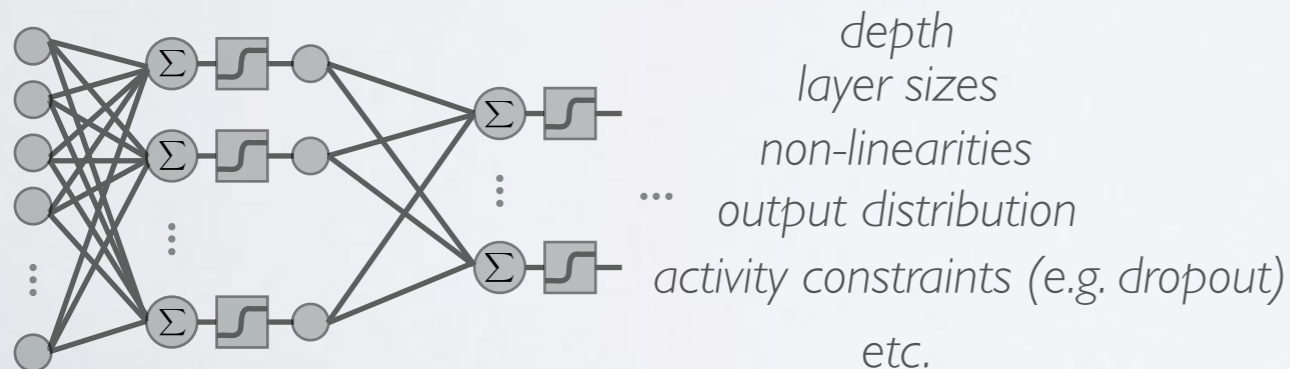
example/class similarity

two sources for improving a model



priors correspond to knowledge (or suspicions) that we already have about the task

a **model prior** is relevant information about the model/task

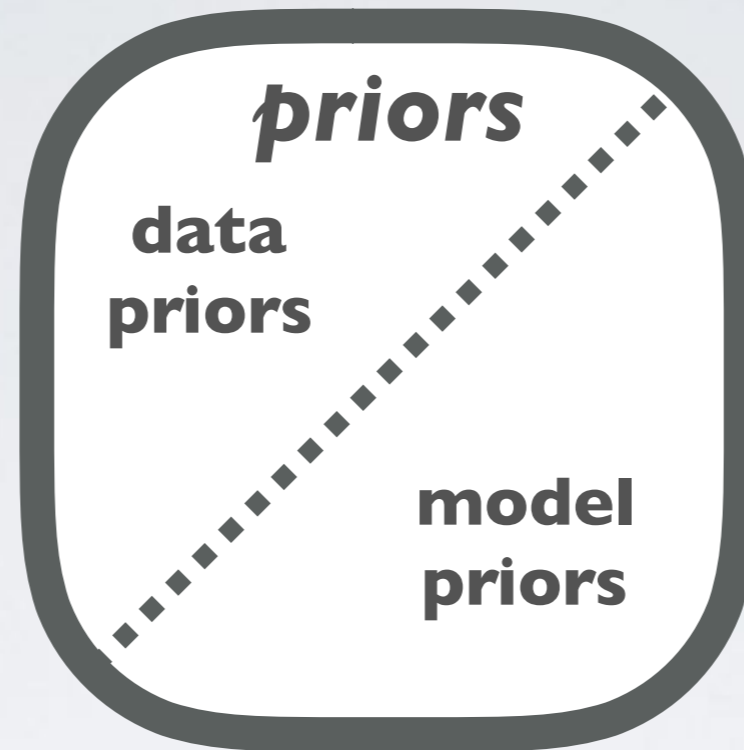


- weight magnitude constraints (L1 /L2)*
- transfer learning from a similar task*
- hand-crafted features*
- etc.*

model class/architecture

parameter constraints/values

two sources for improving a model



priors are necessary for any task

without them, we would have no way of knowing what/how to learn

priors can vary in strength

with *strong* priors, we don't need data (we already *know* the solution)

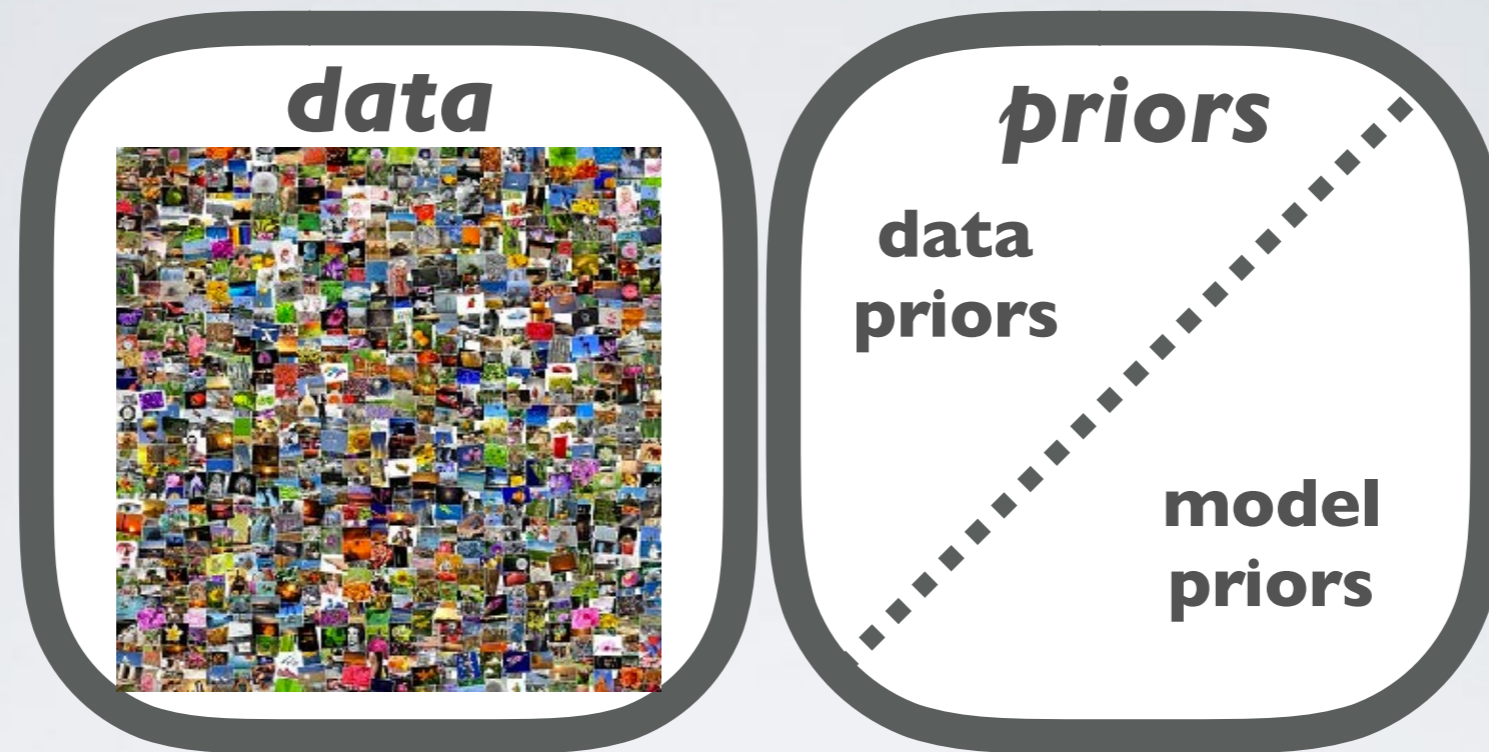
with *weak* priors, we need a lot of data (we mostly *learn* the solution)

priors can be good or bad

good/correct priors make learning easier

bad/incorrect priors make learning more difficult or impossible

two sources for improving a model



our current approach to visual object recognition
relies too heavily on data

need to impose additional/stronger priors
to simplify learning

we'll impose model priors to restrict the model class
for this task

properties of images

images have a notion of **locality**, which operates at *multiple scales*:

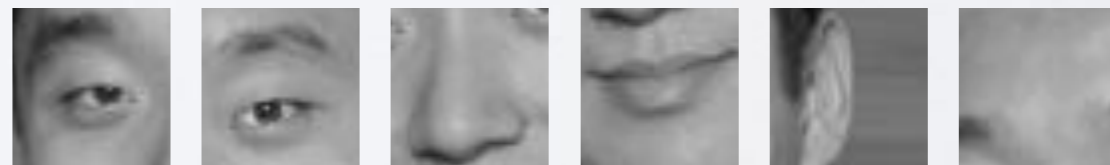
neighboring *pixels* tend to be similar and vary in particular ways



nearby *patches* tend to share characteristics and are combined in particular ways



nearby *regions* (of objects) tend to be found in particular arrangements



properties of images

what does **locality** imply for our model?

more meaningful to work in image space than with reshaped vectors



units should restrict their inputs to areas of nearby units in the previous layer



properties of images

objects have a notion of **translation invariance**



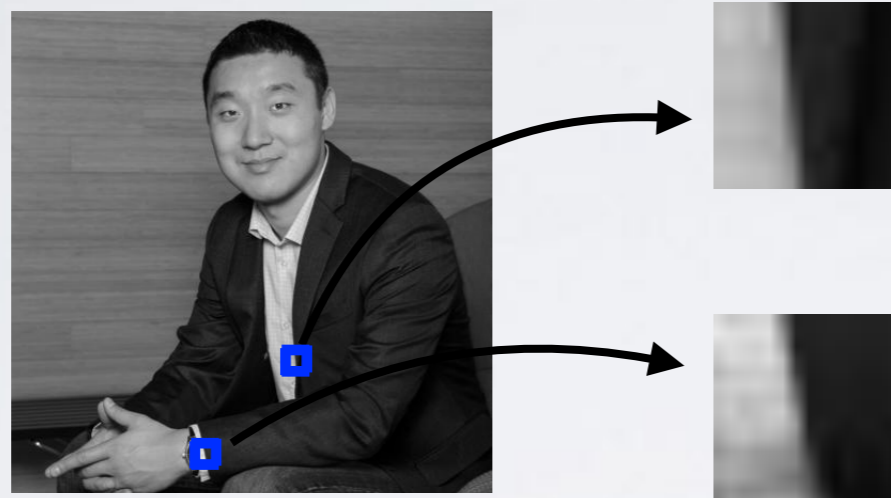
Yisong's identity is independent of his spatial location

similar statistics apply throughout the image

properties of images

what does **translation invariance** imply for our model?

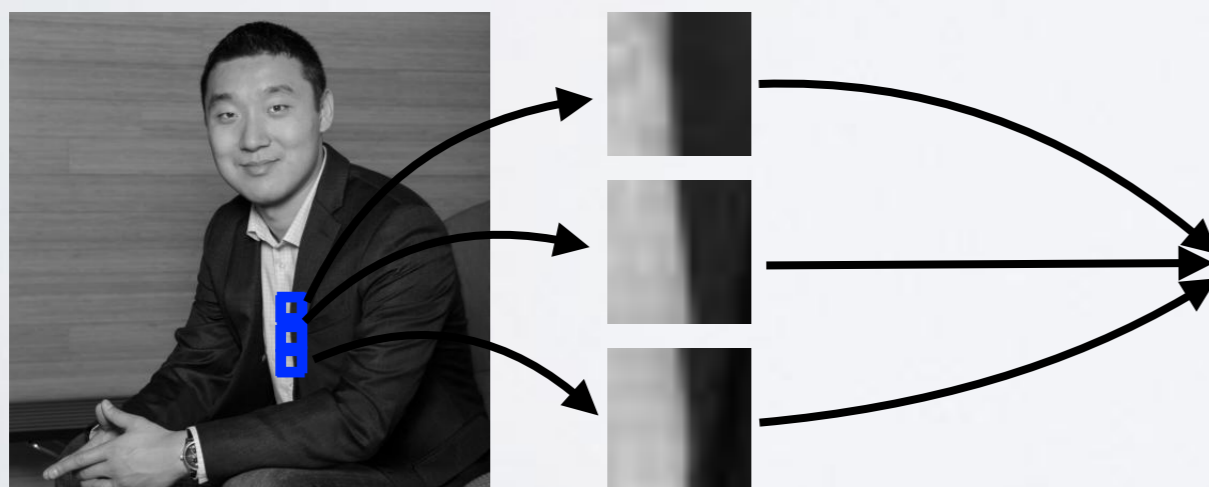
the same weights should apply throughout the input



can use the same weights to detect both edges

only need a single vertical edge detector to find all vertical edges

we can aggregate (*pool*) over a feature to detect whether or not it is present



vertical white/black edge

decrease the spatial size by 'summarizing' the lower-level activations

keep only a relevant summary

builds translation invariance

additional model priors - summary

work in the image-space

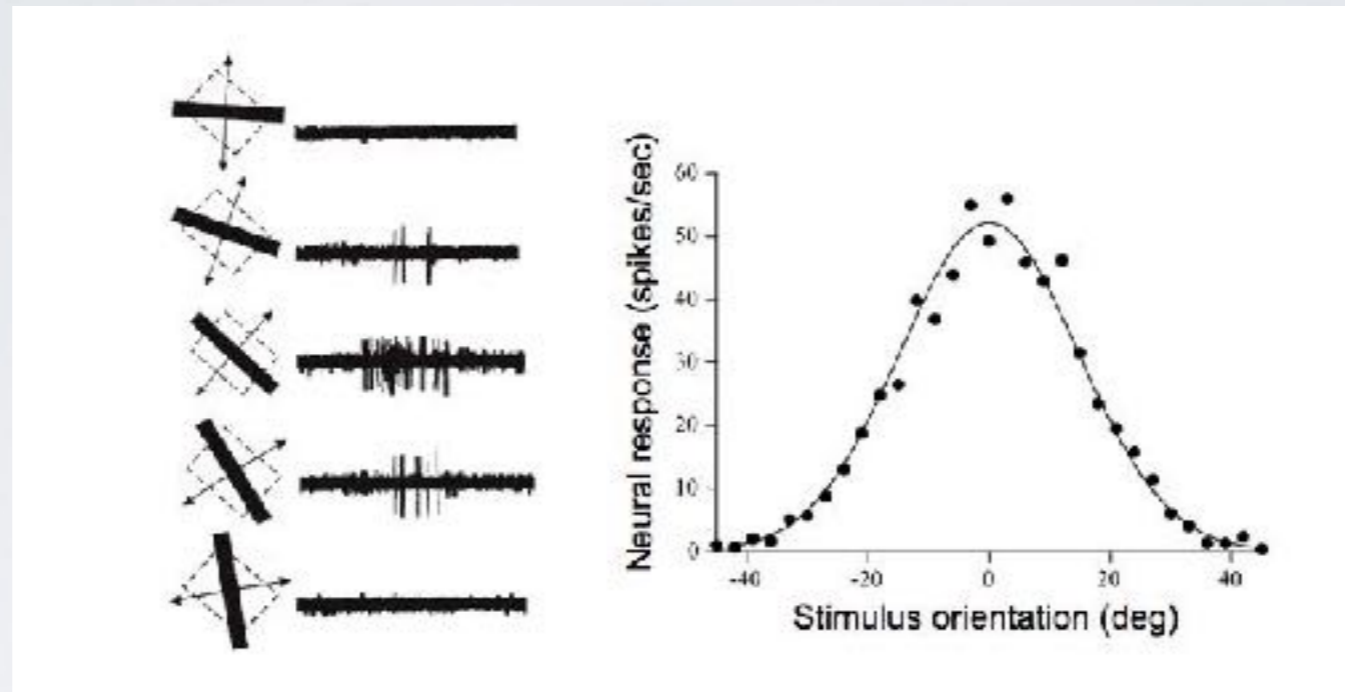
units only take a small window of inputs

weights will be shared across multiple units

pool each feature to create translation invariance

biological inspiration

how do animals recognize visual stimuli?



