

Neural Networks for Visual Recognition

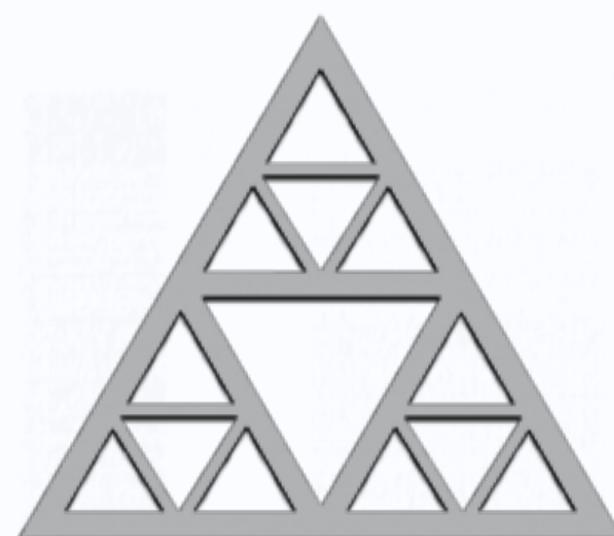
Gül Varol

IMAGINE team, École des Ponts ParisTech

gul.varol@enpc.fr

<http://imagine.enpc.fr/~varolg/>

@RecVis, 05.11.2024



École des Ponts
ParisTech

Announcements

Room for next weeks.

Assignment 2 due Tuesday Nov 12.

Assignment 3 Kaggle competition is out. (due Nov 26)

Custom final project topics: email in advance

Neural Networks

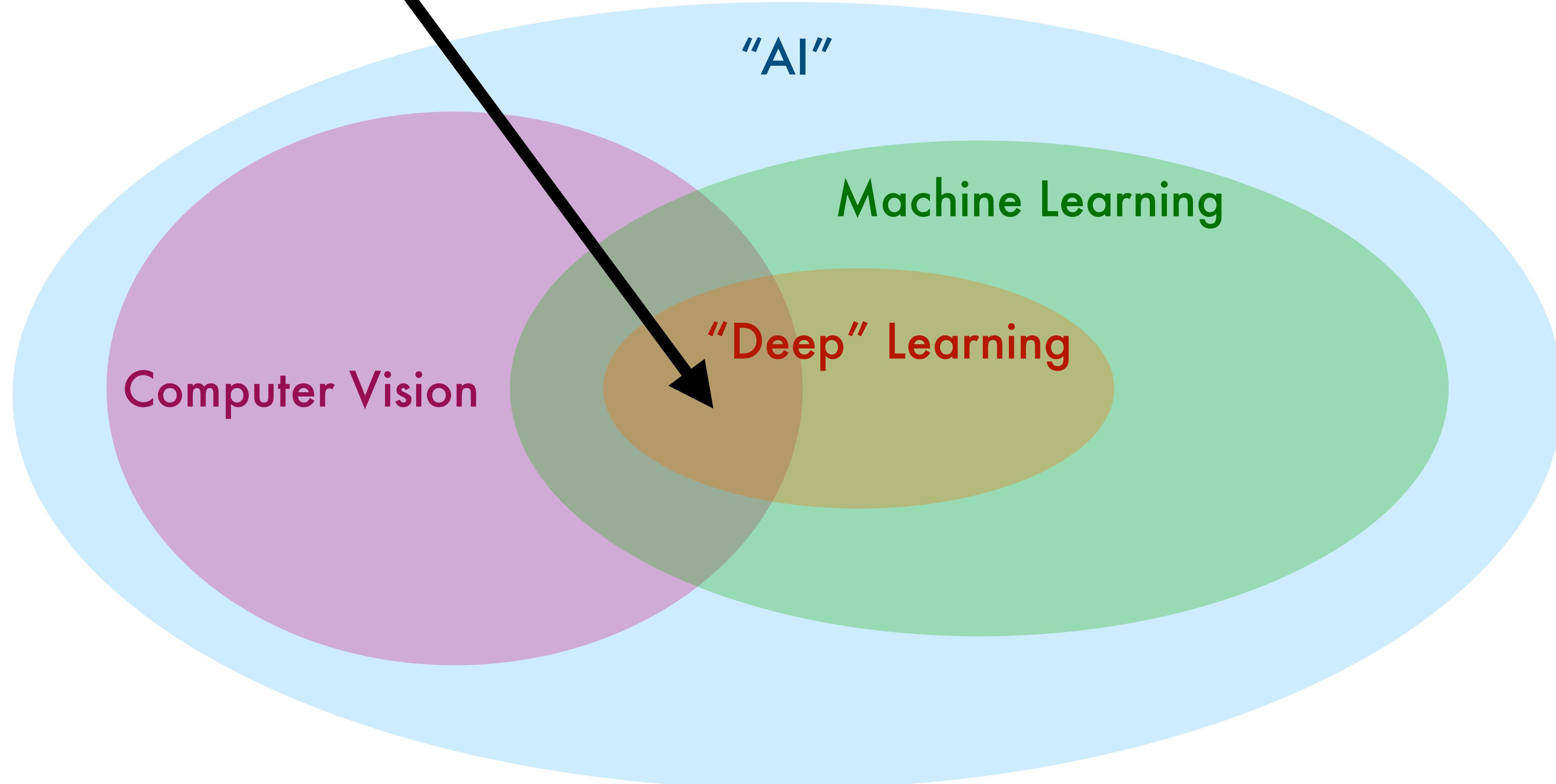
Last week: Introduction to neural networks

This week: Neural networks for **visual recognition**:

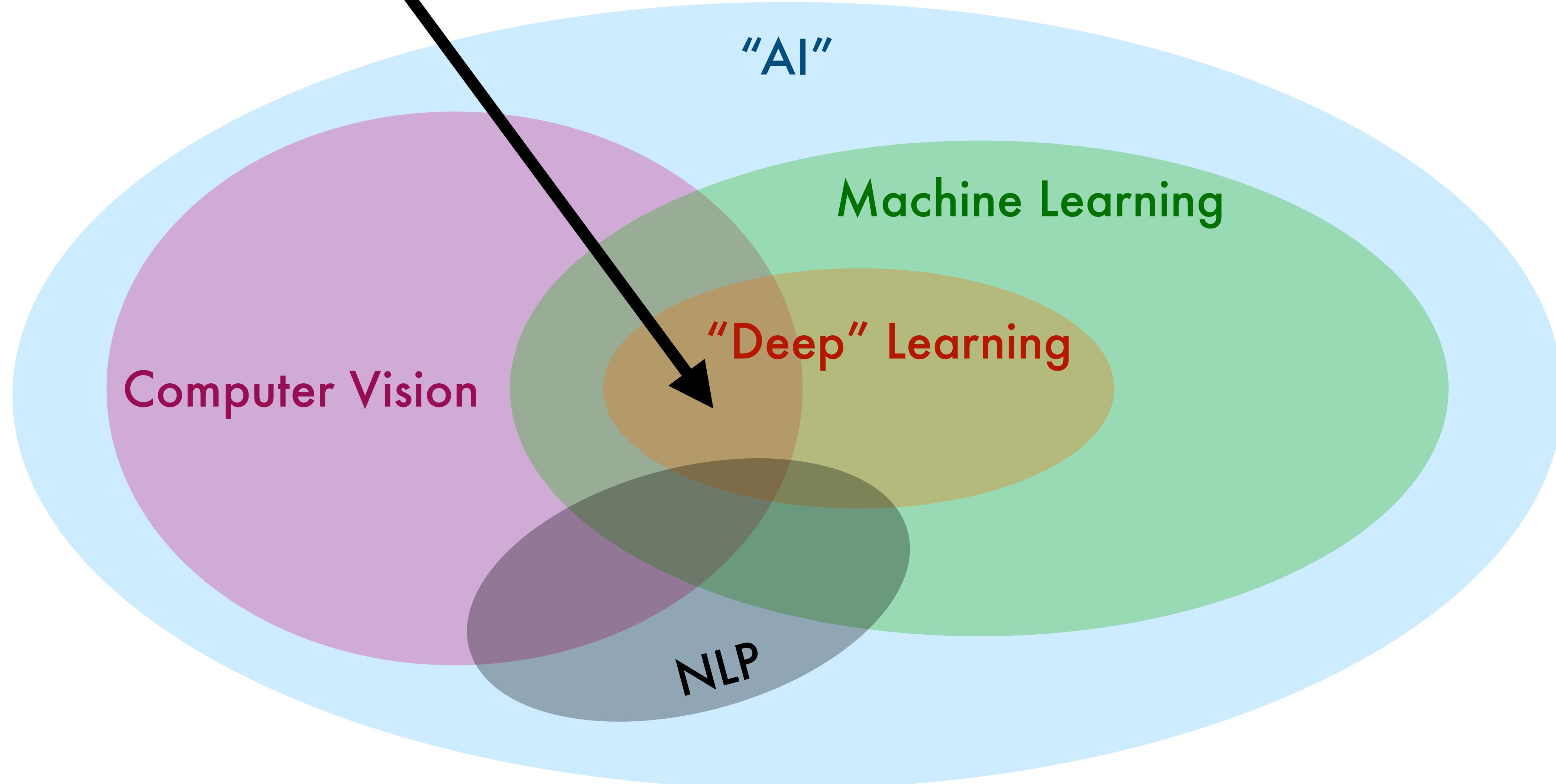
Convolutional Neural Networks (CNNs) for image classification

Next week: Beyond classification: Object detection, Segmentation, Human pose;
Beyond CNNs: Transformers

This lecture



This lecture



Agenda

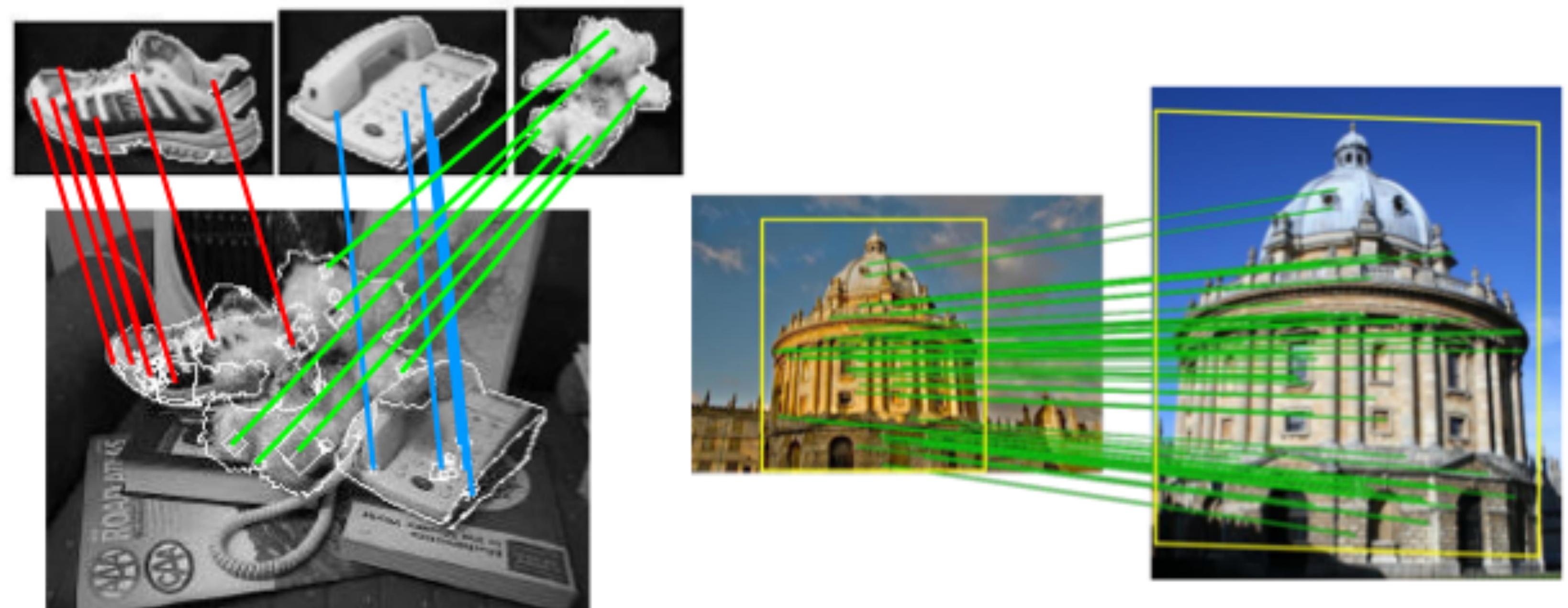
- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Recap: Image recognition so far

*Instance-level
recognition*



Category Recognition

- Image classification: assigning a class label to the image



Car: present
Cow: present
Bike: not present
Horse: not present
...

Category Recognition

- Image classification: assigning a class label to the image



Car: present
Cow: present
Bike: not present
Horse: not present
...

- Object localization: define the location and the category



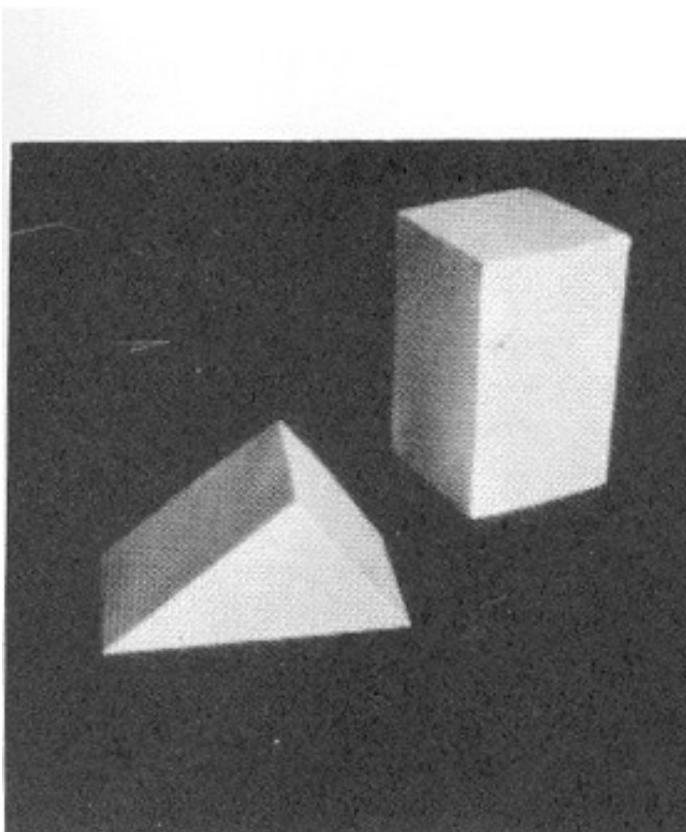
Location
Category

Difficulties: within-class variations

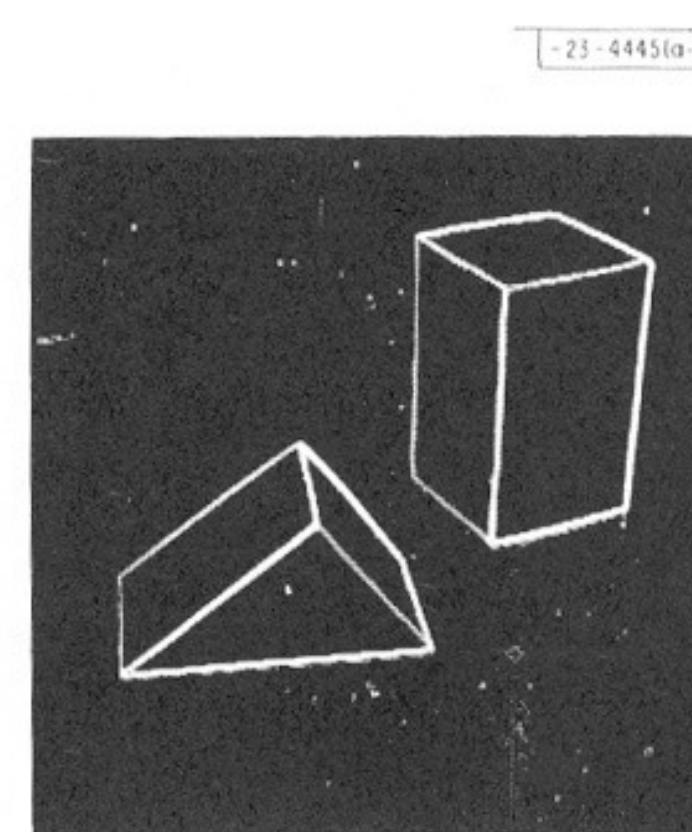


Why machine learning?

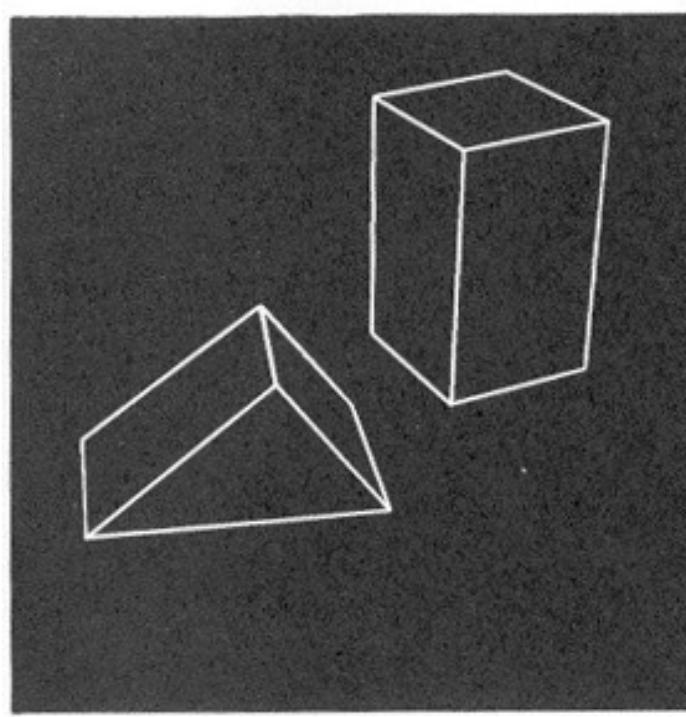
- Early approaches: simple features + handcrafted models
- Can handle only few images, simple tasks



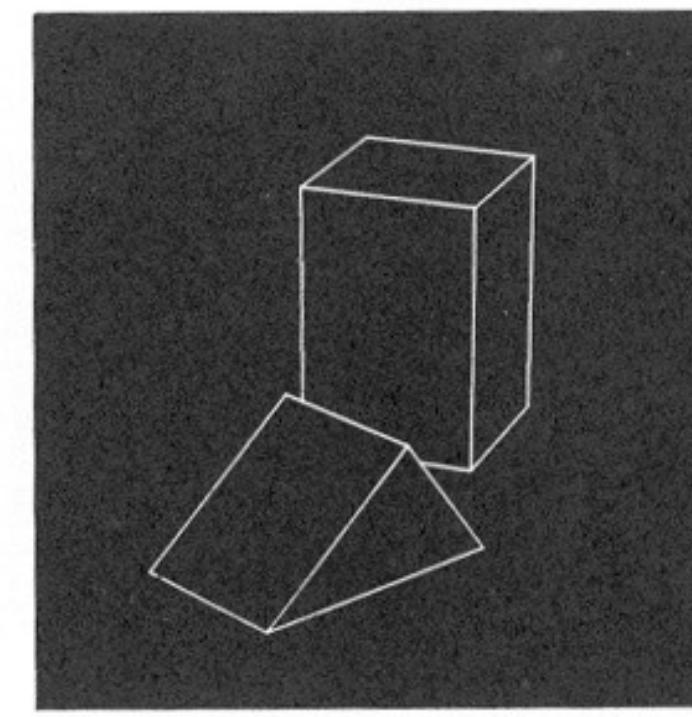
(a) Original picture.



(b) Differentiated picture.



(c) Line drawing.



(d) Rotated view.

L. G. Roberts, *Machine Perception of Three Dimensional Solids*,
Ph.D. thesis, MIT Department of Electrical Engineering, 1963.

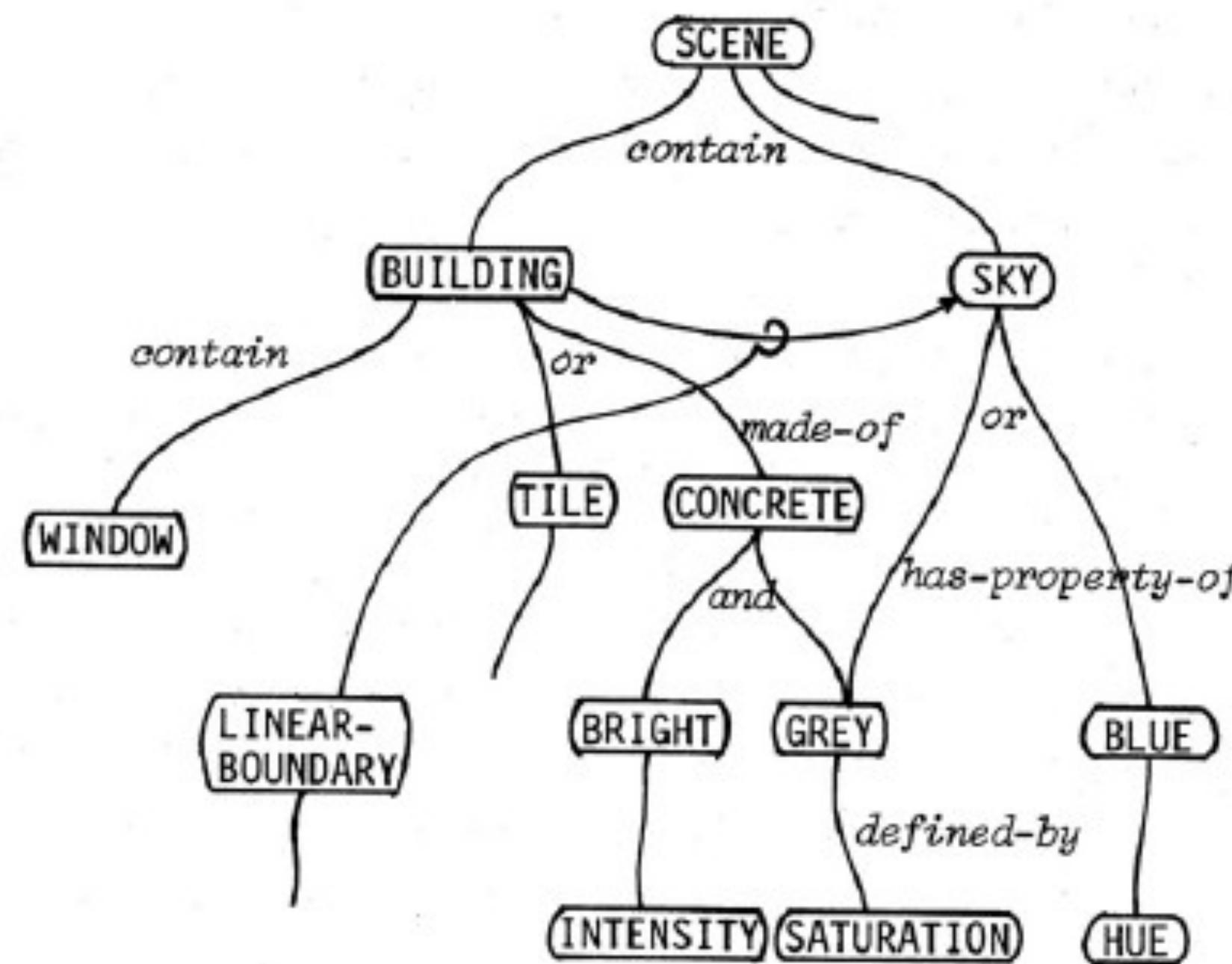
ABSTRACT

In order to make it possible for a computer to construct and display a three-dimensional array of solid objects from a single two-dimensional photograph, the rules and assumptions of depth perception have been carefully **analyzed** and **mechanized**. It is assumed that a photograph is a perspective projection of a set of objects which can be constructed from transformations of known three-dimensional models, and that the objects are supported by other visible objects or by a ground plane. These assumptions enable a computer to obtain a reasonable, three-dimensional description from the edge information in a photograph by means of a topological, mathematical process.

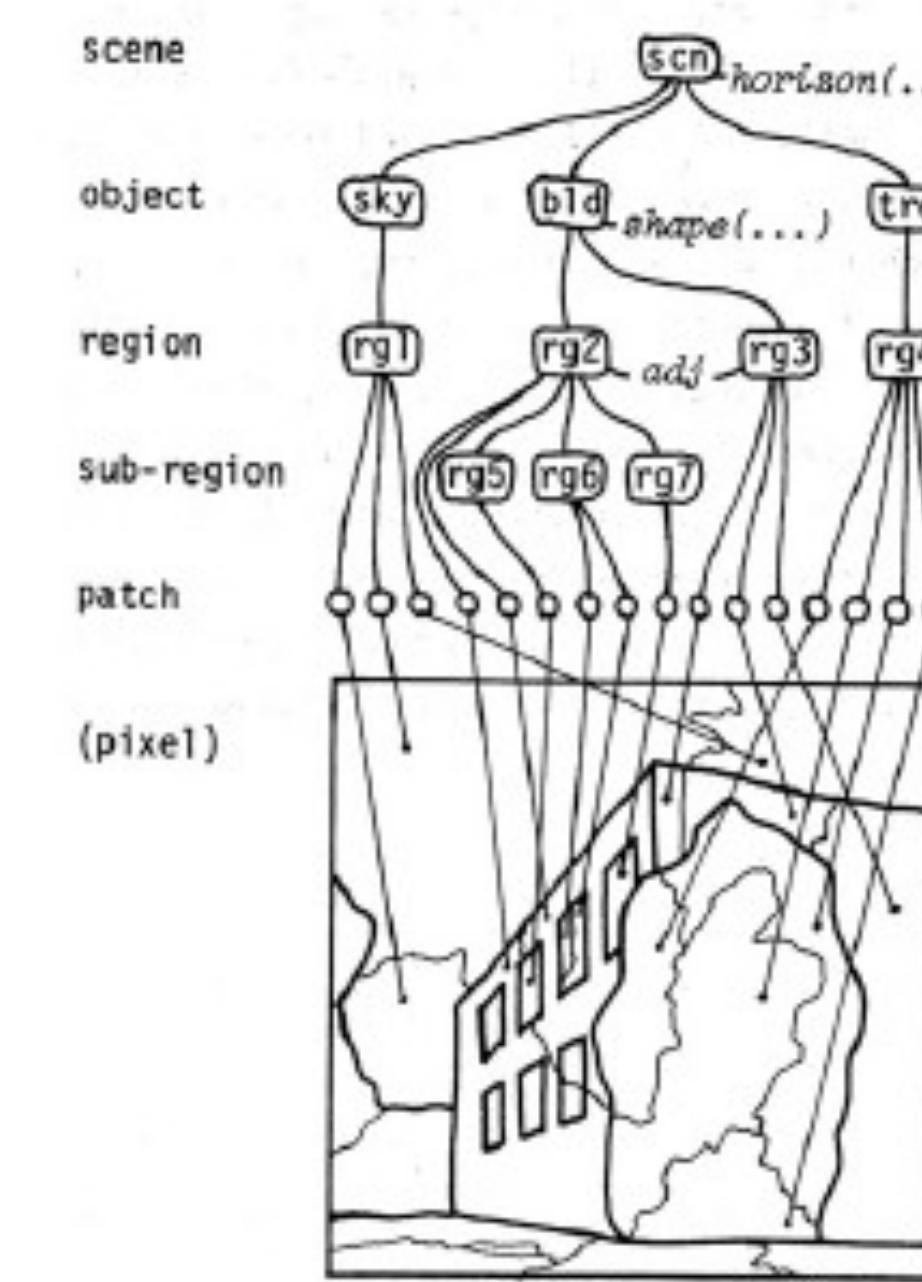
A computer program has been written which can process a photograph into a line drawing, transform the line drawing into a three-dimensional representation, and finally, display the three-dimensional structure with all the hidden lines removed, from any point of view. The 2-D to 3-D construction and 3-D to 2-D display processes are sufficiently general to handle most collections of planar-surfaced objects and provide a valuable starting point for future investigation of computer-aided three-dimensional systems.

Why machine learning?

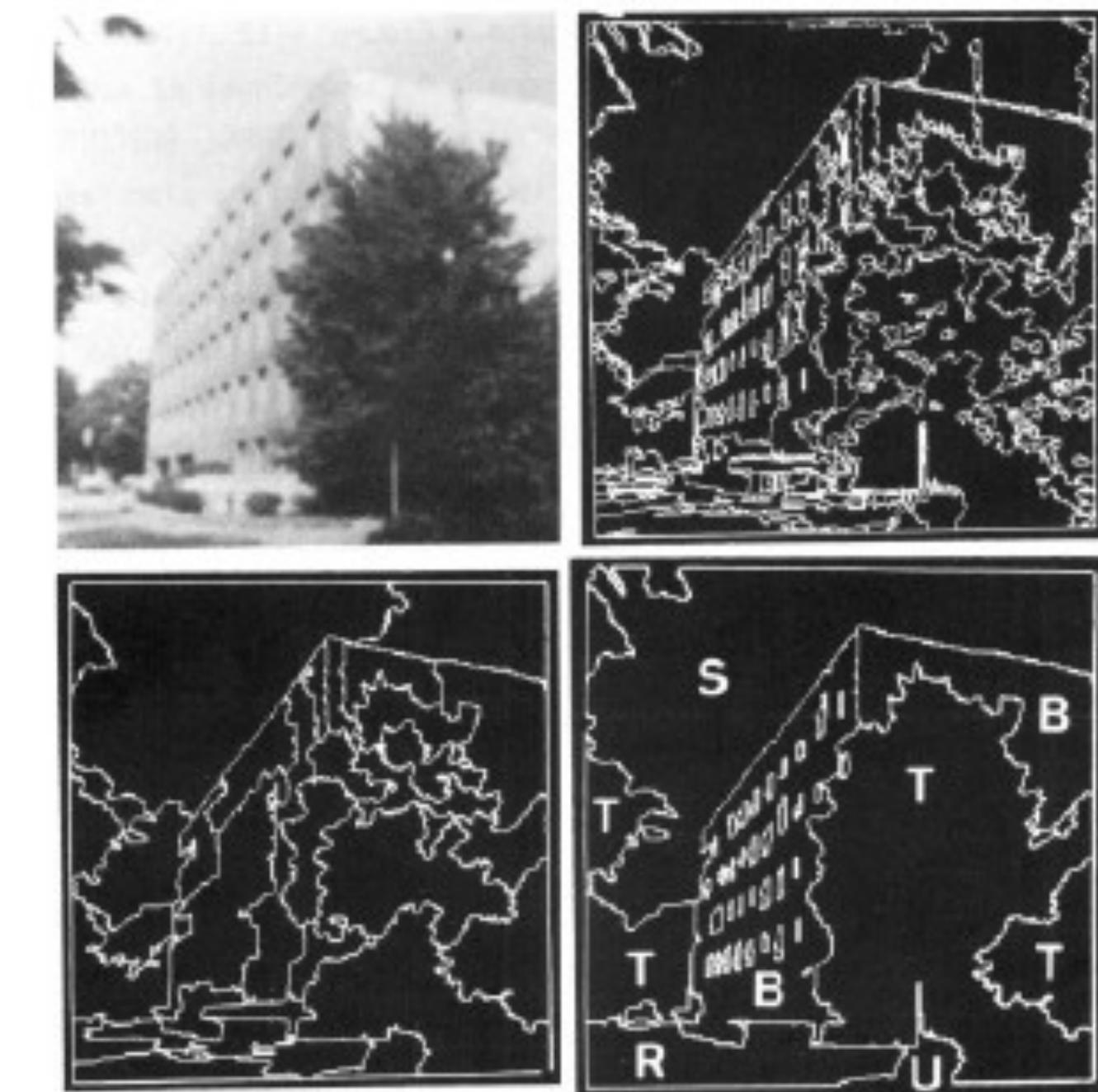
- Early approaches: manual programming of rules
- Tedious. limited. and does not take data into account



(a) Bottom-up process



(b) Top-down process



(c) Result

Figure 3. A system developed in 1978 by Ohta, Kanade and Sakai [33, 32] for knowledge-based interpretation of outdoor natural scenes. The system is able to label an image (c) into semantic classes: S-sky, T-tree, R-road, B-building, U-unknown.

Why machine learning?

- Today lots of data, complex tasks



Internet images,
personal photo albums

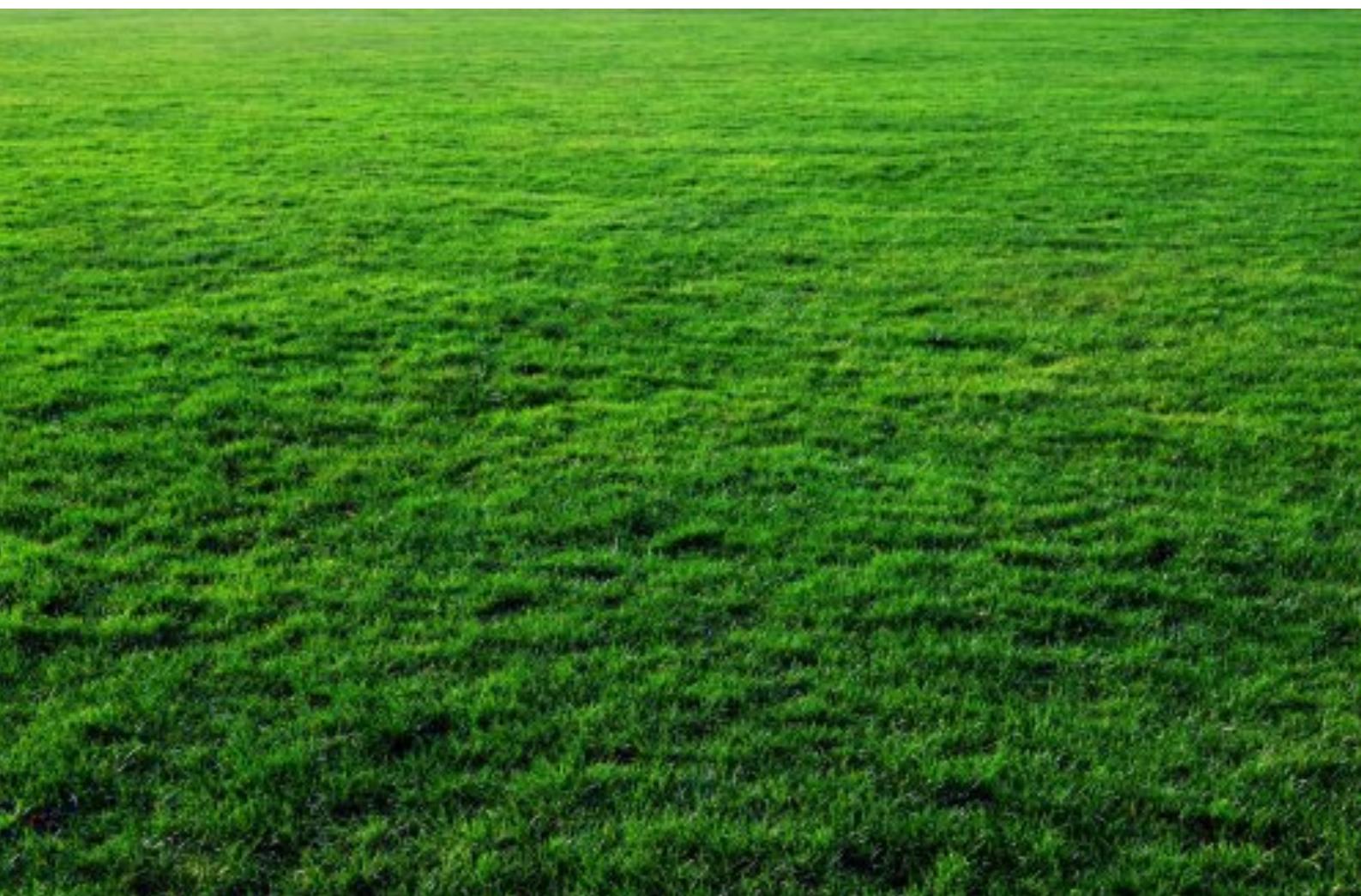


Movies, news, sports

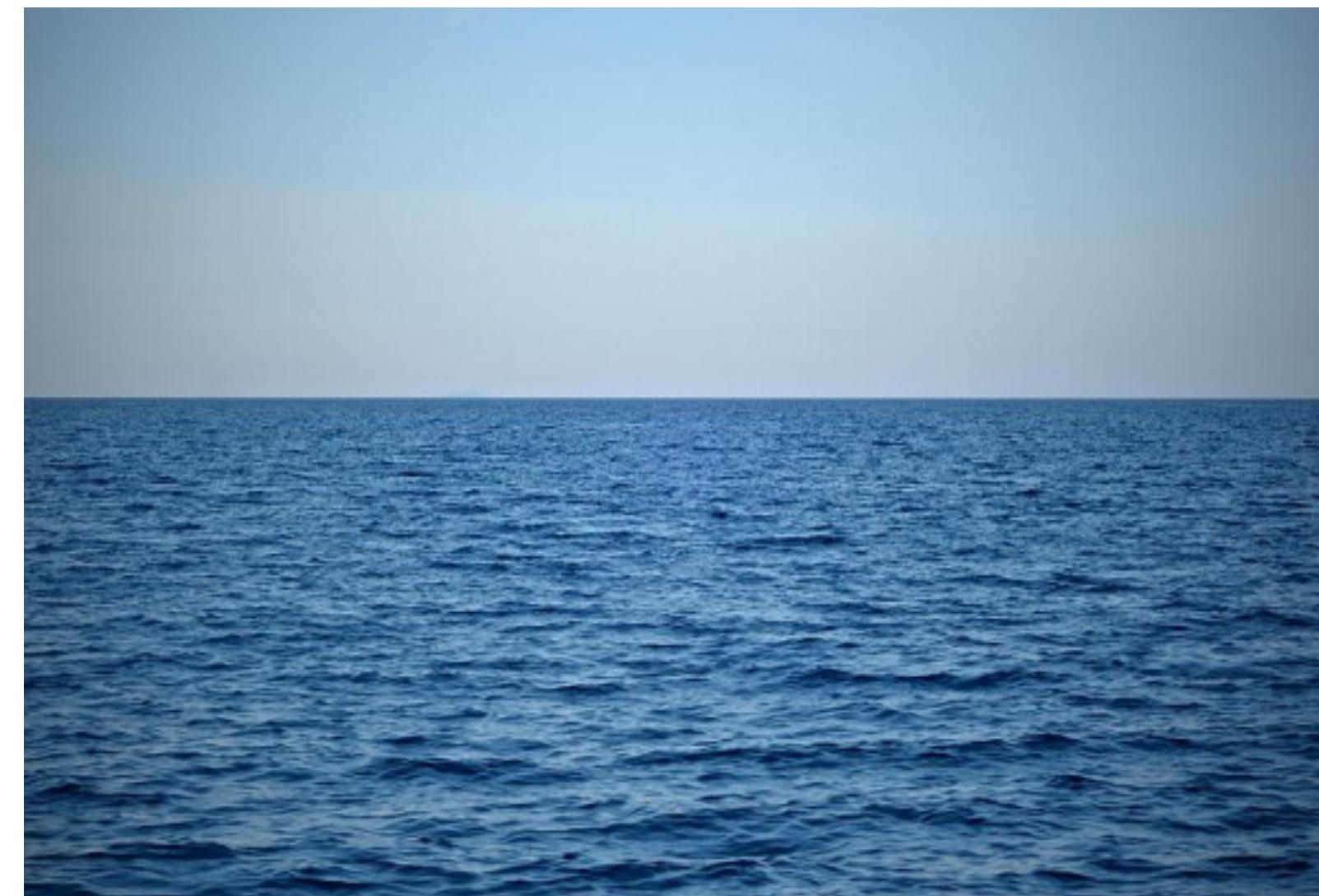
- Instead of trying to encode rules directly, learn them from examples of inputs and desired outputs

Texture Classification

- Profound observation: Grass and sea pictures don't look the same!
- Basic idea: Model the distribution of "texture" over the image (or over a region) and classify in different classes based on the texture models learned from training examples.