Hubel & Wiesel - 1950s

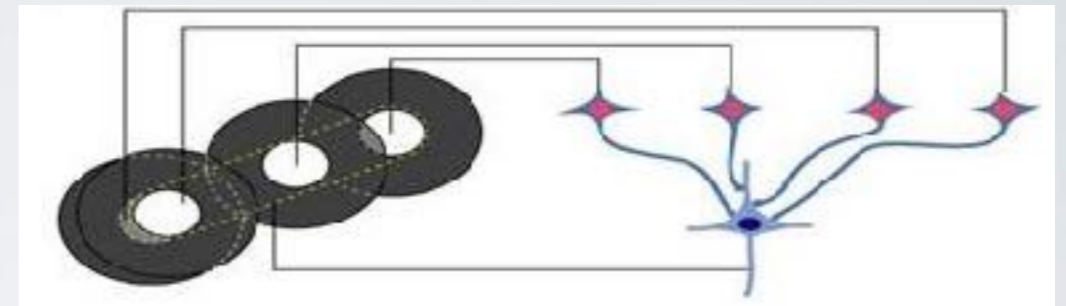
recorded responses of neurons in primary visual cortex (V1) to simple visual stimuli

found neurons selective for bars at a specific orientation at specific locations

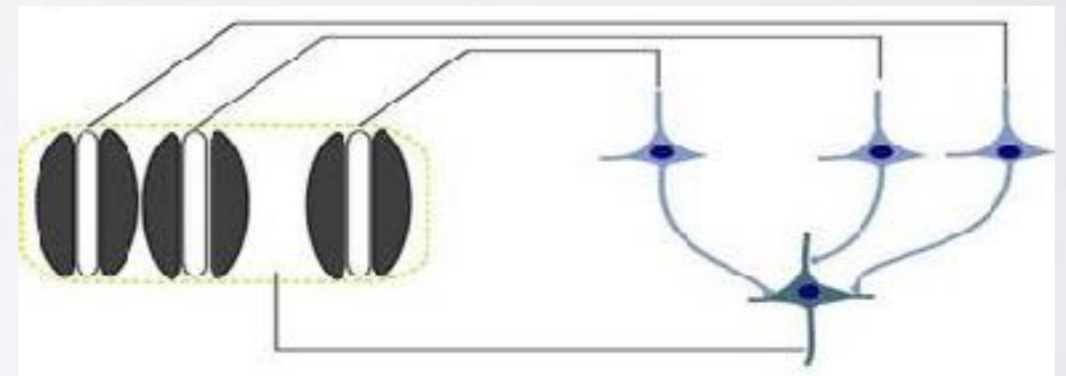
biological inspiration

how do animals recognize visual stimuli?

simple cells combine lower level features (on/off ganglion responses) within a receptive field to select for more complex features

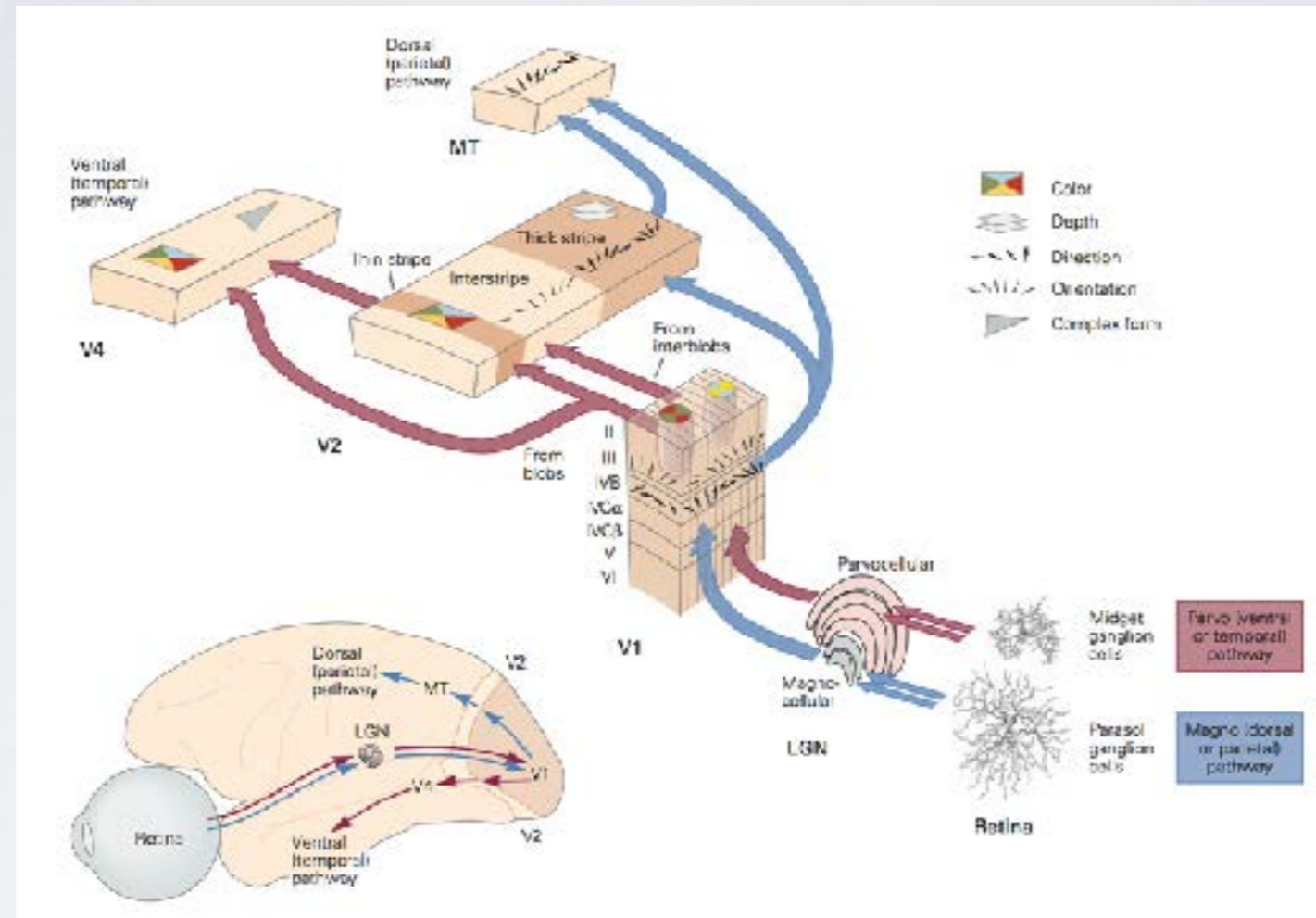


complex cells combine responses from simple cells within a larger receptive field to develop translation invariance



biological inspiration

how do animals recognize visual stimuli?



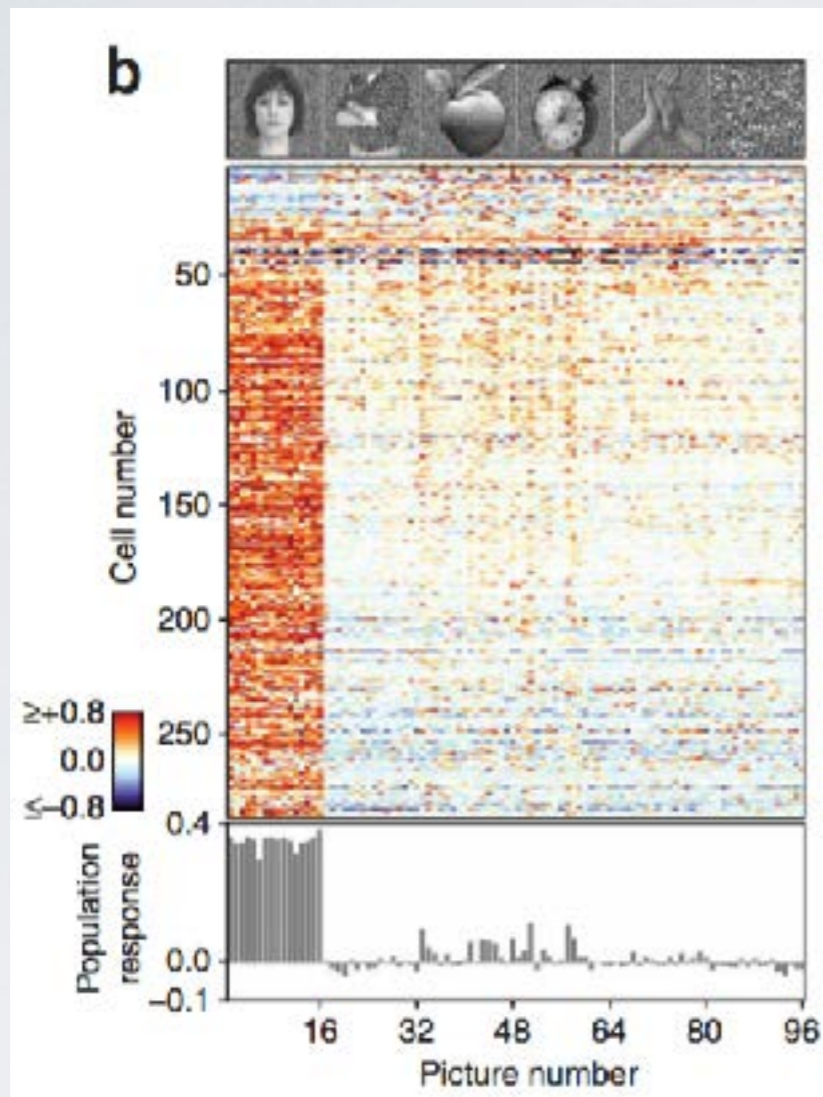
two main pathways:

recognition/what (ventral) pathway & location/where (dorsal) pathway

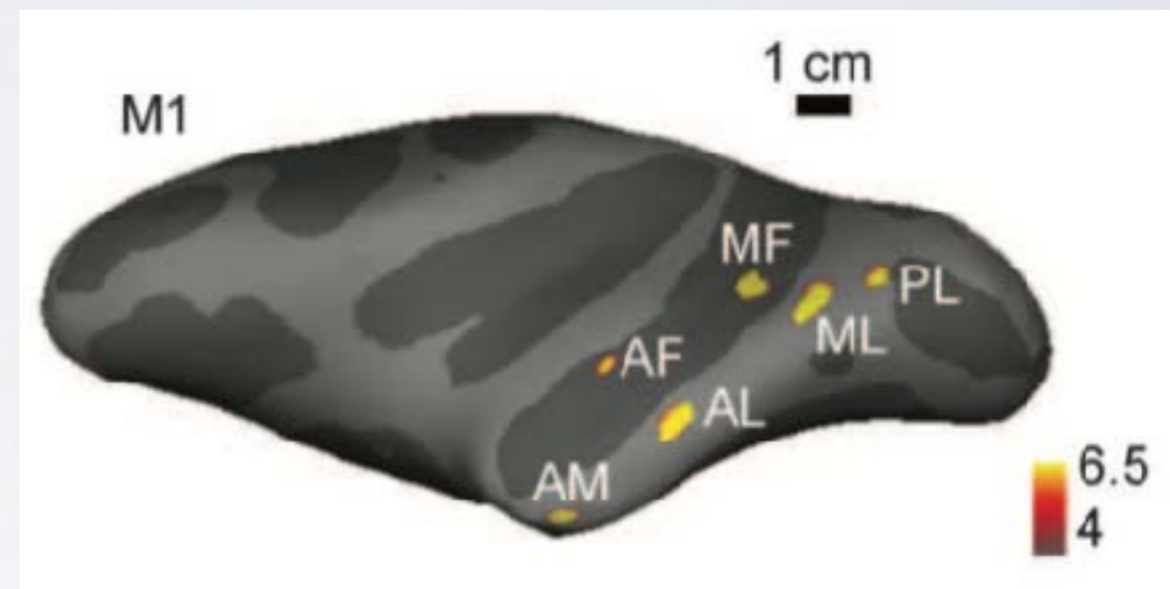
simple visual features are combined hierarchically to select for more complicated visual features

biological inspiration

how do animals recognize visual stimuli?



face 'patches'



areas higher in the recognition hierarchy are selective for highly specific features

these areas tend to be densely interconnected
and relatively invariant to spatial location

CONVOLUTIONAL NEURAL NETWORKS

(discrete) convolution

convolution is a filtering operation

convolve a filter/kernel with the input

Gaussian blur



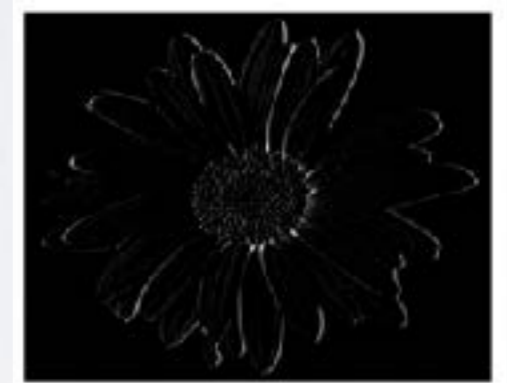
$$* \quad \text{[Gaussian kernel]} \quad =$$



edge detection



$$* \quad \text{[Edge detection kernel]} \quad =$$



sharpening



$$* \quad \text{[Sharpening kernel]} \quad =$$



(discrete) convolution

example:

$$\text{filter weights} = \begin{pmatrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{pmatrix}$$