Grass



Sea

Image categorization

- Profound observation: Cows and buildings don't look the same!
- Basic idea: Model the distribution of "texture" over the image (or over a region) and classify in different classes based on the texture models learned from training examples.



Cow

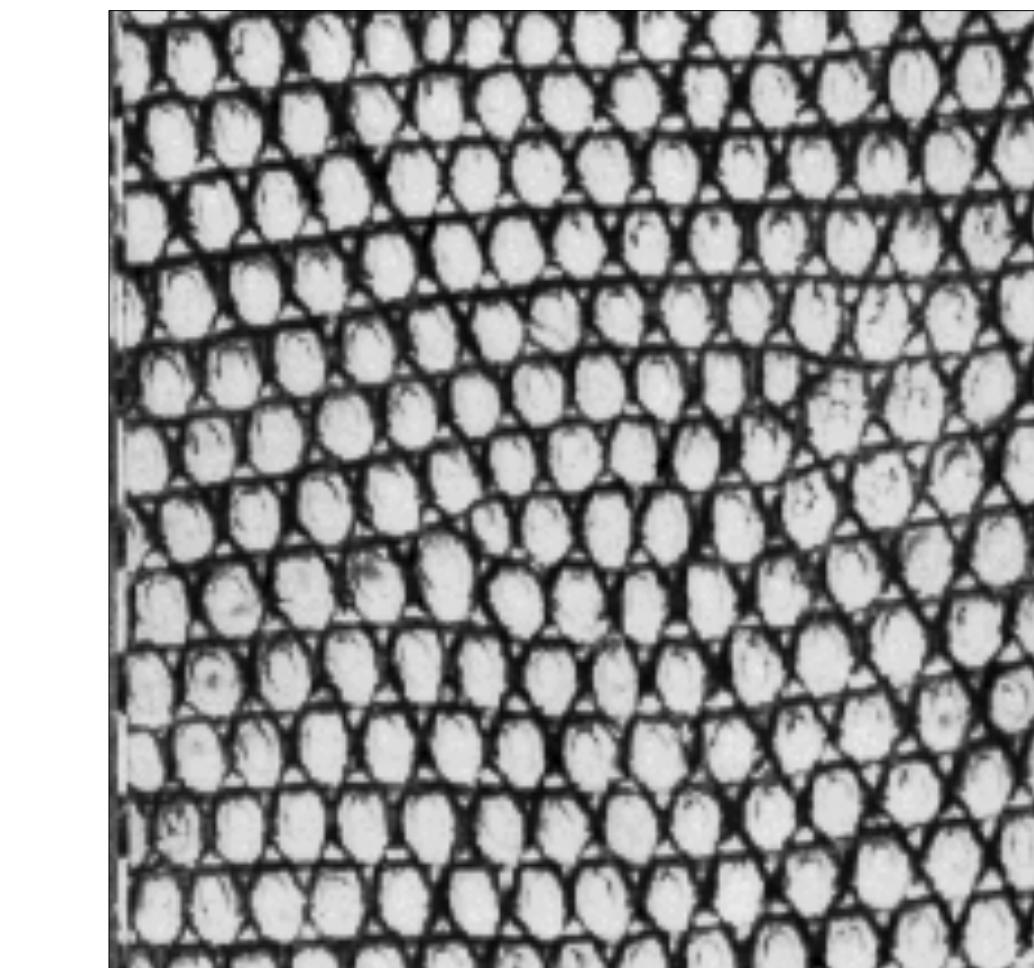
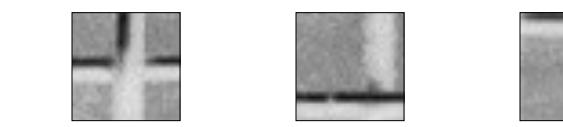
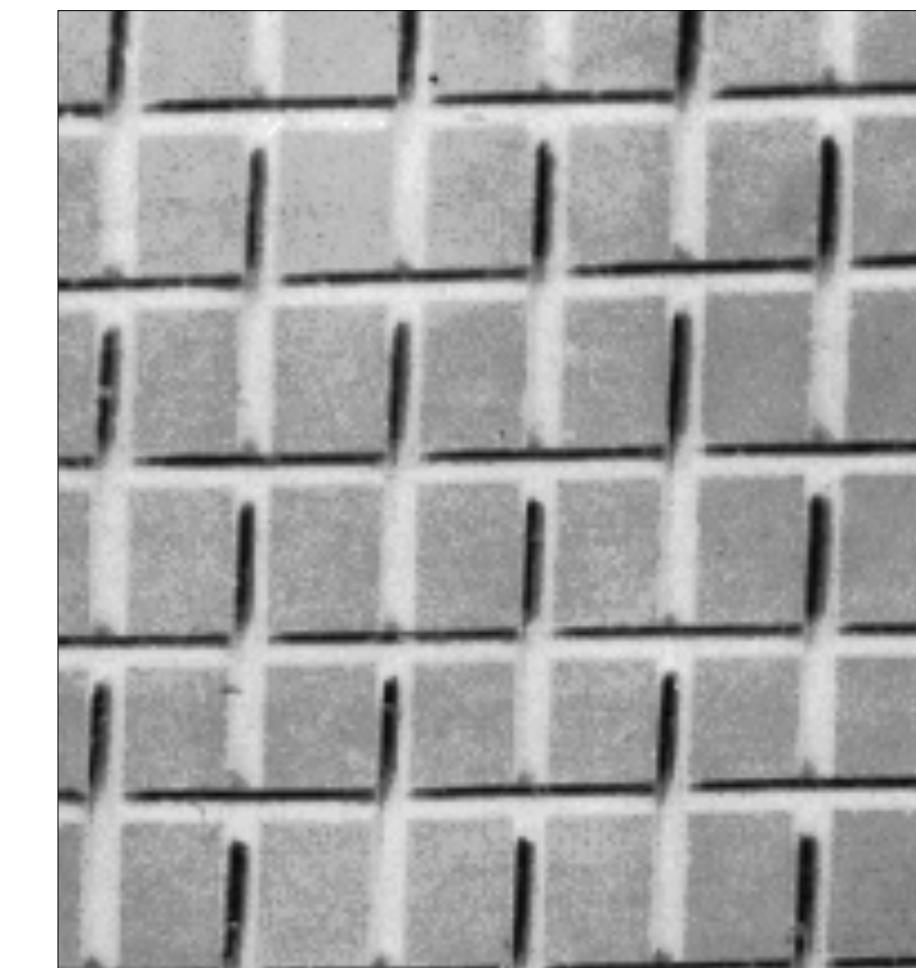
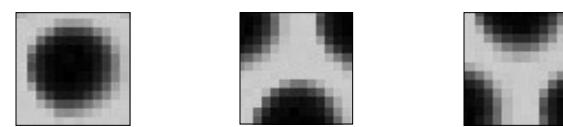
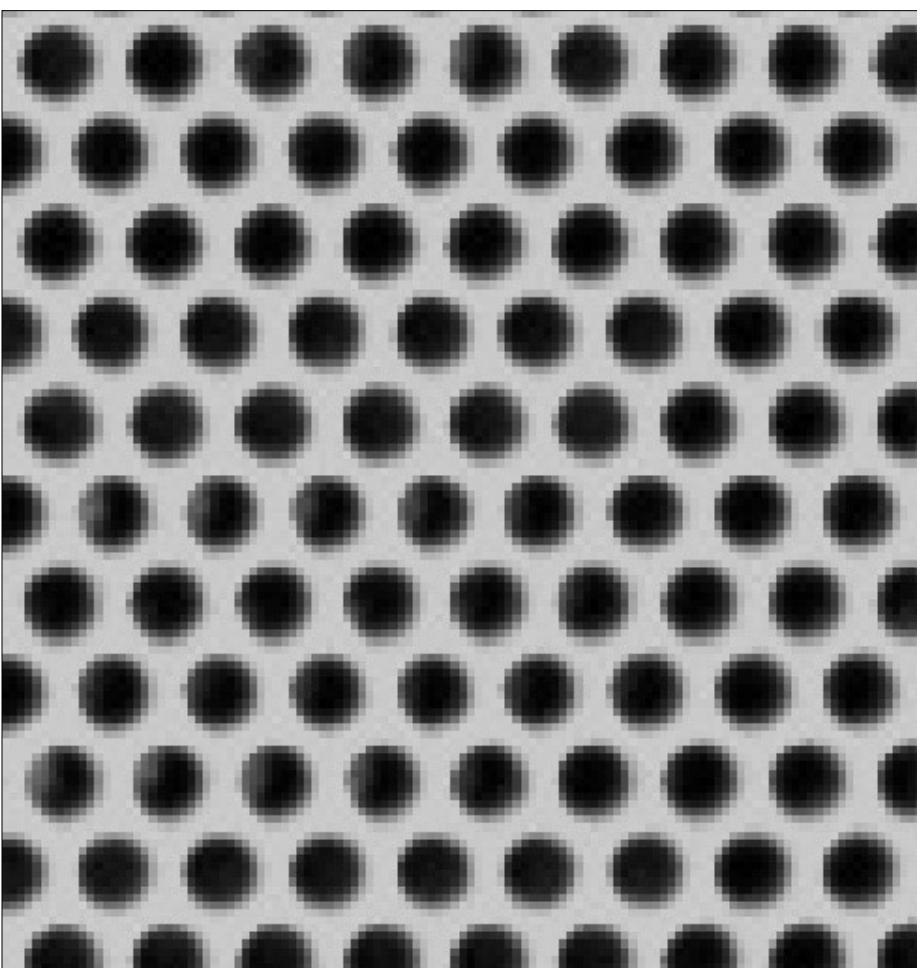


Building

Bag-of-features for image classification

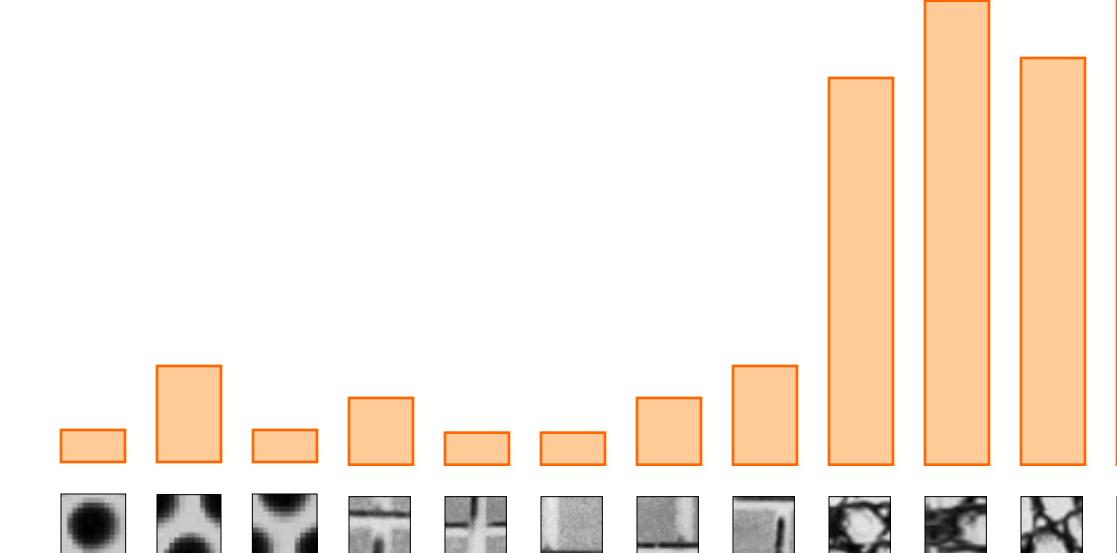
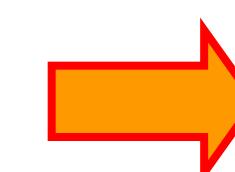
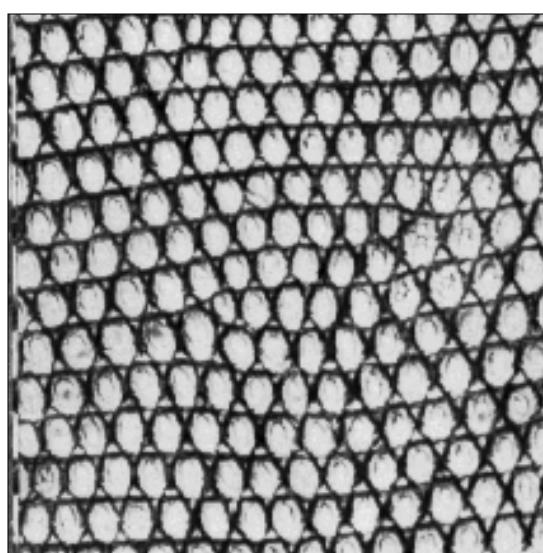
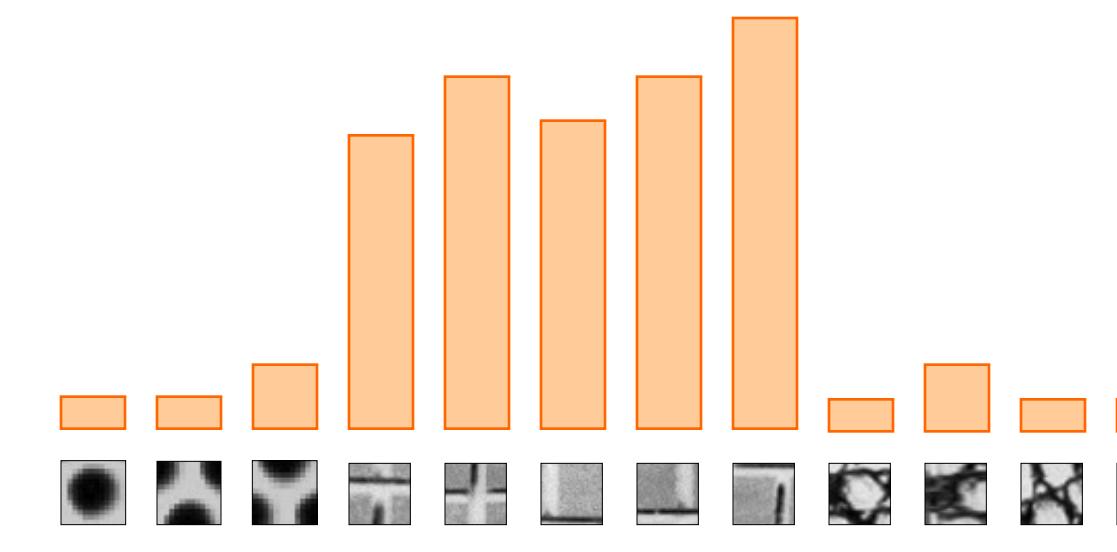
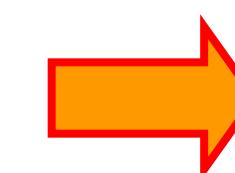
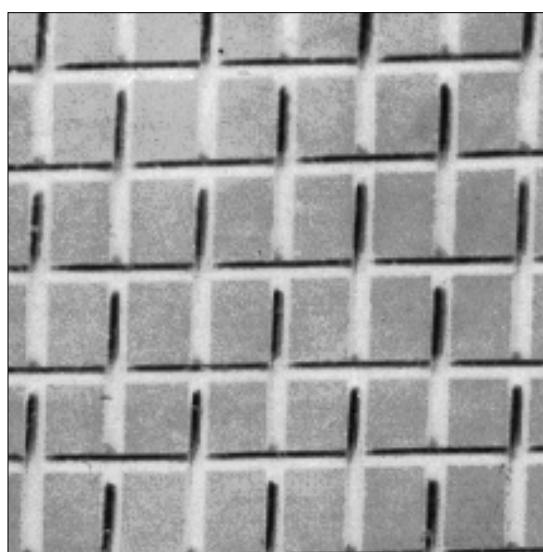
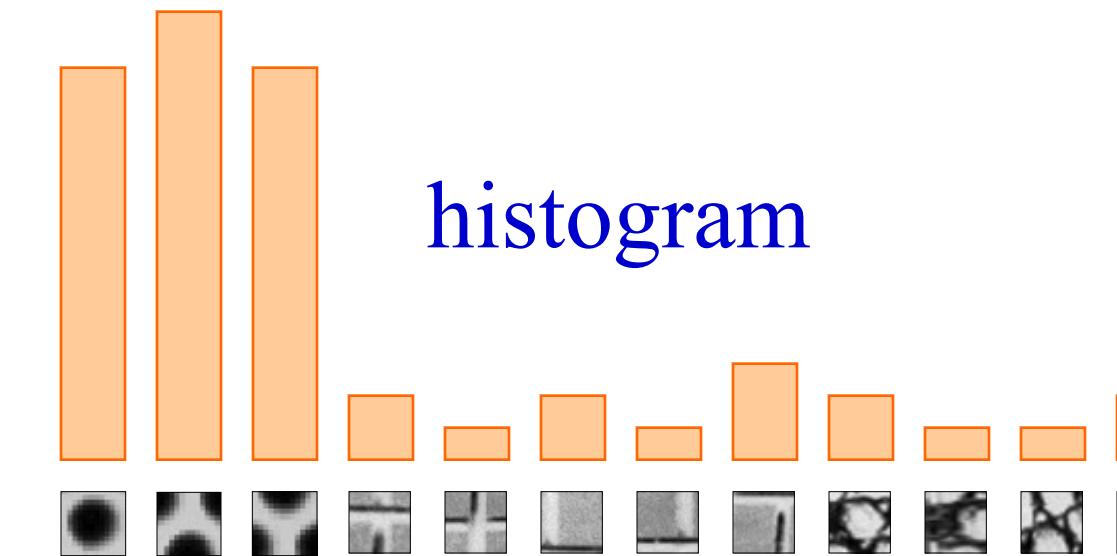
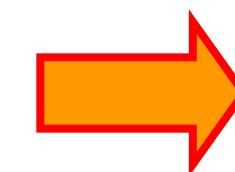
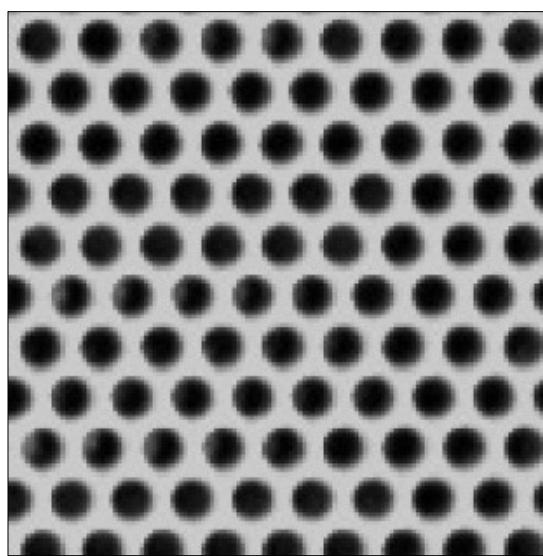
Origin: texture recognition

- Texture is characterized by the repetition of basic elements or *textons*



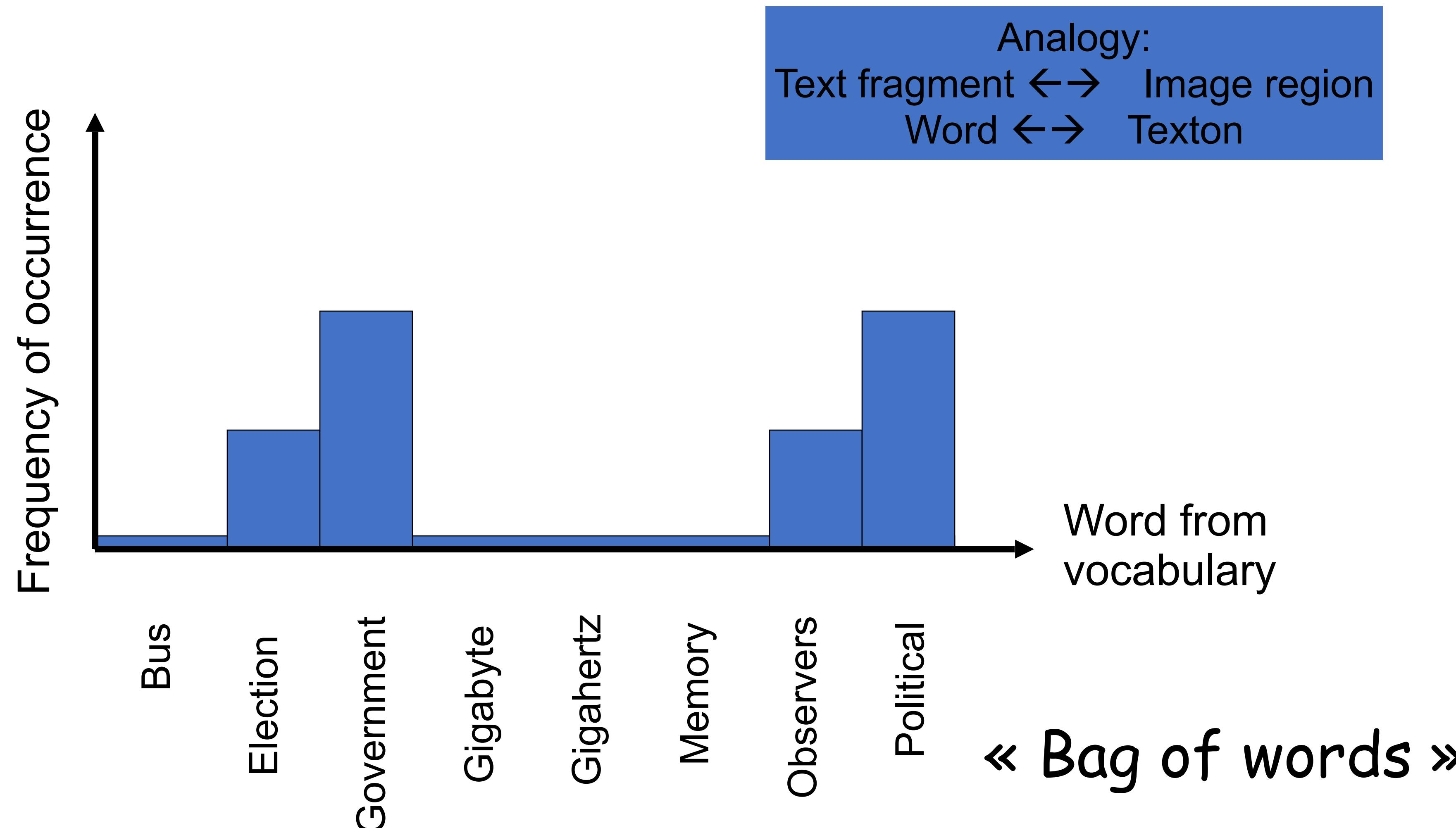
Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001;
Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

Bag-of-features for image classification



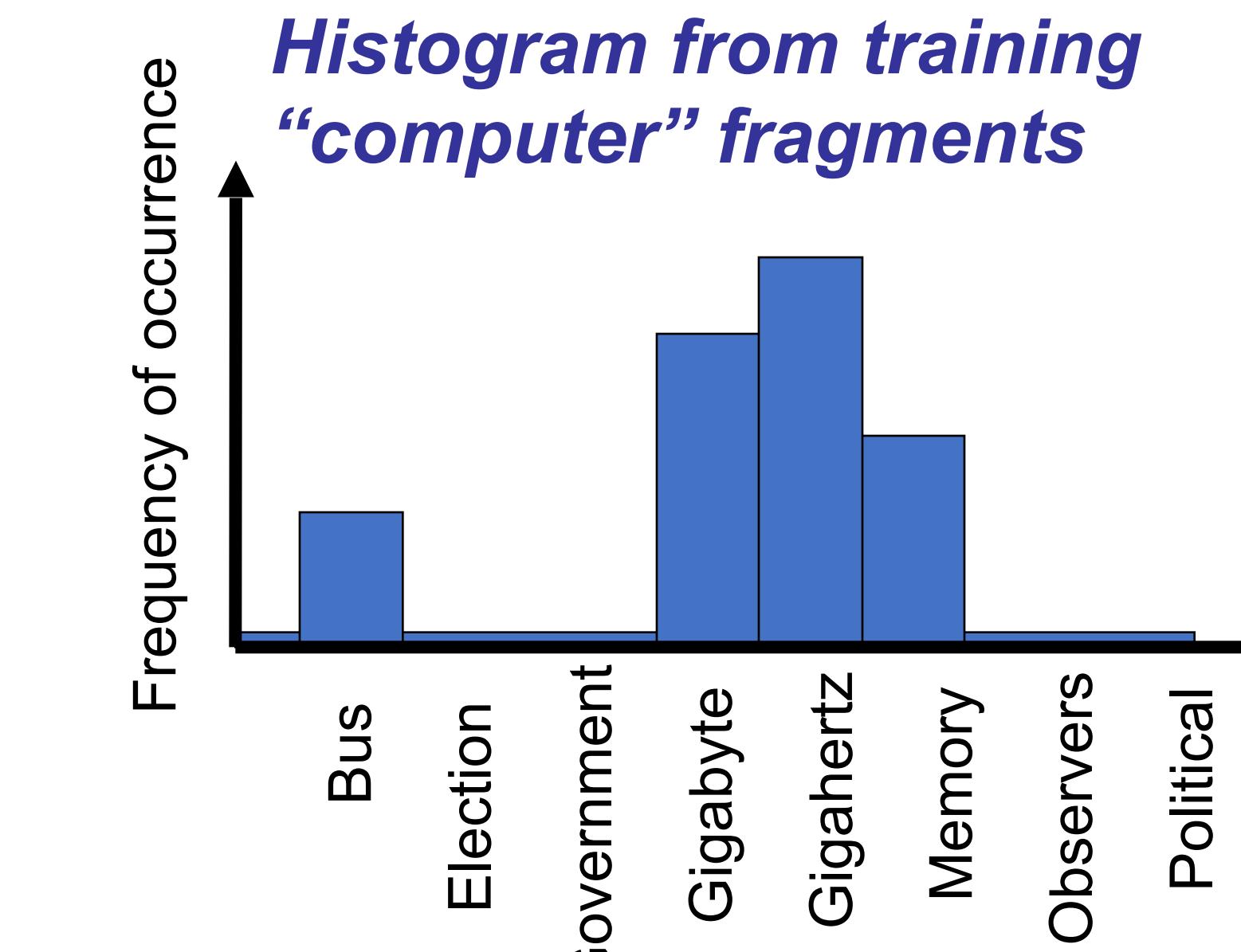
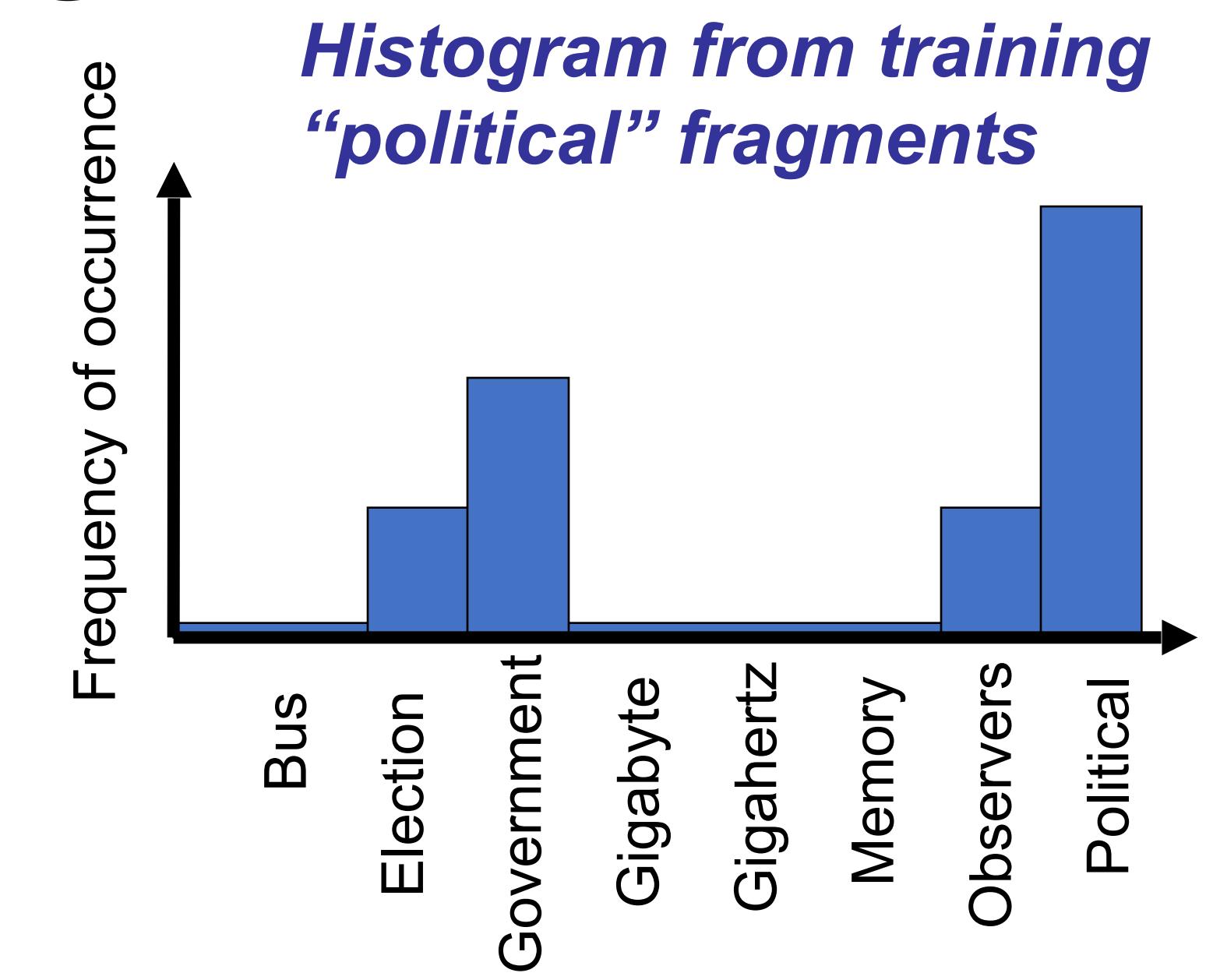
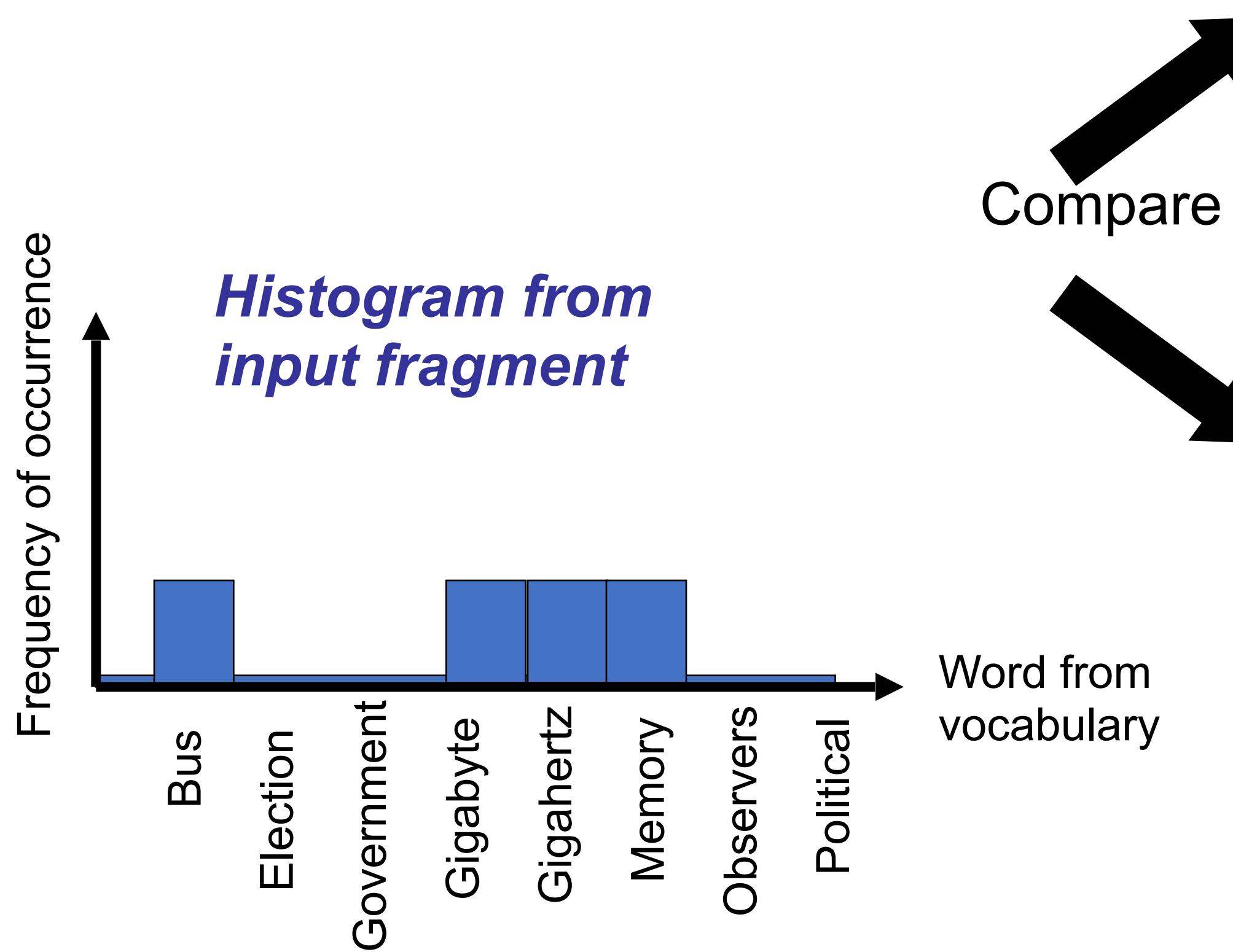
Analogy with Text Analysis

Political observers say that the government of Zorgia does not control the political situation. The government will not hold elections ...

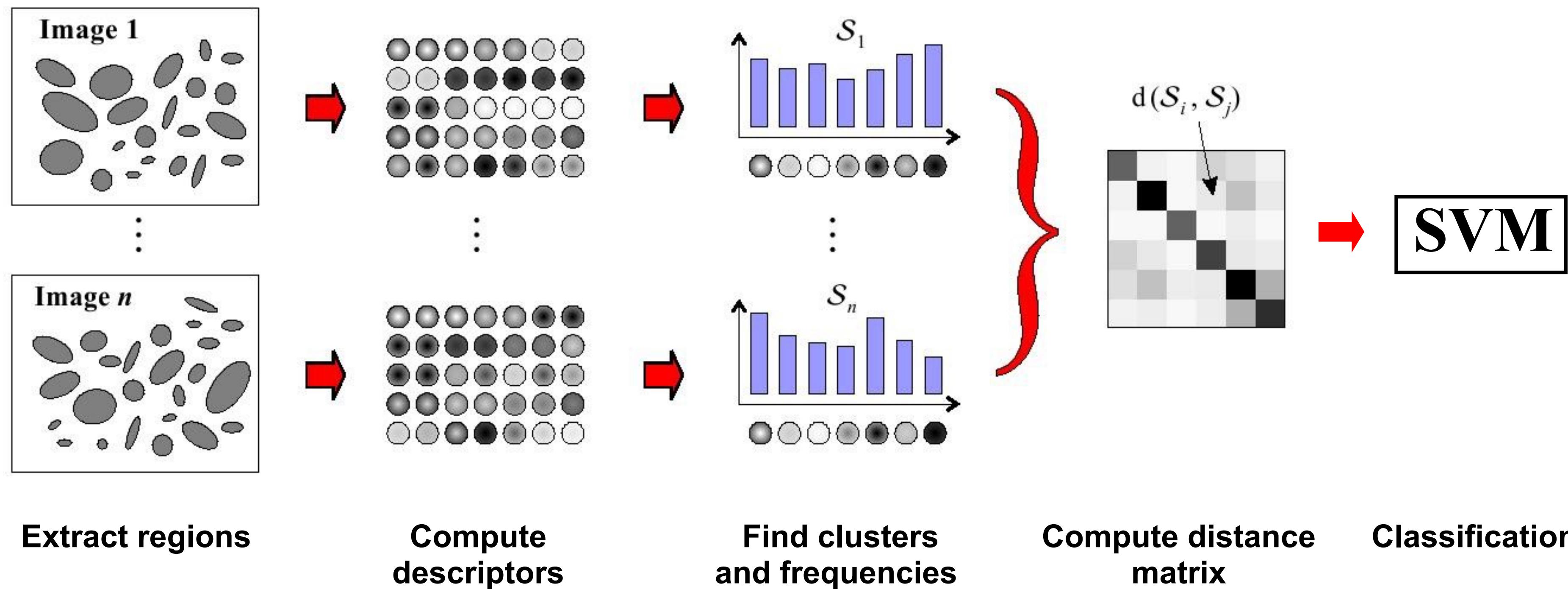


Analogy with Text Analysis

The ZH-20 unit is a 200Gigahertz processor with 2Gigabyte memory. Its strength is its bus and high-speed memory.....