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

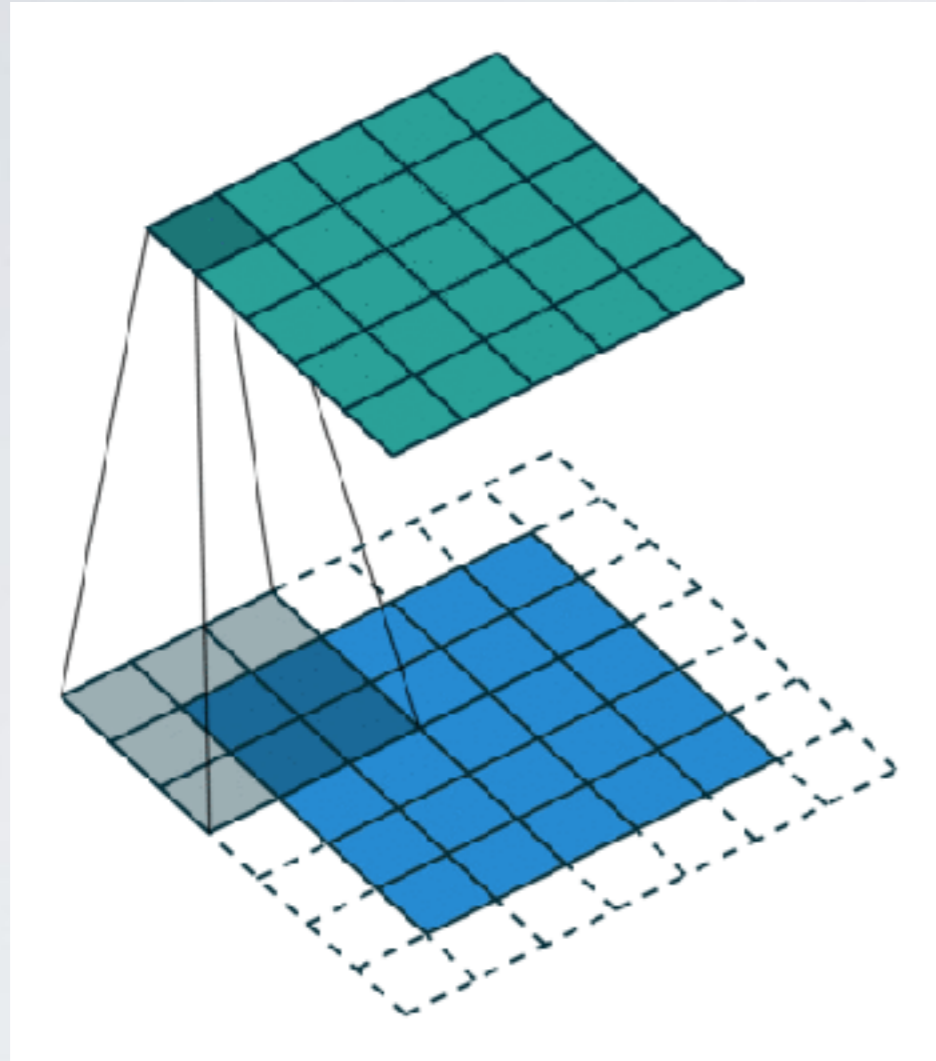
12	12	17
10	17	19
9	6	14

take the *inner-product* of filter weights and input patches across entire input

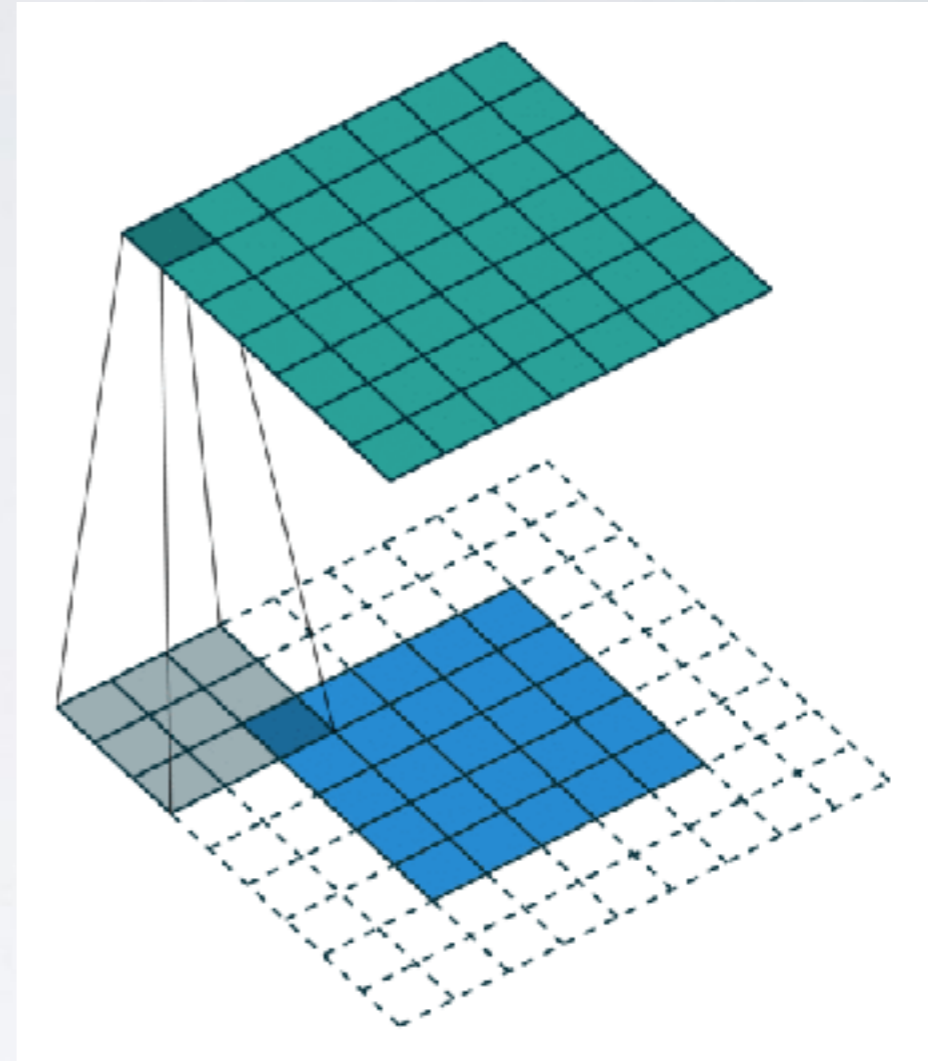
the output is a measure of the degree of the filter feature's presence at each location in the input, known as a **feature map**

padding

half (same) padding



full padding



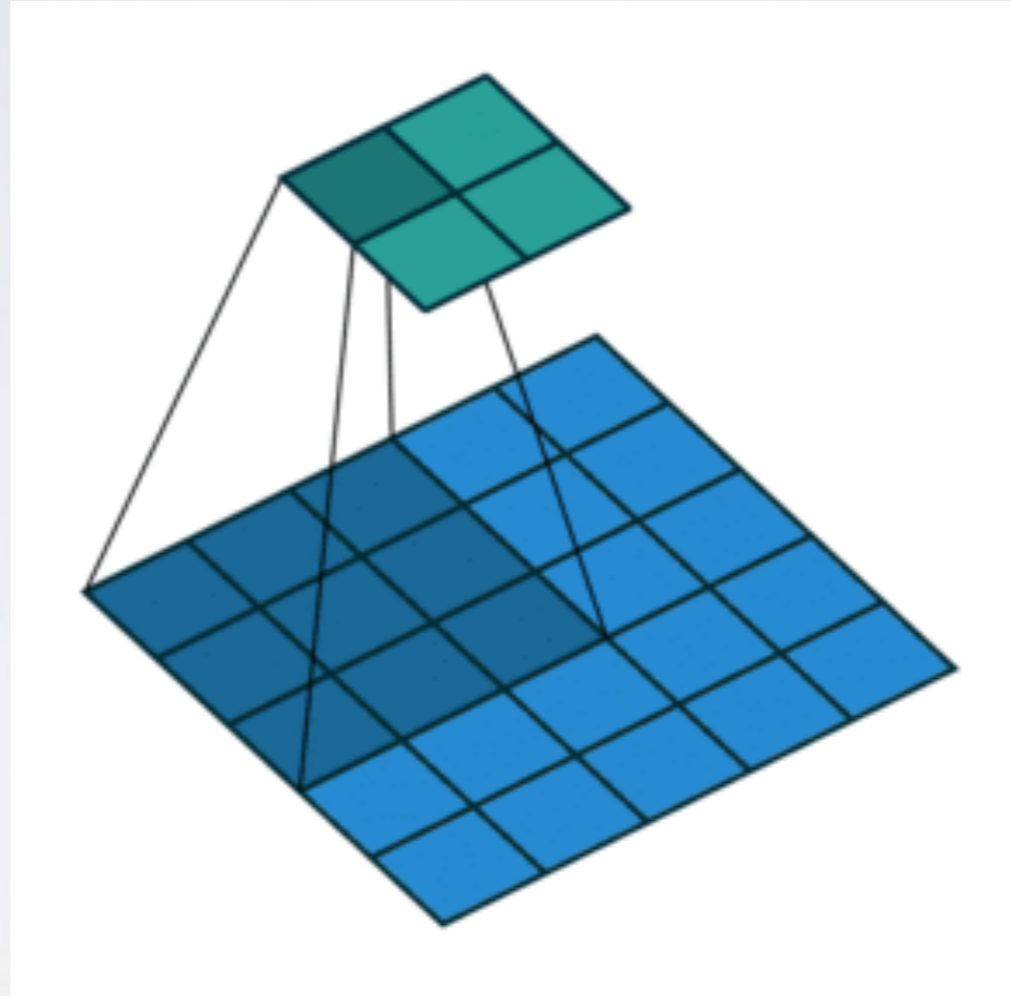
we can **pad** the input with additional values (typically zeros) around the perimeter

this increases the output spatial size in comparison with non-padded convolutions

'same' padding maintains the input size
note that 'valid' padding refers to no padding

stride

vertical stride = horizontal stride = 2

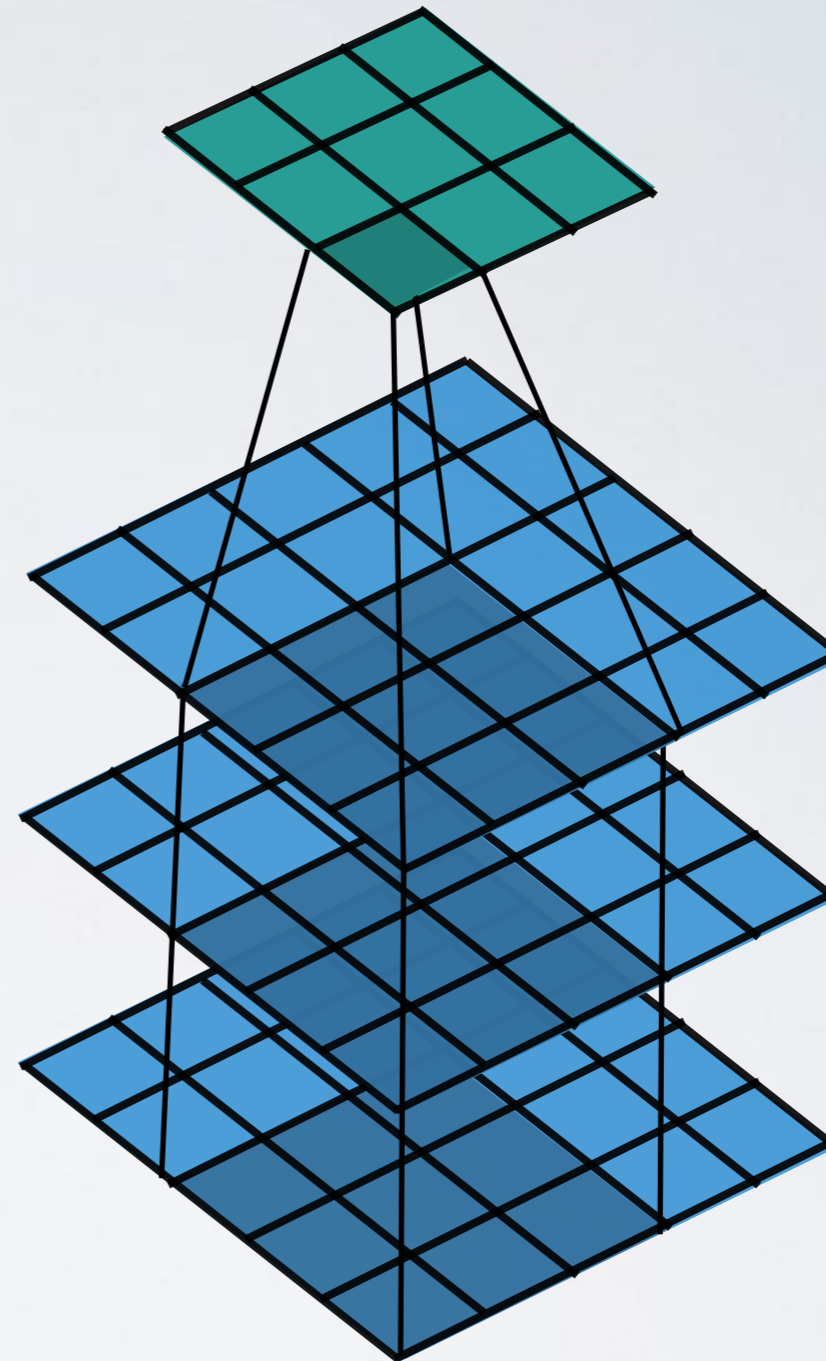


we can apply the kernel with a **stride**, only computing the output at certain integer intervals

this decreases the output spatial size in comparison with non-strided convolutions,
where the stride is one

(discrete) convolutions

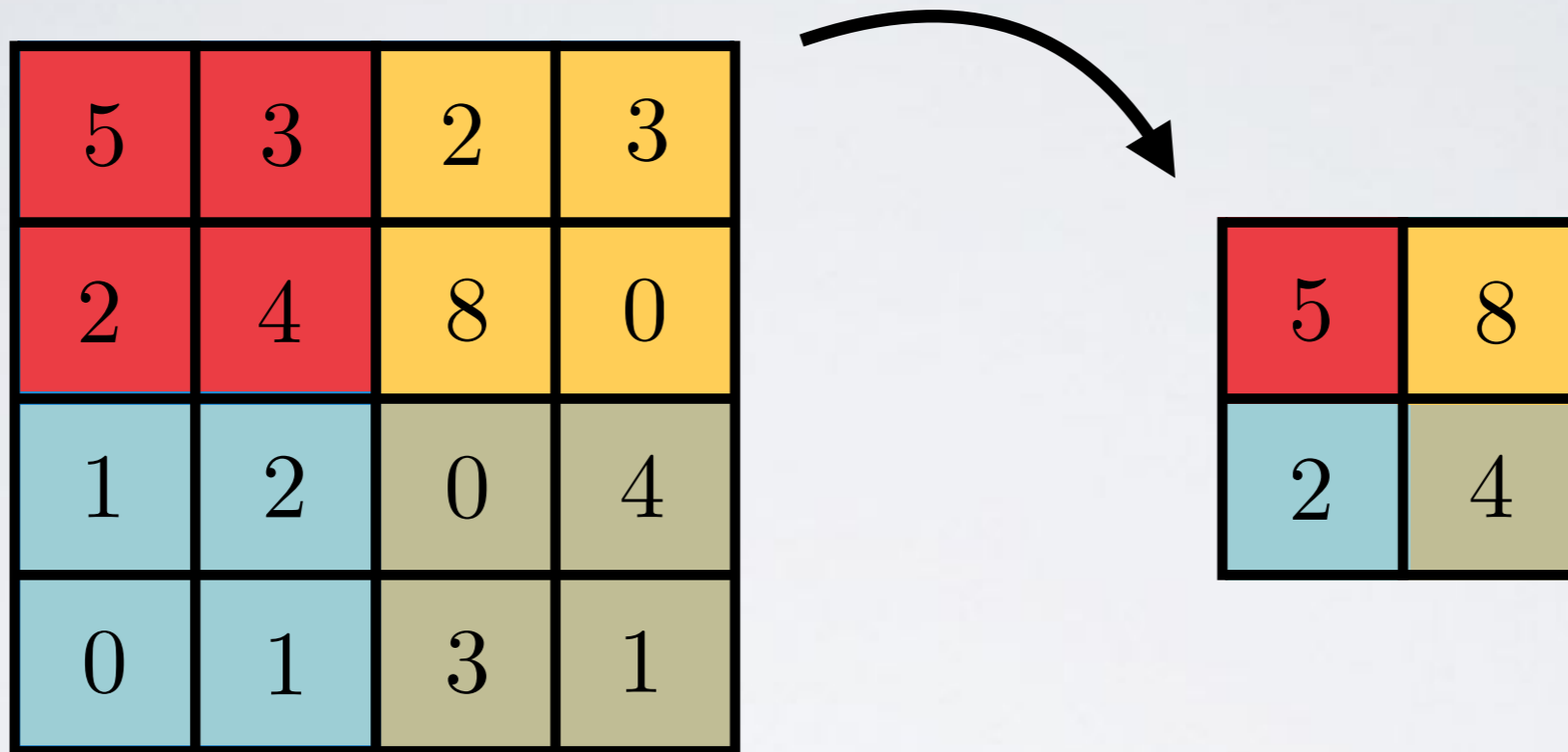
$3 \times 3 \times 3$ filter tensor



convolutions can be applied over *multiple* input feature maps
in this case, instead of being a matrix, *the kernel/filter is a tensor*
can handle RGB images

pooling

example: *max pooling, 2 x 2 window, stride 2, no padding*



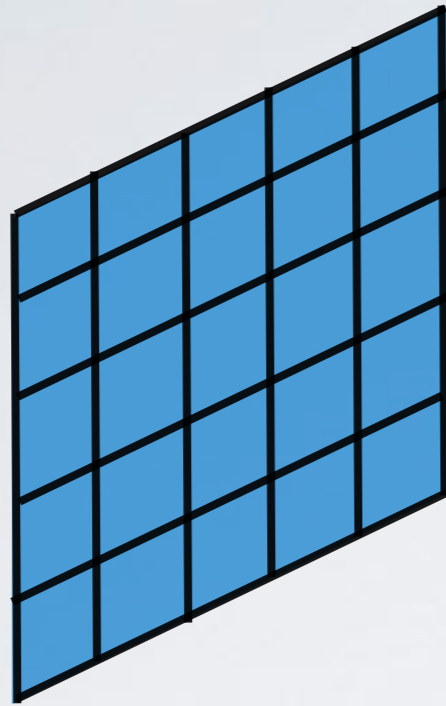
aggregate (pool) over each feature map

pooling is performed over a window

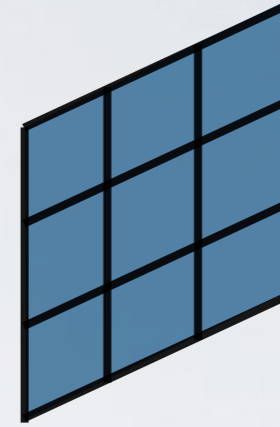
pad and *stride* can also be applied

can use different operations, e.g. max, average, etc.