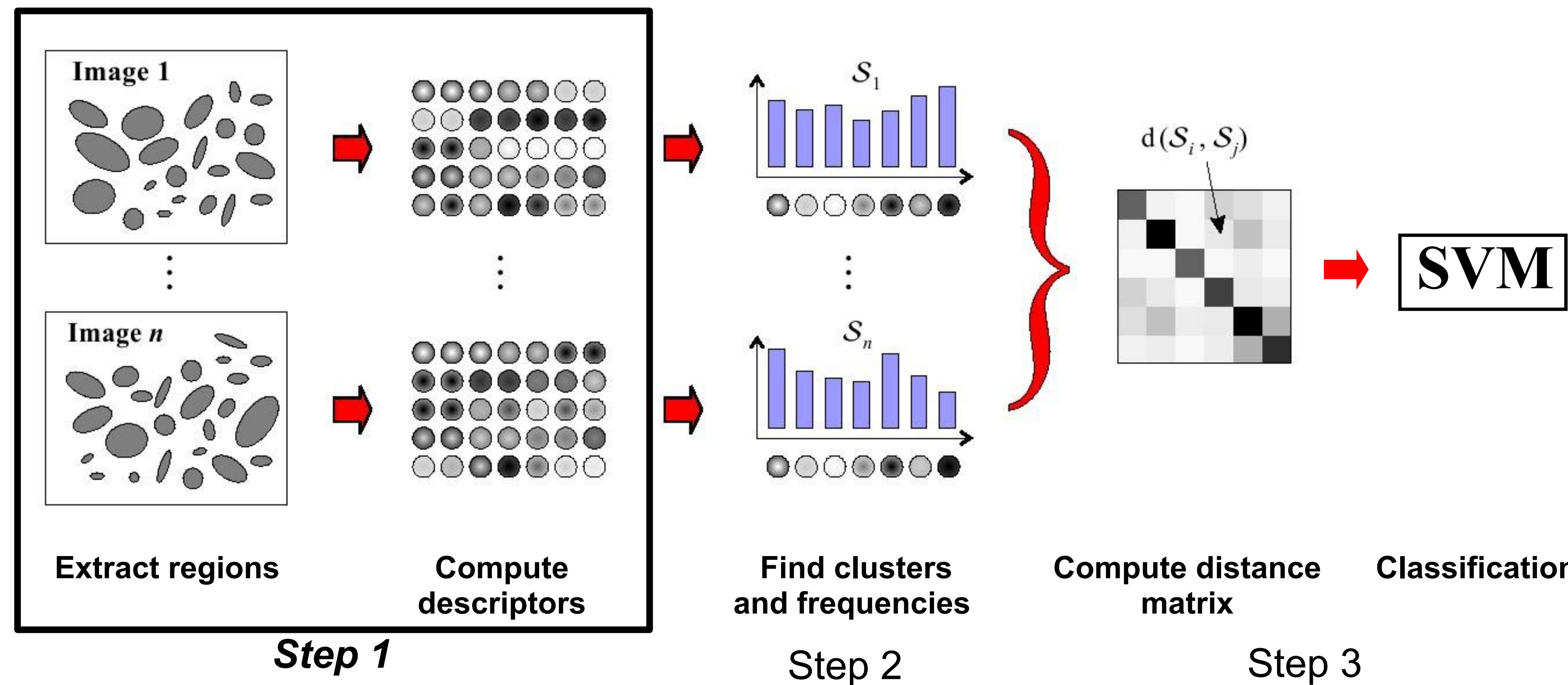


Bag-of-features for image classification



[Csurka et al. WS'2004], [Nowak et al. ECCV'06], [Zhang et al. IJCV'07]

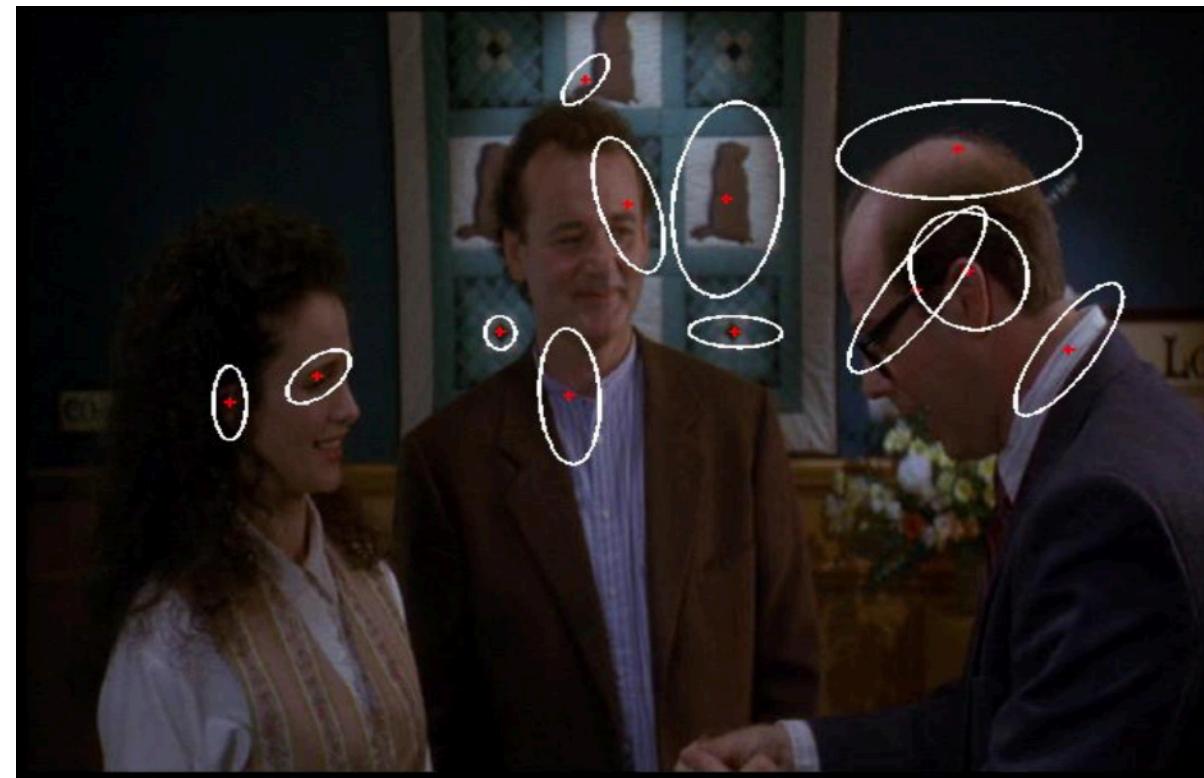
Bag-of-features for image classification



Step 1: feature extraction

Sparse sampling

- SIFT as interest point detector



Dense sampling

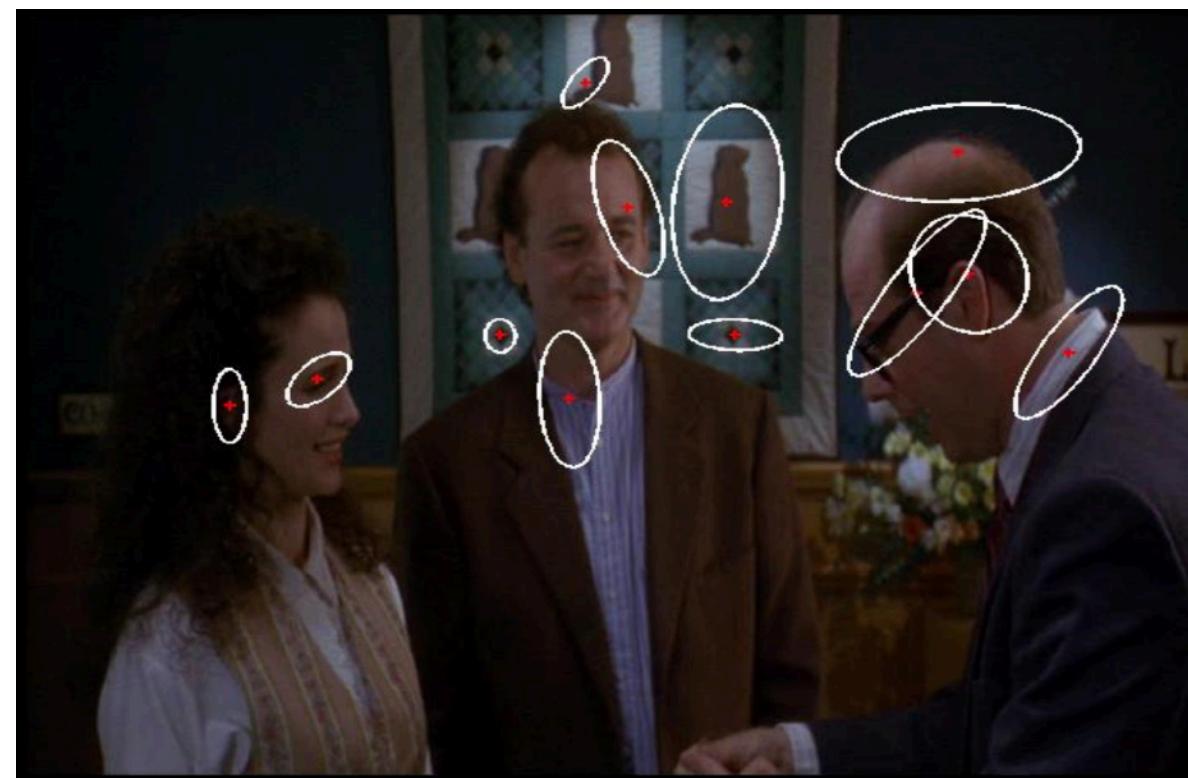
- Interest points do not necessarily capture “all” features



Step 1: feature extraction

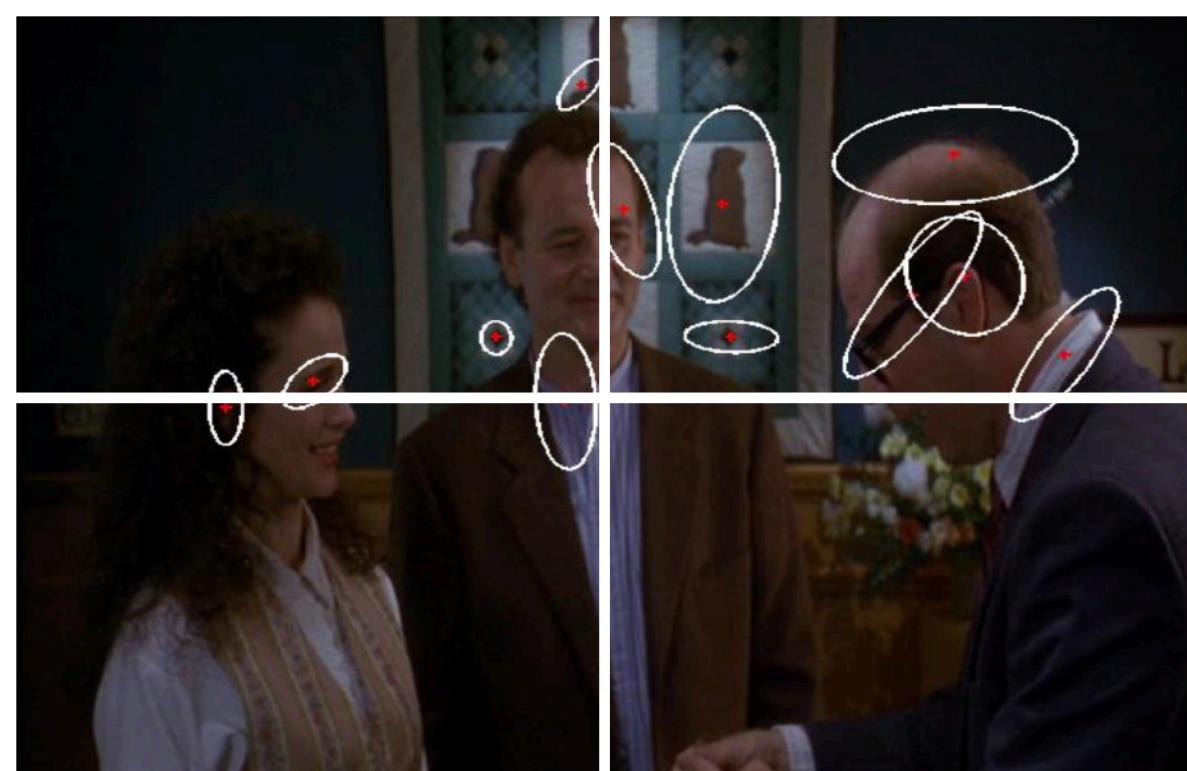
Sparse sampling

- SIFT as interest point detector

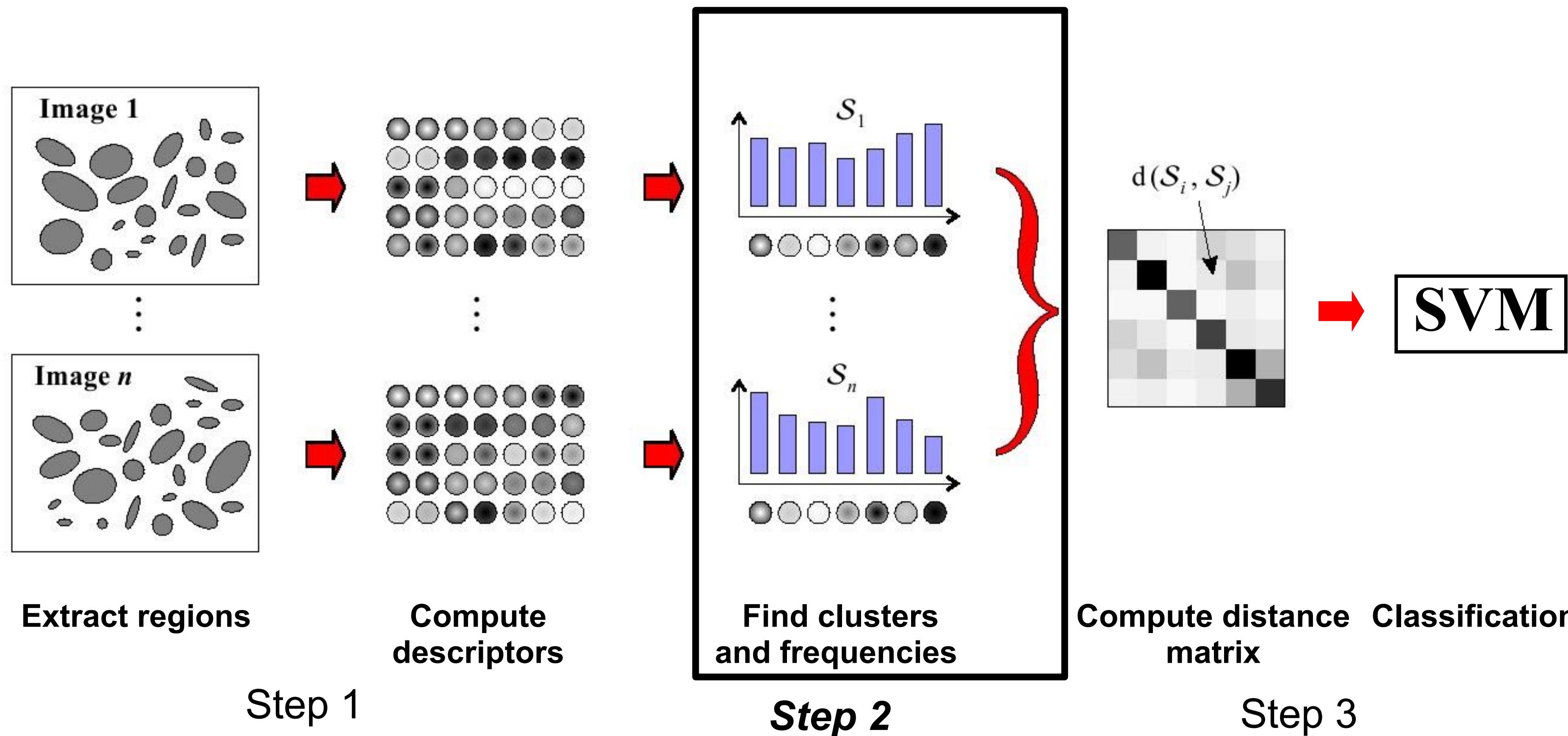


Dense sampling

- Interest points do not necessarily capture “all” features
- Spatial pyramid (Lazebnik, Schmid & Ponce, CVPR 2006)



Bag-of-features for image classification



Step 2: Quantization

Cluster descriptors

- K-means
- Gaussian mixture model

Assign each visual word to a cluster

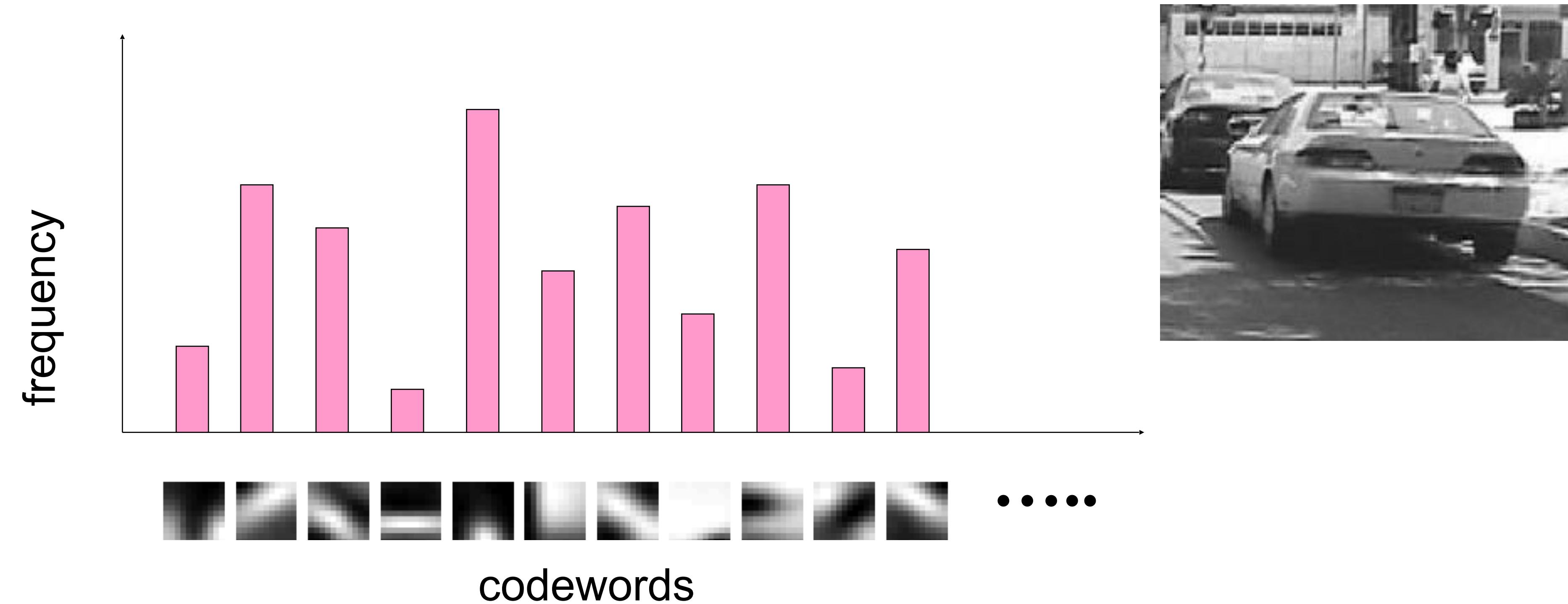
- Hard or soft assignment

Build frequency histogram

Examples for visual words

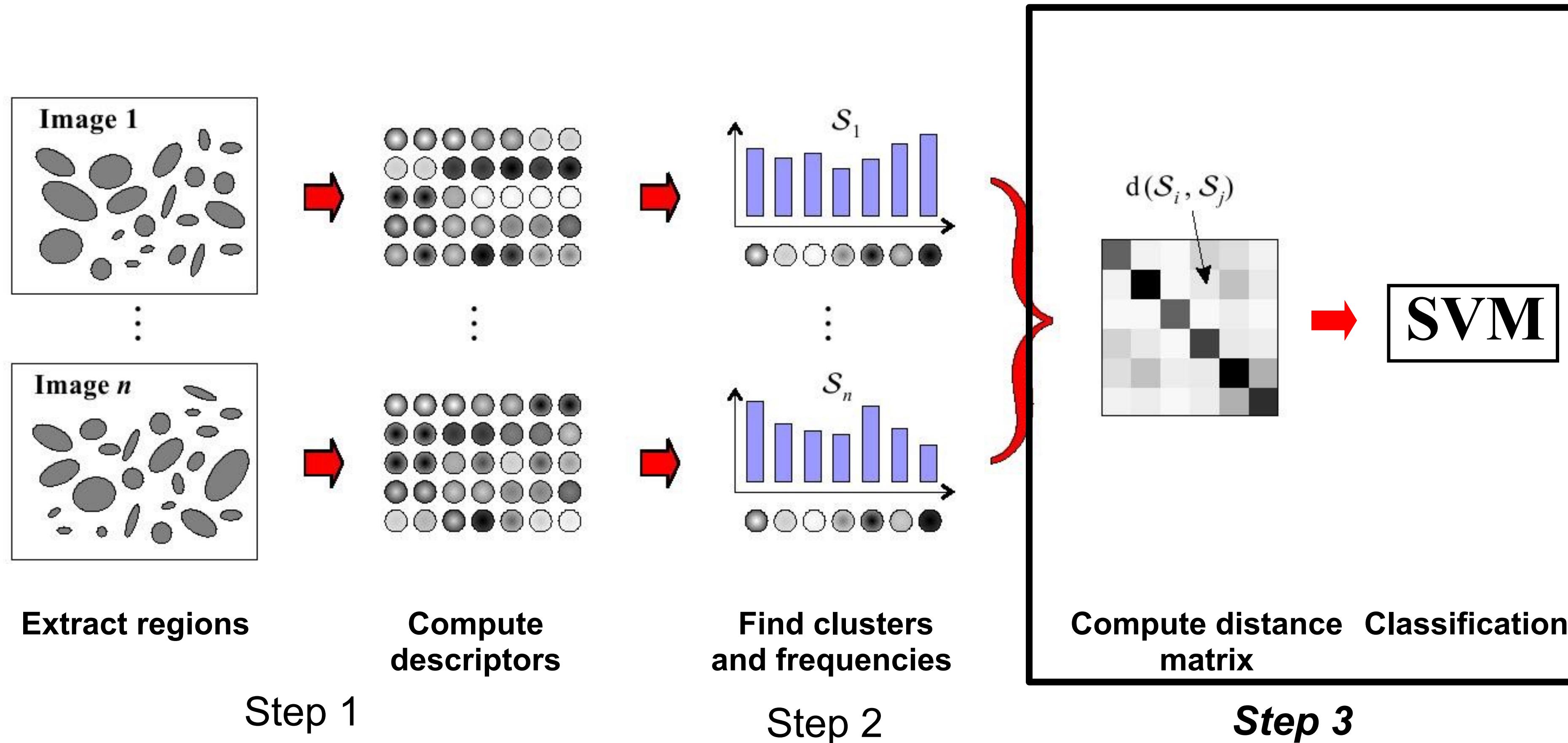
Airplanes		
Motorbikes		
Faces		
Wild Cats		
Leaves		
People		
Bikes		

Image representation



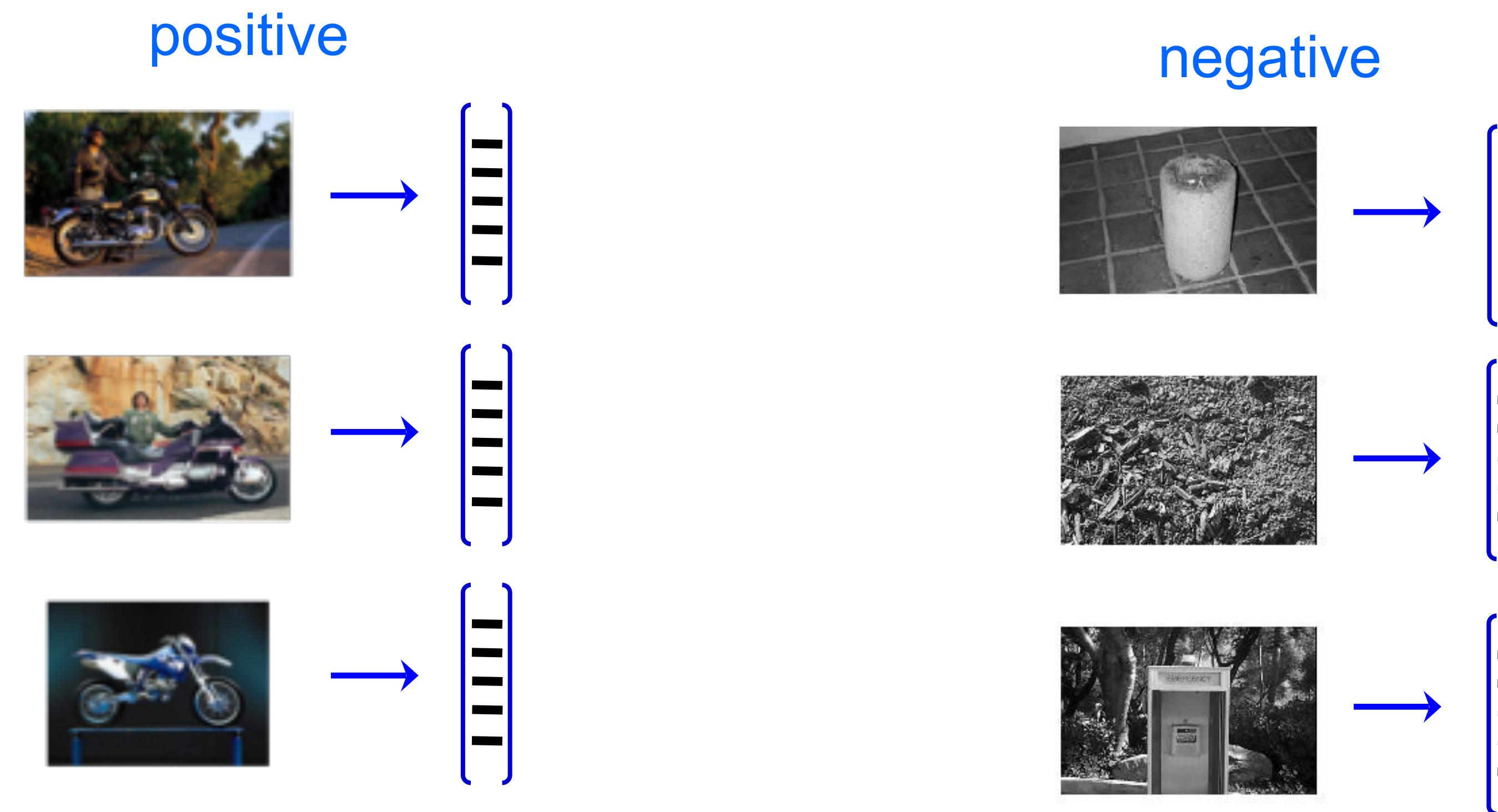
- Each image is represented by an aggregated histogram vector, typically 1000-4000 dimensional
- Normalized with L2 norm
- Fisher Vectors [Perronnin et al. ECCV'10]: improvements over Bag of Features

Bag-of-features for image classification



Step 3: Classification

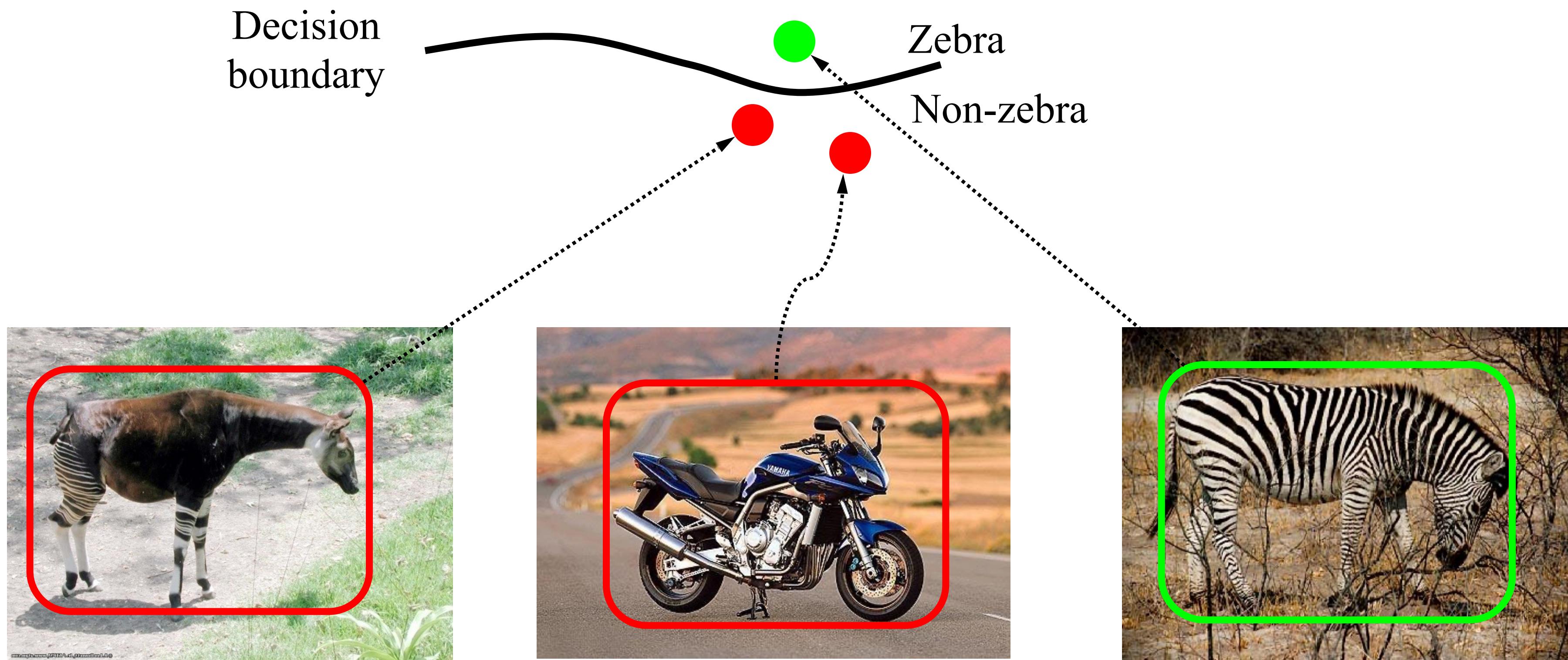
Training data: Vectors are histograms, one from each image



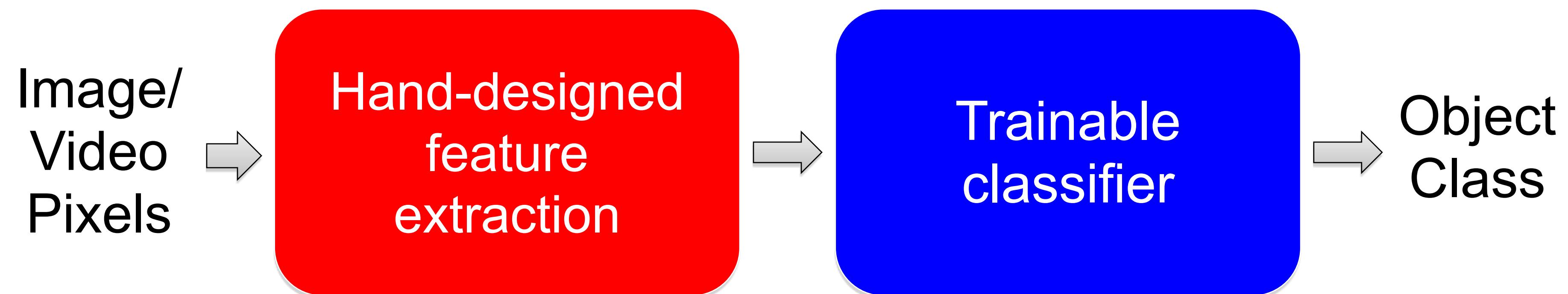
Train classifier, e.g. SVM

Step 3: Classification

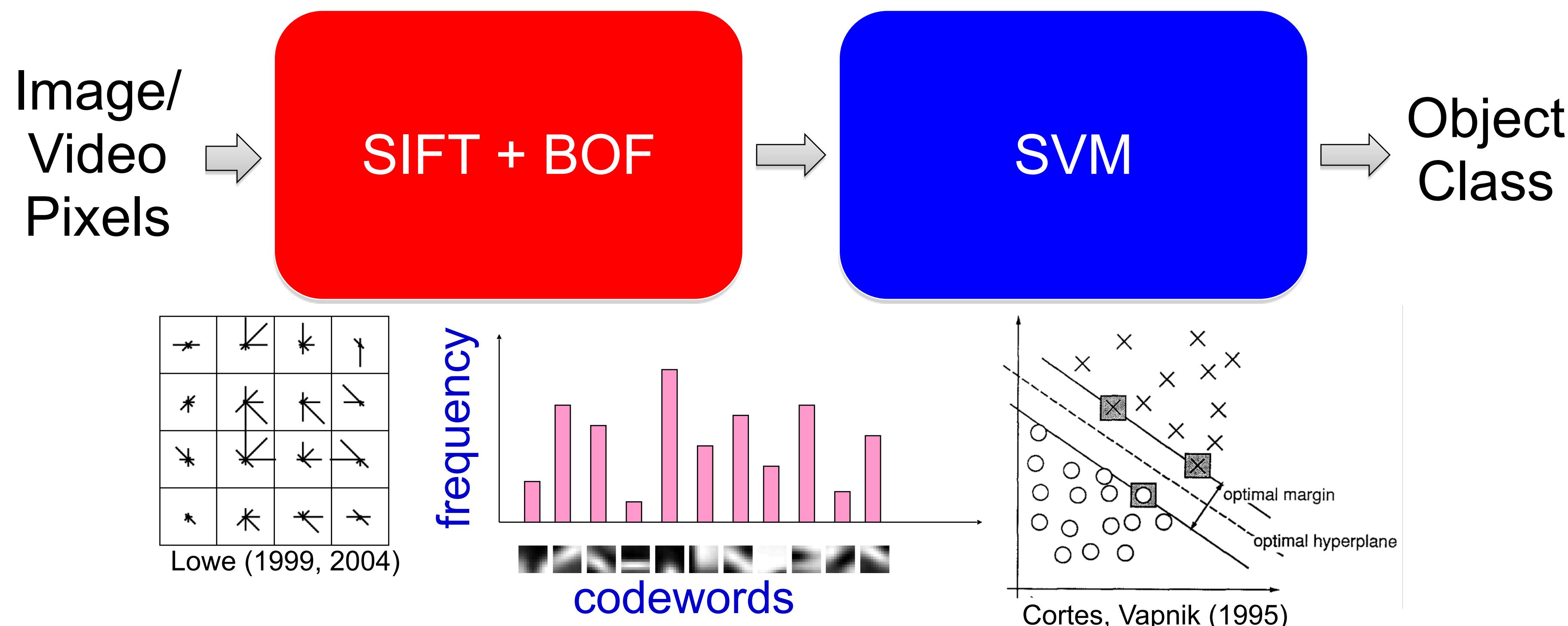
Learn a decision rule (classifier) assigning bag-of-features representations of images to different classes



Traditional Recognition Approach

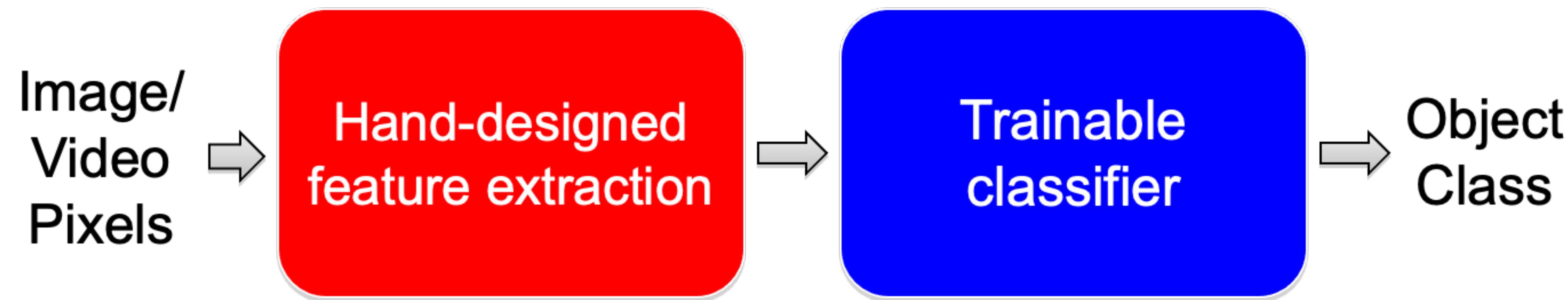


Traditional Recognition Example



- SIFT features
- BOF: Bag of Features / Visual Words (inspired by Bag of Words in NLP)
- SVM: Support Vector Machines for classification

Analogy to the traditional visual recognition pipeline



- Features are not learned (e.g., HOG, SIFT, Bag of Features)
- Trainable classifier is often generic (e.g., SVM, Random Forest)

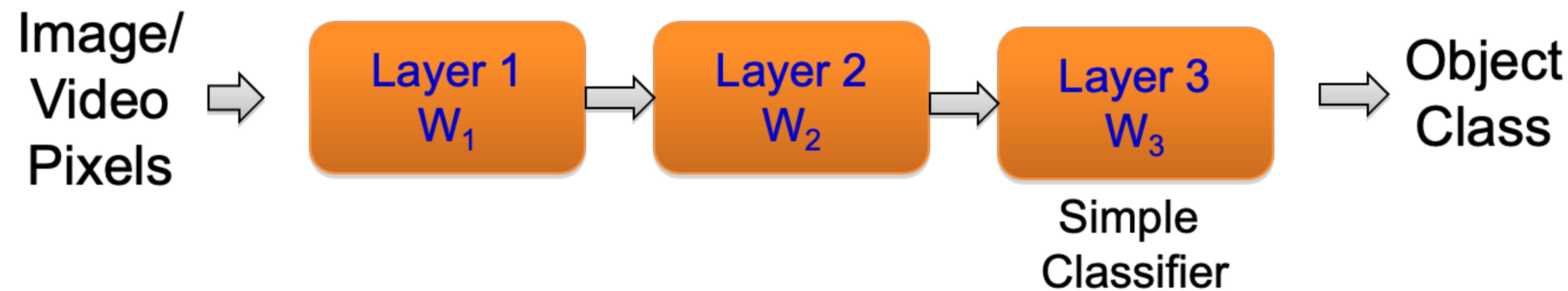
Analogy to the traditional visual recognition pipeline

What about learning the features?