convolutional pop quiz



5 × 5 input feature map



3 × 3 filter

a 3 × 3 filter is convolved with a 5 × 5 input feature map

if we use *unit strides* and *no padding* ('valid'), what is the output spatial size? 3 × 3

if we use a *stride of 2* and *no padding* ('valid'), what is the output spatial size? 2 × 2

if we use *unit strides* and *half padding* ('same'), what is the output spatial size? 5 × 5

if we use a *stride of 2* and *half padding* ('same'), what is the output spatial size? 3 × 3

fully-connected vs. convolutional

example calculation

convolutional layer with 64 3×3 filters operating on 3 input channels

$$64 \times \text{[3D cube icon]} + 64$$

→ $64 \text{ filters} \times 3 \times 3 \times 3 + 64 \text{ biases} = \underline{1792 \text{ parameters}}$

note that the number of parameters is independent of the spatial size

if the input dimension is 128×128 and the convolution is applied with unit strides and 'same' padding, the convolutional layer will have an output size of $64 \times 128 \times 128 = \underline{1,048,576 \text{ units}}$

to get the same number of output units from a **fully-connected layer** would require $(128 \times 128 \times 3) \times 1,048,576 + 1,048,576 = \underline{51,540,656,128 \text{ parameters}}$

if the assumptions we made about images (locality, translation invariance) are valid, then we have reduced the number of parameters by a factor of over 10 million

convolutional layers allow us to trade off *flexibility* for an *increased representation size*

image datasets



ImageNet

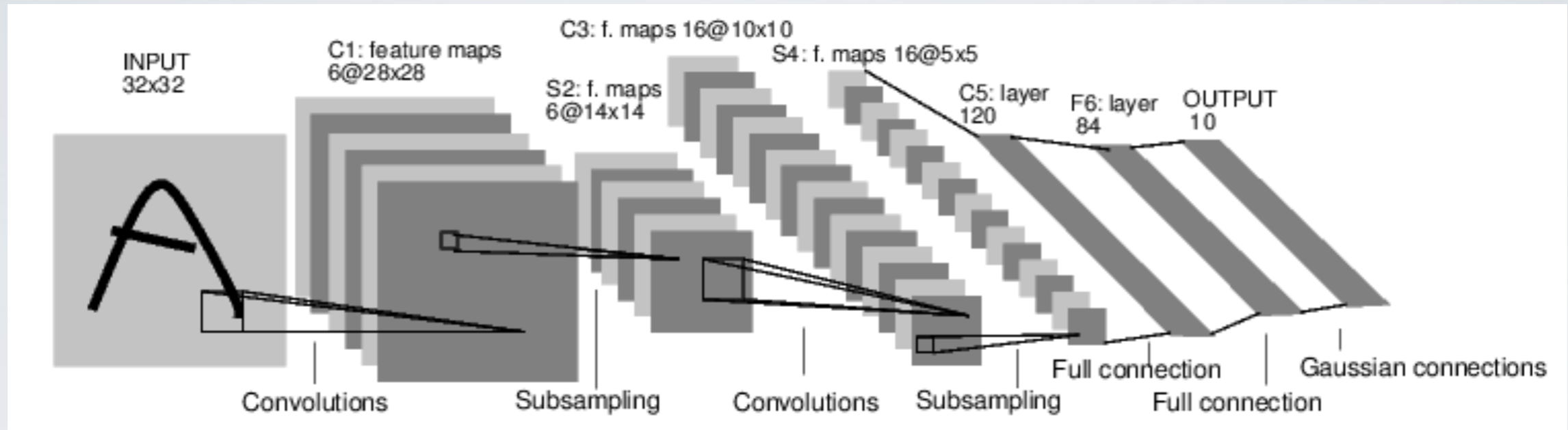
Competition (ILSVRC) Dataset

1,000 classes,
1.2 million images

Full Dataset

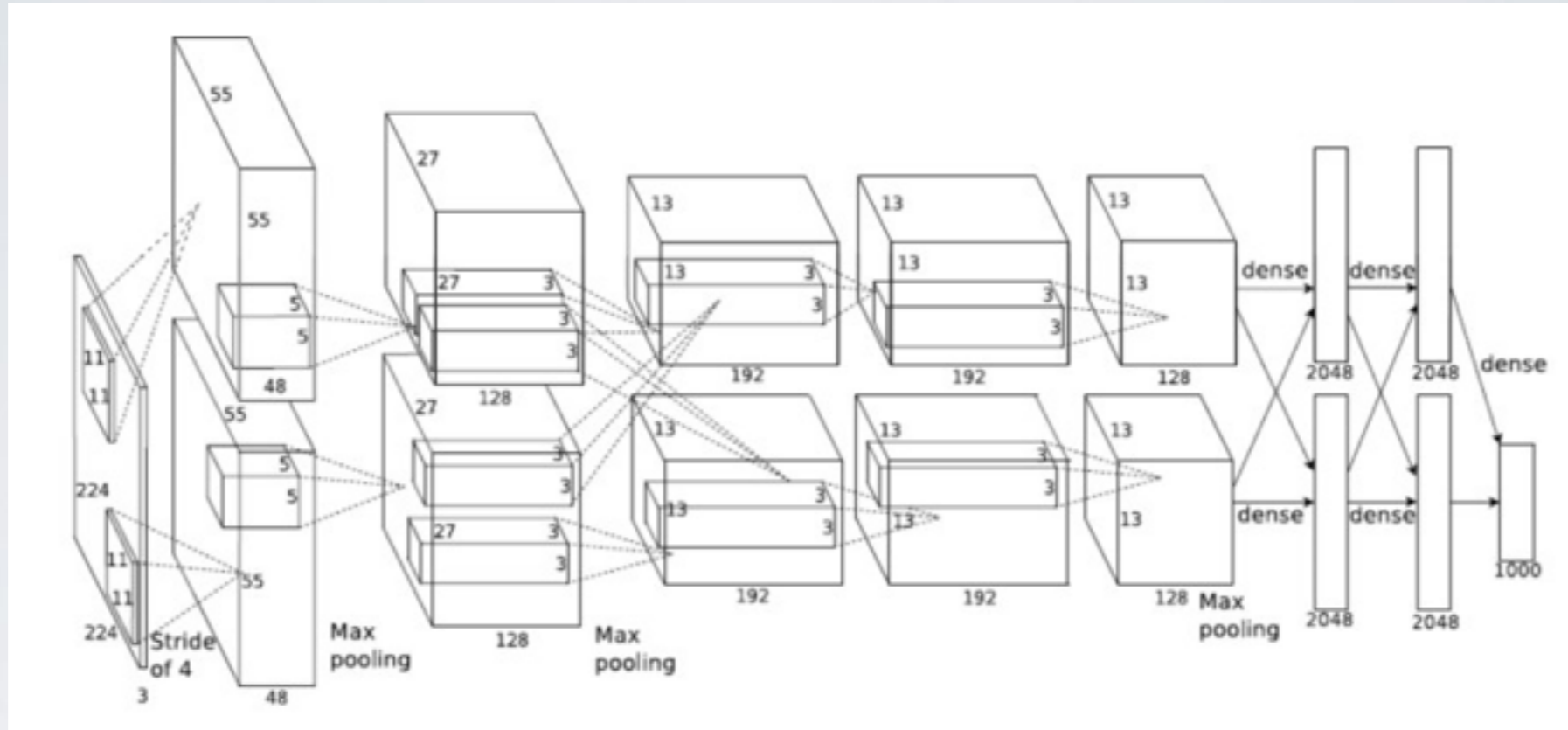
21,841 classes,
14 million images

models



LeNet - 1998

models



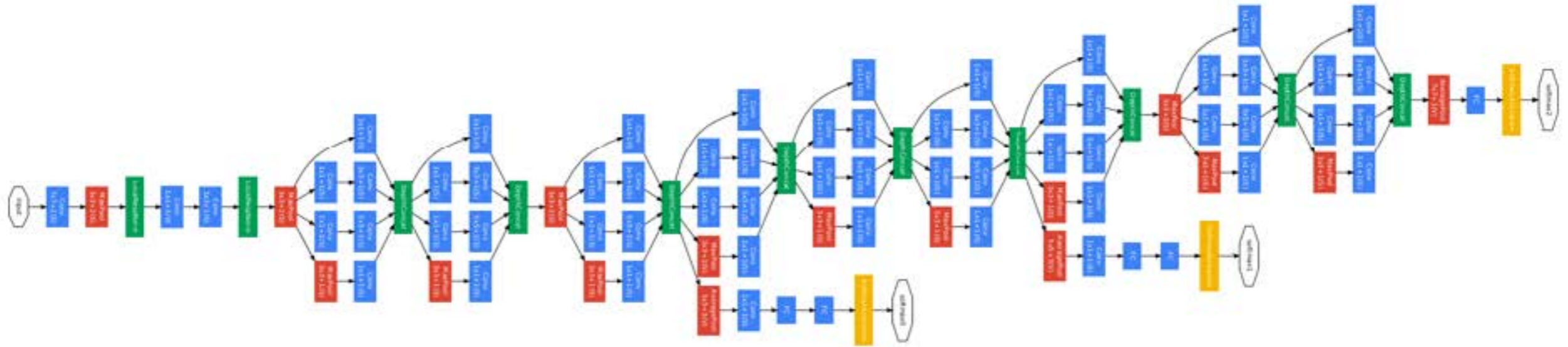
AlexNet - 2012

models



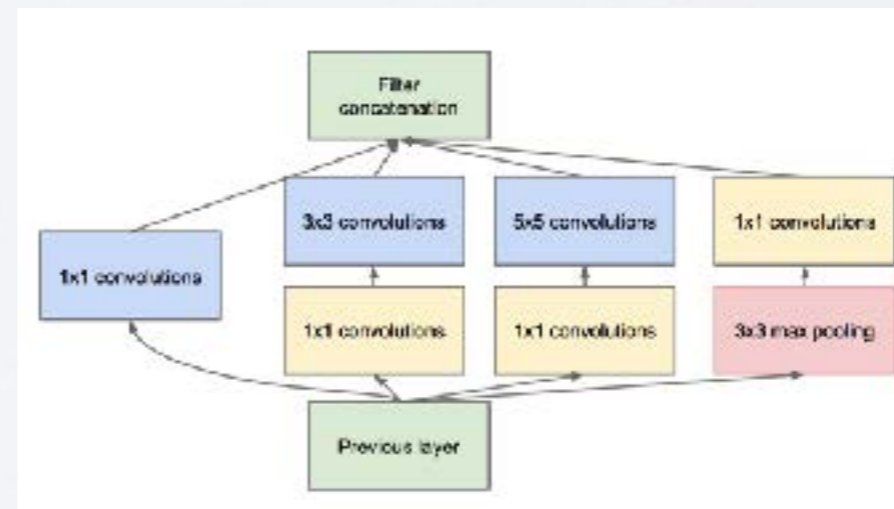
VGG - 2014

models

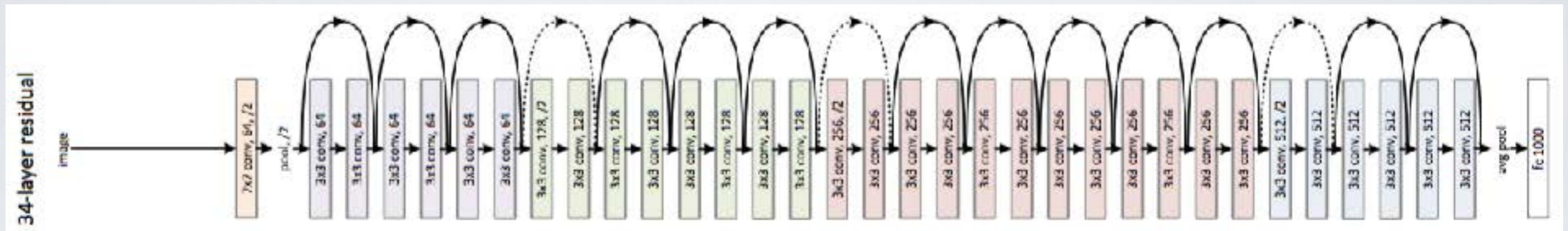


GoogLeNet - 2014

Inception Block:

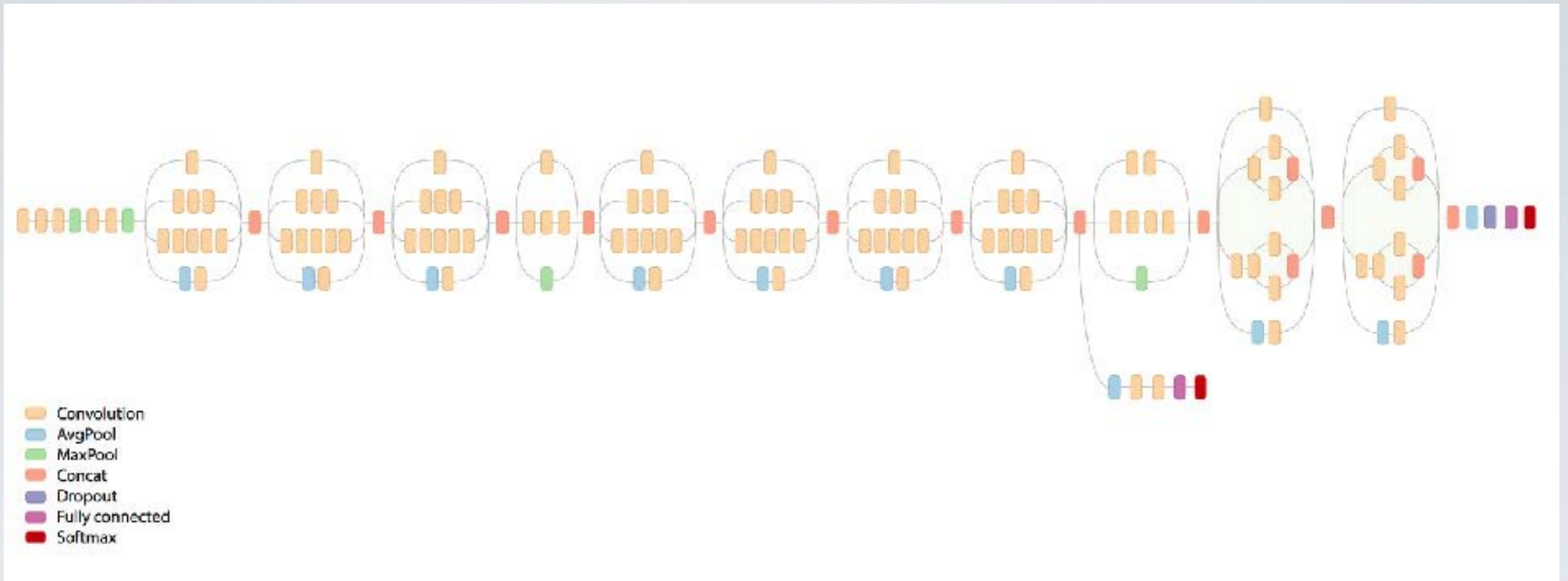


models



ResNet - 2015

models



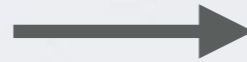
Inception v4 - 2016

top-5 error on ILSVRC

top-5 error denotes the percentage of examples for which the ground truth label is not in the model's top 5 predictions