- Features are learned “end-to-end” (i.e., pixels are input)
- “Feature hierarchy” all the way from pixels to classifier
- Each layer extracts features from the output of previous layer
- Train all layers jointly

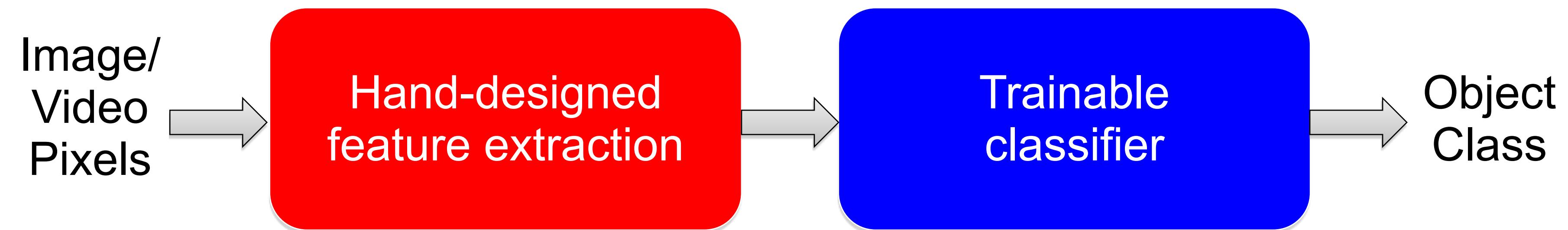
Analogy to the traditional visual recognition pipeline



- Features are learned “end-to-end” (i.e., pixels are input)
- “Feature hierarchy” all the way from pixels to classifier
- Each layer extracts features from the output of previous layer
- Train all layers jointly

“Shallow” vs. “deep” models

Traditional recognition: “Shallow” architecture



Deep learning: “Deep” architecture



Agenda

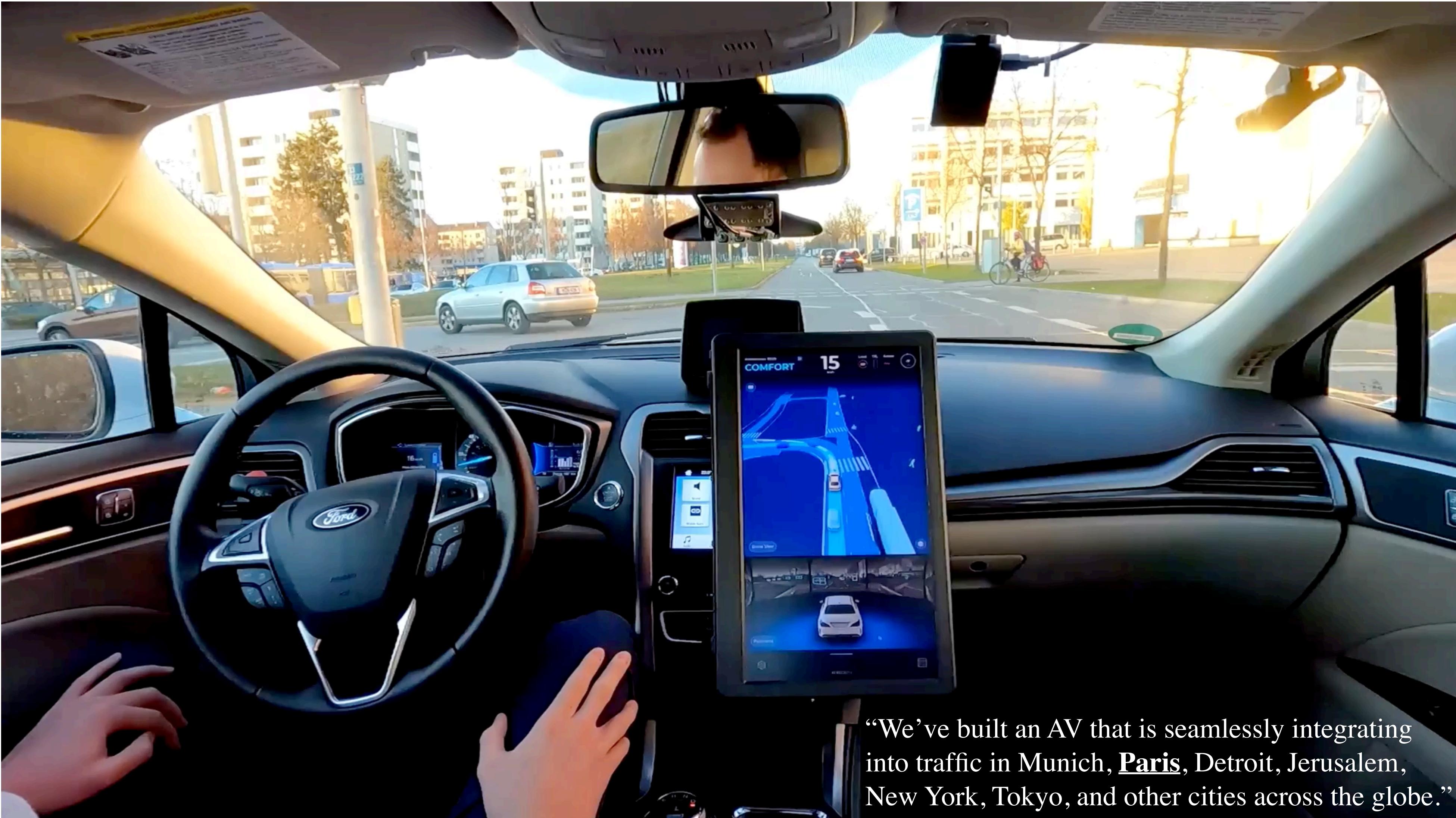
- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Neural Networks in Production

Face detection



Self-driving cars / Autonomous vehicles



“We’ve built an AV that is seamlessly integrating into traffic in Munich, Paris, Detroit, Jerusalem, New York, Tokyo, and other cities across the globe.”

Shopping



Google Translate

≡ Google Translate

Text

Documents

Websites

DETECT LANGUAGE

ENGLISH

FRENCH

SF ▾



FRENCH

TURKISH

ENGLISH



this course is so interesting



ce cours est tellement intéressant



29 / 5,000



What is “Deep” Learning?

Recap: Basics of supervised learning

- n training data pairs

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$$

- Learn a predictor/decision function

$$\hat{f} : \mathcal{X} \rightarrow \mathcal{A}$$

- By minimizing

$$\sum_{i=1}^n l(f(x_i), y_i)$$

Recap: Basics of supervised learning

- n training data pairs
- Learn a predictor/decision function
- By minimizing

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$$

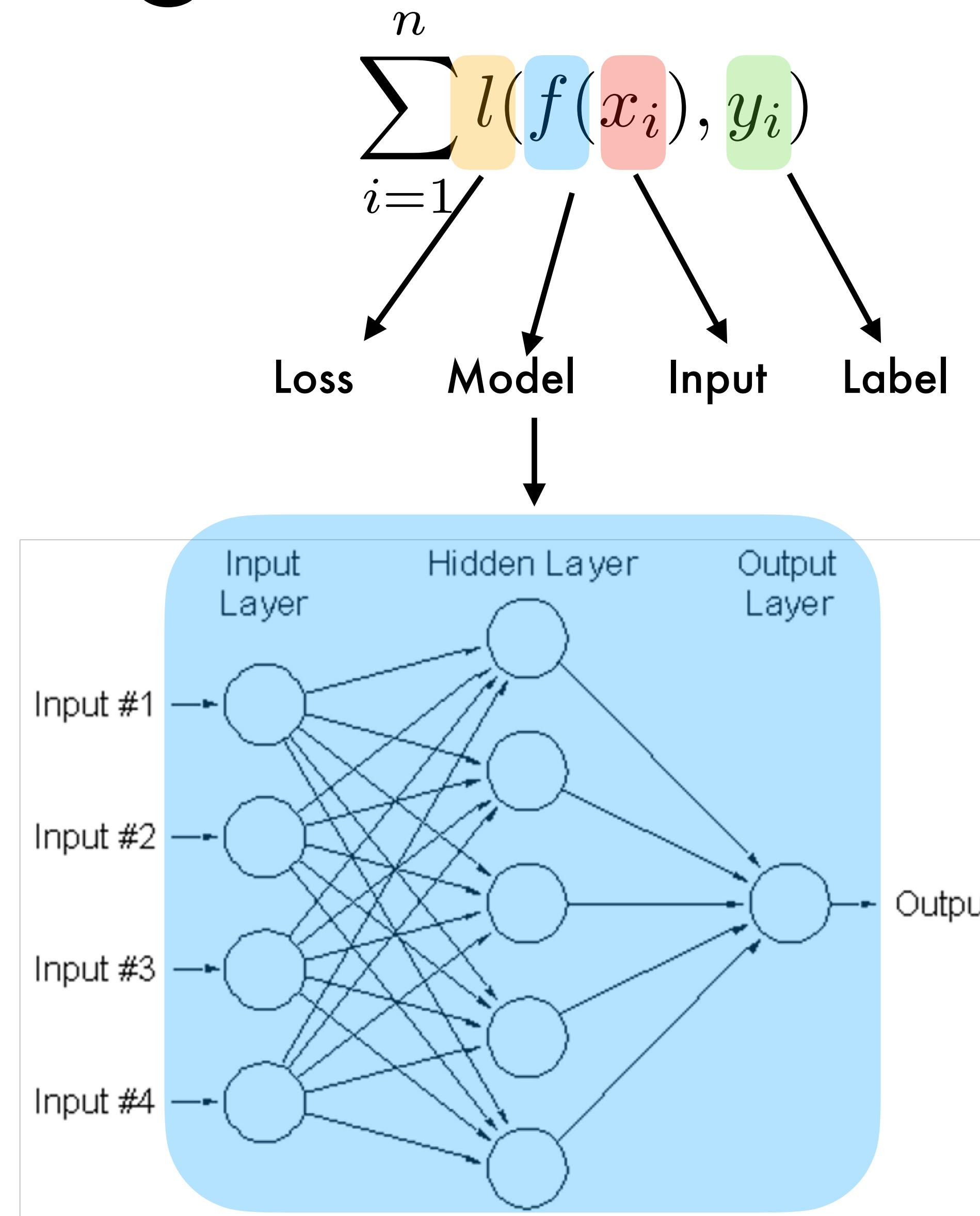
$$\hat{f}: \mathcal{X} \rightarrow \mathcal{A}$$

$$\sum_{i=1}^n l(f(x_i), y_i)$$

The diagram illustrates the components of a loss function. It shows a summation expression $\sum_{i=1}^n l(f(x_i), y_i)$. The term l is highlighted in yellow. The term $f(x_i)$ is highlighted in blue. The terms x_i and y_i are highlighted in red and green respectively. Arrows point from these highlighted terms to the words "Model", "Input", and "Label" below the equation.

Loss Model Input Label

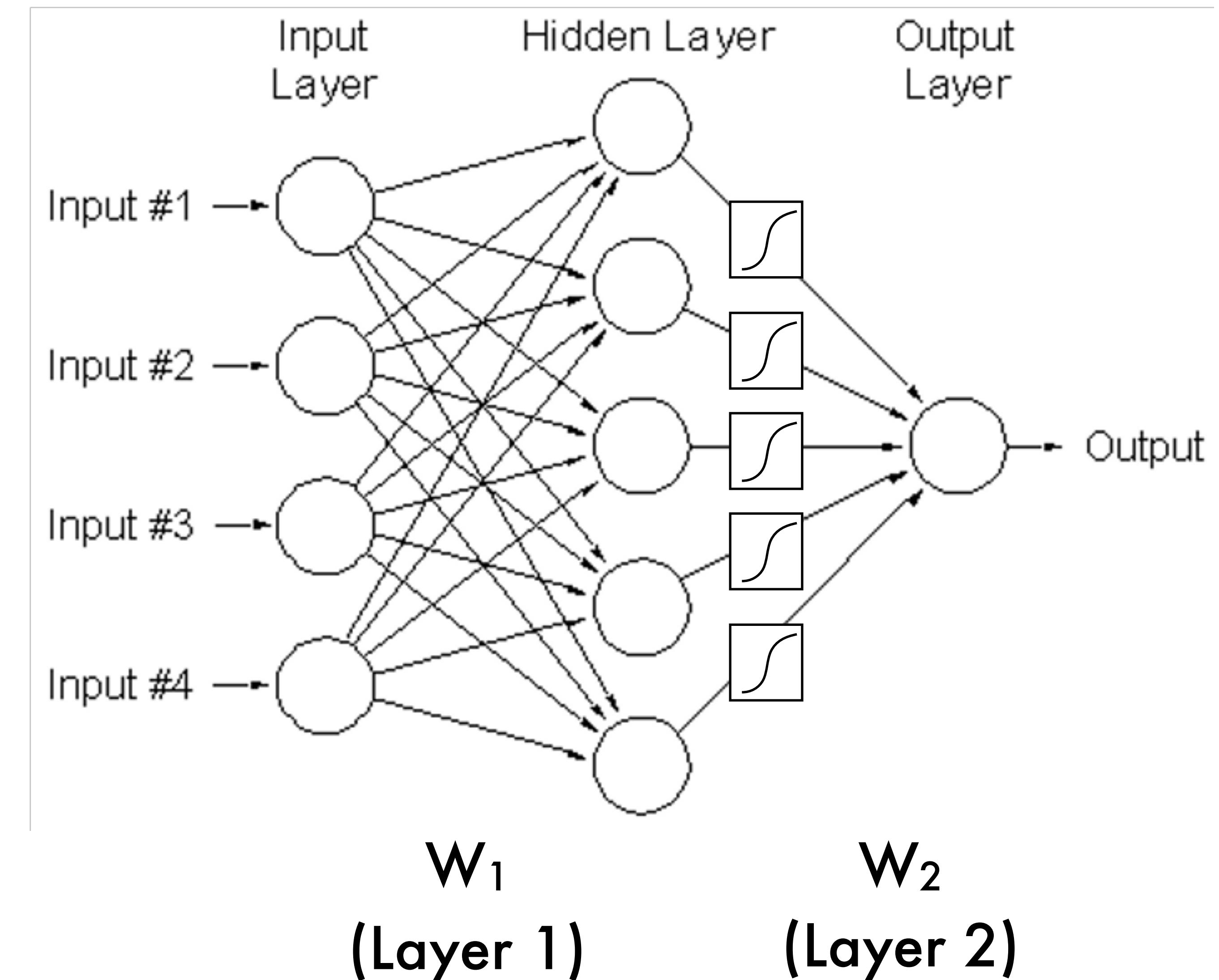
Deep learning



**Deep learning:
Model = neural network**

What is a “deep” neural network?

Stacking more than one **layer**



What is a layer?

Typically matrix multiplication! (But the function can take many forms*)

- **Fully-connected** layer
- **Convolution** layer
- **Pooling** layer (e.g., Max-pooling)
- **Non-linearity** layer (e.g., ReLU)
- **Attention** layer
- ...

*requirement to be differentiable if optimized with gradient descent algorithm variants

Recap: Perceptrons

[Rosenblatt, 1957]

Most basic form of a neural network

Input

Weights

x_1

w_1

x_2

w_2

x_3

w_3

.

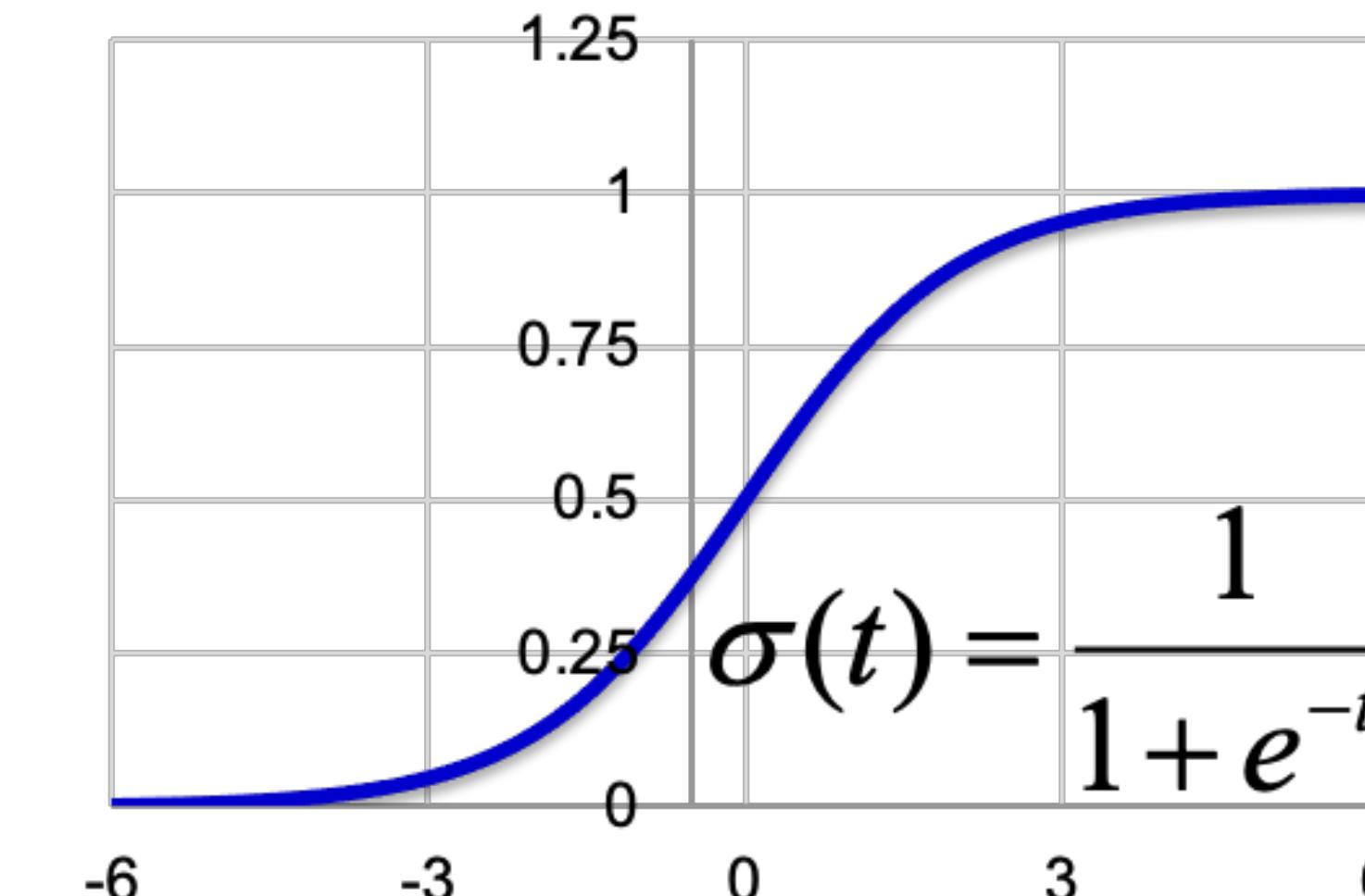
.

.

x_d

w_d

Sigmoid function:



Output: $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

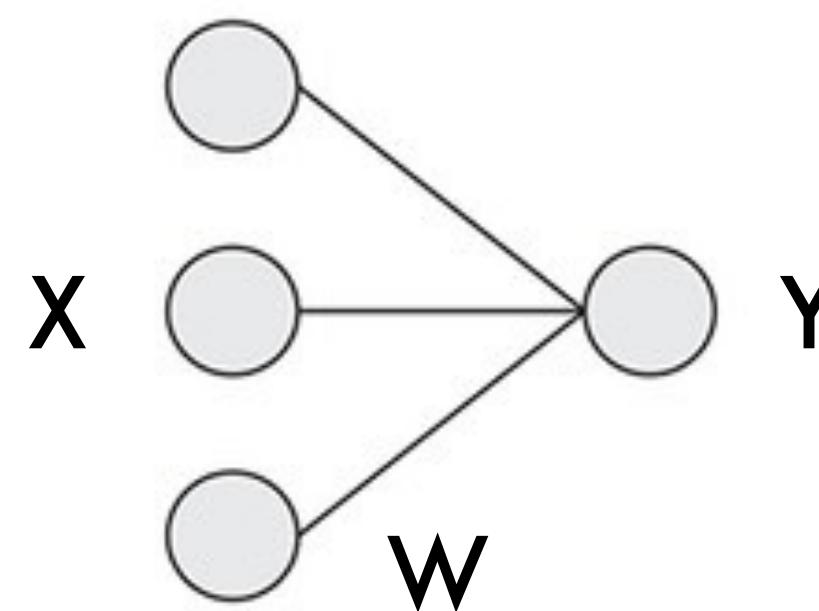
Non-linearity

Bias

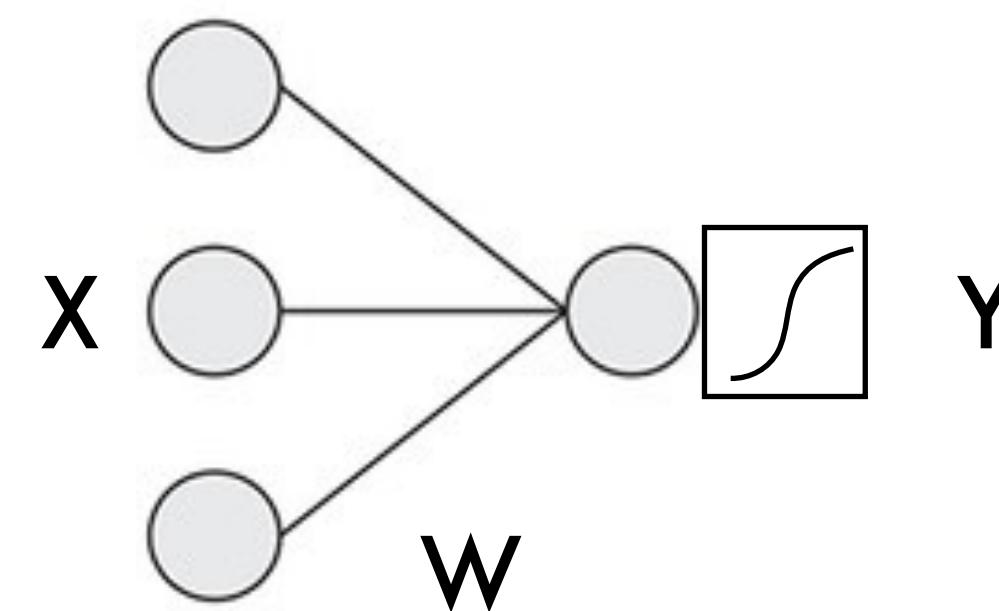
Linear combination
of inputs

Recap: Multi-Layer Perceptron (MLP)

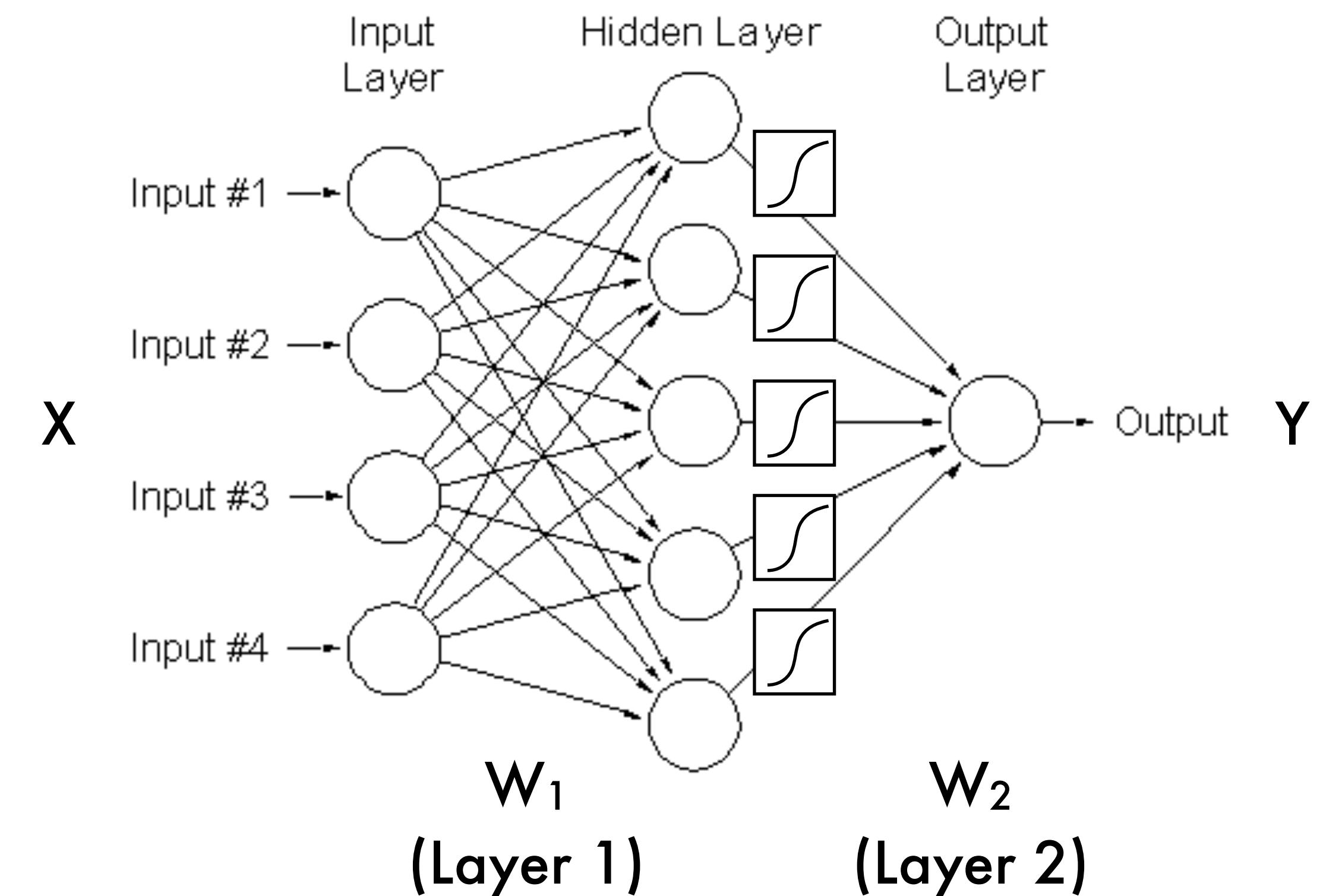
Linear regression:



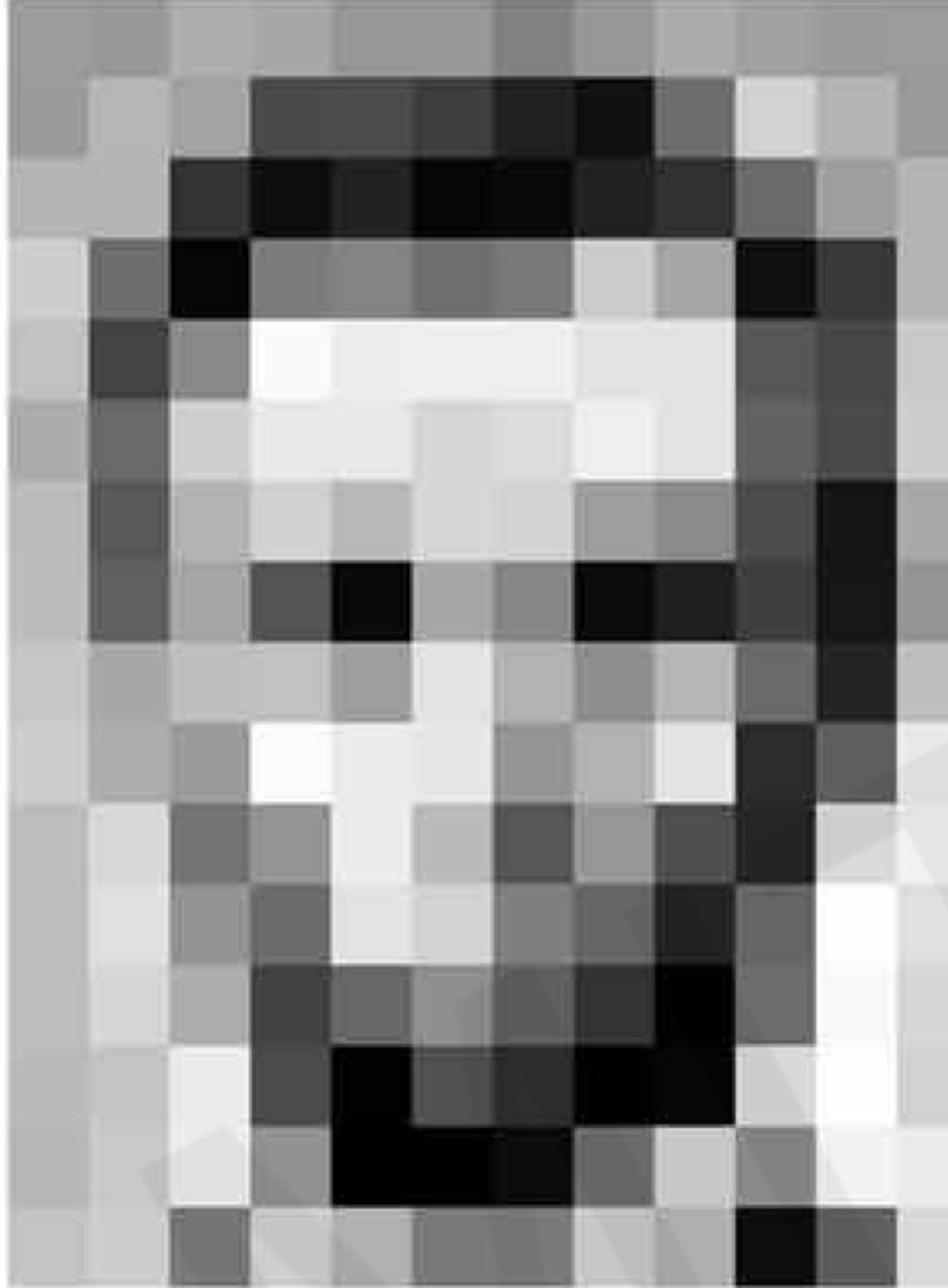
Perceptron:



MLP:



Images are numbers



157	153	174	168	150	152	129	151	172	161	156	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	63	179	209	185	215	211	158	139	15	20	169
189	97	165	64	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	118	149	236	187	86	150	79	36	218	241
190	224	147	168	227	210	127	102	56	101	255	224
190	214	173	65	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

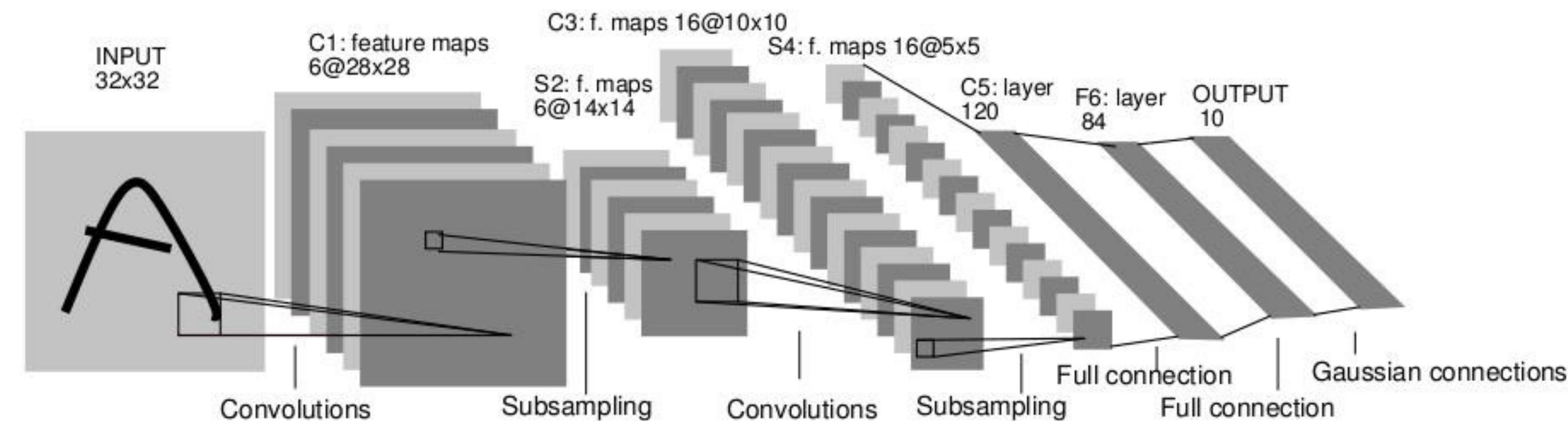
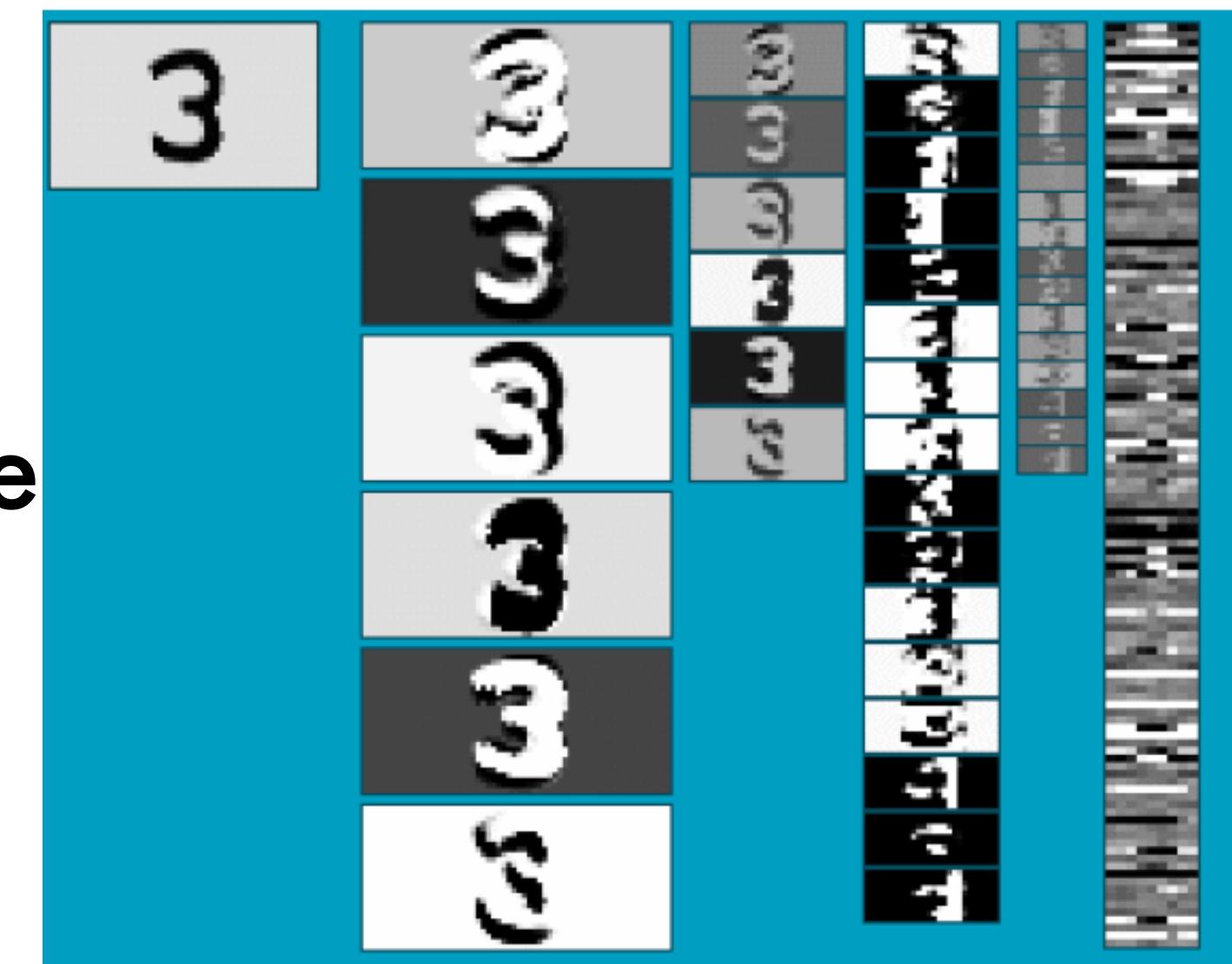
What the computer sees

157	153	174	168	150	152	129	151	172	161	156	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	15	20	169
189	97	165	64	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	118	149	236	187	86	150	79	36	218	241
190	224	147	168	227	210	127	102	56	101	255	224
190	214	173	65	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

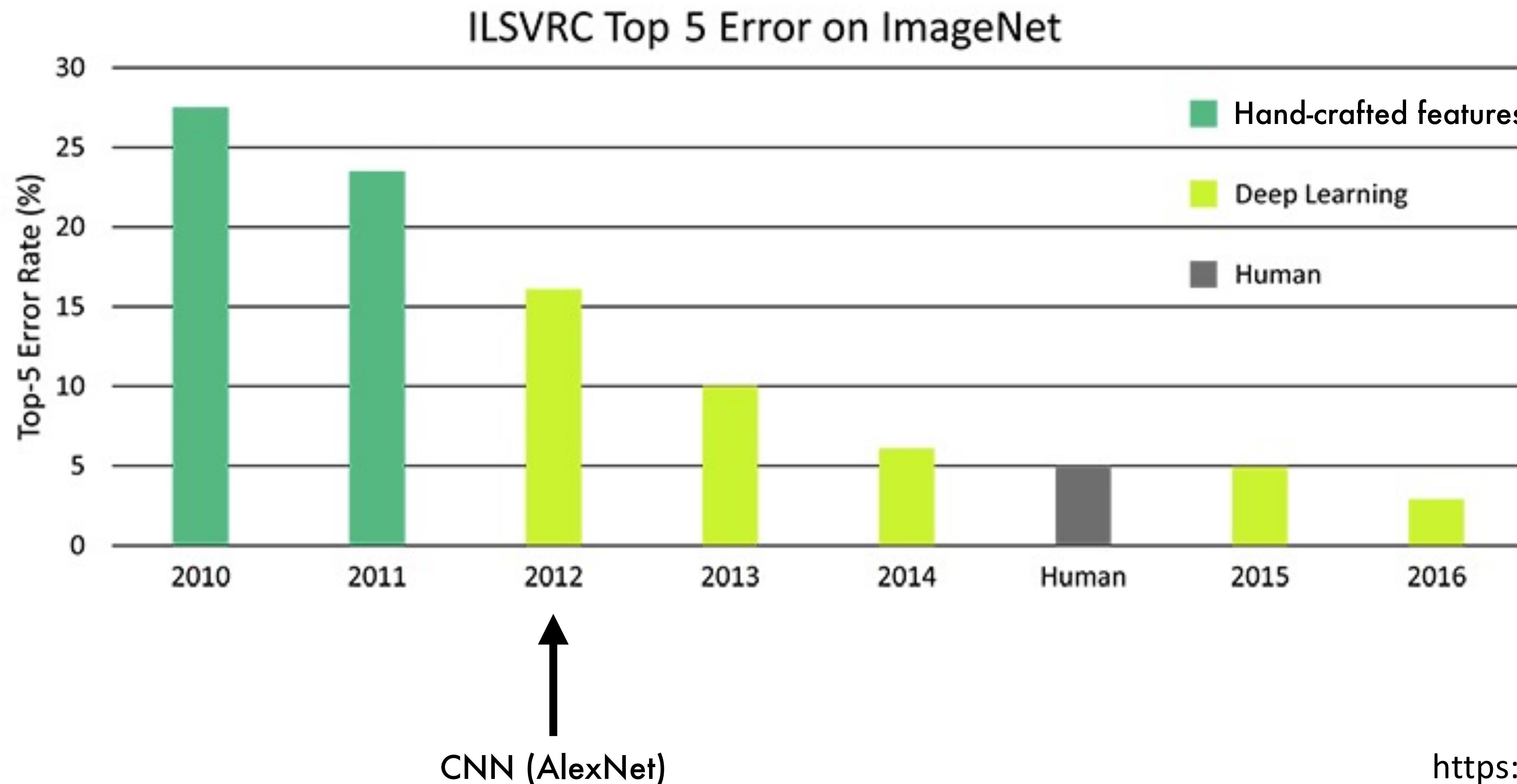
Review: Convolutional Neural Networks (CNN)

- Neural network with specialized **connectivity** structure
- Stack multiple stages of feature extractors
- Higher stages compute more global, more invariant features
- Classification layer at the end



Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278–2324, 1998.

Progress on ImageNet

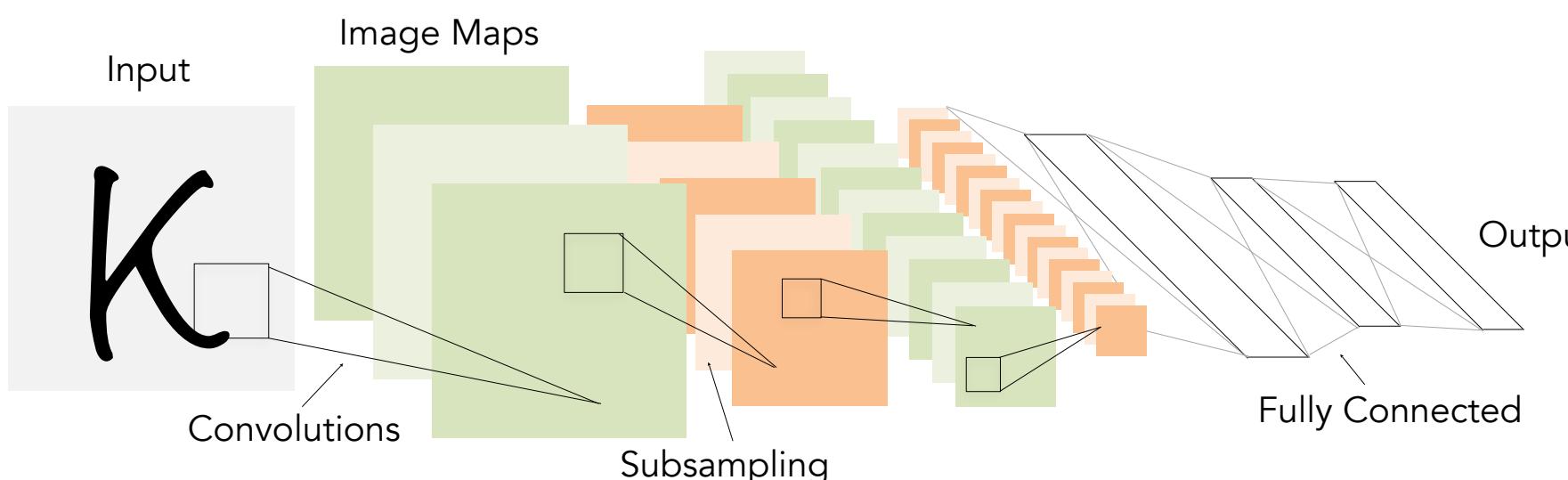


<https://www.dsiac.org>

CNNs were not invented overnight

1998

LeCun et al.

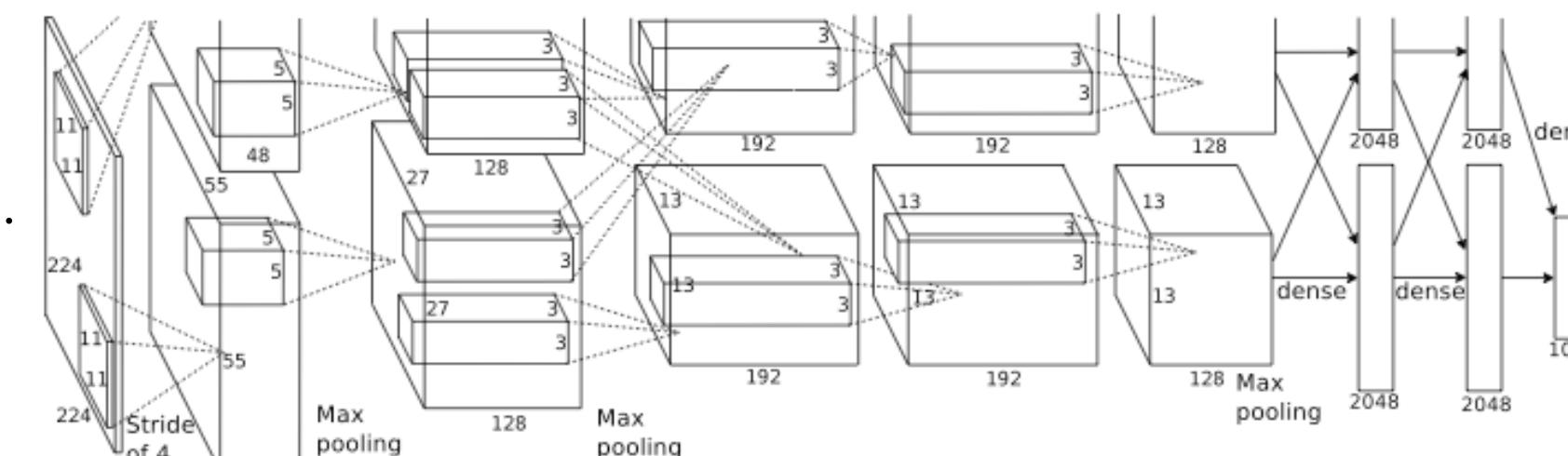


of transistors
10⁶
pentium® II

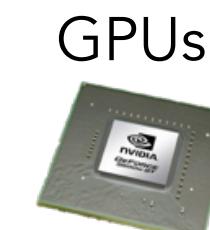
of pixels used in training
10⁷ NIST

2012

Krizhevsky et al.



of transistors
10⁹
Intel Xeon processor



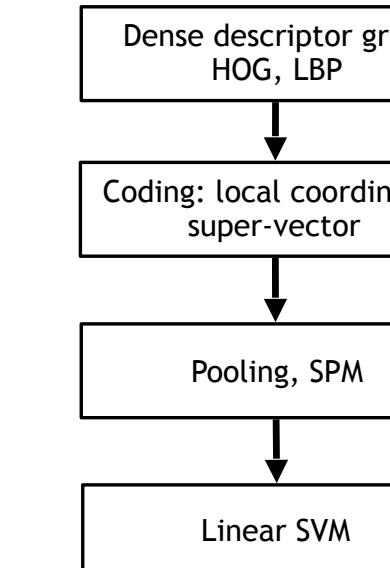
of pixels used in training
10¹⁴ IMAGENET

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012.
Reproduced with permission.

IMAGENET Large Scale Visual Recognition Challenge

Year 2010

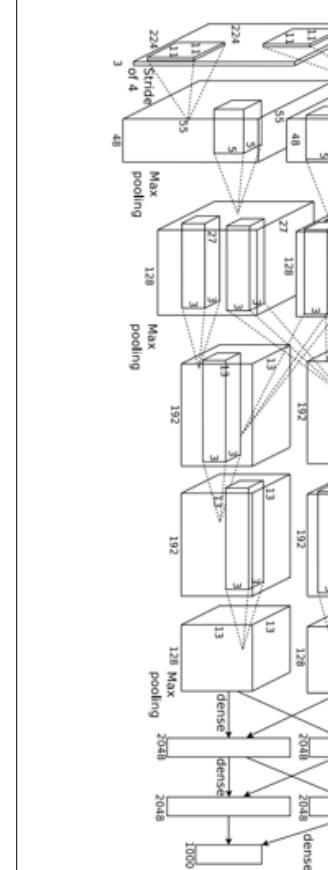
NEC-UIUC



[Lin CVPR 2011]

Year 2012

SuperVision

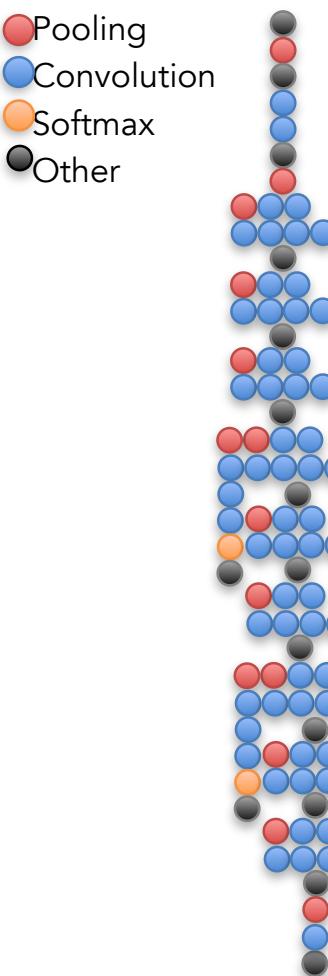


[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012.
Reproduced with permission.

Year 2014

GoogLeNet

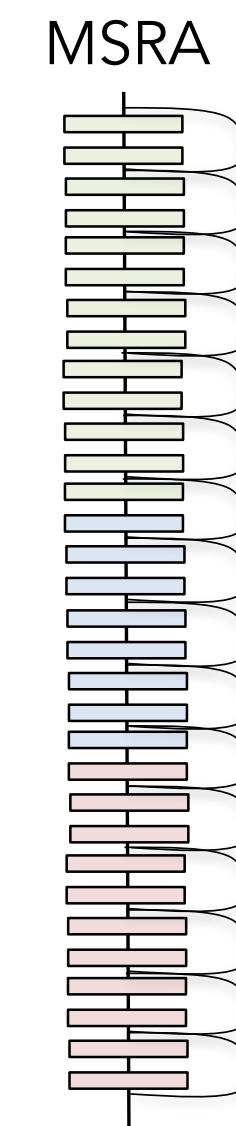
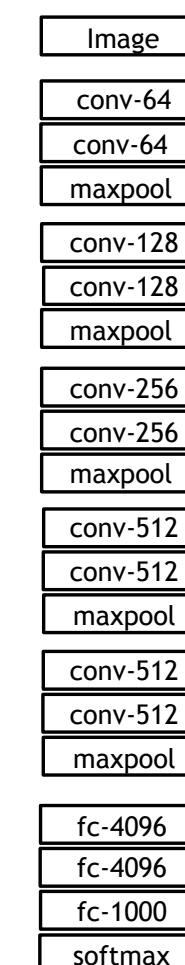


[Szegedy arxiv 2014]

[Simonyan arxiv 2014]

Year 2015

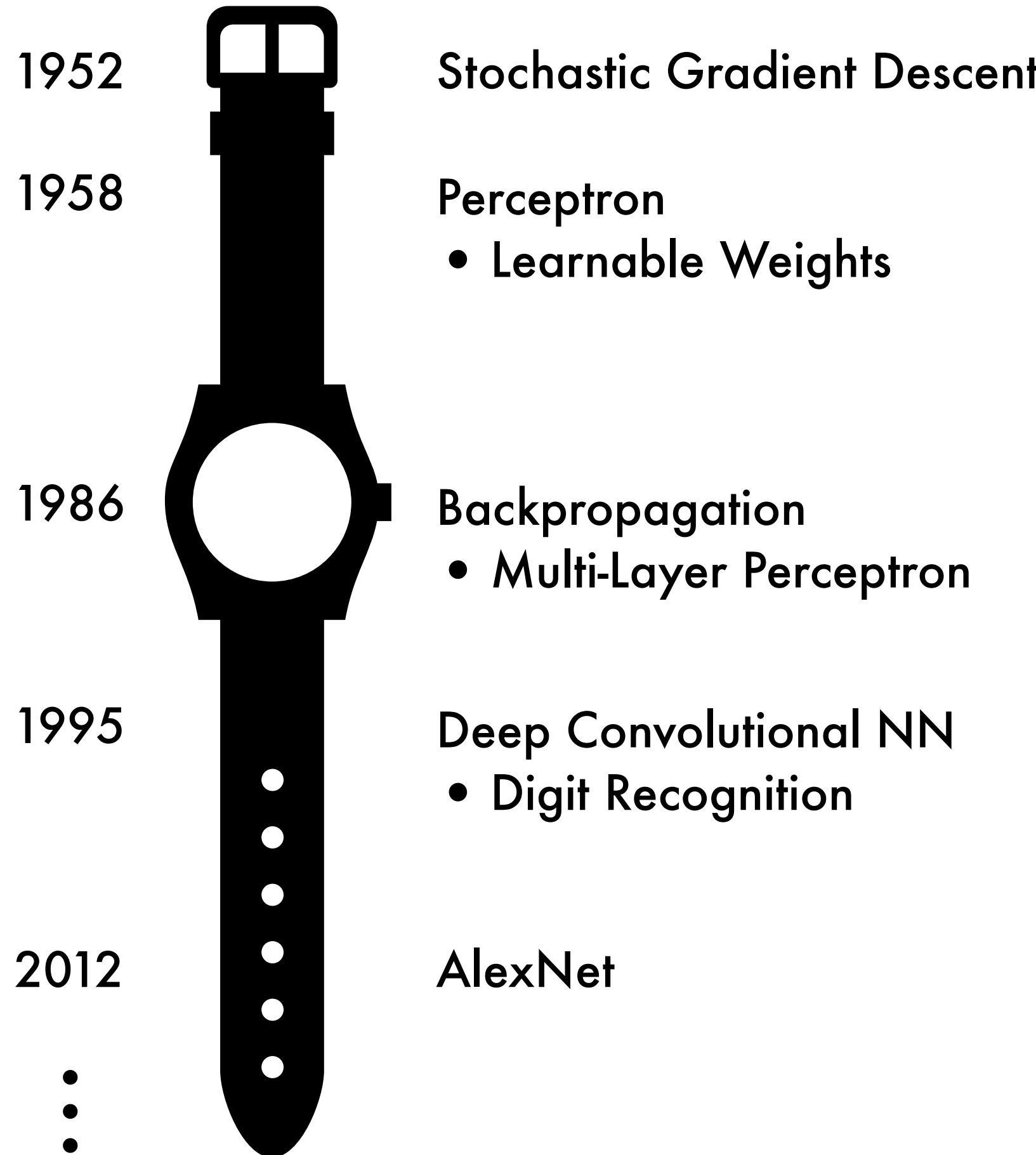
VGG



[He ICCV 2015]

Why now?

Neural Networks date back decades.



1. Big Data

- Larger **datasets**
- Easier collection & storage

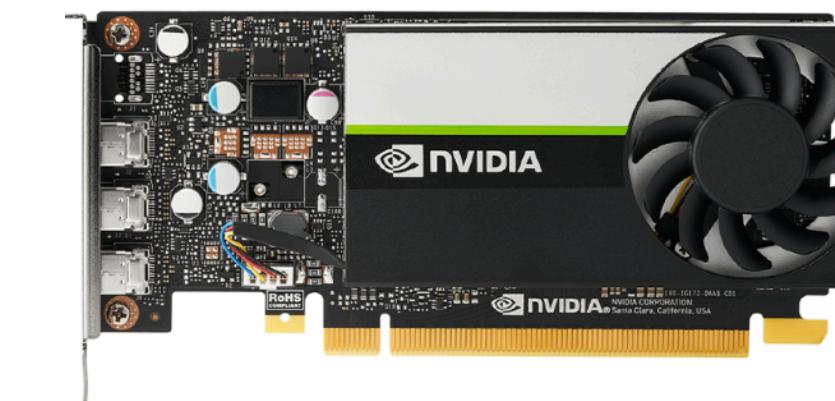


LAION-400-MILLION OPEN DATASET

...

2. Hardware

- Graphics Processing Units (**GPUs**)
- Massively Parallelizable



3. Software

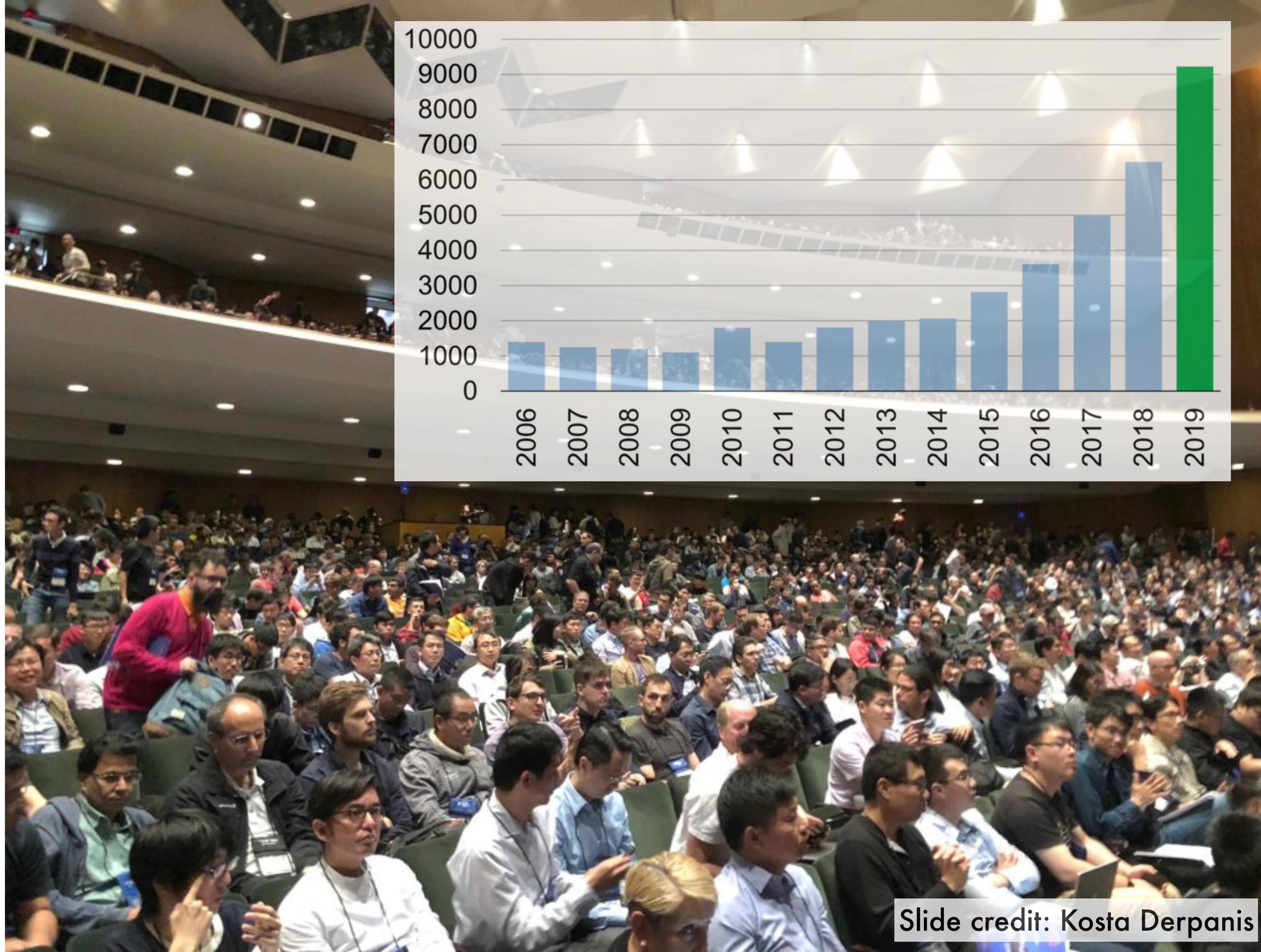
- Improved Techniques
- New Models
- **Toolboxes**

 PyTorch

 TensorFlow

...

CVPR:
**(Computer Vision Pattern
Recognition Conference)**



Slide credit: Kosta Derpanis

Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Standard layers

- 1. **Fully-connected** layer
- 2. **Convolution** layer
- 3. **Pooling** layer (e.g., Max-pooling)
- 4. **Non-linearity** layer (e.g., ReLU)
- 5. **Normalization** layer (e.g., BatchNorm)
- ...

Input
image

[$224 \times 224 \times 3$]

Conv1

Conv2

Conv3

Conv4

Conv5

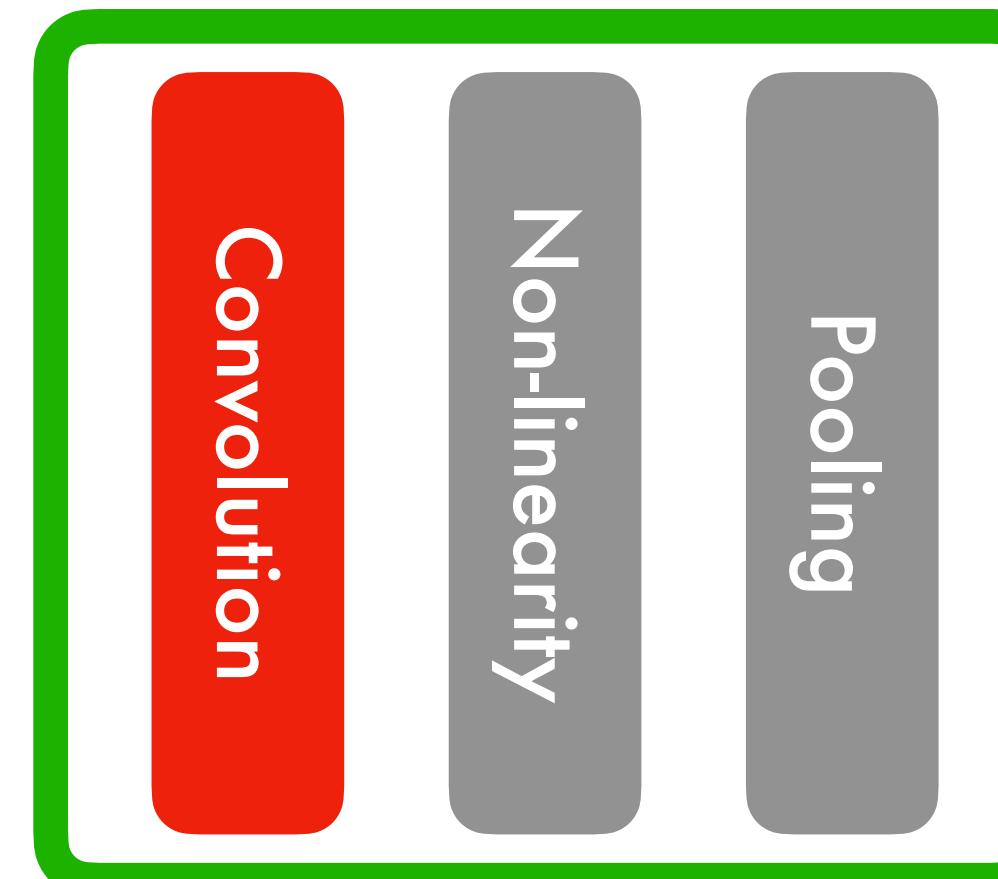
FC6

FC7

FC8

Convolutional block Fully-connected block

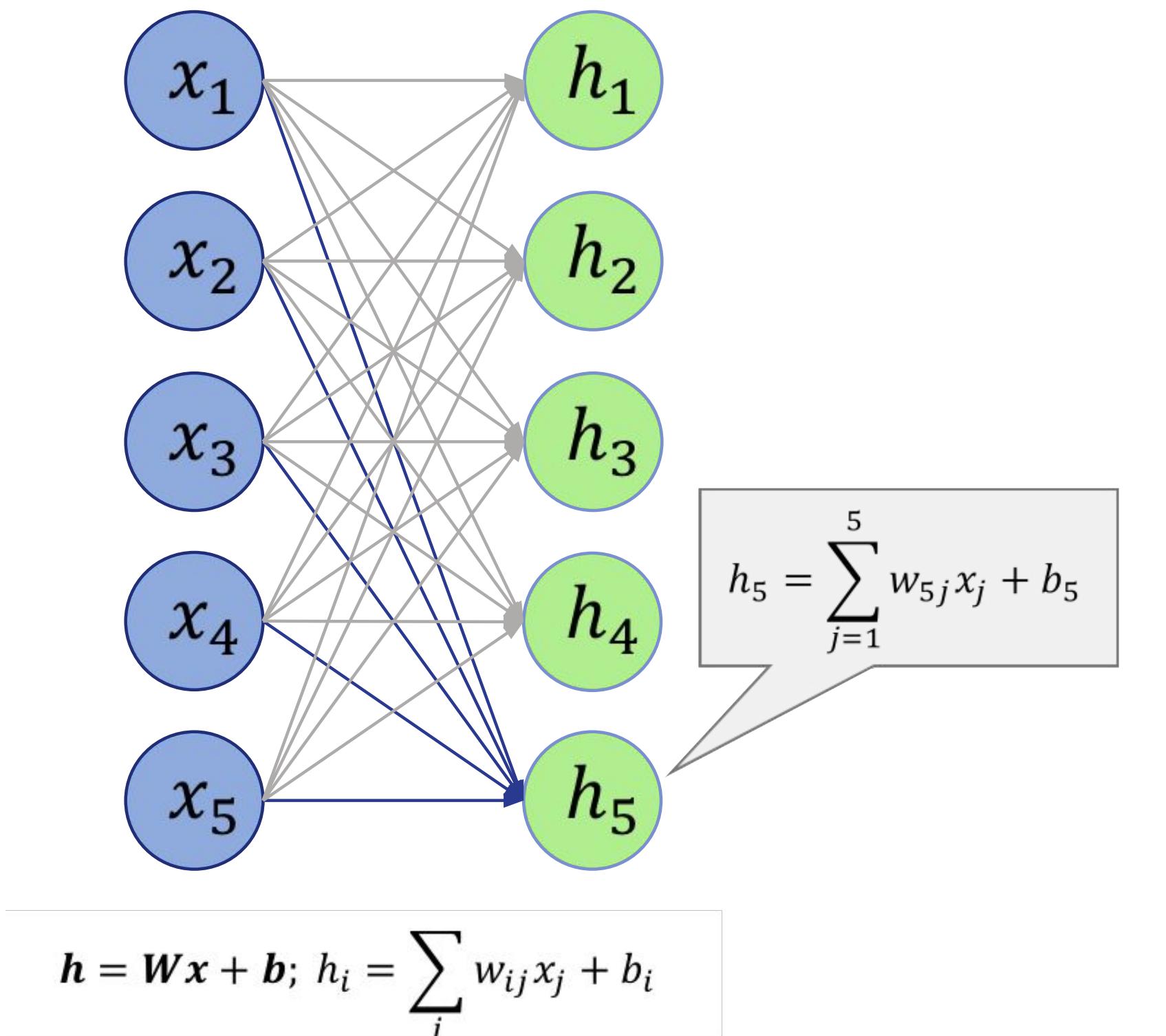
} Learnable
parameters



Class
probability

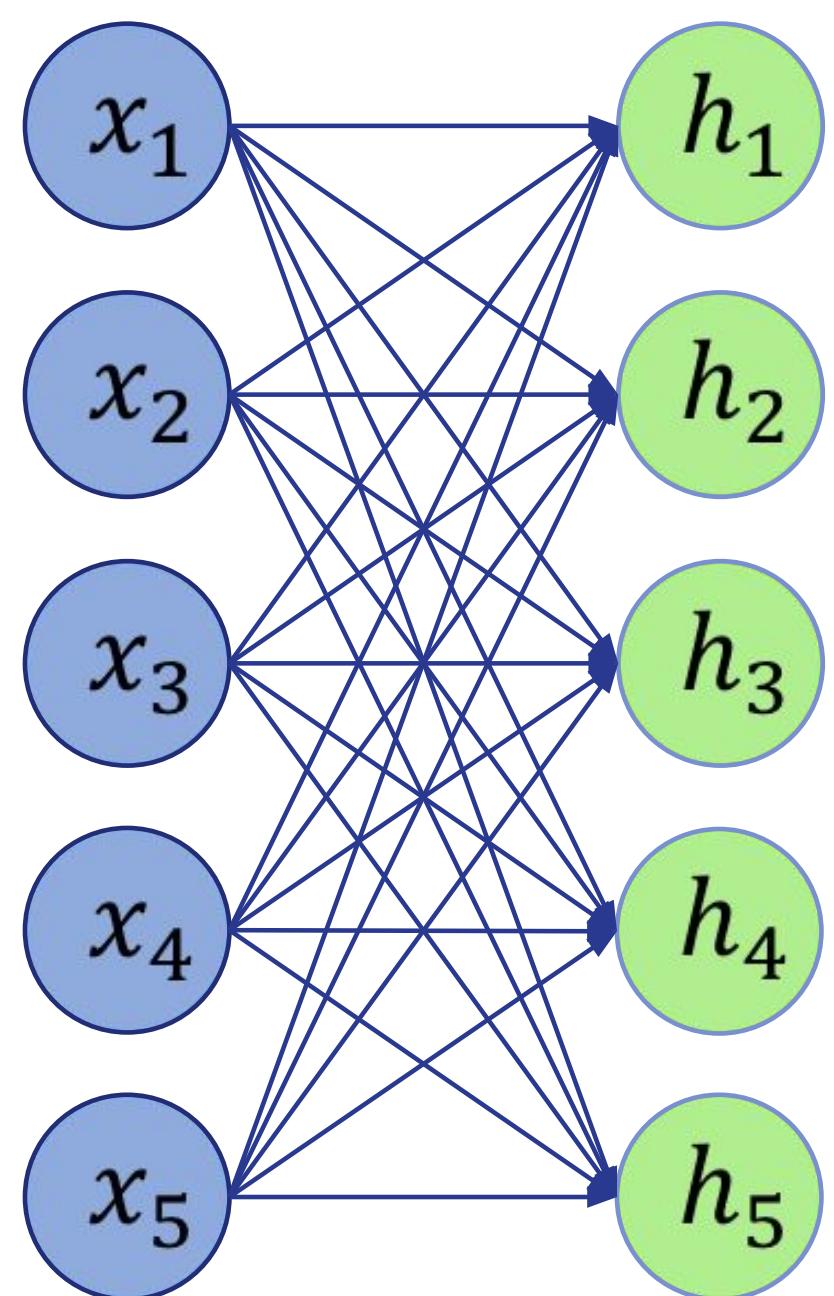
[1000]

1. Fully-connected layer



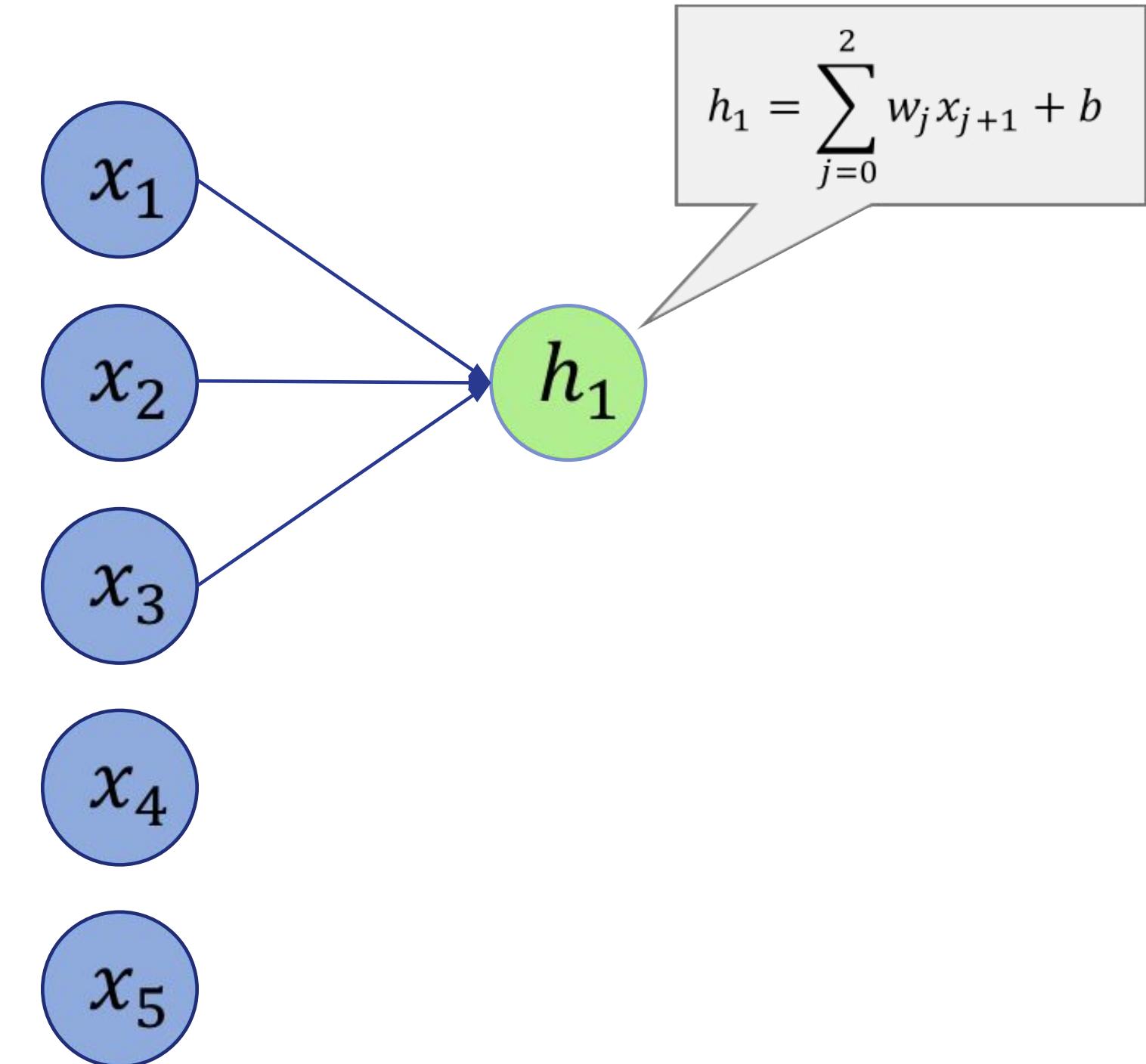
2. Convolution layer

Fully-connected



$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}; h_i = \sum_j w_{ij}x_j + b_i$$

1D Convolutional

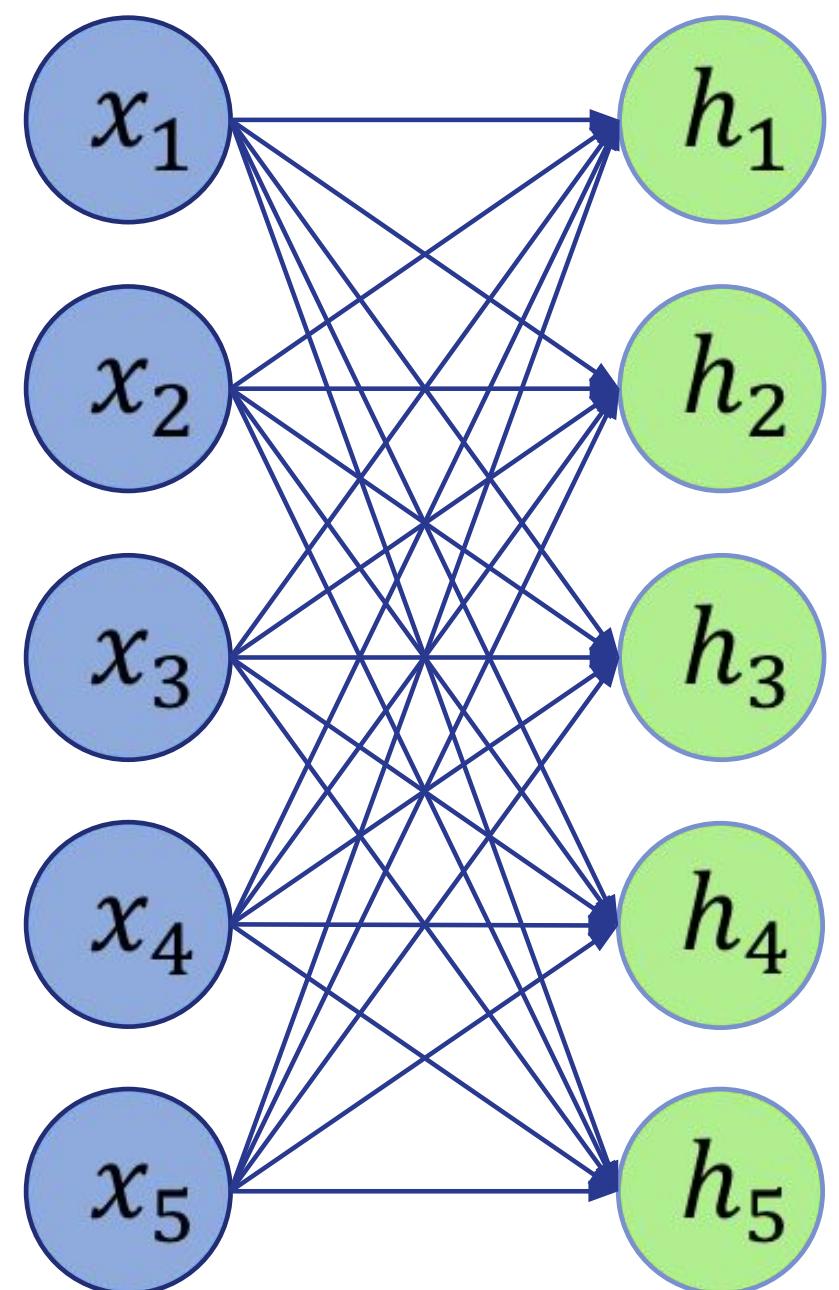


$$h_i = \sum_j w_j x_{j+i} + b$$

- Layer with a special connectivity structure
- Dependencies are local
- Translation invariance

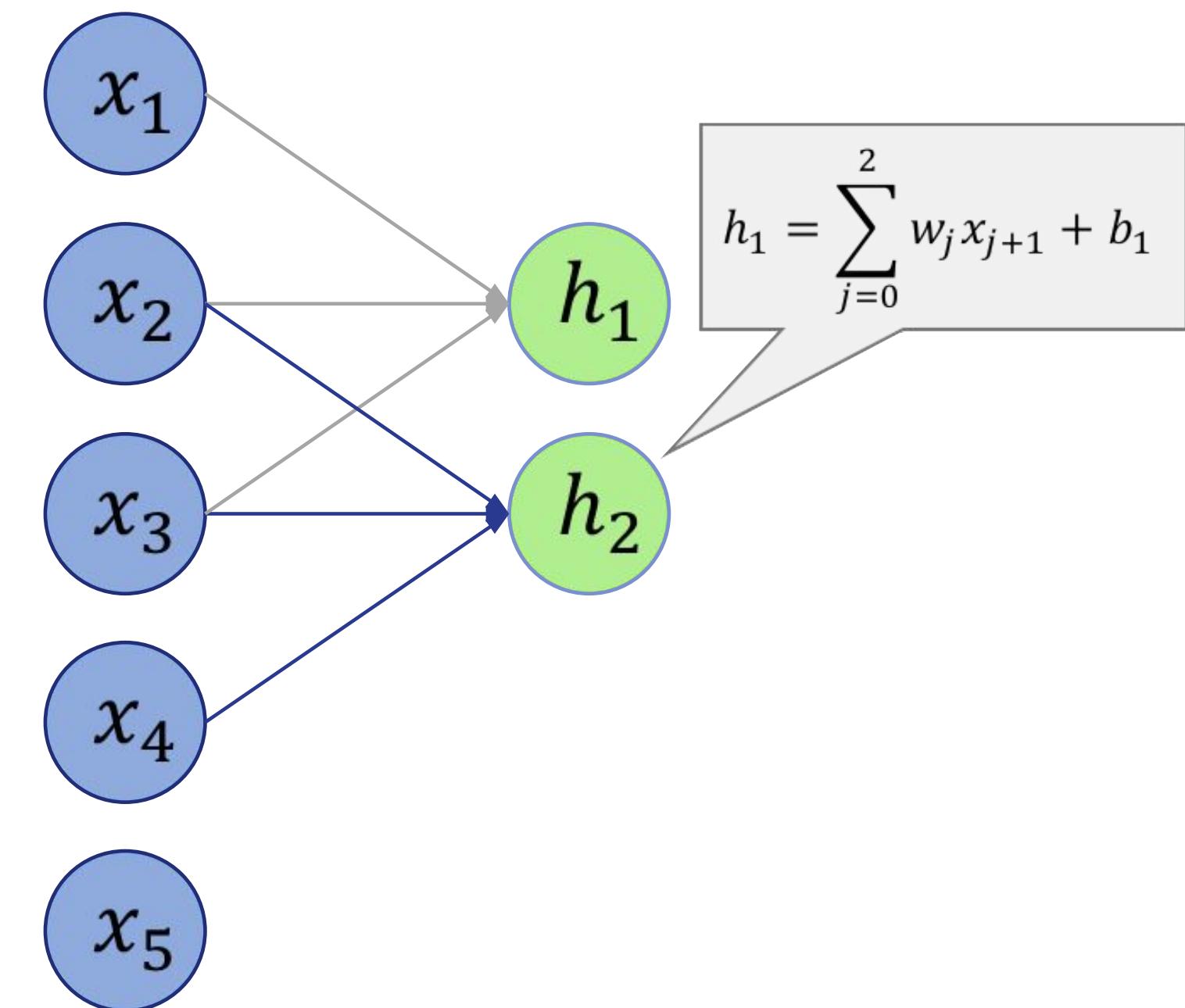
2. Convolution layer

Fully-connected



$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}; h_i = \sum_j w_{ij}x_j + b_i$$