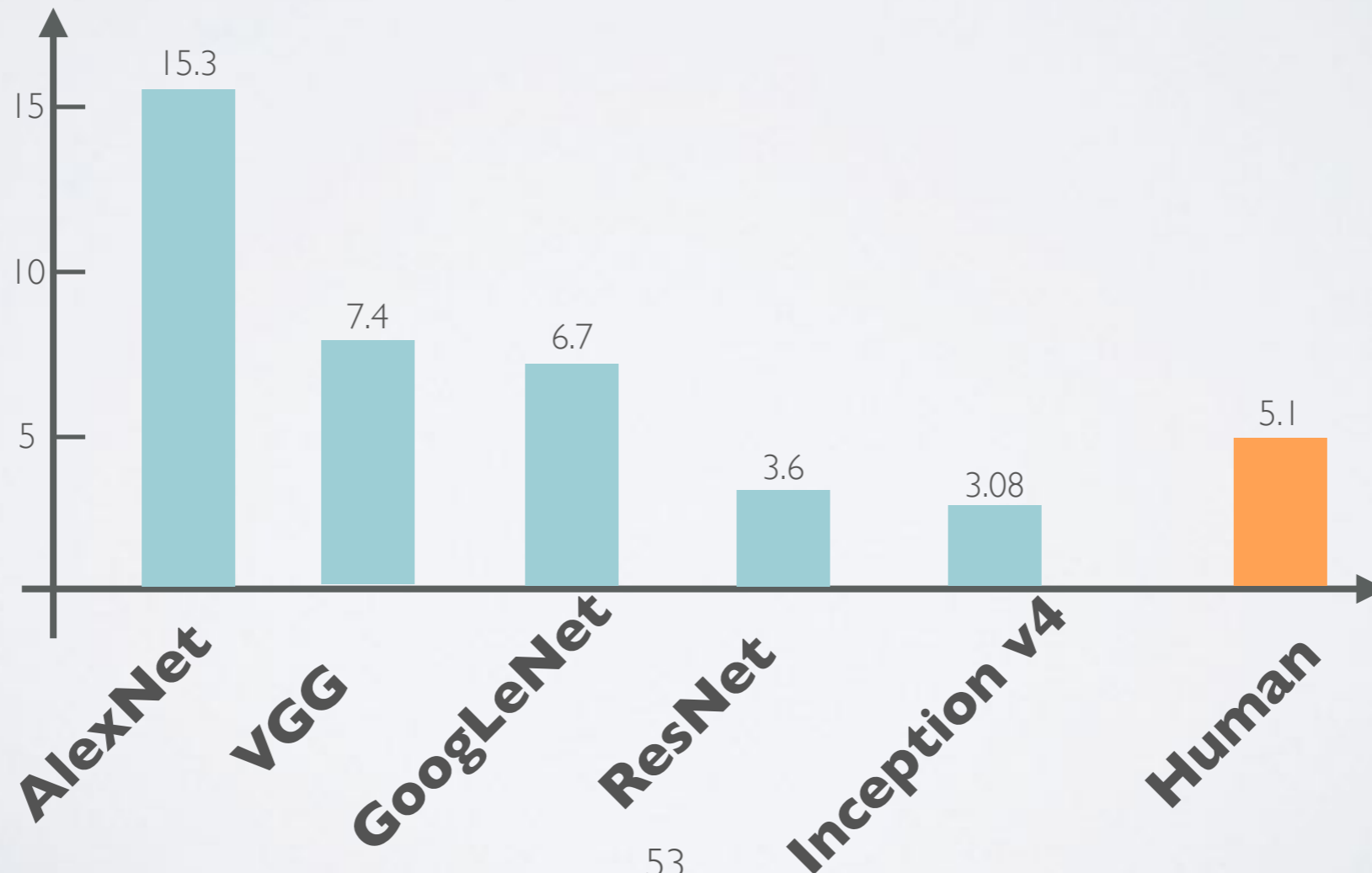


1. television
2. radio
3. walkie talkie
4. pager
- 5. cell phone**



correct top-5 prediction

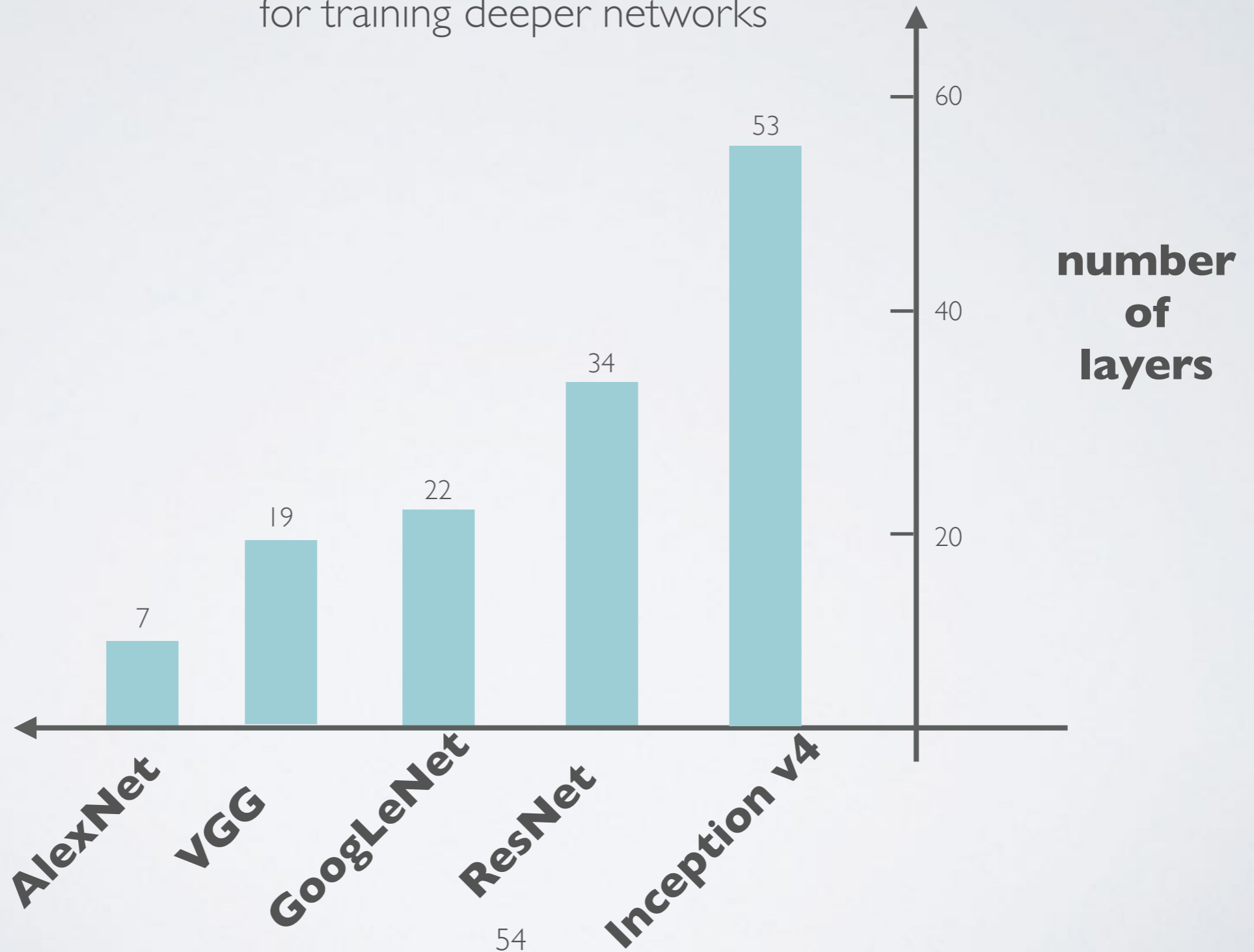
top-5 error (%)



number of layers

the number of layers with weights in state-of-the-art networks has grown

this is primarily due to new tricks that have been developed for training deeper networks



number of layers

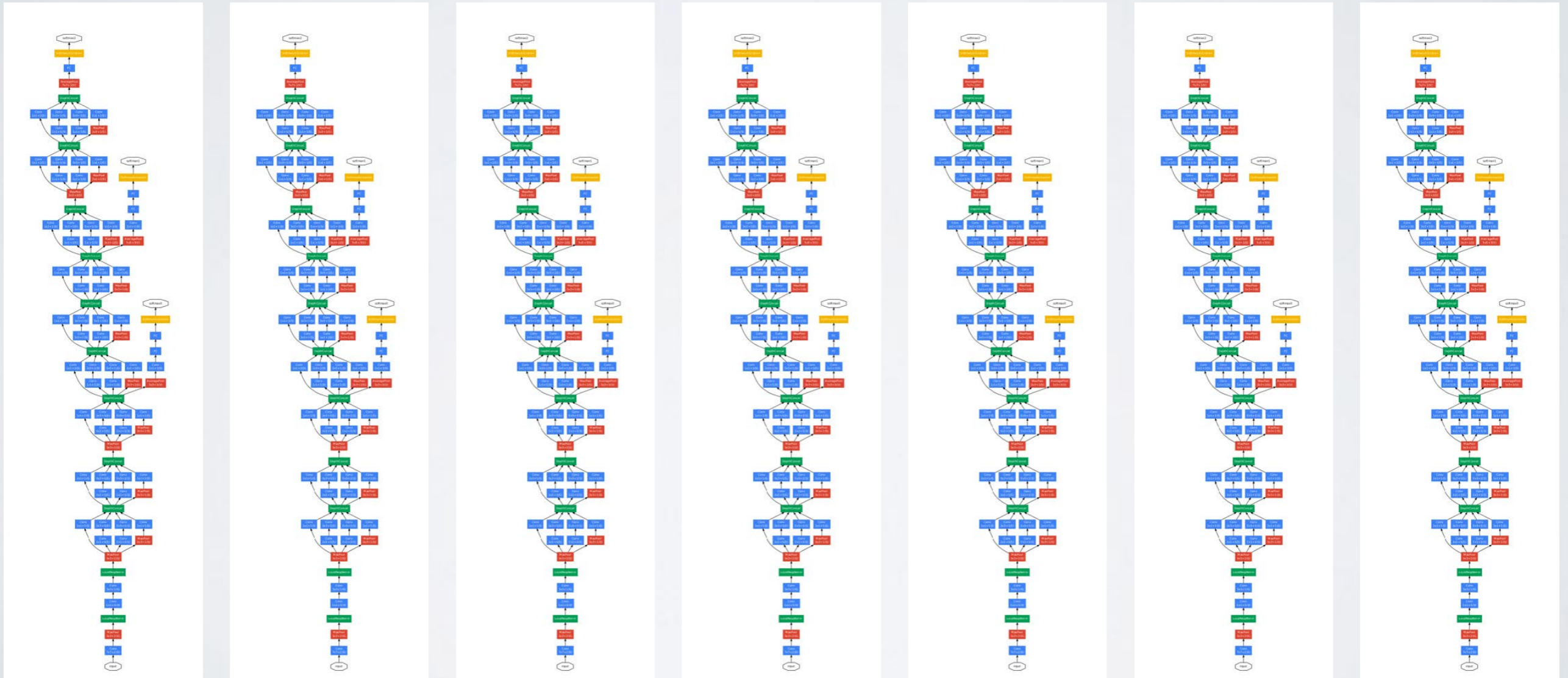


ensembles

as a side note:

the winning entries are typically *ensembles* of networks

since each network likely converges to a different local minimum,
averaging their predictions helps in generalization



data augmentation

in training these networks, people often use *data augmentation* to effectively boost the number of examples in the training set

we use priors on the nature of the data to create additional examples

example data augmentation:

original image



crops



rotations



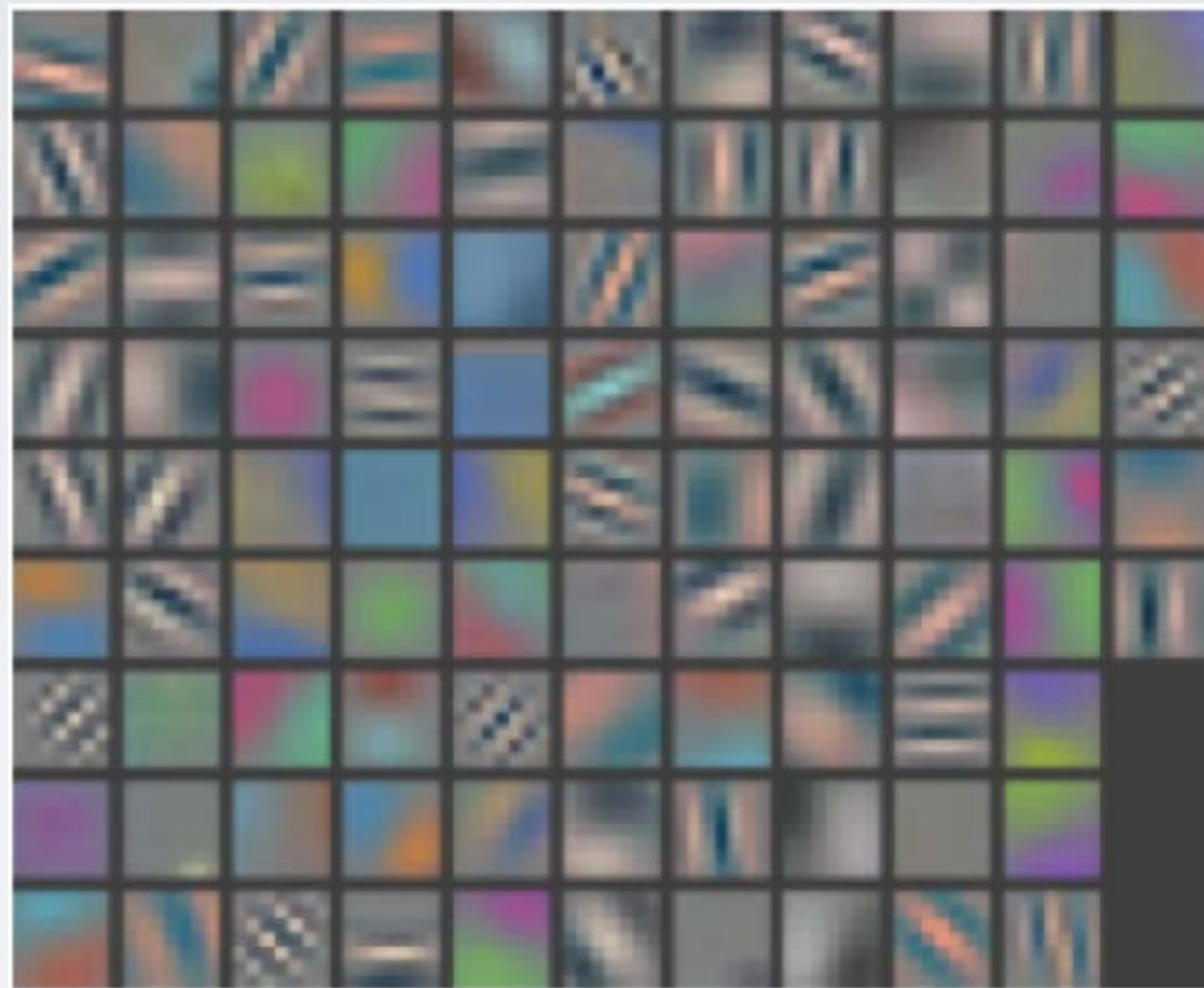
left-right flip



what are these models learning?

filter visualization

at the first layer, we can visualize the filters as RGB images



filter visualization

for later layers, we can no longer visualize them directly in the image domain

there are three main ways in which to visualize them:

maximal images from dataset

- feed in all images from the dataset and see which images make the filter activate maximally

deconvolution

- run the network in reverse from an intermediate layer to try to convert its activation back to the image domain

optimized image

- backpropagate from a filter to the image itself to find an image that would make it fire maximally

filter visualization



filter visualization



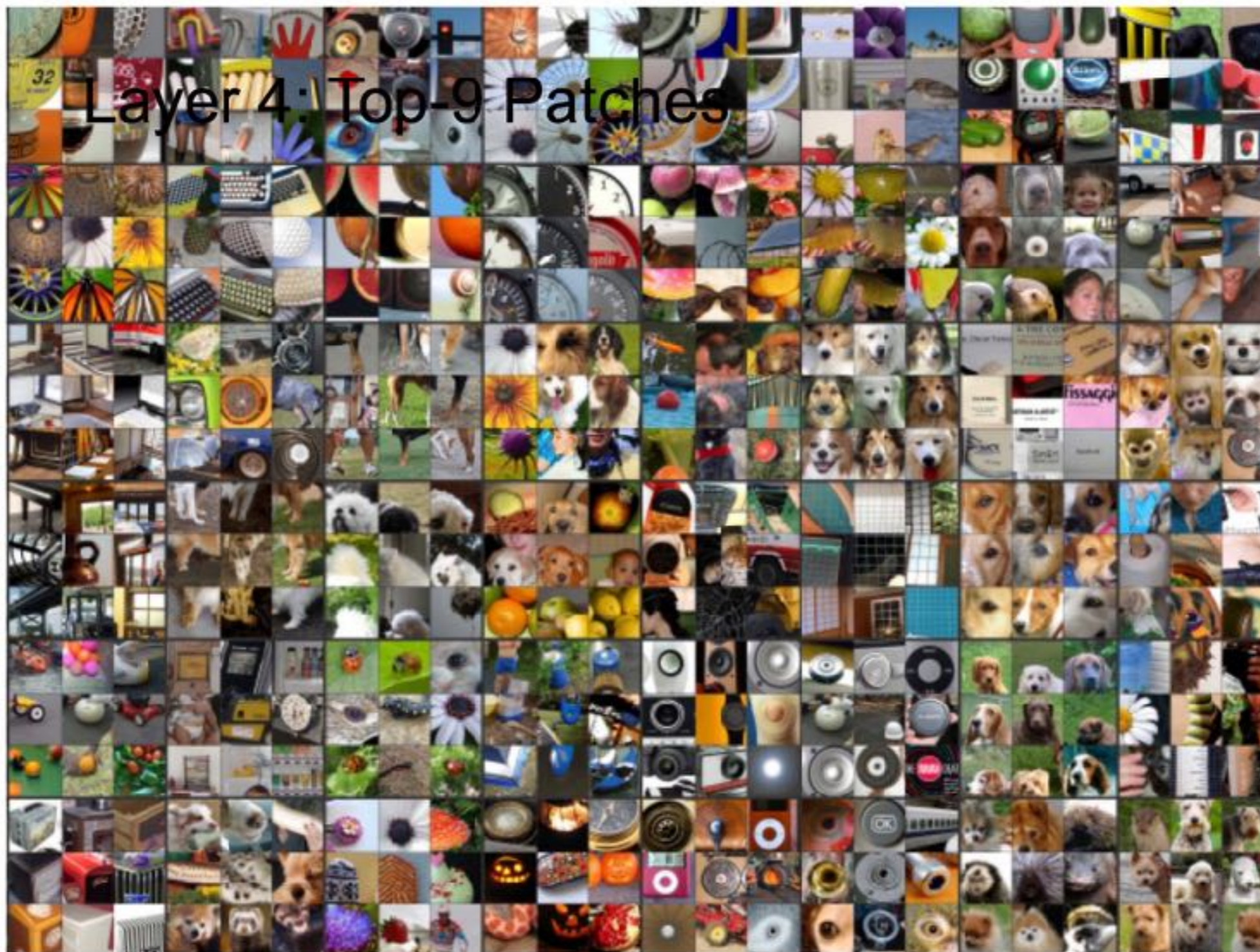
filter visualization



filter visualization



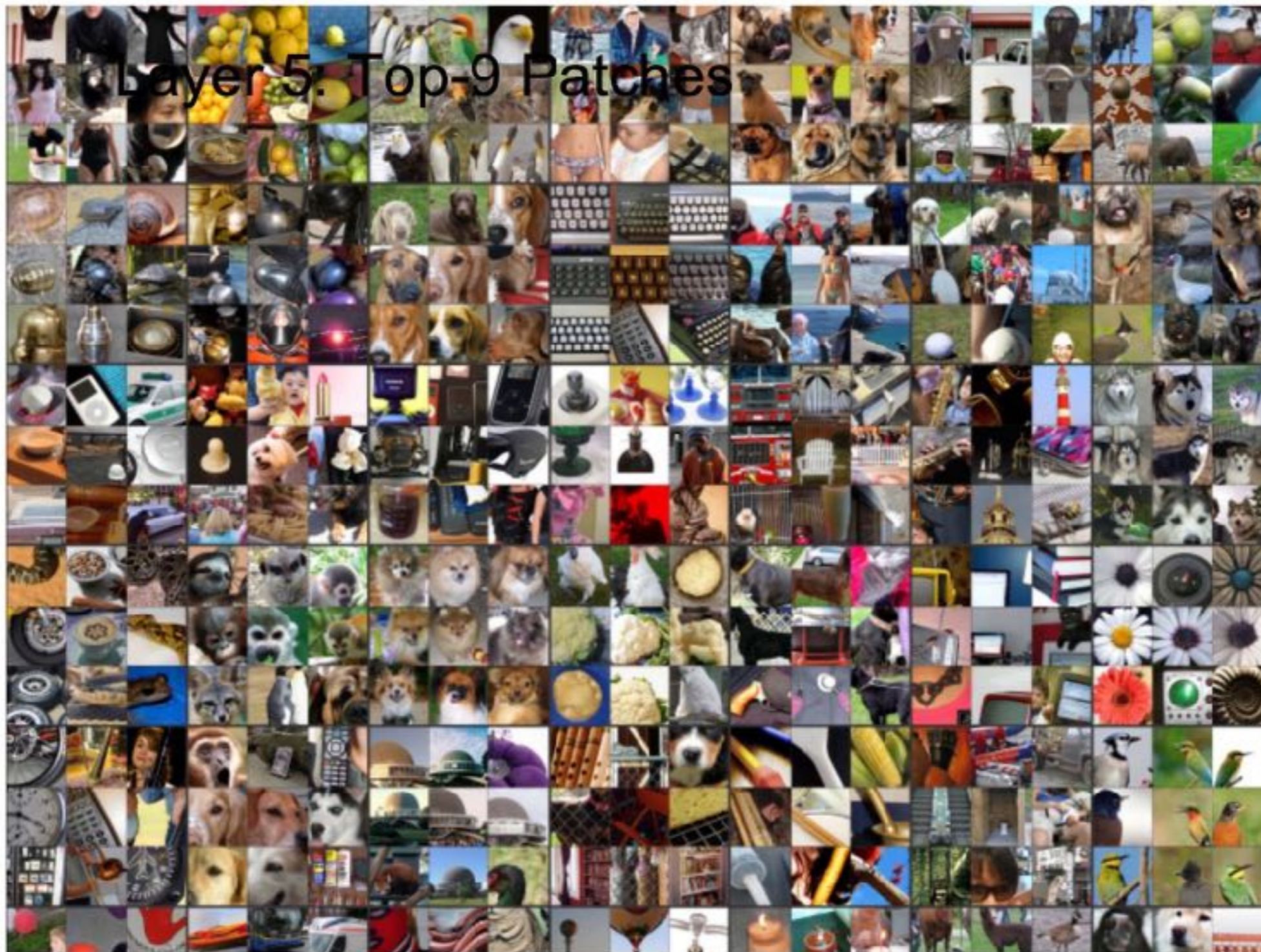
filter visualization



filter visualization



filter visualization

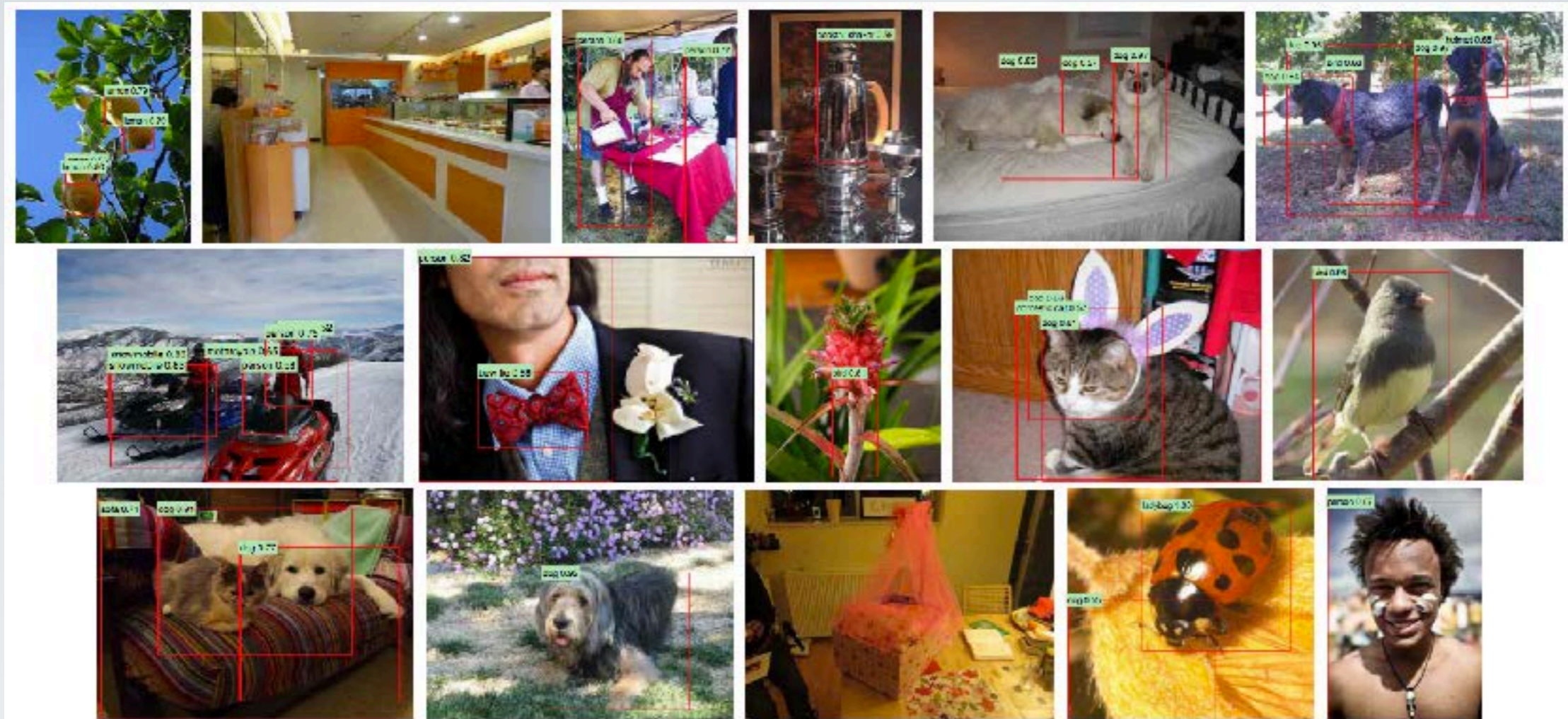


filter visualization

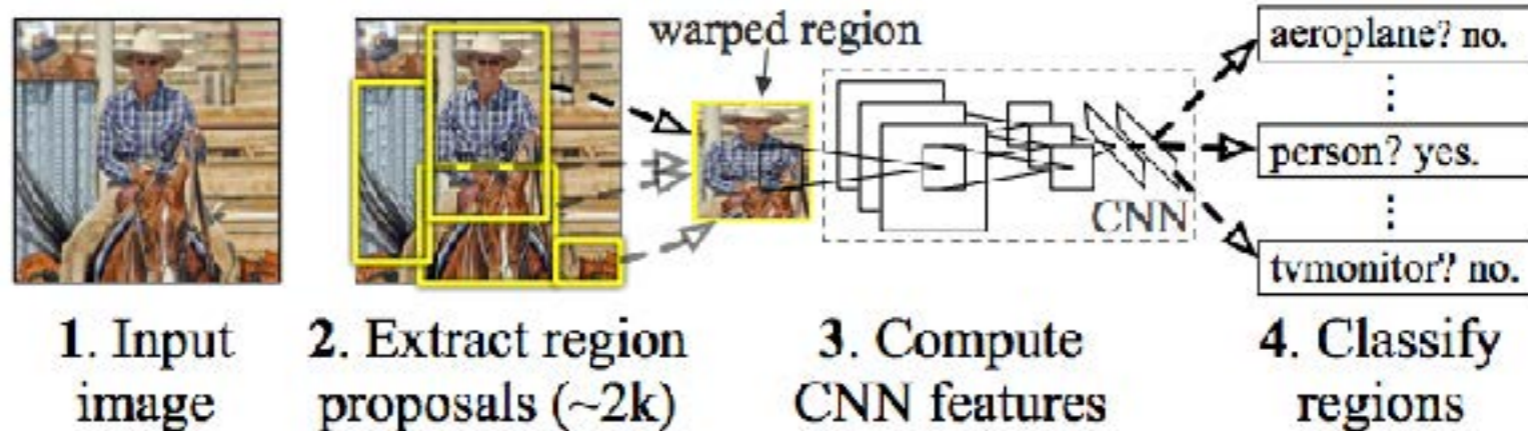


moving beyond object recognition

object detection

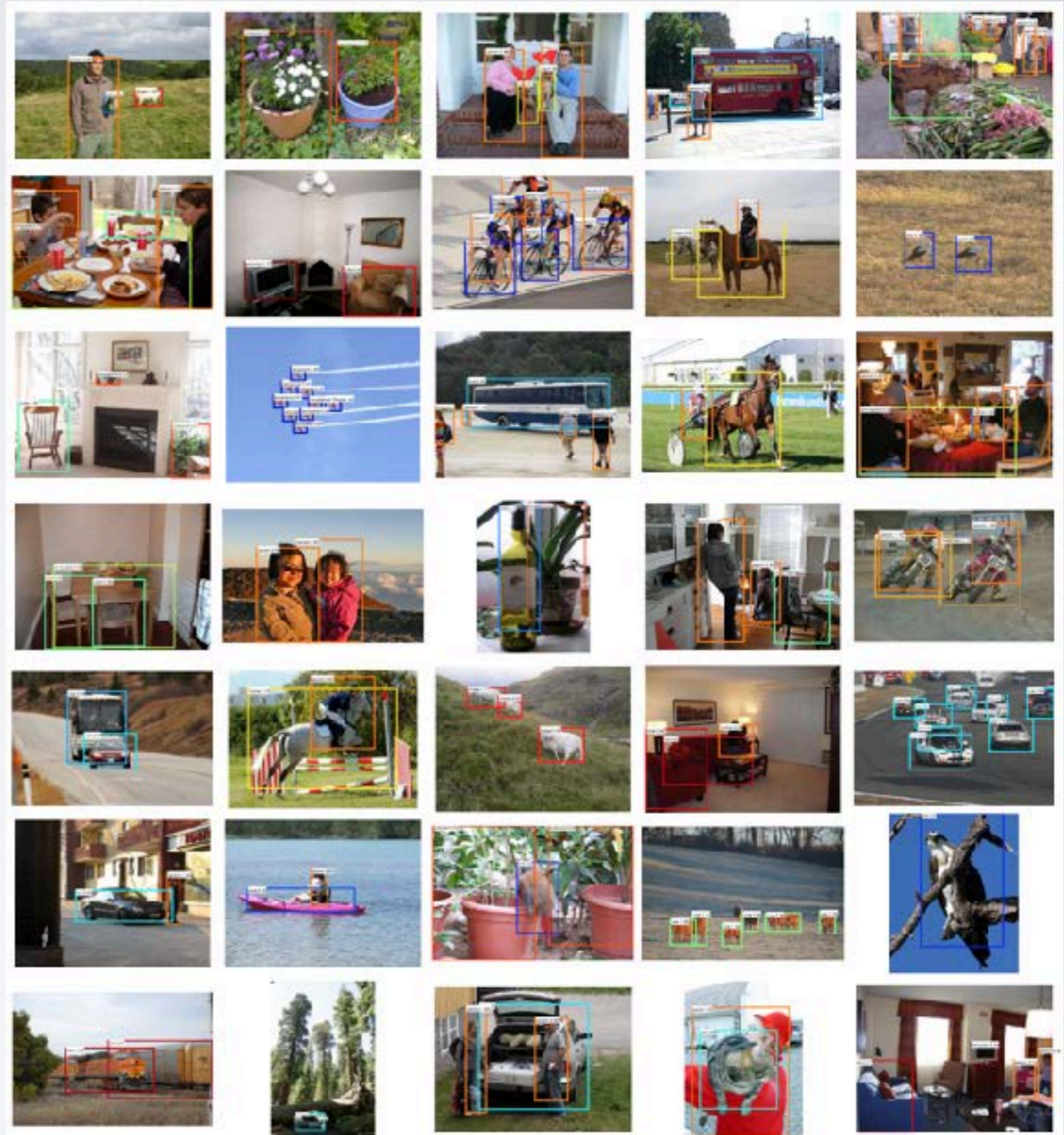
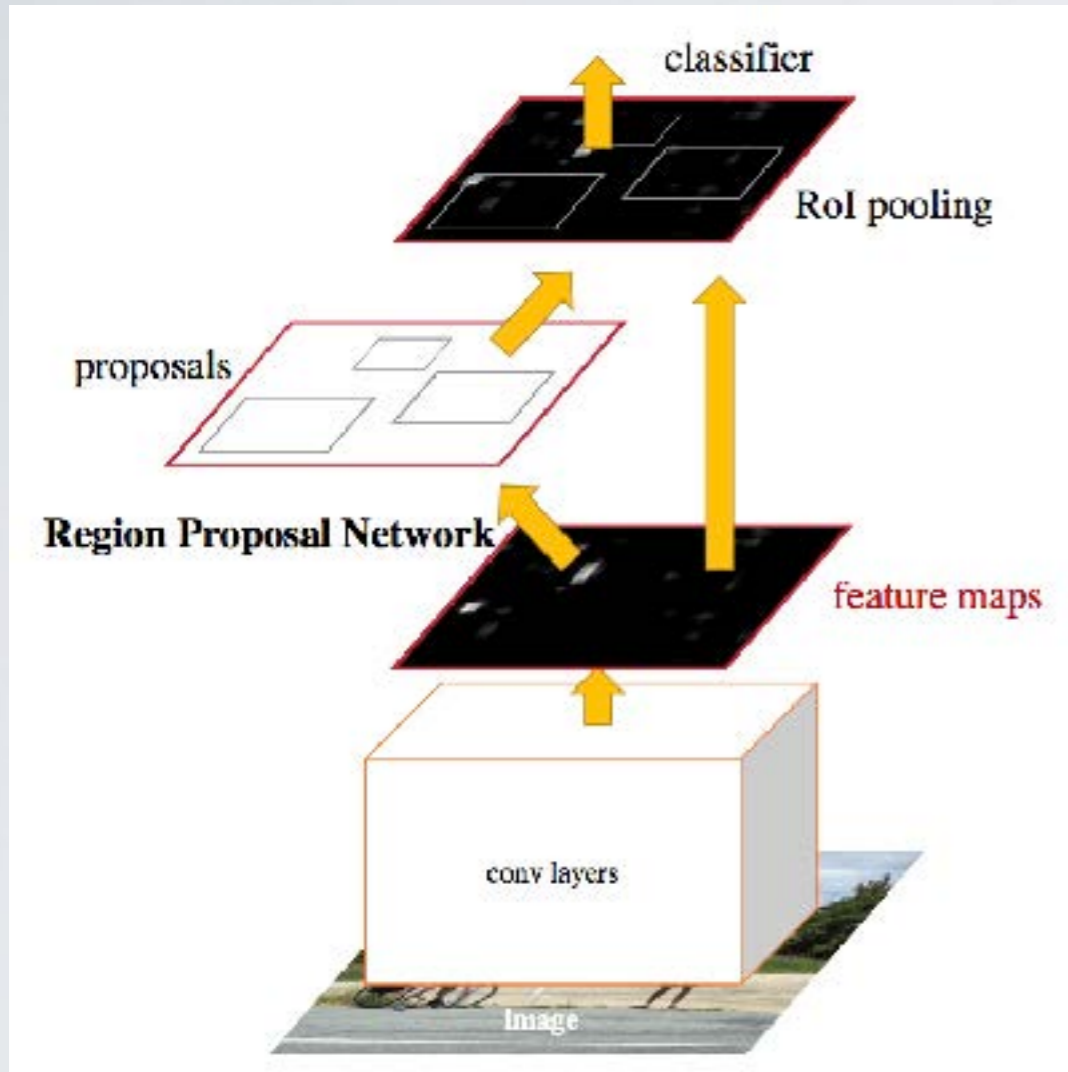