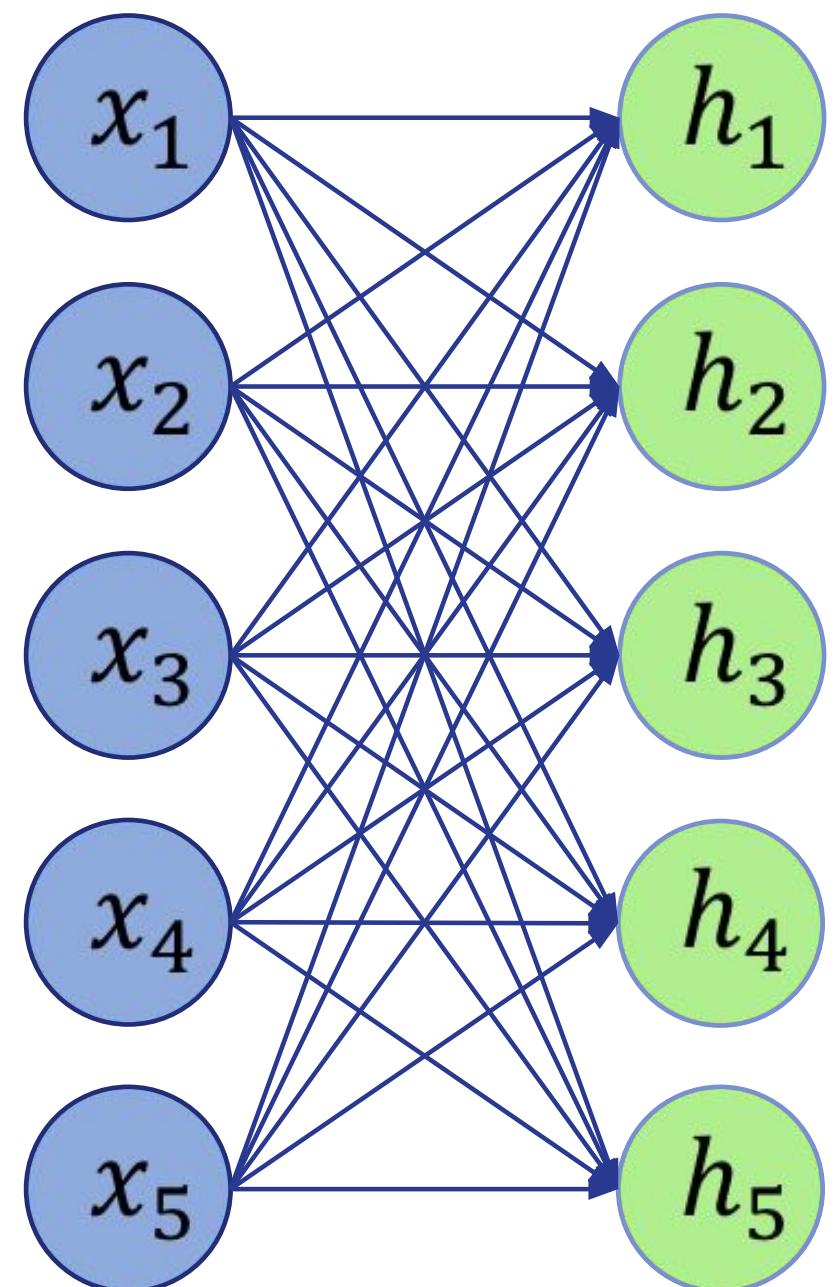
1D Convolutional



$$h_i = \sum_j w_j x_{j+i} + b$$

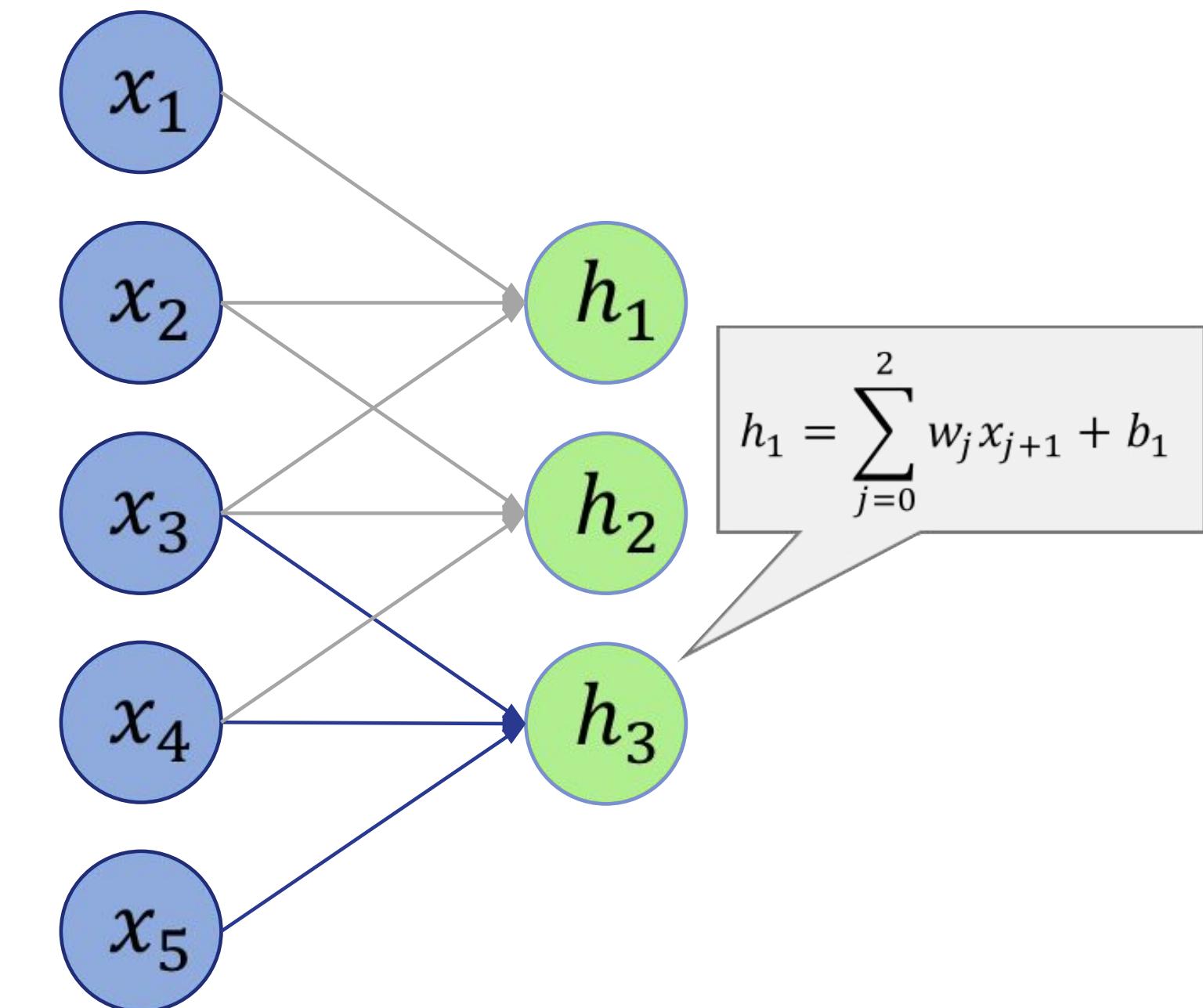
2. Convolution layer

Fully-connected



$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}; h_i = \sum_j w_{ij}x_j + b_i$$

1D Convolutional



$$h_i = \sum_j w_j x_{j+i} + b$$

2. Convolution layer

2D Convolutions

$$I \star K$$

The diagram shows the convolution operation between two matrices, I and K . The input matrix I is a 5x5 grid with values: [1, 0, 0, 1, 2], [0, 0, 0, 3, 0], [0, 1, 2, 1, 1], [1, 1, 3, 0, 0], [3, 0, 0, 0, 1]. The kernel matrix K is a 3x3 grid with values: [0, 0, 1], [0, 2, 0], [1, 1, 0]. The result of the convolution is shown as a 3x3 grid of values: [0, 0, 1], [0, 2, 0], [1, 1, 0]. The convolution is performed using zero padding.

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

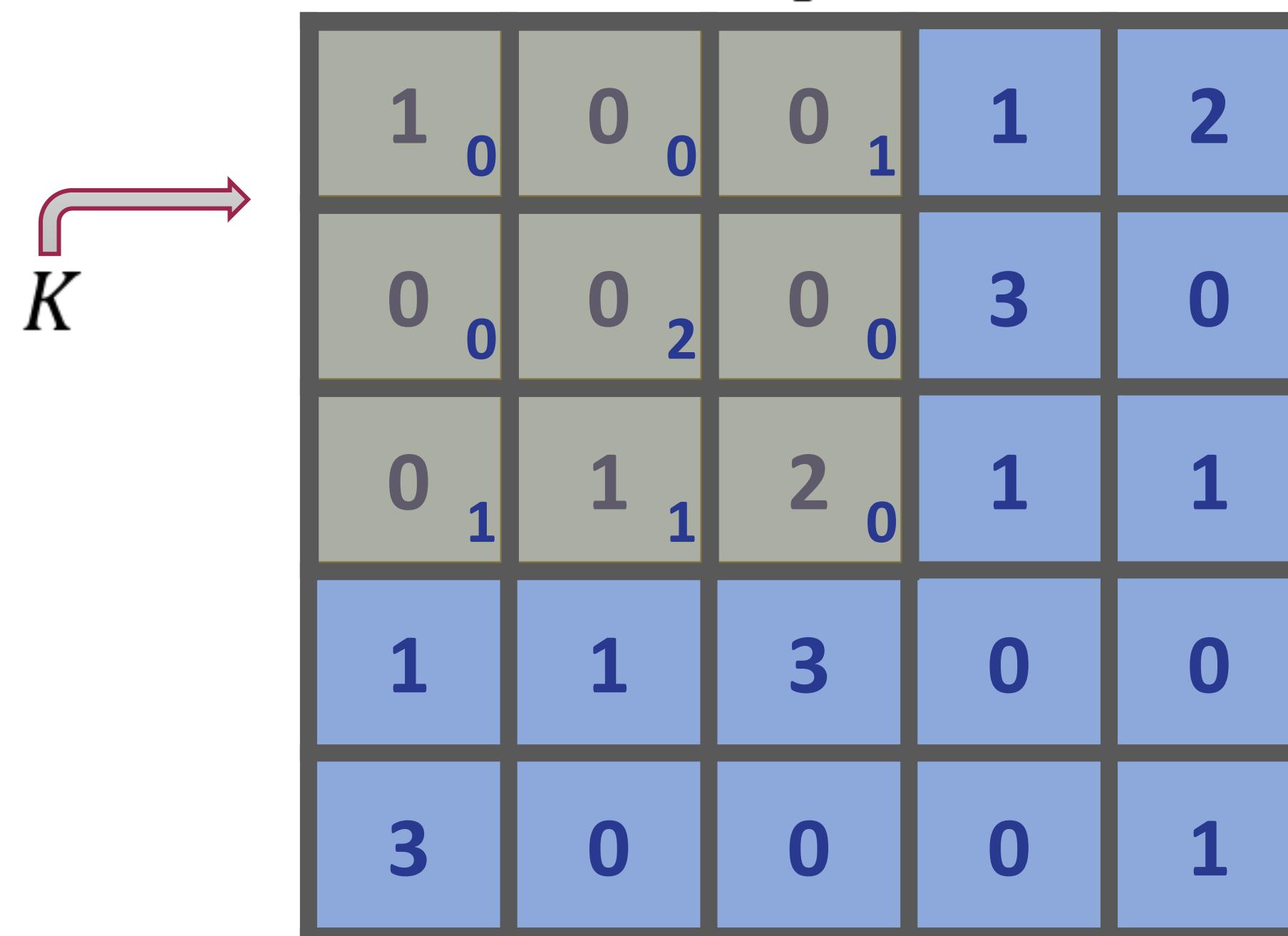
\star

0	0	1
0	2	0
1	1	0

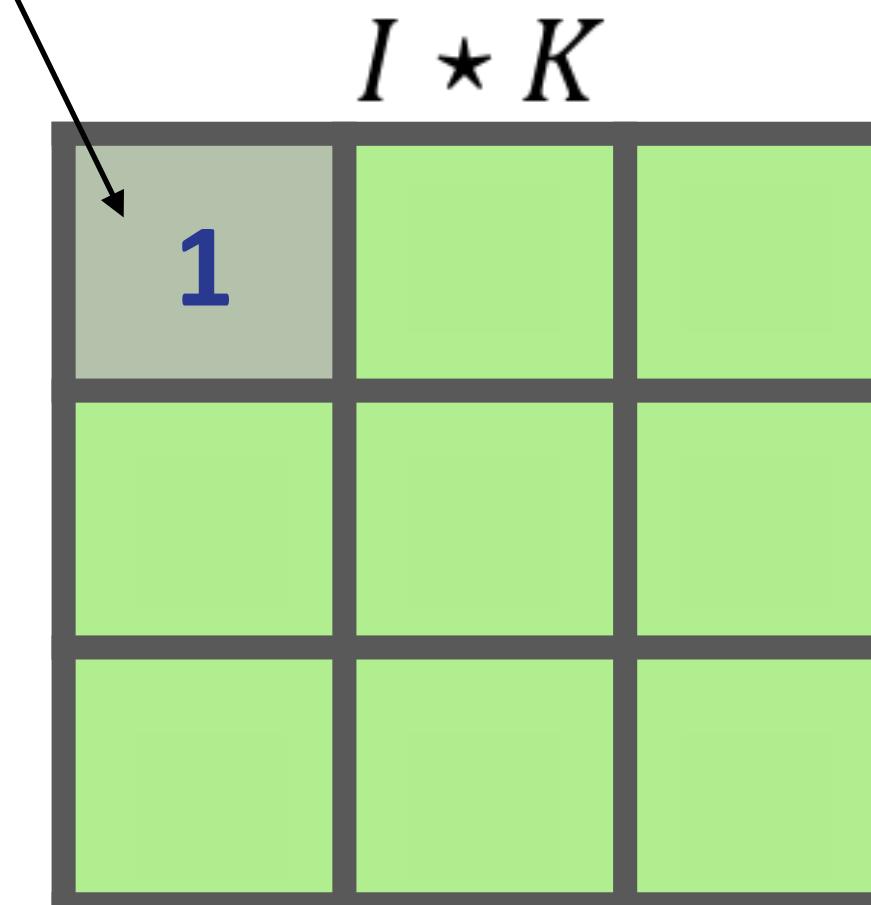
$$(I \star K)(i, j) == \sum_m \sum_n I(m, n)K(i + m, j + n)$$

2. Convolution layer

2D Convolutions



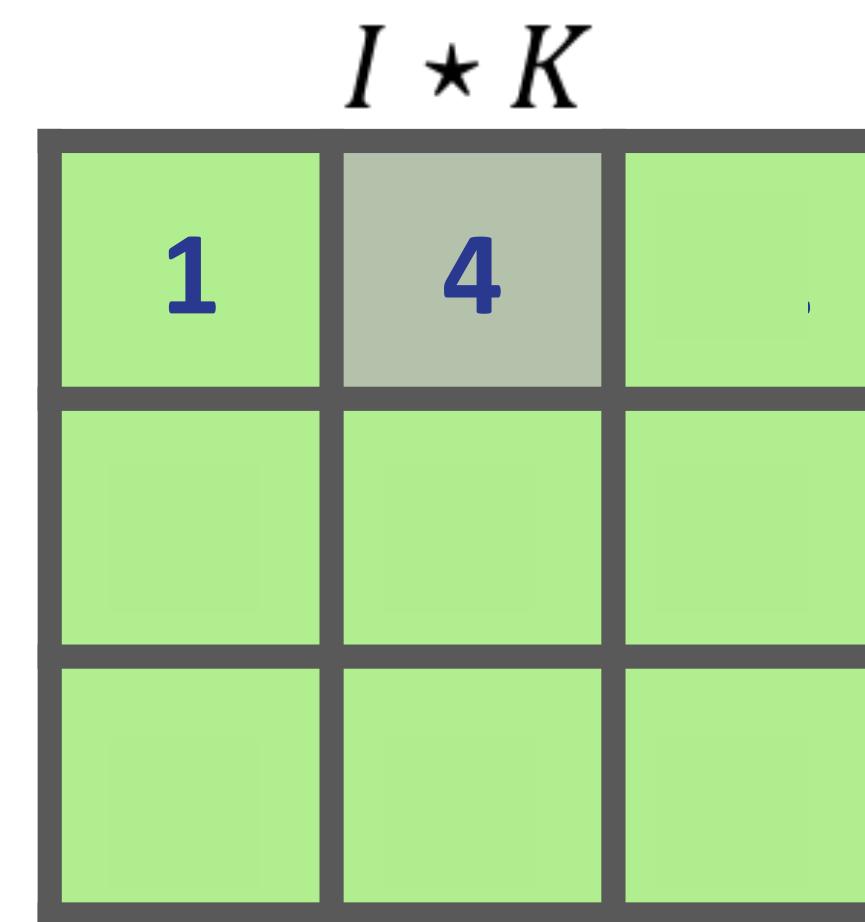
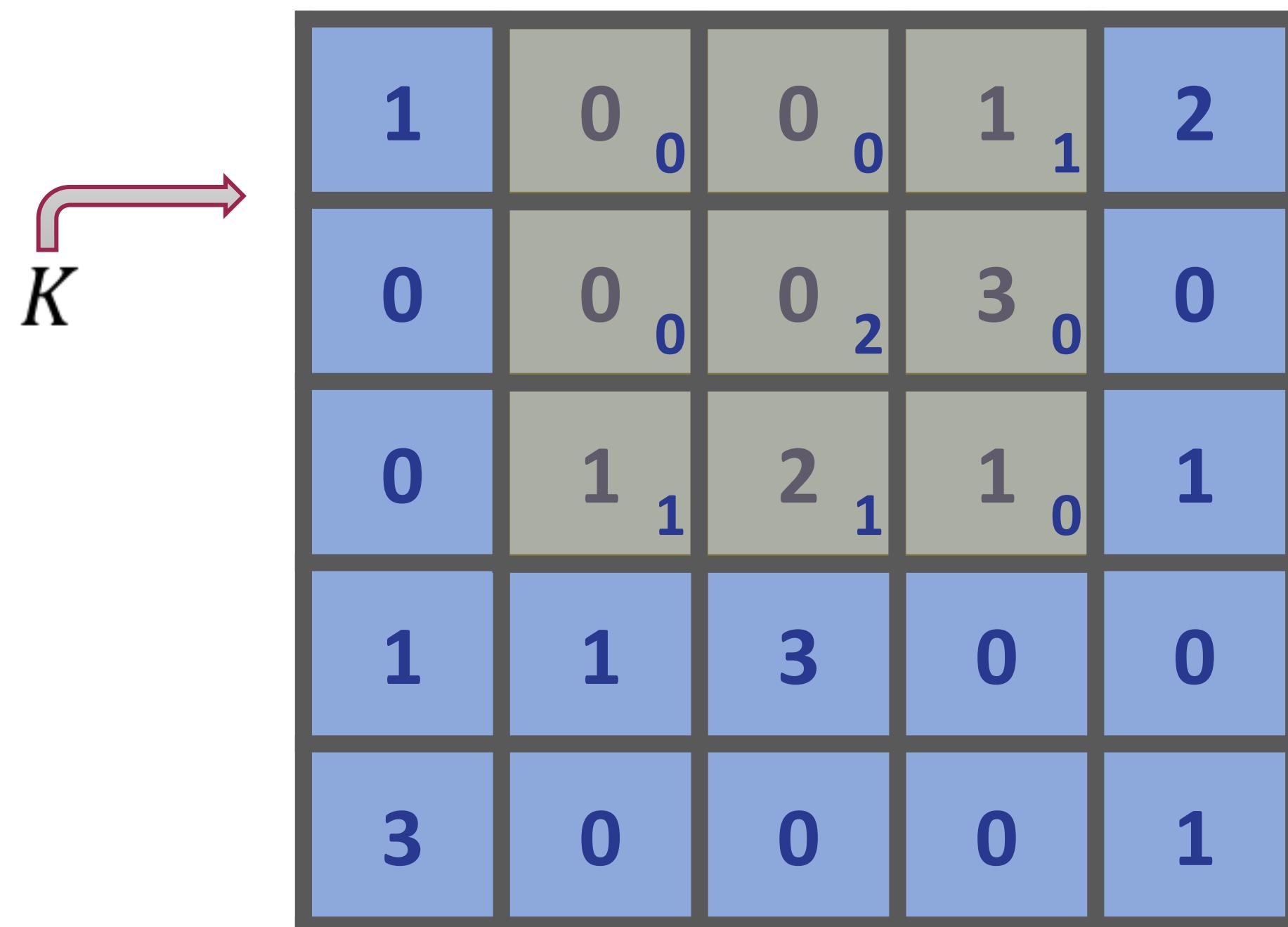
$$(1 \times 0) + (0 \times 0) + (0 \times 1) + \\ (0 \times 0) + (0 \times 2) + (0 \times 0) + \\ (0 \times 1) + (1 \times 1) + (2 \times 0)$$



$$(I \star K)(i, j) == \sum_m \sum_n I(m, n)K(i + m, j + n)$$

2. Convolution layer

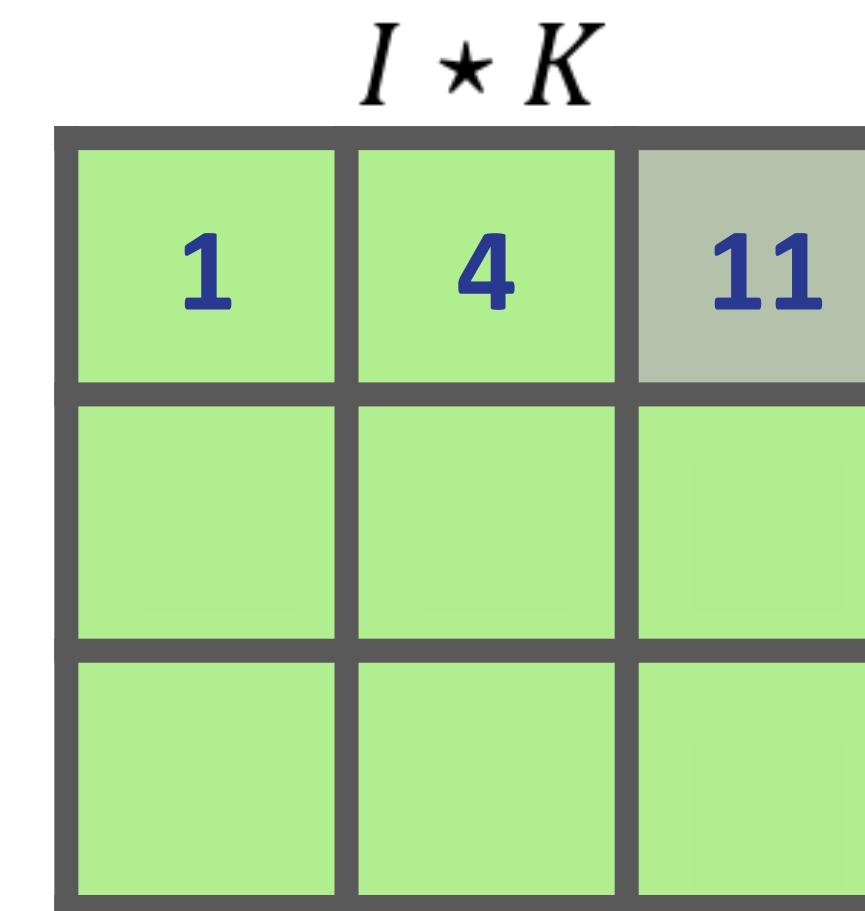
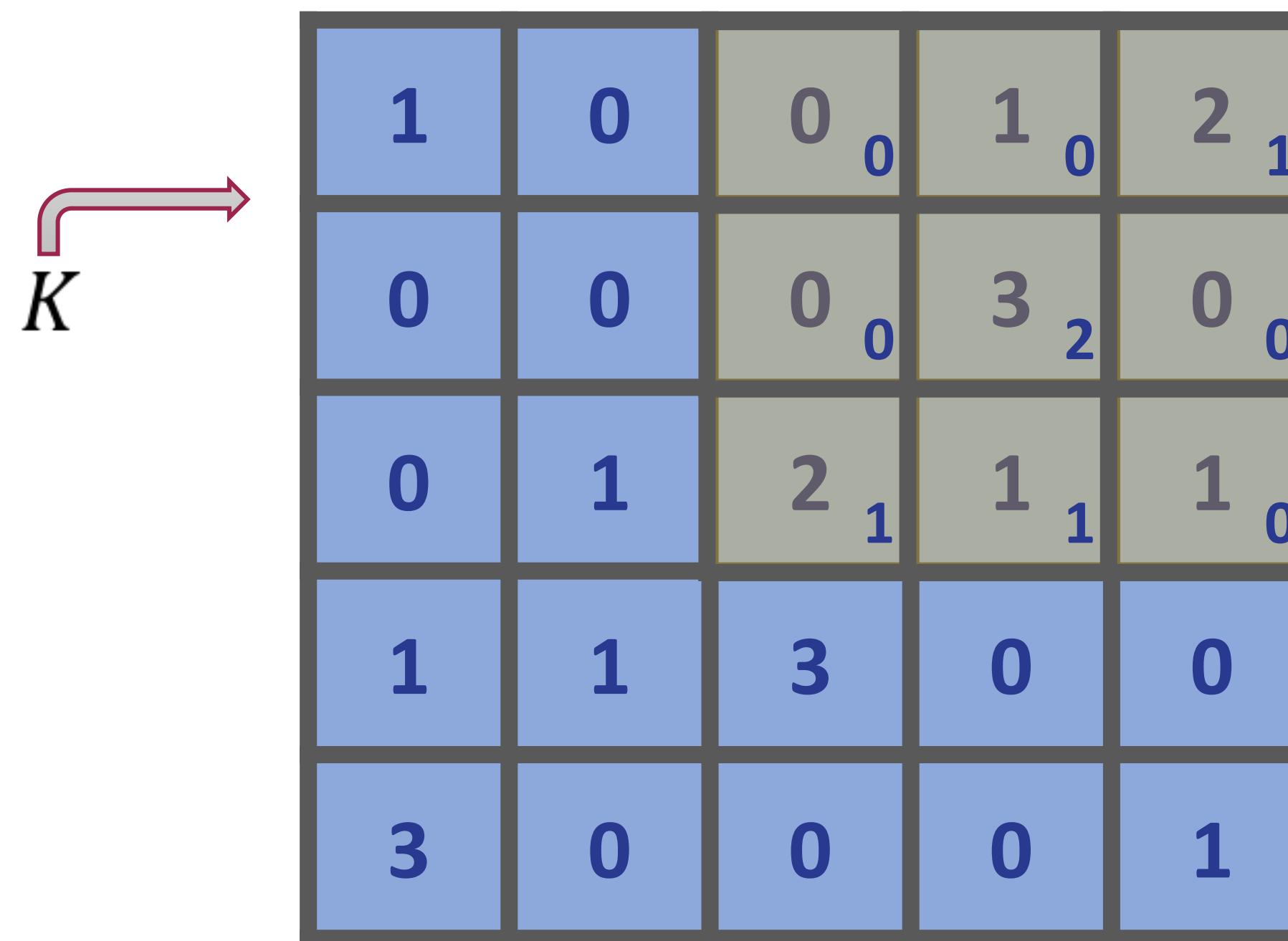
2D Convolutions



$$(I \star K)(i, j) == \sum_m \sum_n I(m, n)K(i + m, j + n)$$

2. Convolution layer

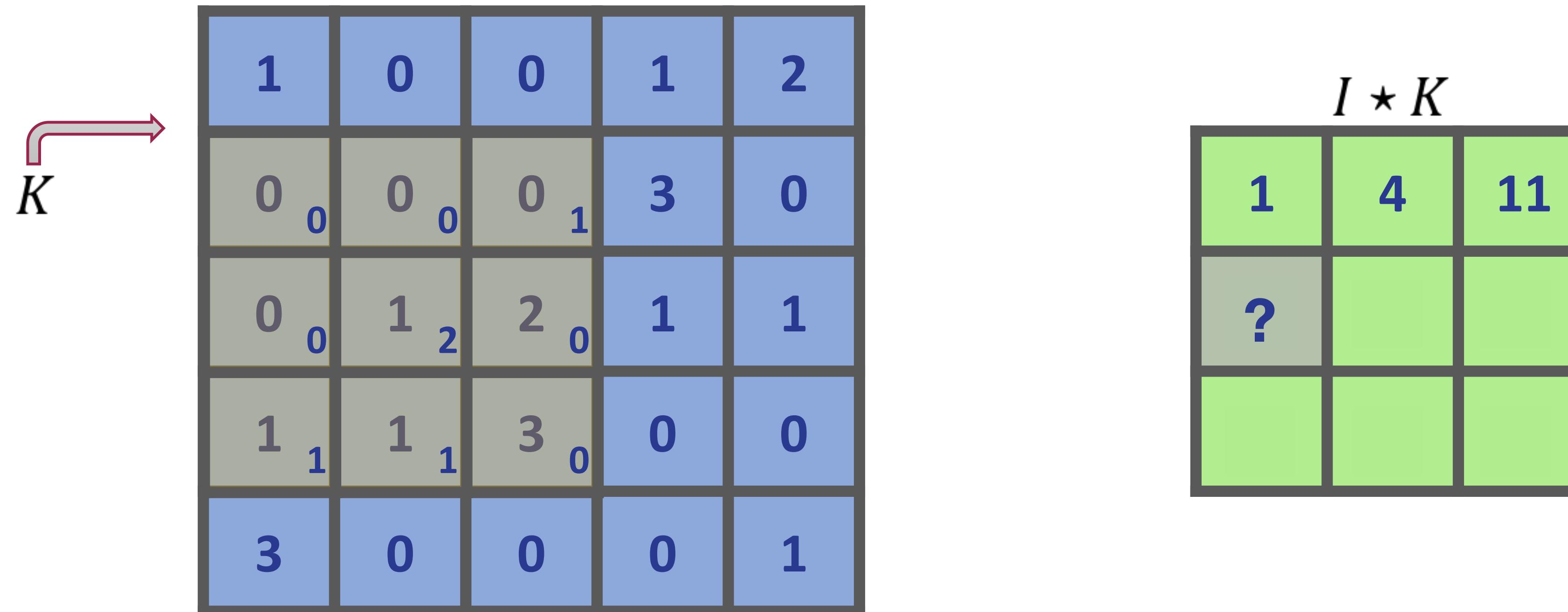
2D Convolutions



$$(I \star K)(i, j) == \sum_m \sum_n I(m, n)K(i + m, j + n)$$

2. Convolution layer

2D Convolutions

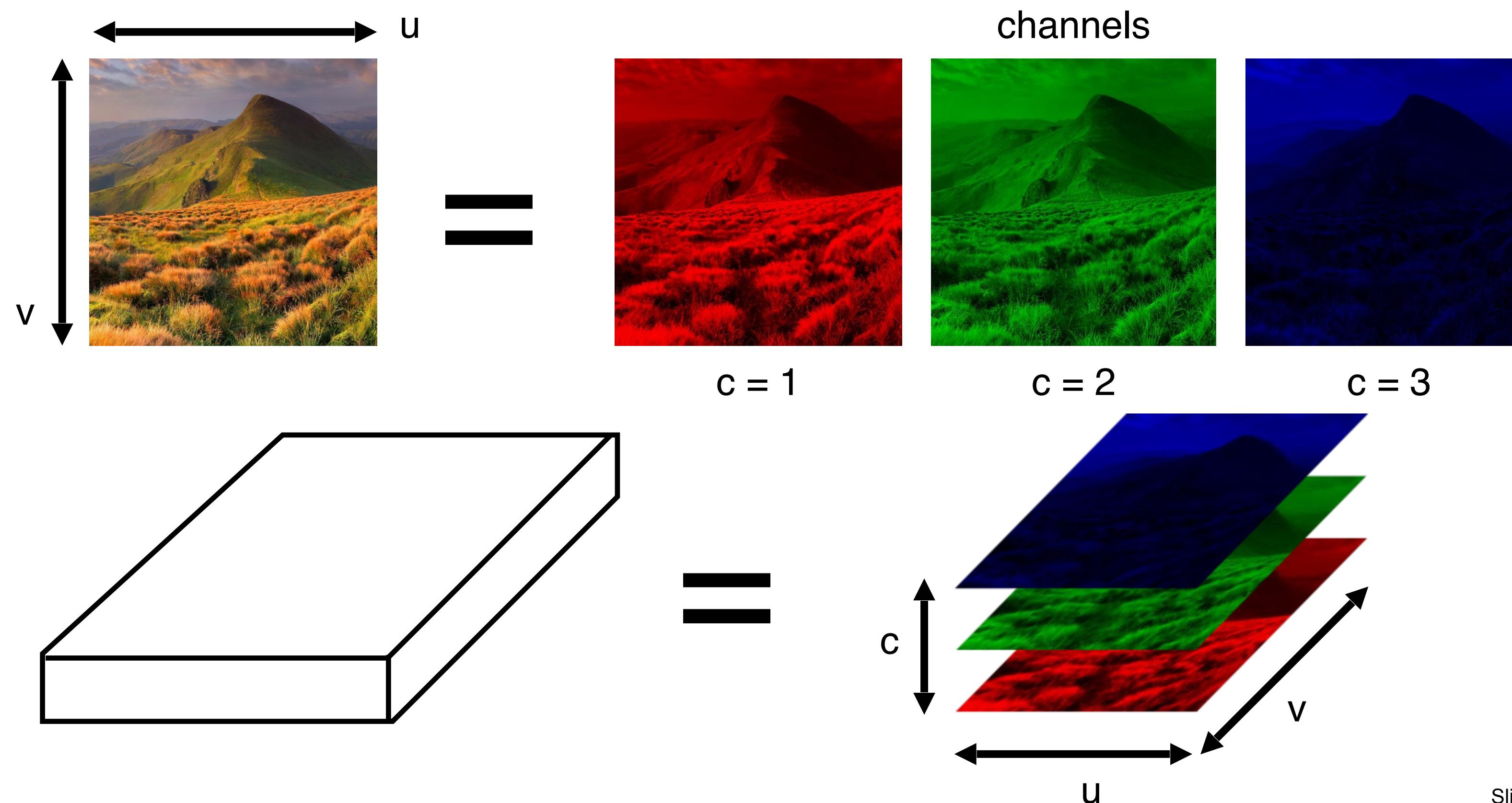


$$(I \star K)(i, j) == \sum_m \sum_n I(m, n)K(i + m, j + n)$$

2. Convolution layer

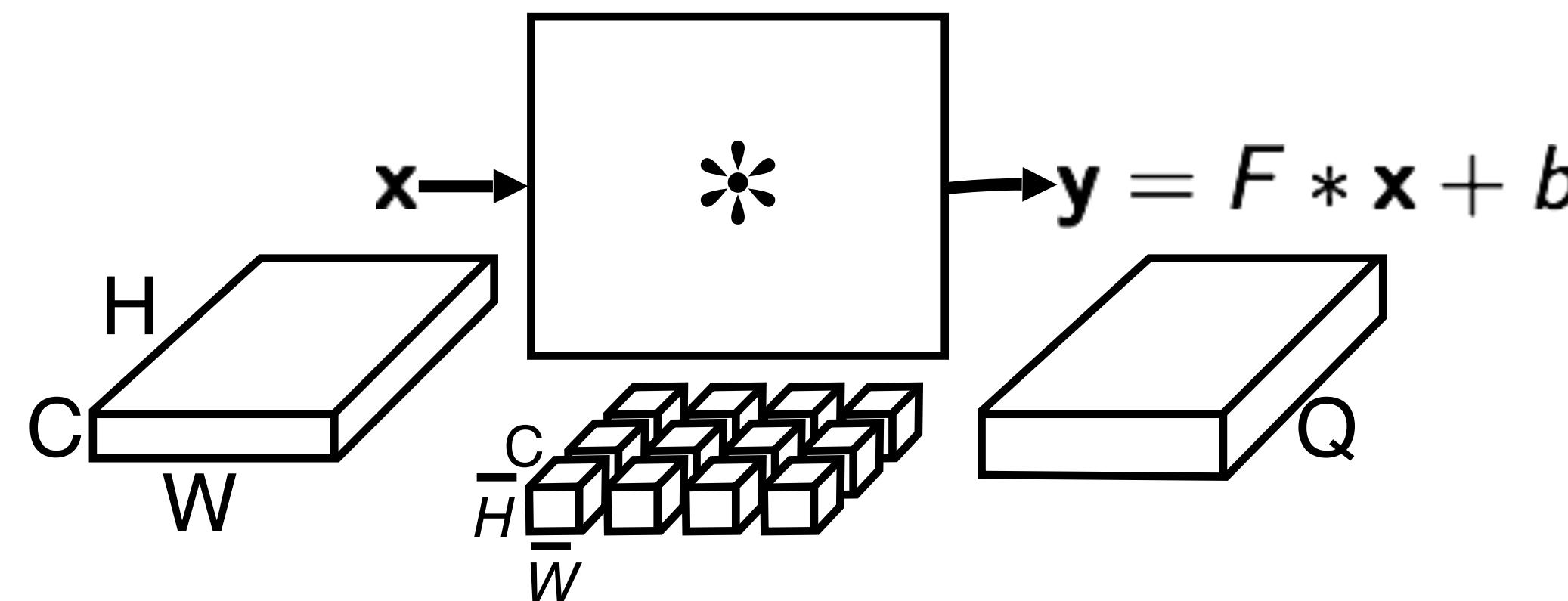
The data manipulated by a CNN has the form of 3D tensors. These are interpreted as **discrete vector fields \mathbf{x}** , assigning a feature vector $(x_{uv1}, \dots, x_{uvC})$ at each spatial location (v,u) .

A colour image is a simple example of a vector field with 3D features (RGB):



2. Convolution layer

With a bank of 3D filters

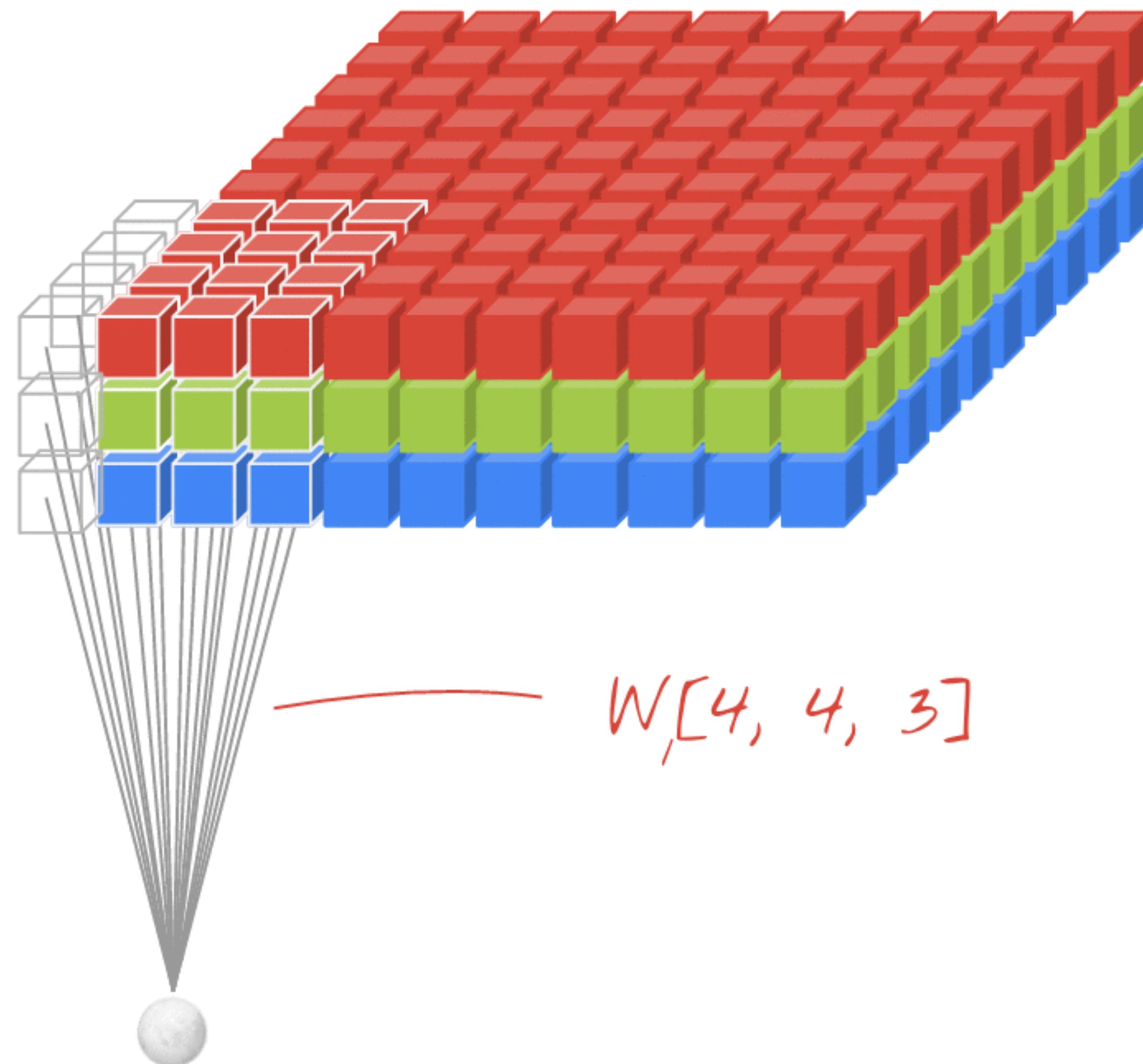


$$y_{v' u' q'} = b_{q'} + \sum_{\bar{v}=1}^{\bar{H}} \sum_{\bar{u}=1}^{\bar{W}} \sum_{c=1}^C x_{v'+\bar{v}-1, u'+\bar{u}-1, c} f_{\bar{v}\bar{u}c q'}$$

Linear convolution applies a bank of linear filters \mathcal{F} to the input tensor \mathbf{x} .

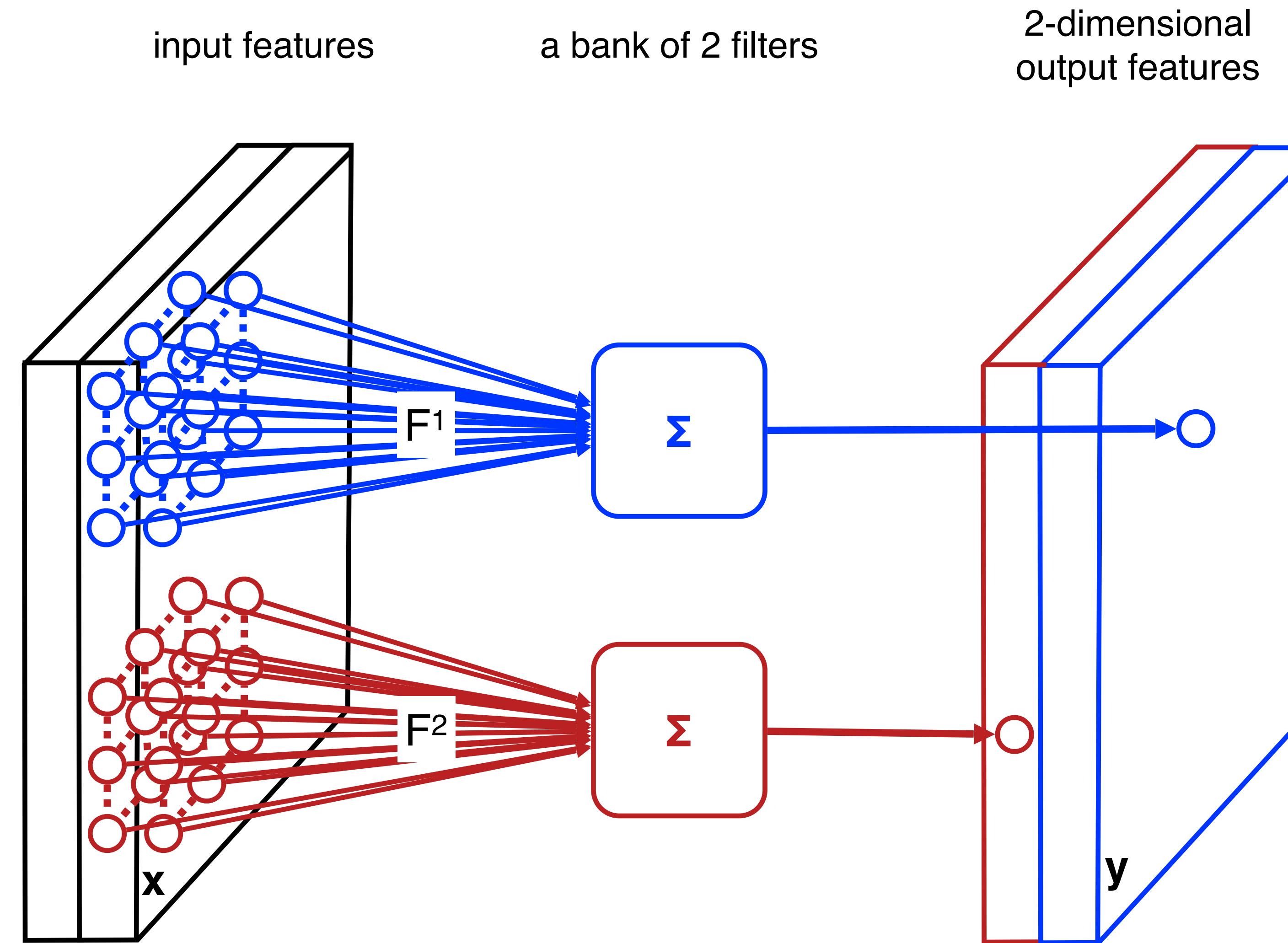
- Input tensor $\mathbf{x} = H \times W \times C$ array
- Filter bank $\mathcal{F} = \bar{H} \times \bar{W} \times C \times Q$ array
- Output tensor $\mathbf{y} = (H - \bar{H} + 1) \times (W - \bar{W} + 1) \times Q$ array

2. Convolution layer



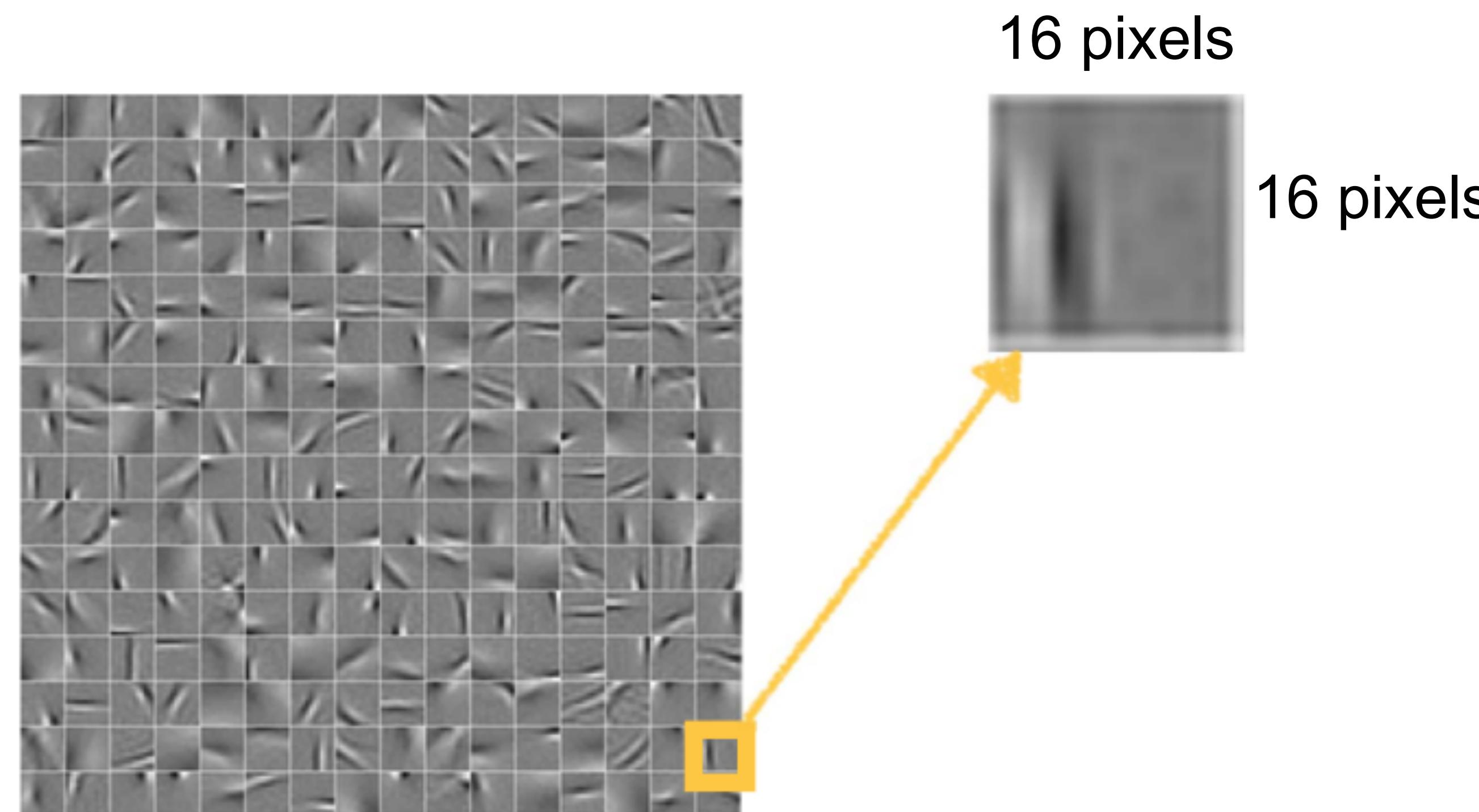
2. Convolution layer

As a neural network



Filter bank example

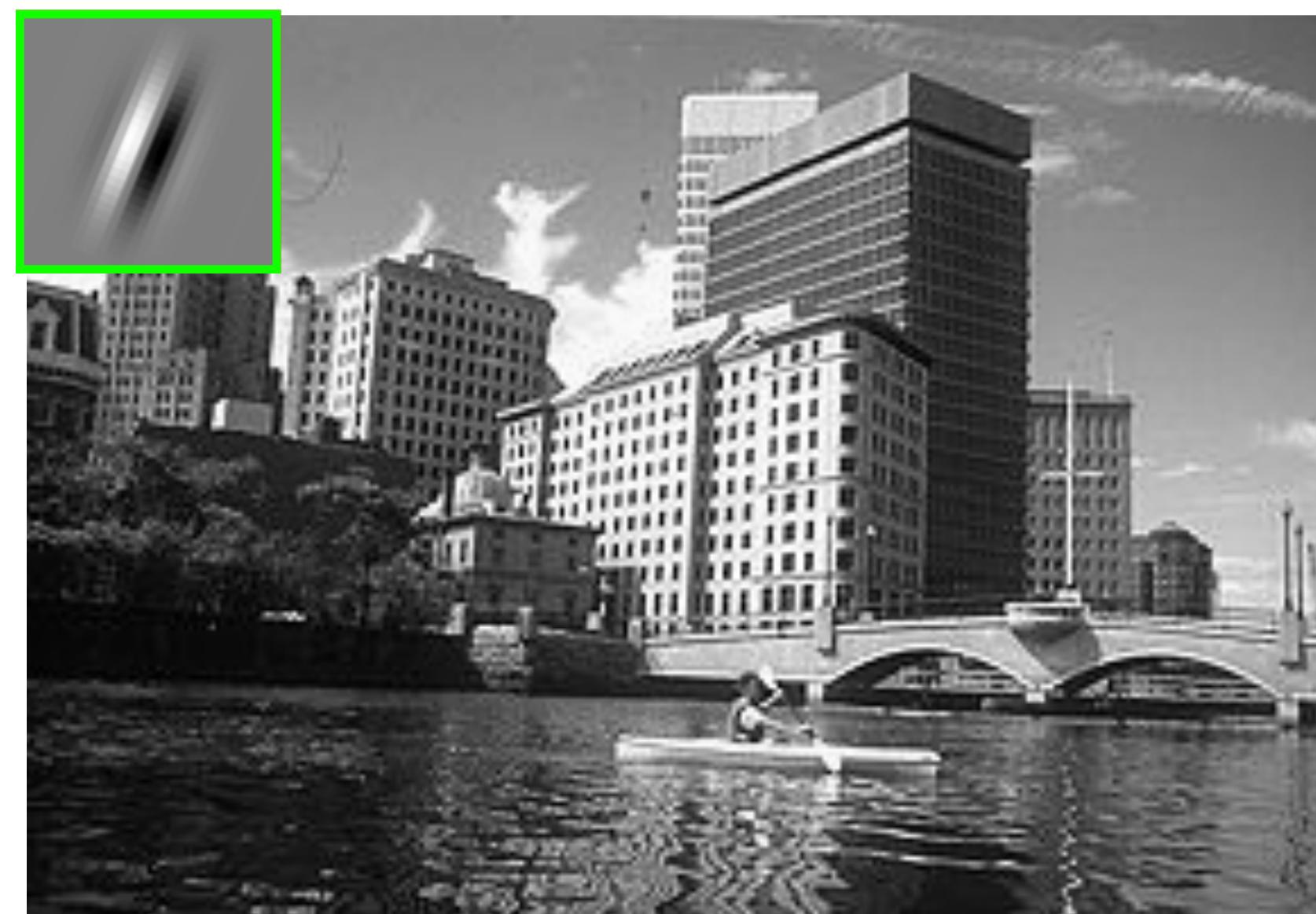
- A bank of 256 filters (learned from data)
- Each filter has 1 channel (it applies to a grayscale image)
- Each filter is 16x16 pixels



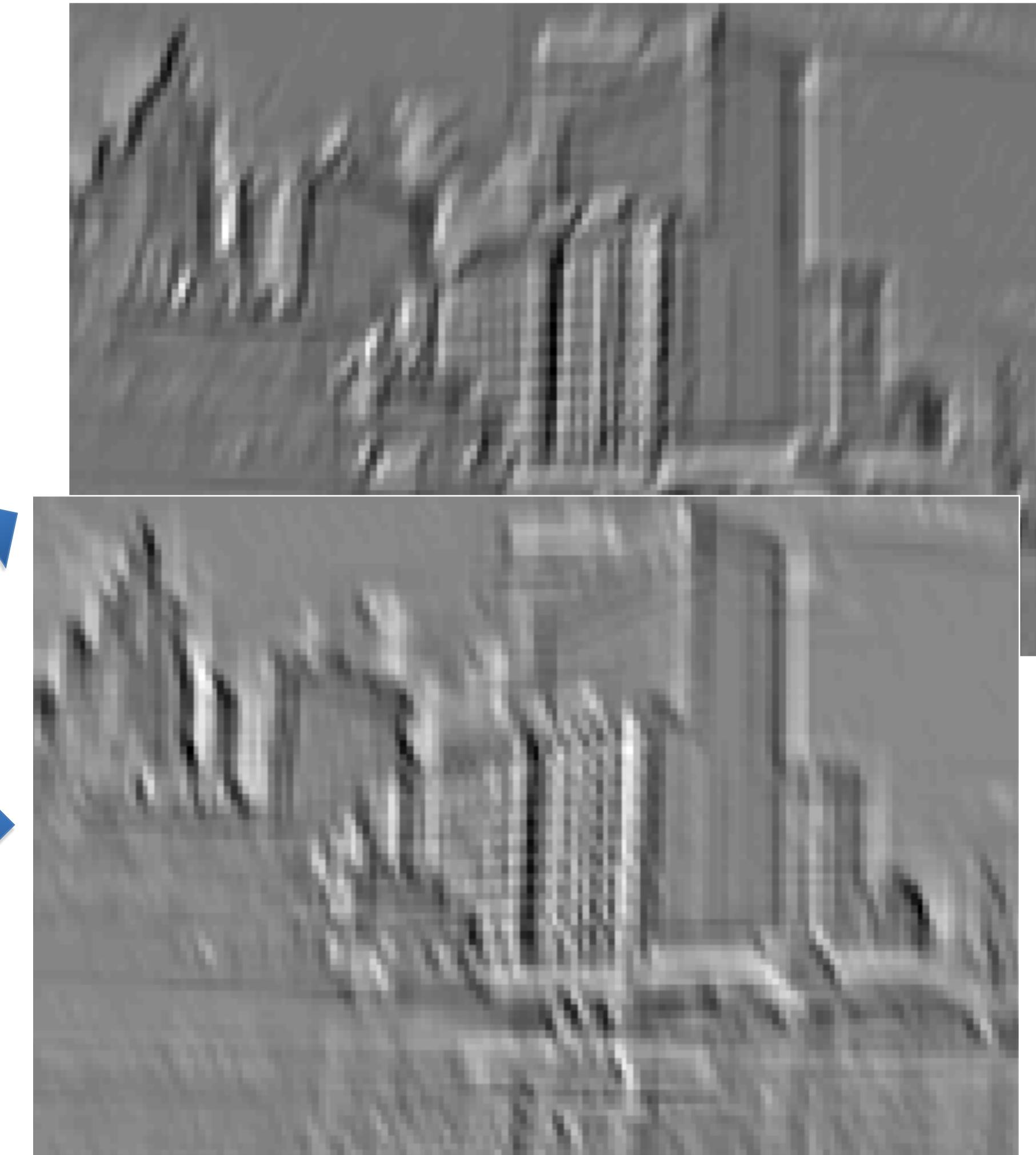
Filtering

Each filter generates a “feature map”

Maximum response when
filter matches signal



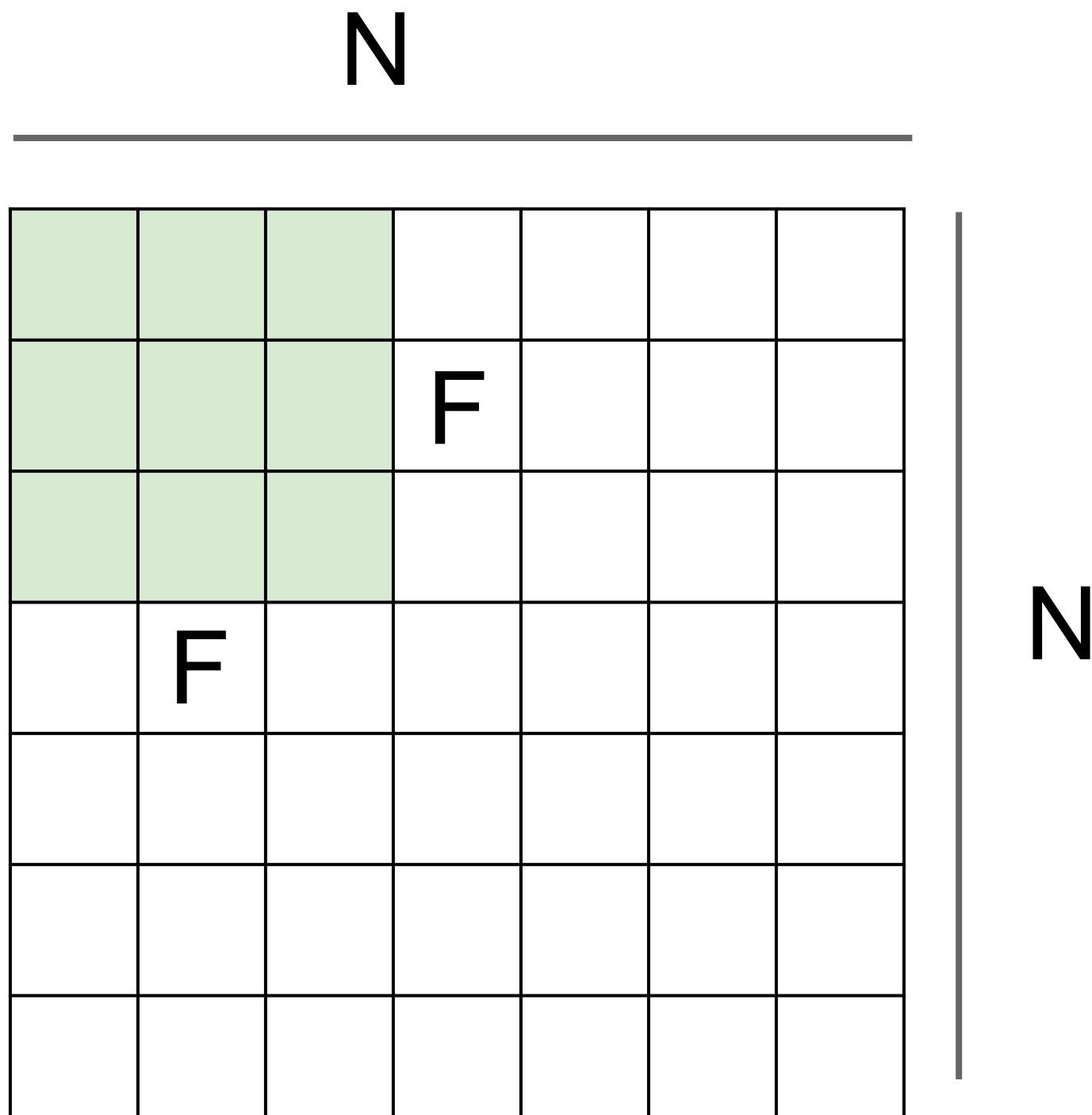
Input



Feature Map

Convolution details

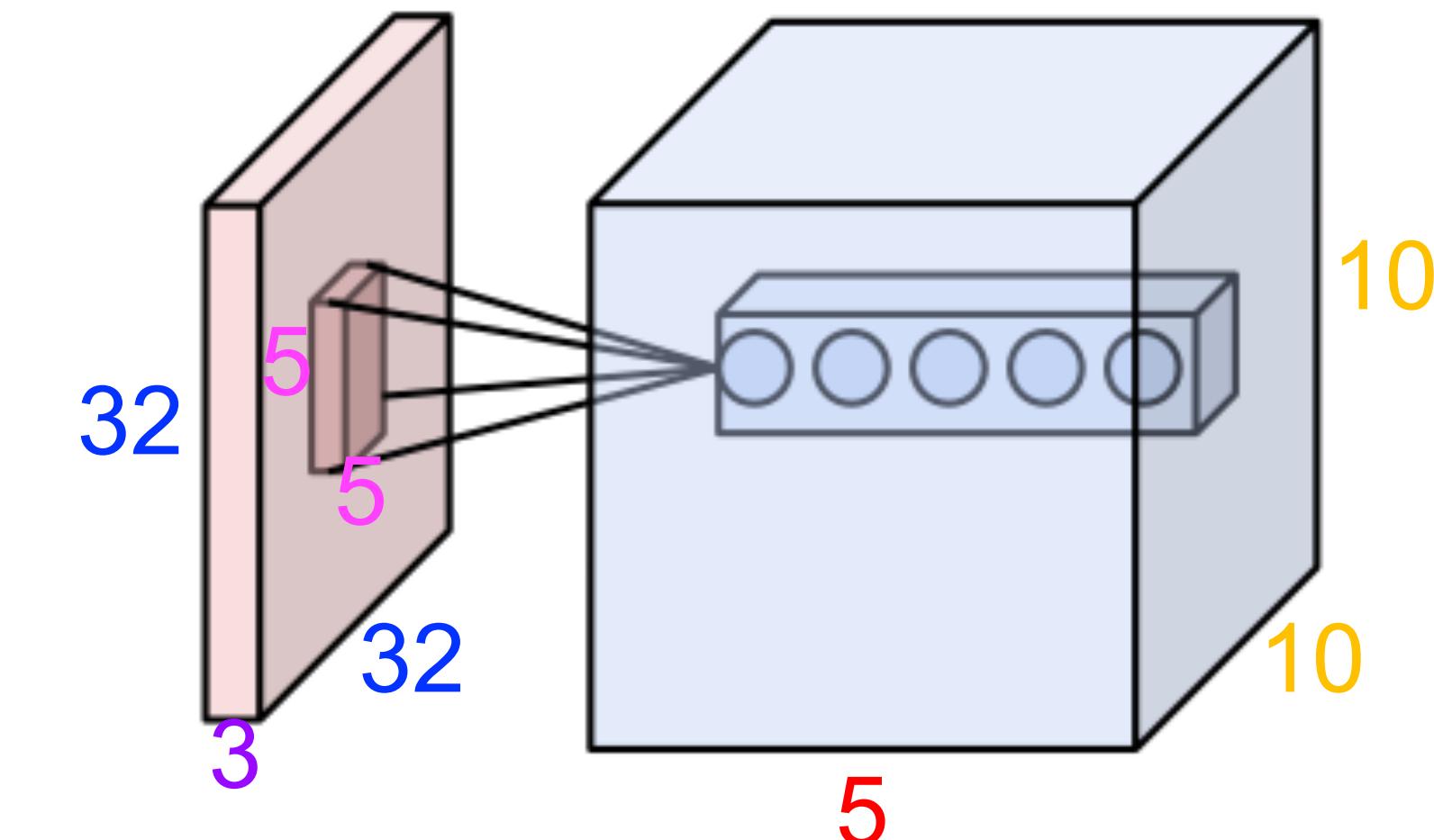
What is the output size?



Output size:
 $(N - F) / \text{stride} + 1$

e.g. $N = 7$, $F = 3$:
stride 1 $\Rightarrow (7 - 3)/1 + 1 = 5$
stride 2 $\Rightarrow (7 - 3)/2 + 1 = 3$
stride 3 $\Rightarrow (7 - 3)/3 + 1 = \dots :\backslash$

Example: What is the output volume?



Input volume: $32 \times 32 \times 3$

Receptive fields: 5×5 , stride 3

Number of neurons: 5

Output volume: $(32 - 5) / 3 + 1 = 10$, so: $10 \times 10 \times 5$

Zero padding (in each channel)

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7
neuron with receptive field 3x3, stride 1
pad with 1 pixel border => what is the output?

7x7 => preserved size!

in general, common to see stride 1, size F, and
zero-padding with $(F-1)/2$.
(Will preserve input size spatially)

What is the number of parameters?

- Consider an input gray-scale image of 1000x1000 pixels.
- What is the number of parameters of a filter bank of 100 7x7 filters?
- How does it compare to a fully connected layer that considers the entire input image?

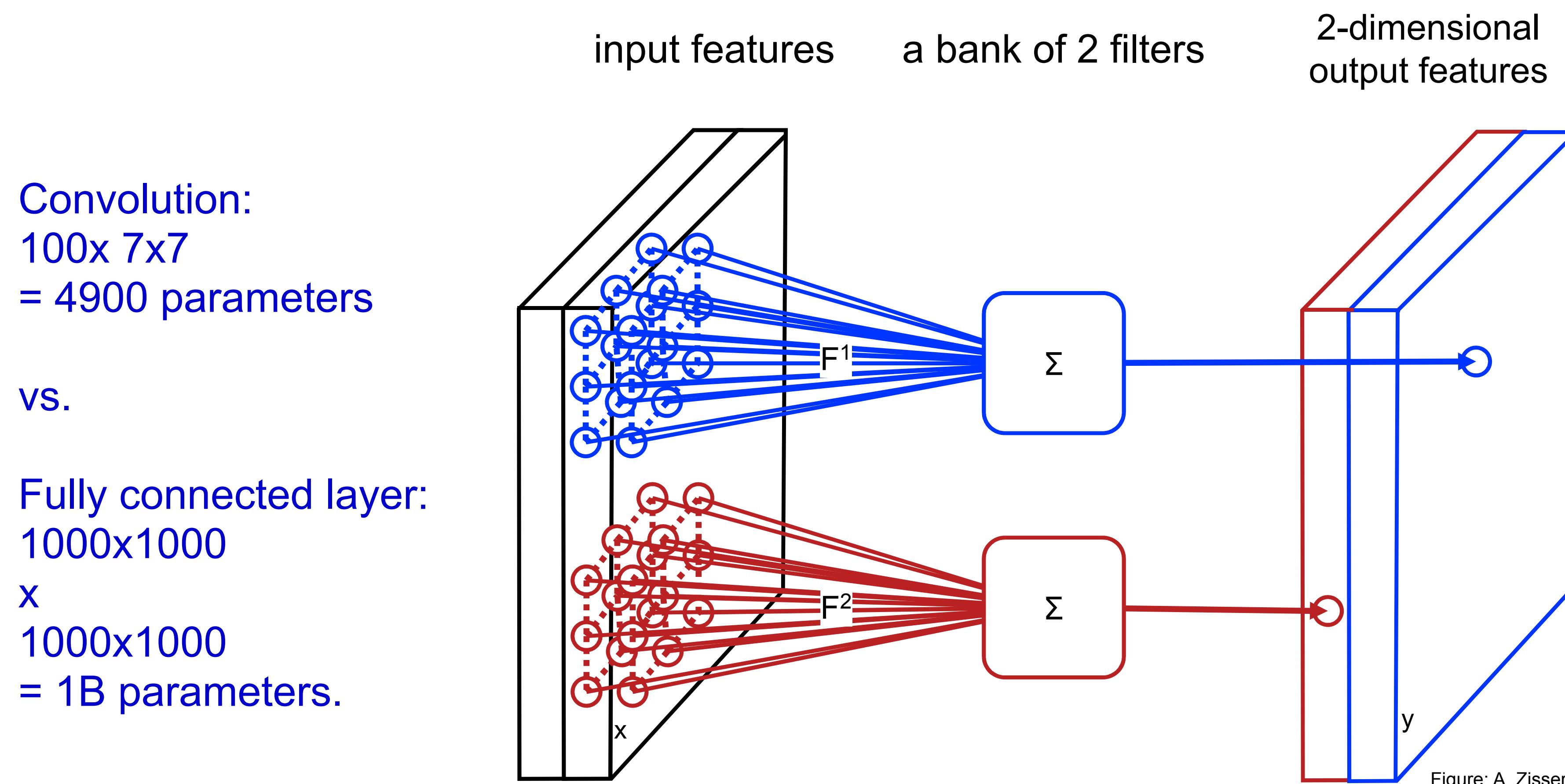
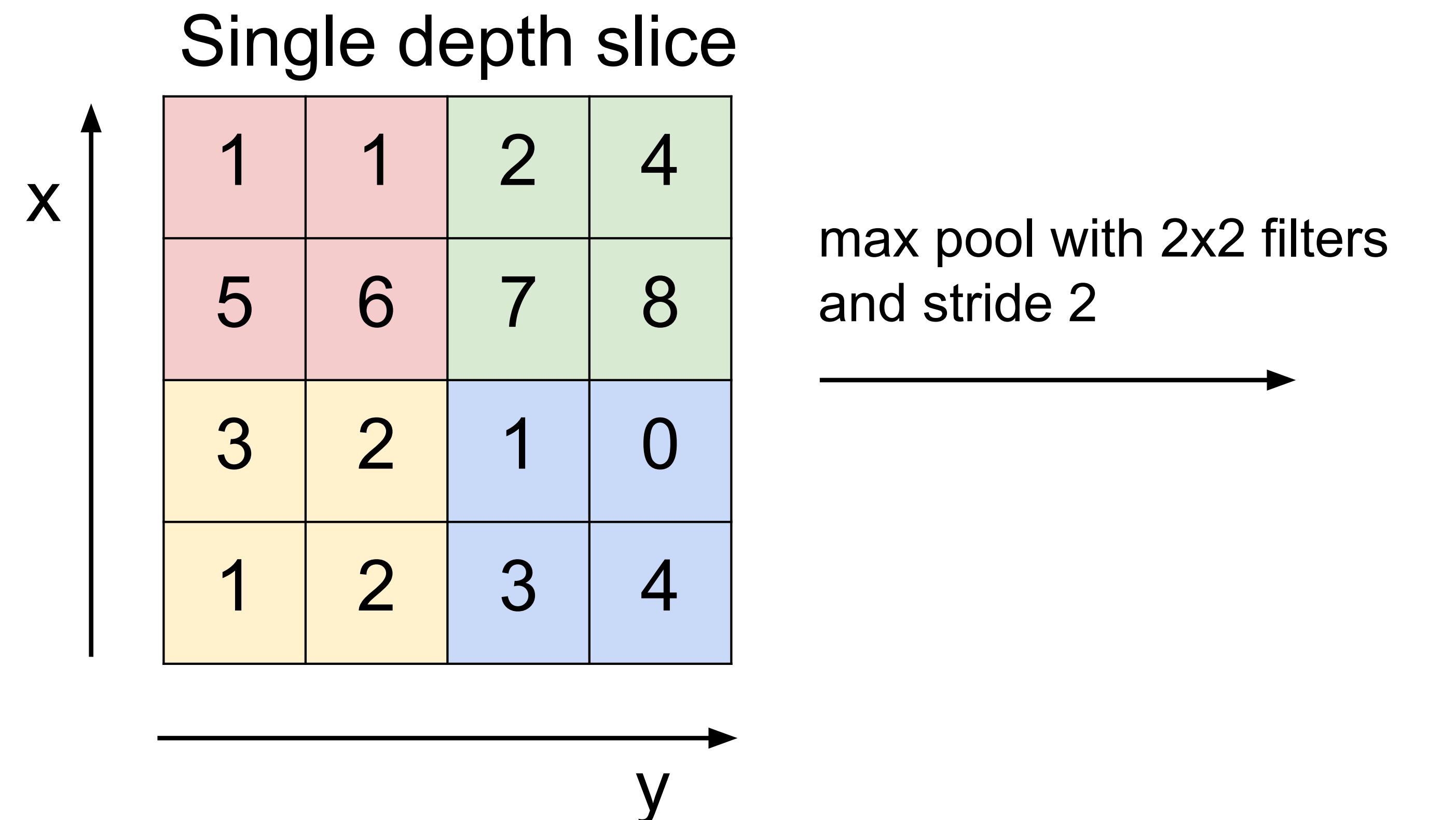


Figure: A. Zisserman

3. Spatial Max Pooling



3. Spatial Max Pooling

Dimensions of pooling outputs

Input volume of size $[W_1 \times H_1 \times D_1]$

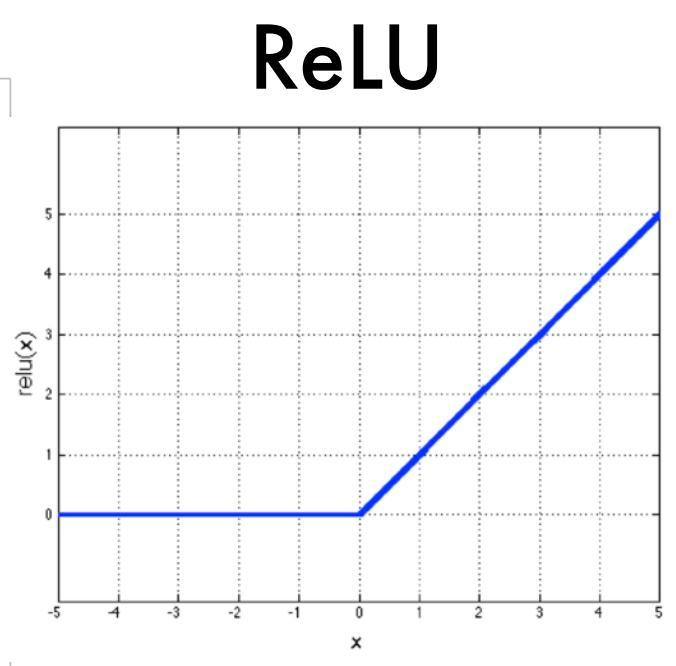
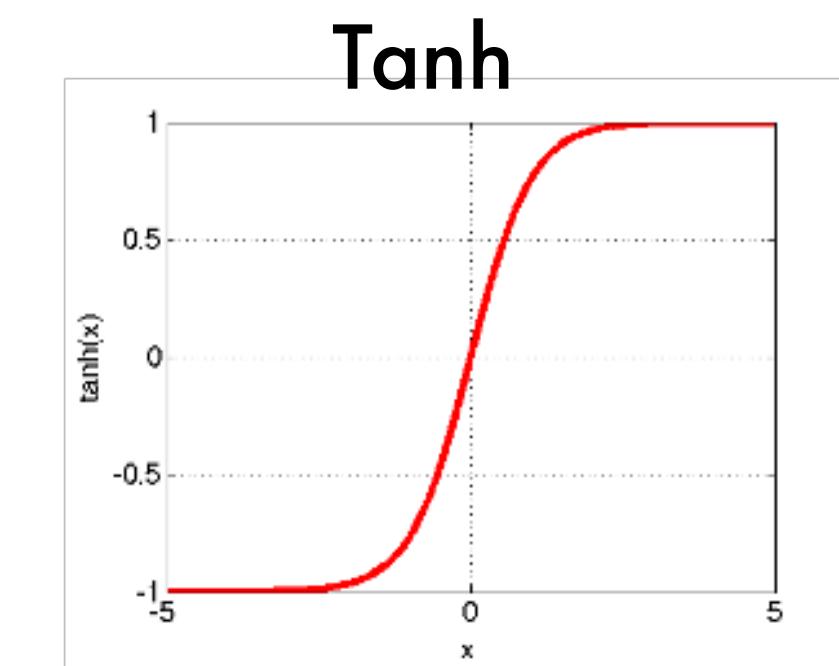
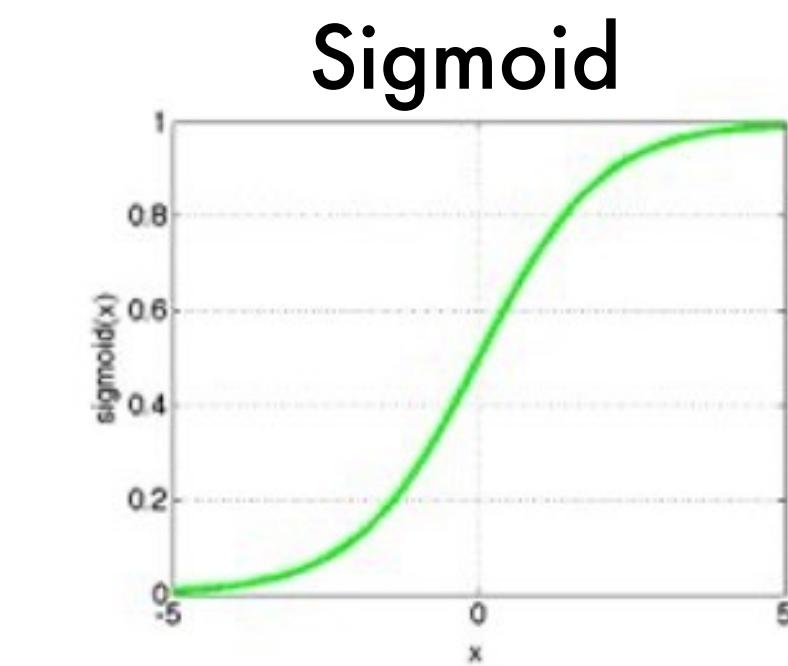
Pooling unit receptive fields $F \times F$ and applying them at strides of S gives

Output volume: $[W_2, H_2, D_1]$

$$W_2 = (W_1 - F)/S + 1, H_2 = (H_1 - F)/S + 1$$

Note: pooling happens independently in each channel/slice

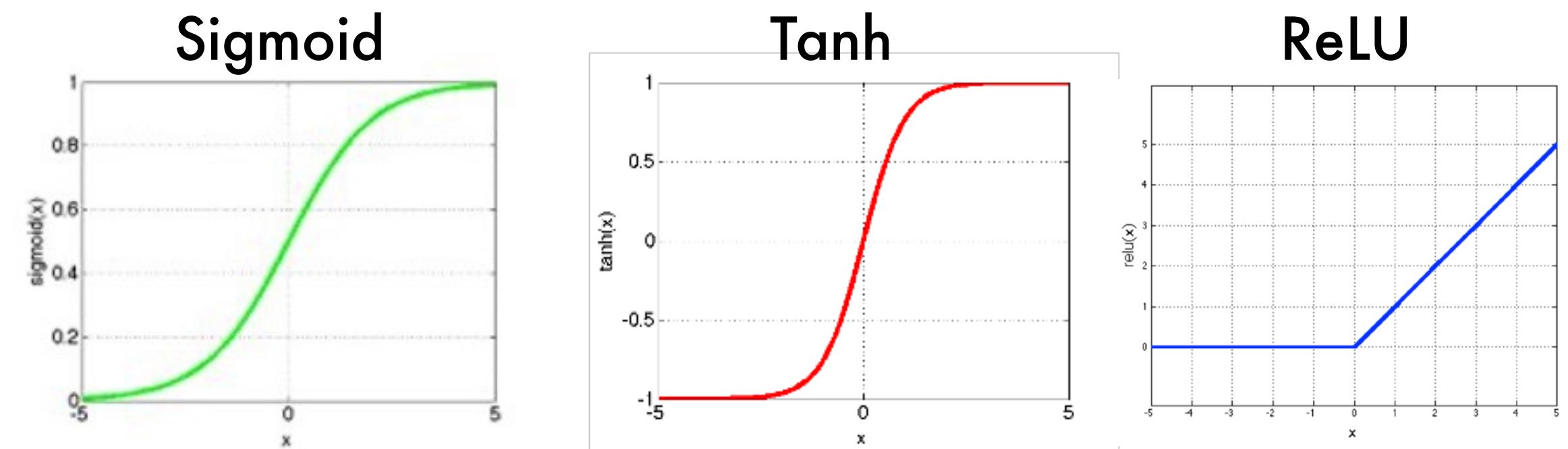
4. Non-linearity



- The non-linear activation functions are essential. Why?