R-CNN: *Regions with CNN features*



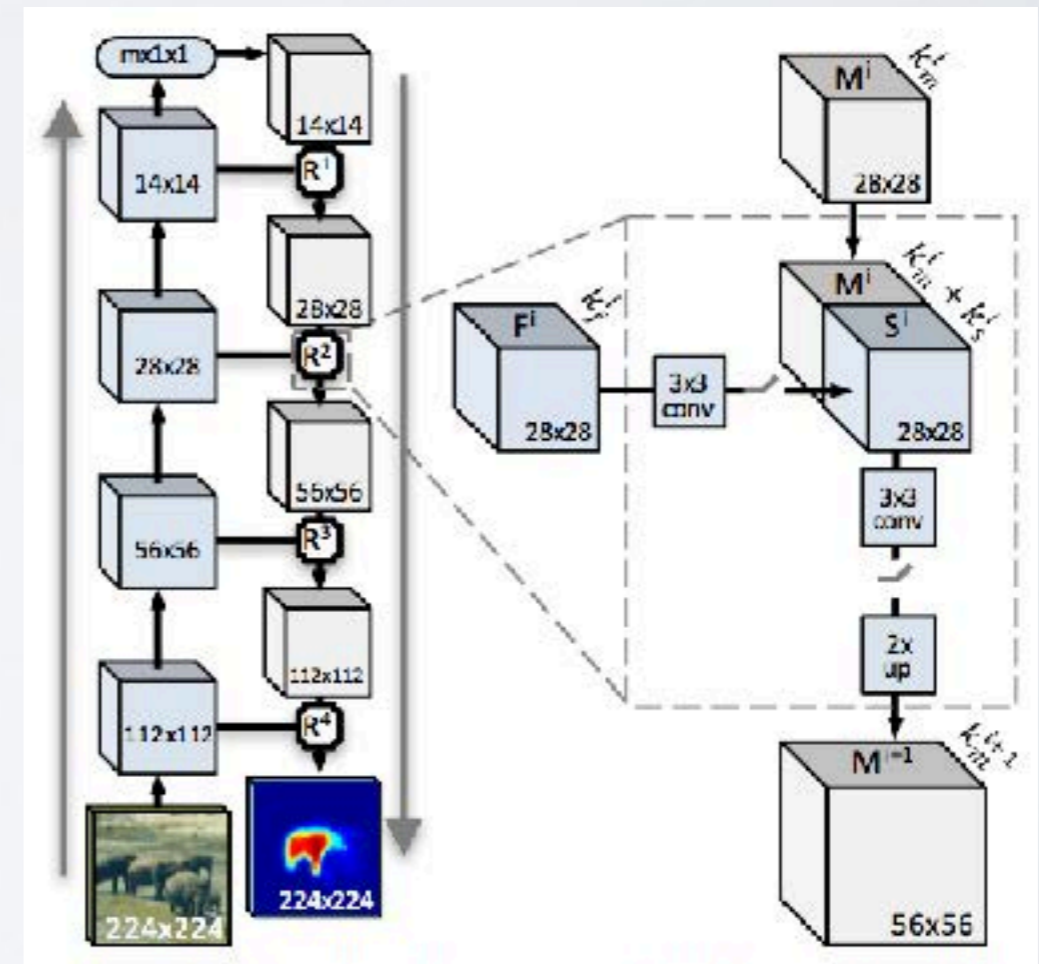
object detection

Faster R-CNN

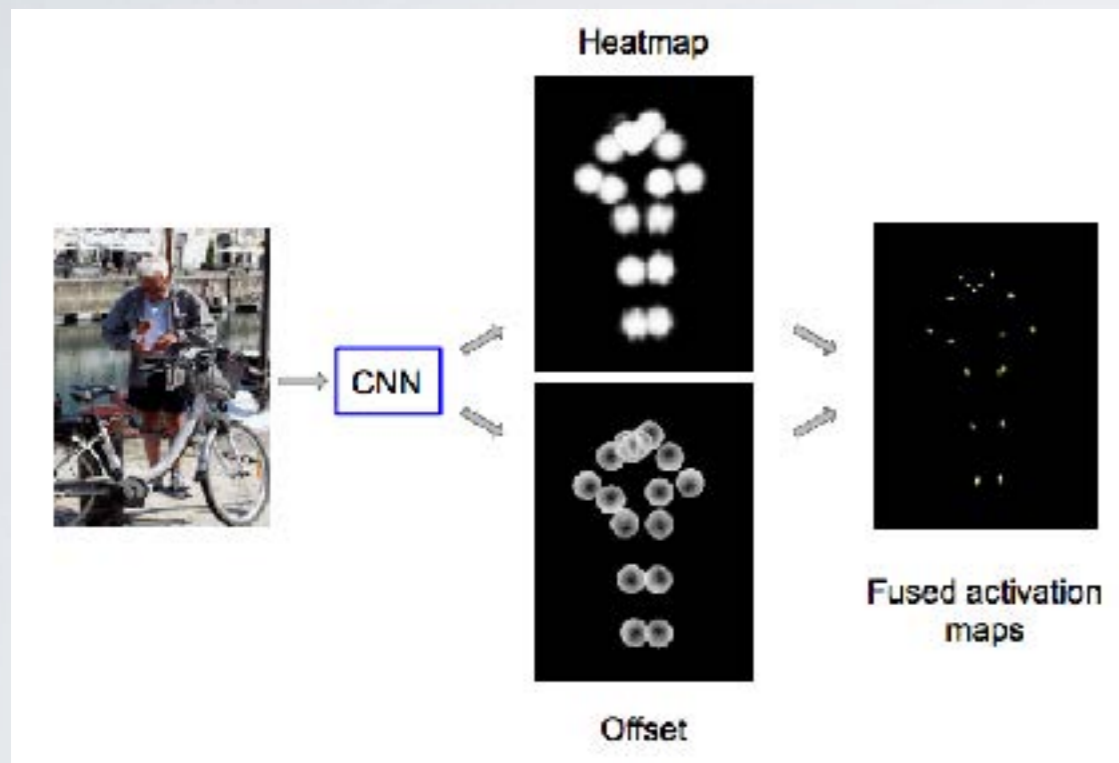


object segmentation

The collage displays 24 different images, each with colored masks identifying objects. The objects include zebra-like patterns, a red bus, a beach scene, a field with animals, a fruit basket, a beach with umbrellas, a person at a bar, a public restroom, a field with animals, a person in a field, a street scene with a stop sign, a field with animals, a public restroom, a street scene with a car, a field with animals, a street scene with a car, a field with animals, a street scene with a car, a field with animals, a street scene with a car, a field with animals, and a street scene with a car.



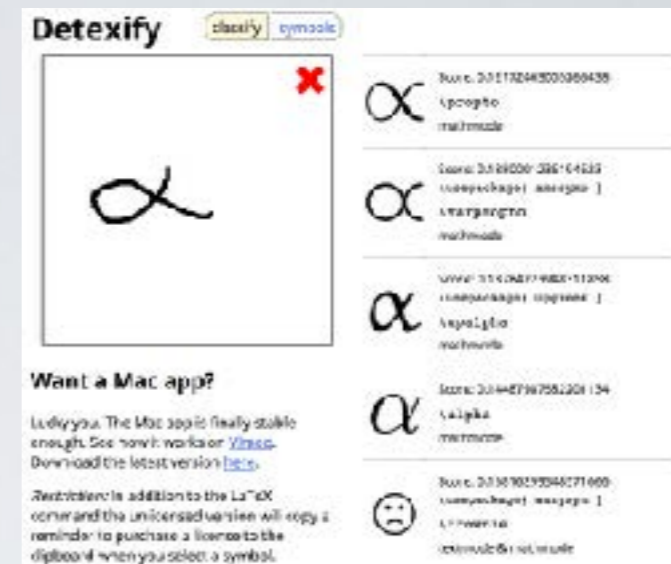
key point estimation



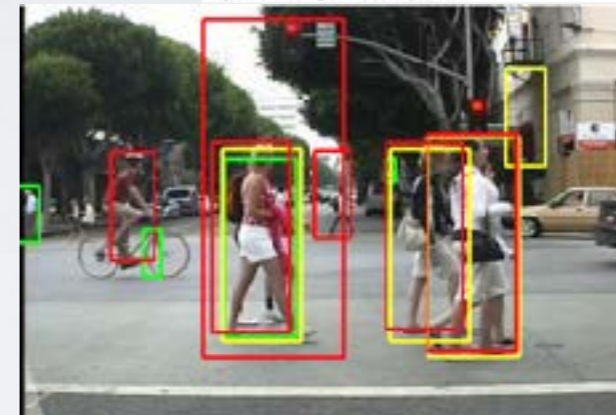
applications

handwriting recognition

- ATM
- note taking

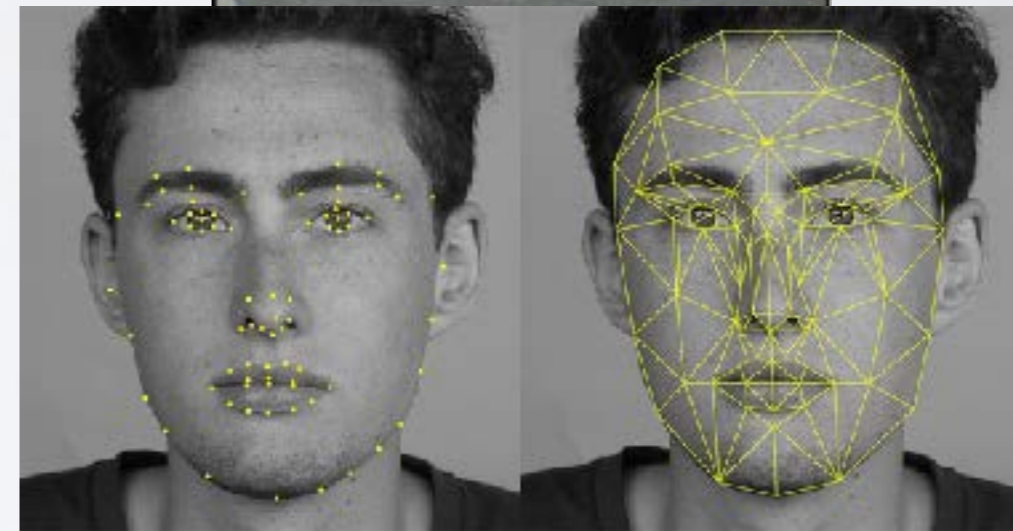


pedestrian/object detection for self-driving cars



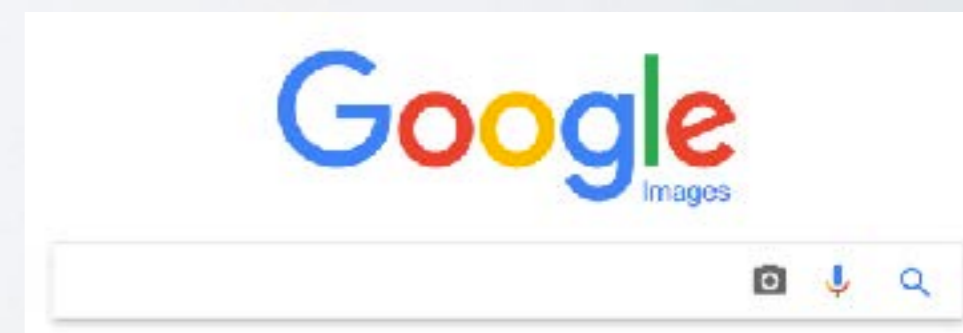
face detection/recognition

- Facebook
- Microsoft
- Snapchat



search by image

- Google image search
- search by image for products on Amazon



applications

fine-grained object recognition

e.g. recognizing any bird species



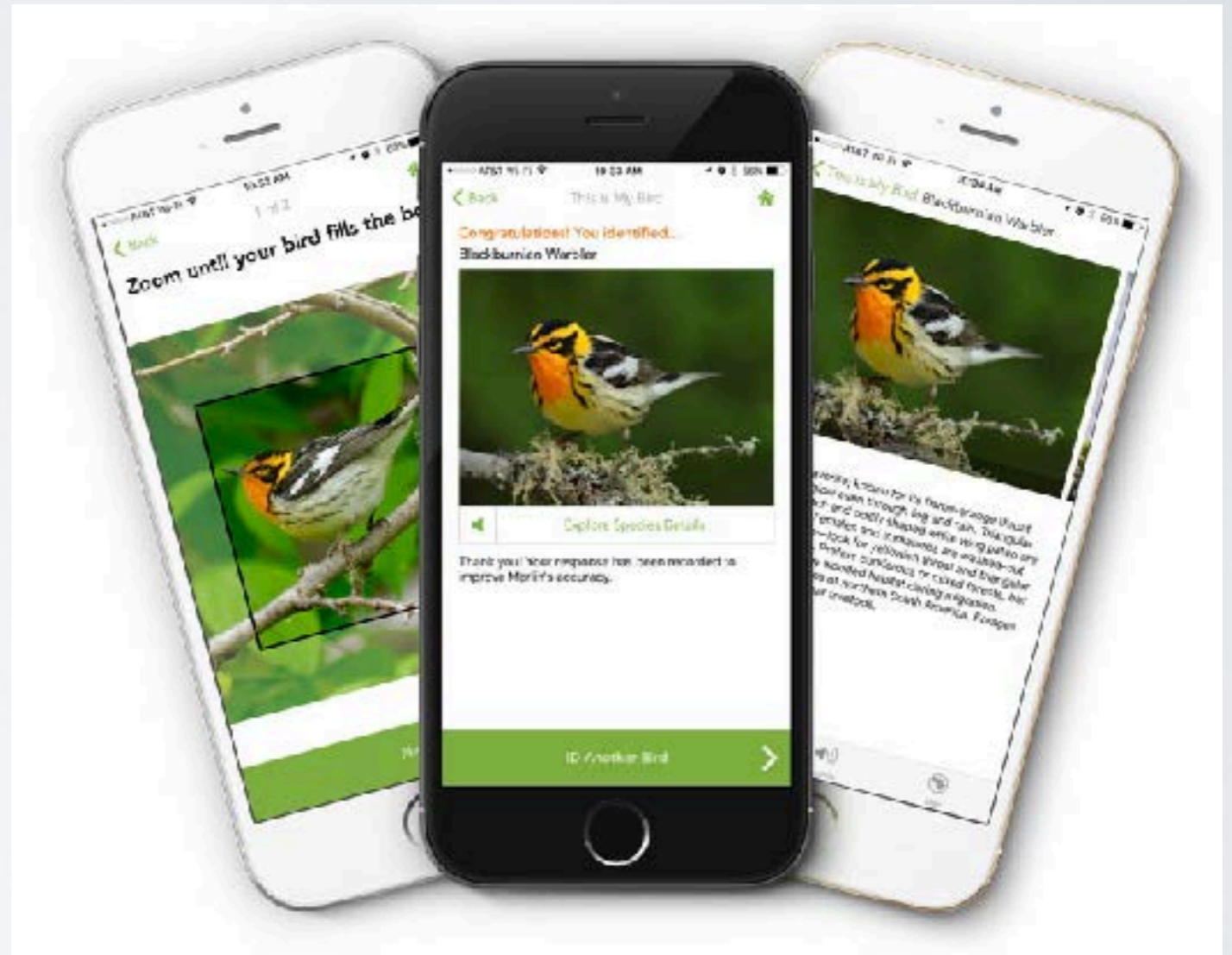
applications

fine-grained object recognition

an expert birder on your phone



Merlin



recognition component developed at Caltech

other cool stuff

start with an input image or random noise

randomly activate filters within the network, backpropagate to the image



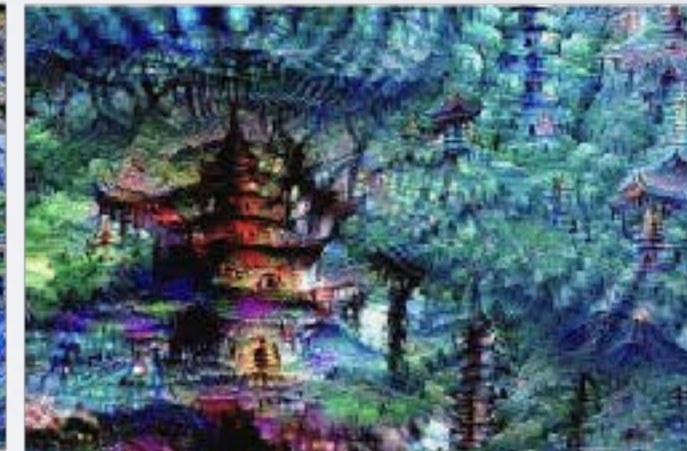
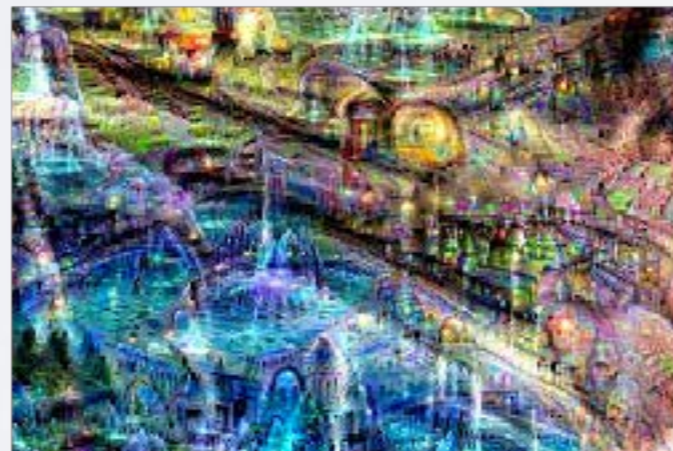
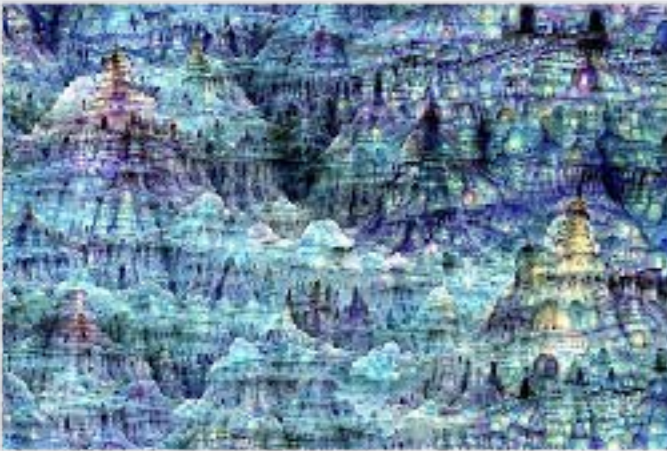
depending on which filters are chosen, has different interesting effects

it's as though the network is dreaming, hence, *deep dream*

<https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>

other cool stuff

deep dream masterpieces



other cool stuff

neural style transfer

start with a 'content image' and a 'style image'

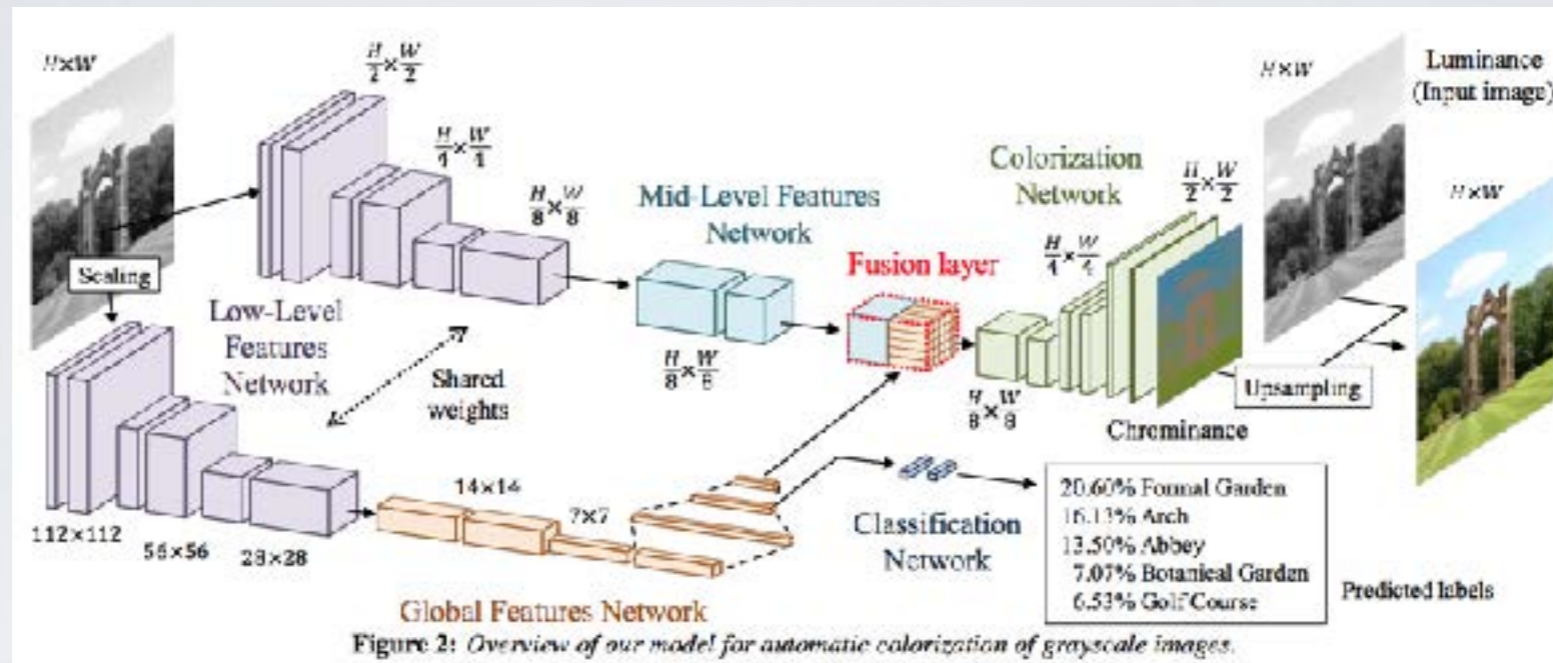
backpropagate the style image's high-level statistics to the content image



other cool stuff

colorization

learn a mapping from black and white images to color images based on visual features



discrimination

*there are patterns in images
that allow us to infer latent
properties*



Yisong

generation

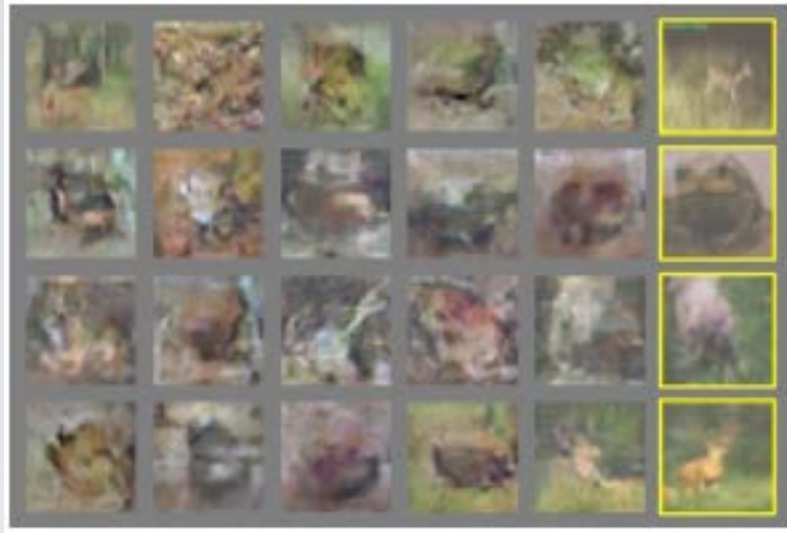
*there are latent properties
that result in specific patterns
in images*

Yisong



other cool stuff

generative modeling of images



Goodfellow et al., 2014

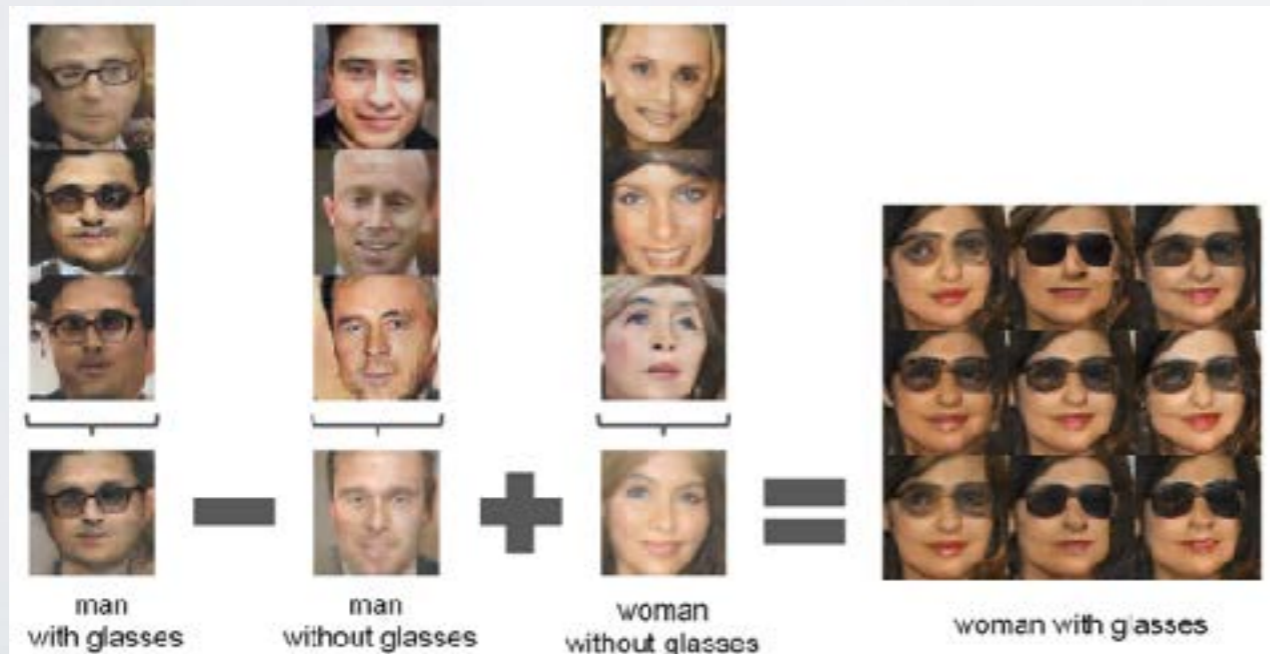


Goodfellow et al., 2016

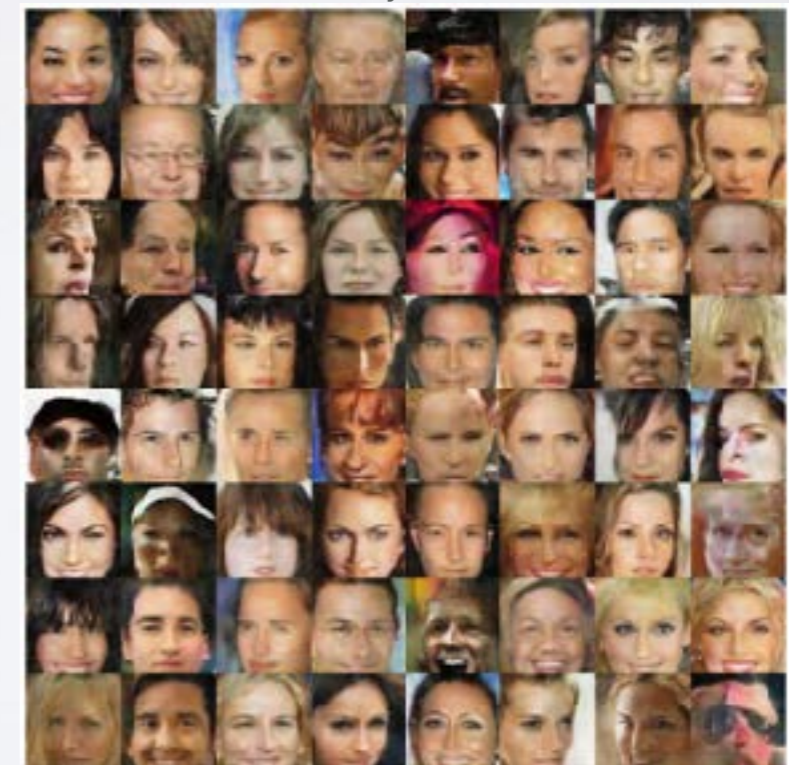


Kingma et al., 2016

Radford et al., 2015



Volkhonskiy et al., 2017



conclusion

convolutions are an ideal choice when working with data with invariances

- increases parameter efficiency through model priors
- can also be applied to 1-D and 3-D data

convolutional neural networks have enabled rapid gains in nearly all areas of computer vision and other fields

- object recognition/detection/segmentation
- success depends heavily on *data* and *hardware*

many interesting new areas to explore

- work with images at multiple levels of abstraction, not just pixel level

open problems

- choosing network architecture
- limited reasoning abilities, just an input-output mapping
- learning from few examples
- understanding what/how these networks are learning, improving training
- etc.