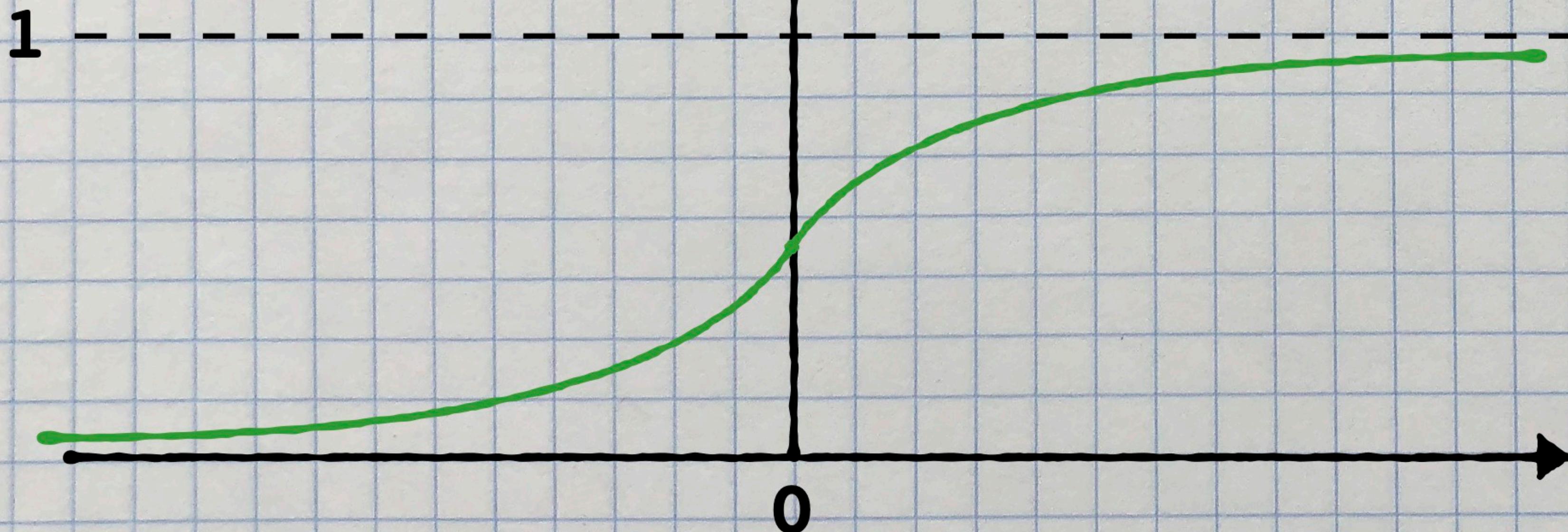
4. Non-linearity

Why?



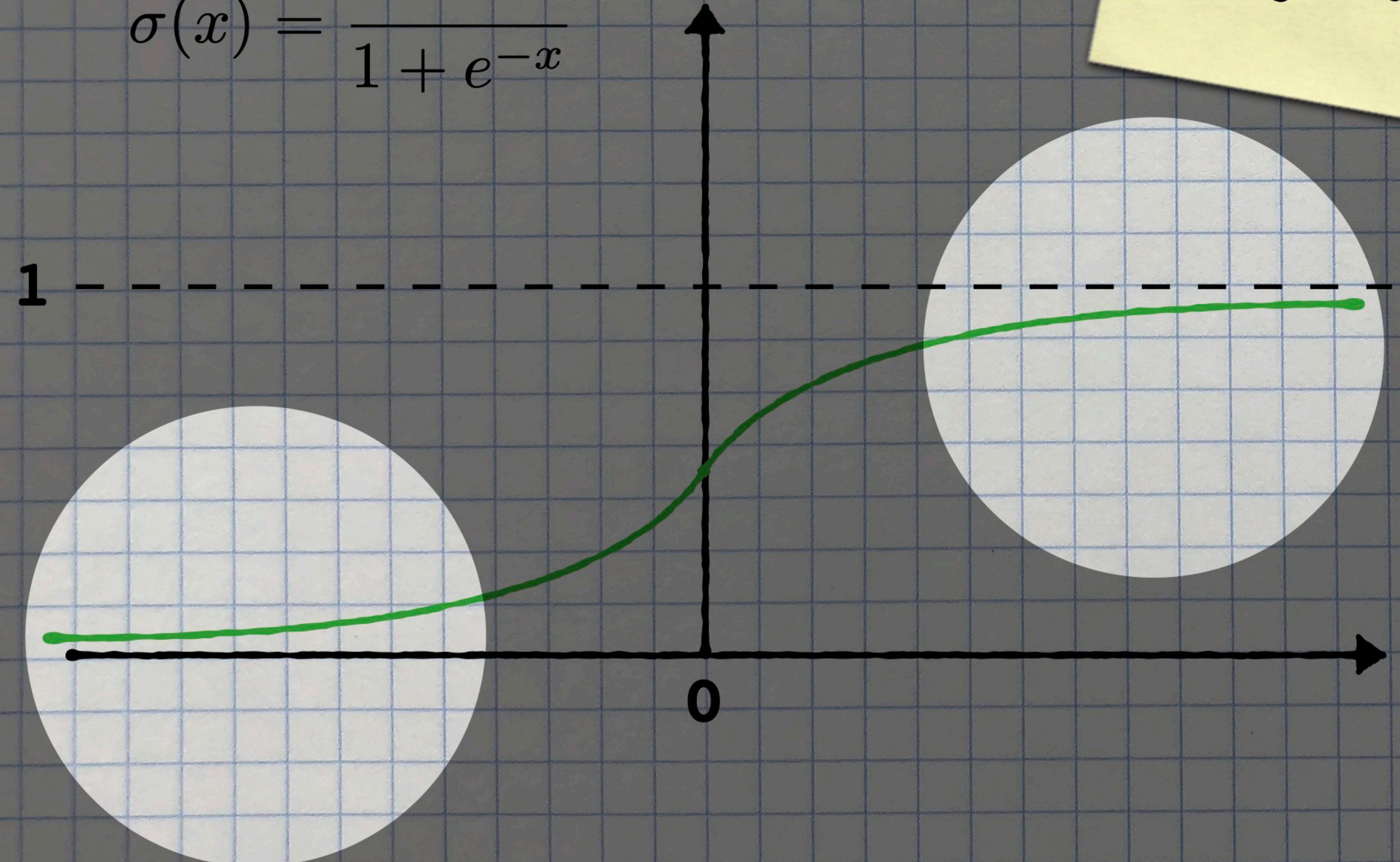
- Non-linearities allow us to approximate arbitrarily **complex functions**.
- **Universal approximation theorem:** A two-layer multilayer perceptron (MLP) with increasing continuous and bounded non-linearity can approximate any continuous function on a compact given enough hidden neurons. [Cybenko 1989]
- Linear activation functions produce **linear decisions no matter what the model size**, i.e., stacking multiple linear functions can be expressed with a single linear function.

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



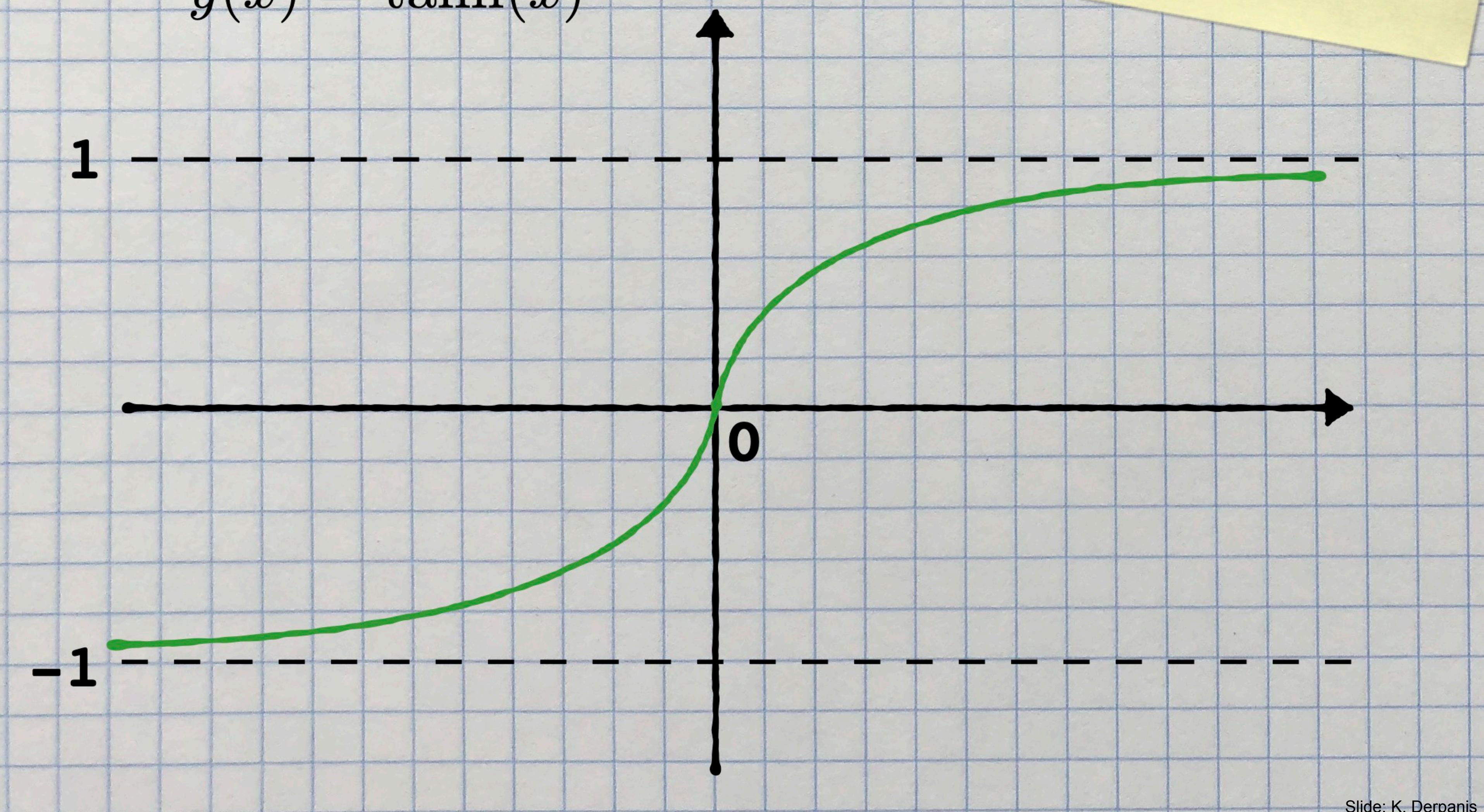
sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



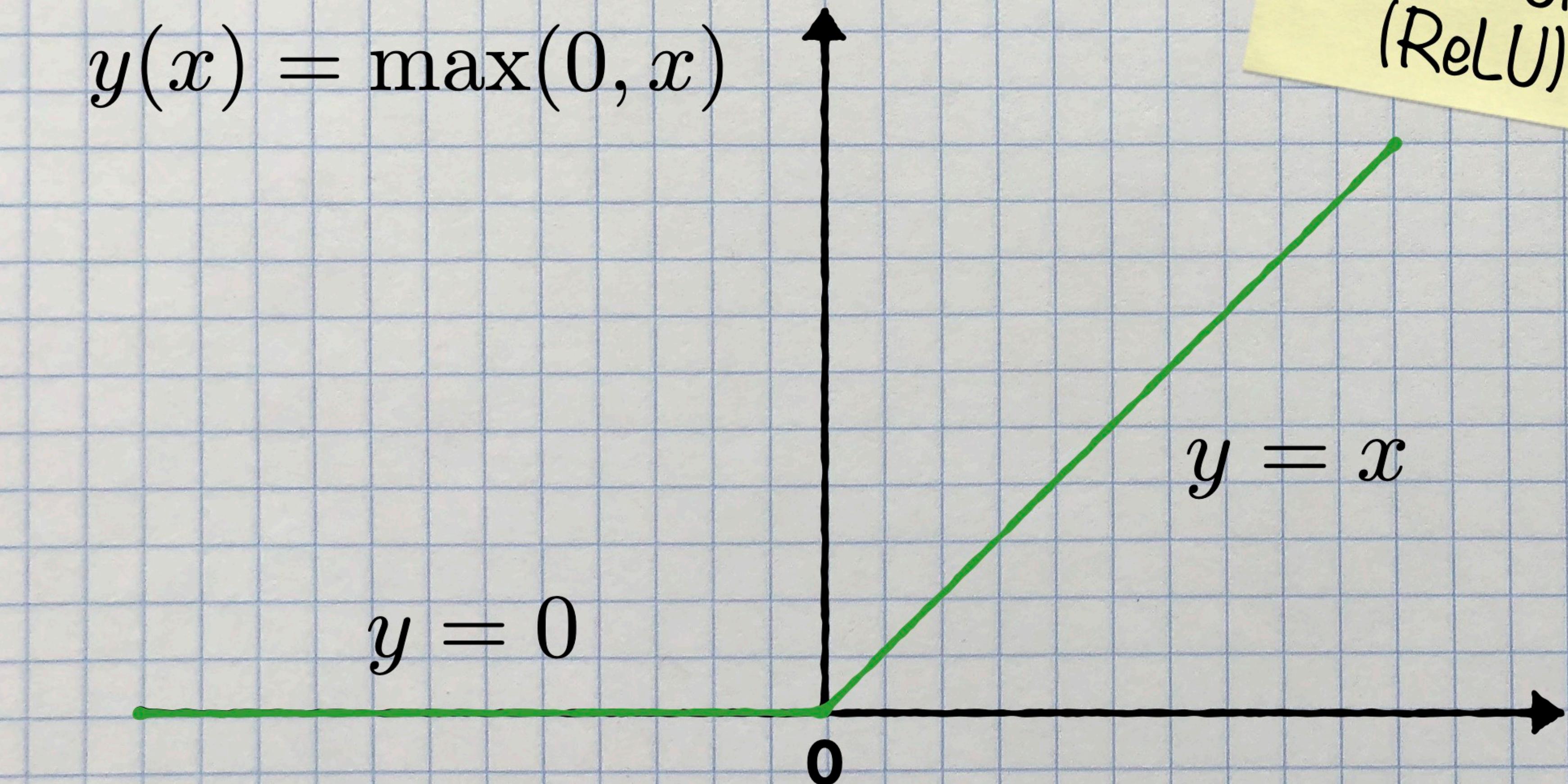
hyperbolic
tangent

$$y(x) = \tanh(x)$$



Rectified
Linear Unit
(ReLU)

$$y(x) = \max(0, x)$$



Rectified
Linear Unit
(ReLU)

$$y(x) = \max(0, x)$$

What is the derivative of the ReLU?

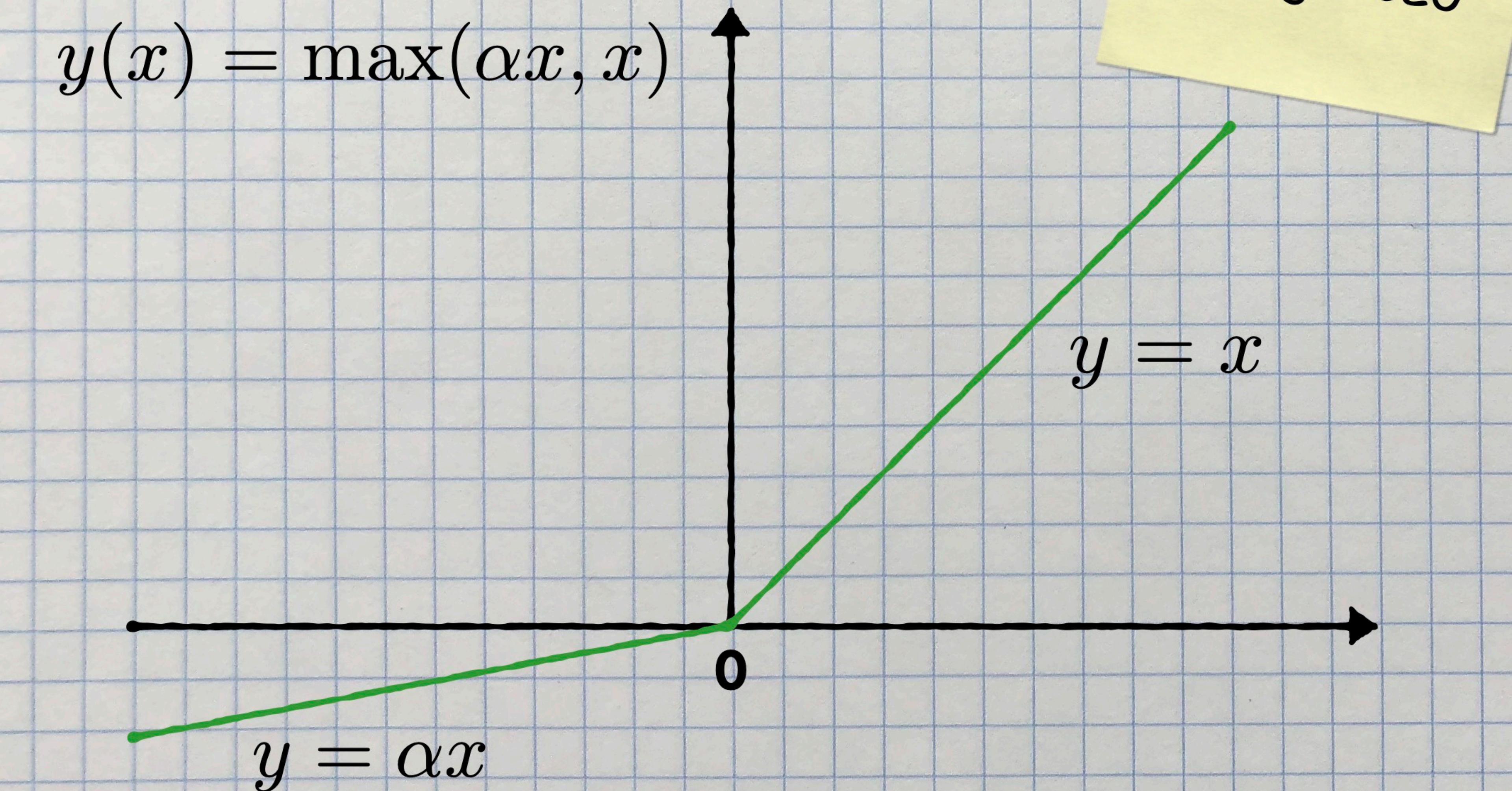
$$y = 0$$

0

$$y = x$$

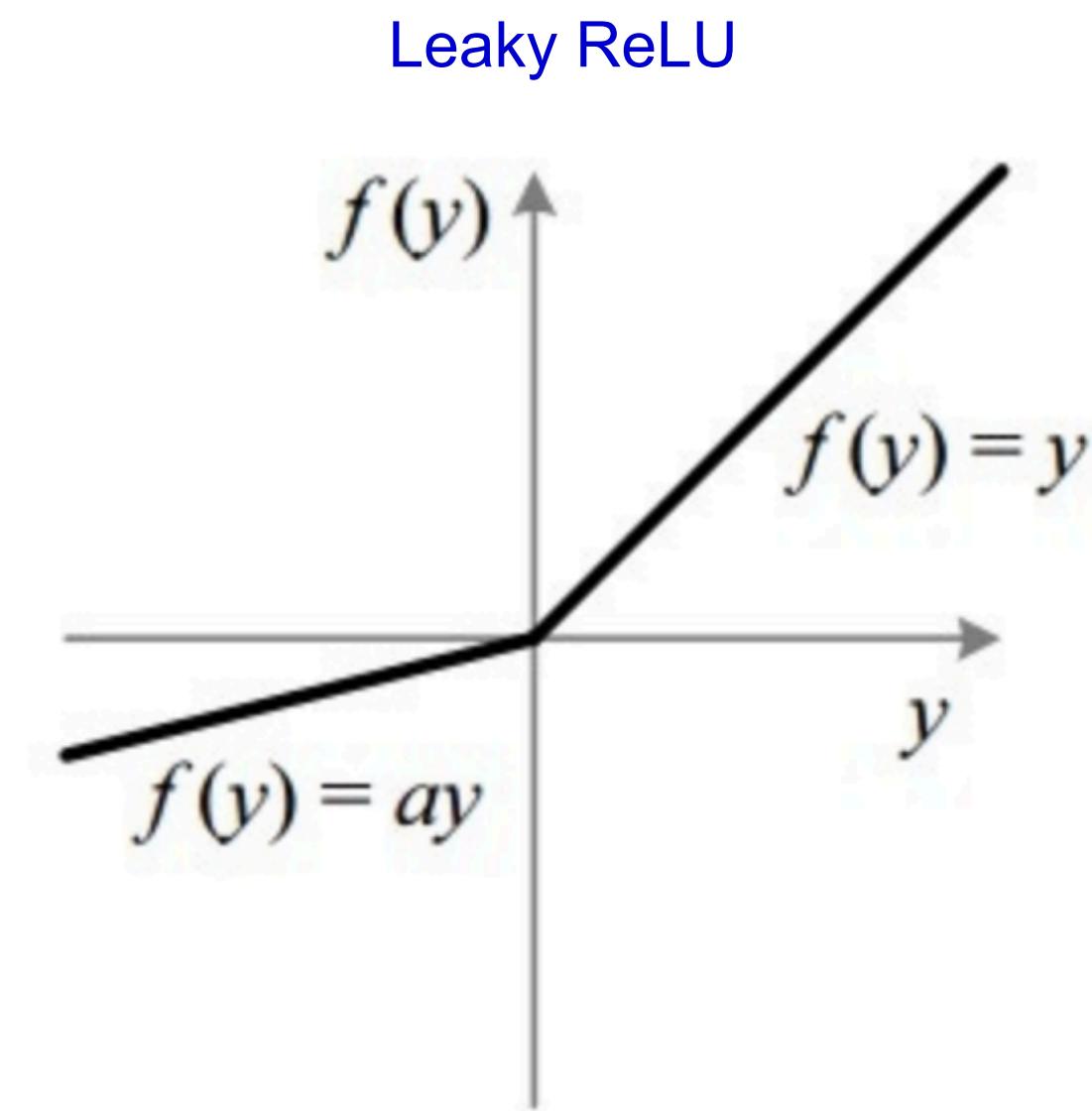
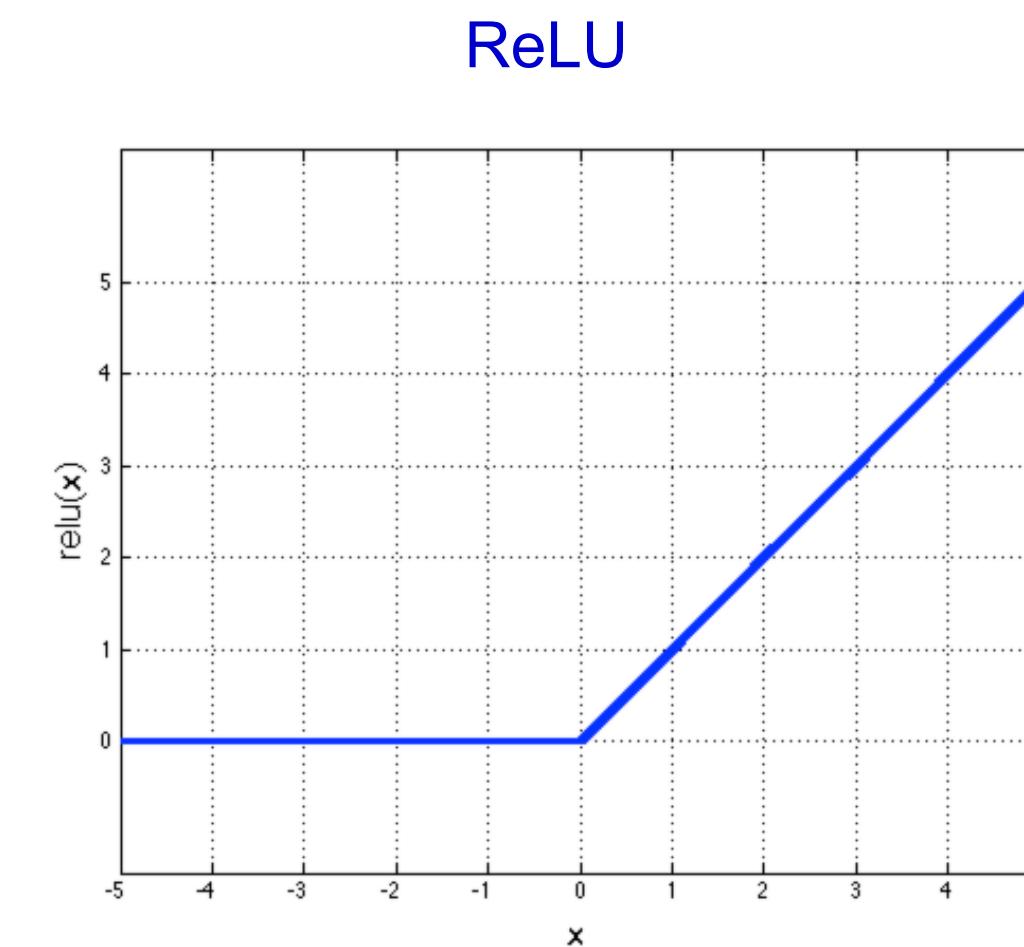
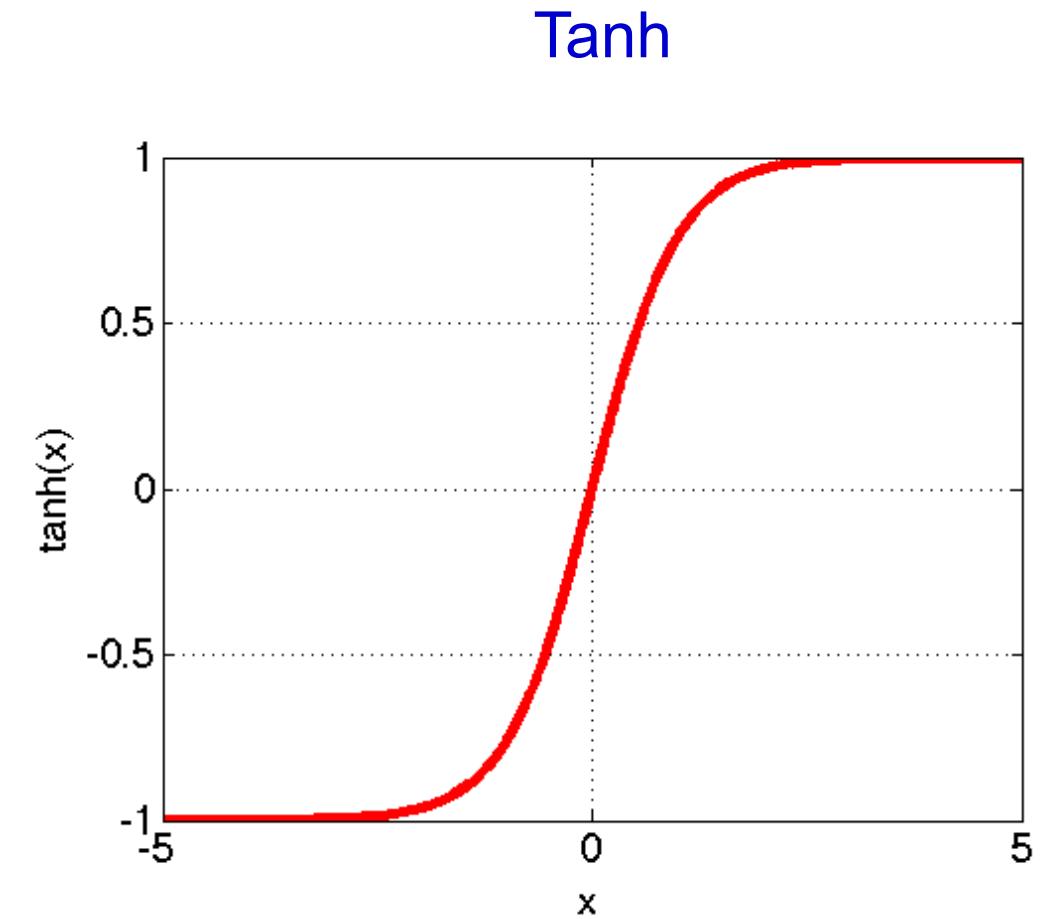
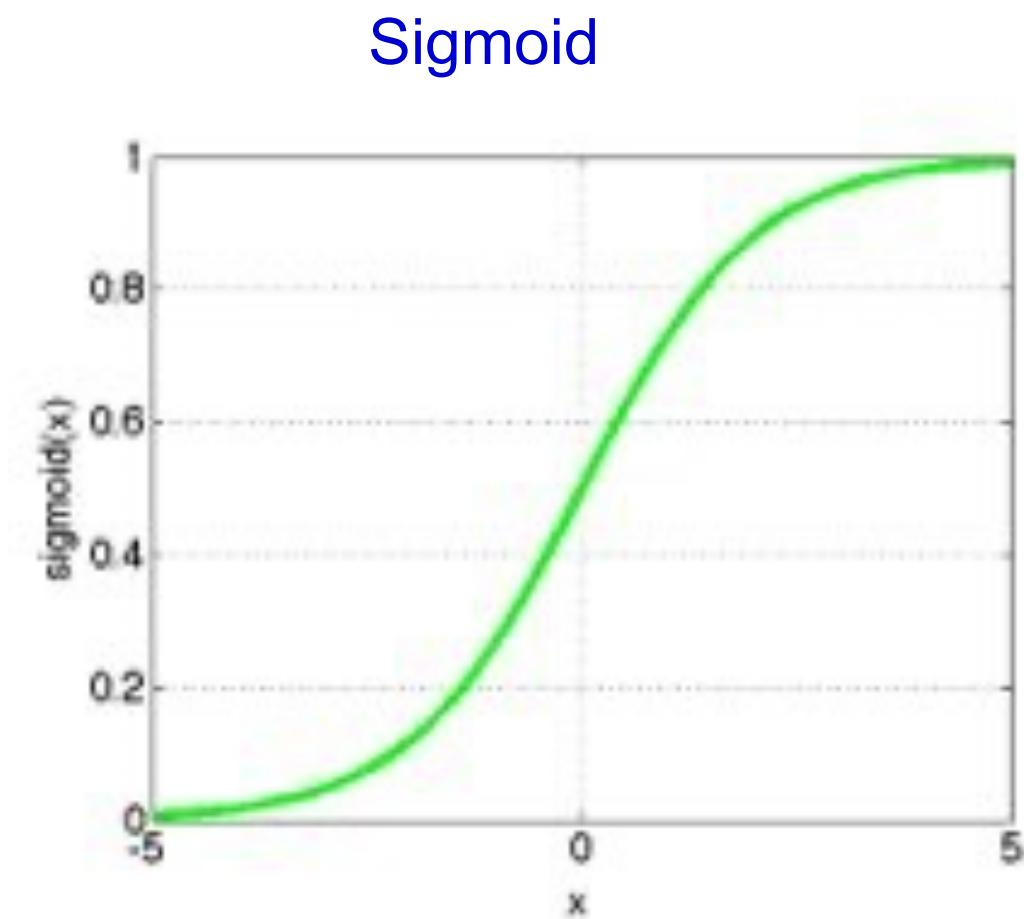
Leaky-ReLU

$$y(x) = \max(\alpha x, x)$$



4. Non-Linearity

- Per-element (independent)
- Options:
 - Sigmoid: $1/(1+\exp(-x))$
 - Tanh
 - Rectified linear unit (ReLU)
 - Simplifies backpropagation
 - Makes learning faster
 - Avoids saturation issues
 - Variants of ReLU, e.g. Leaky ReLU



5. Normalization

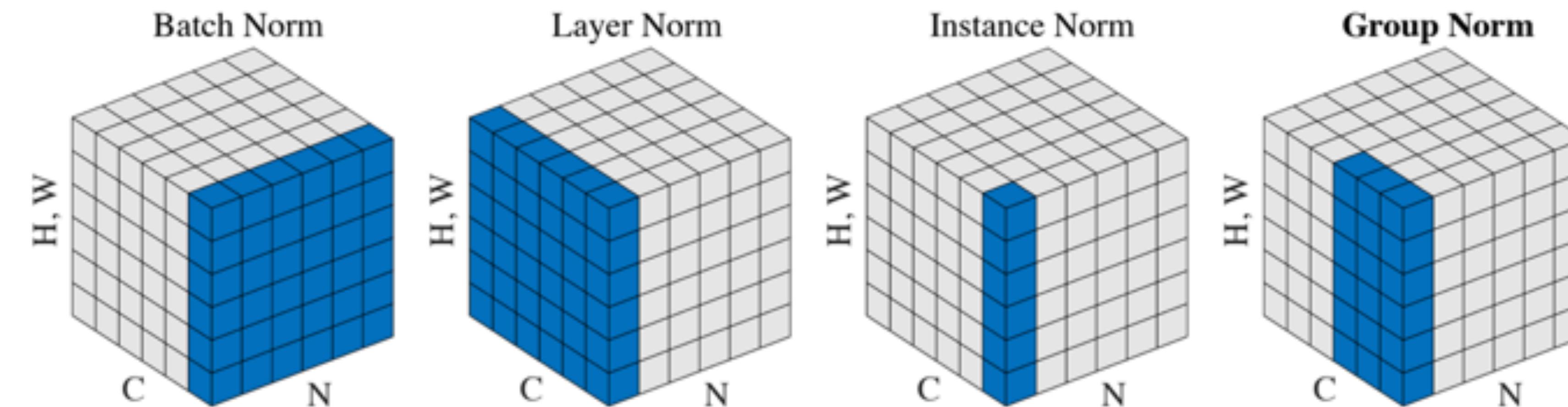
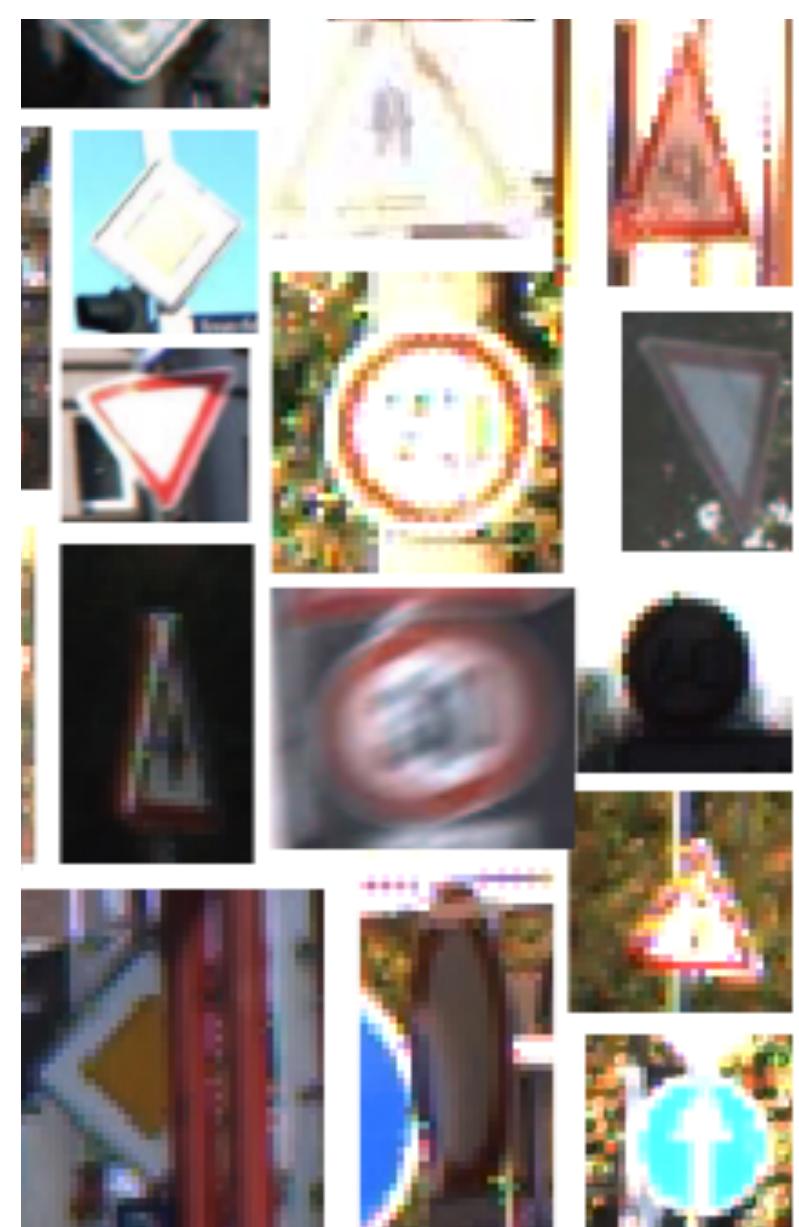


Figure 2. **Normalization methods.** Each subplot shows a feature map tensor, with N as the batch axis, C as the channel axis, and (H, W) as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

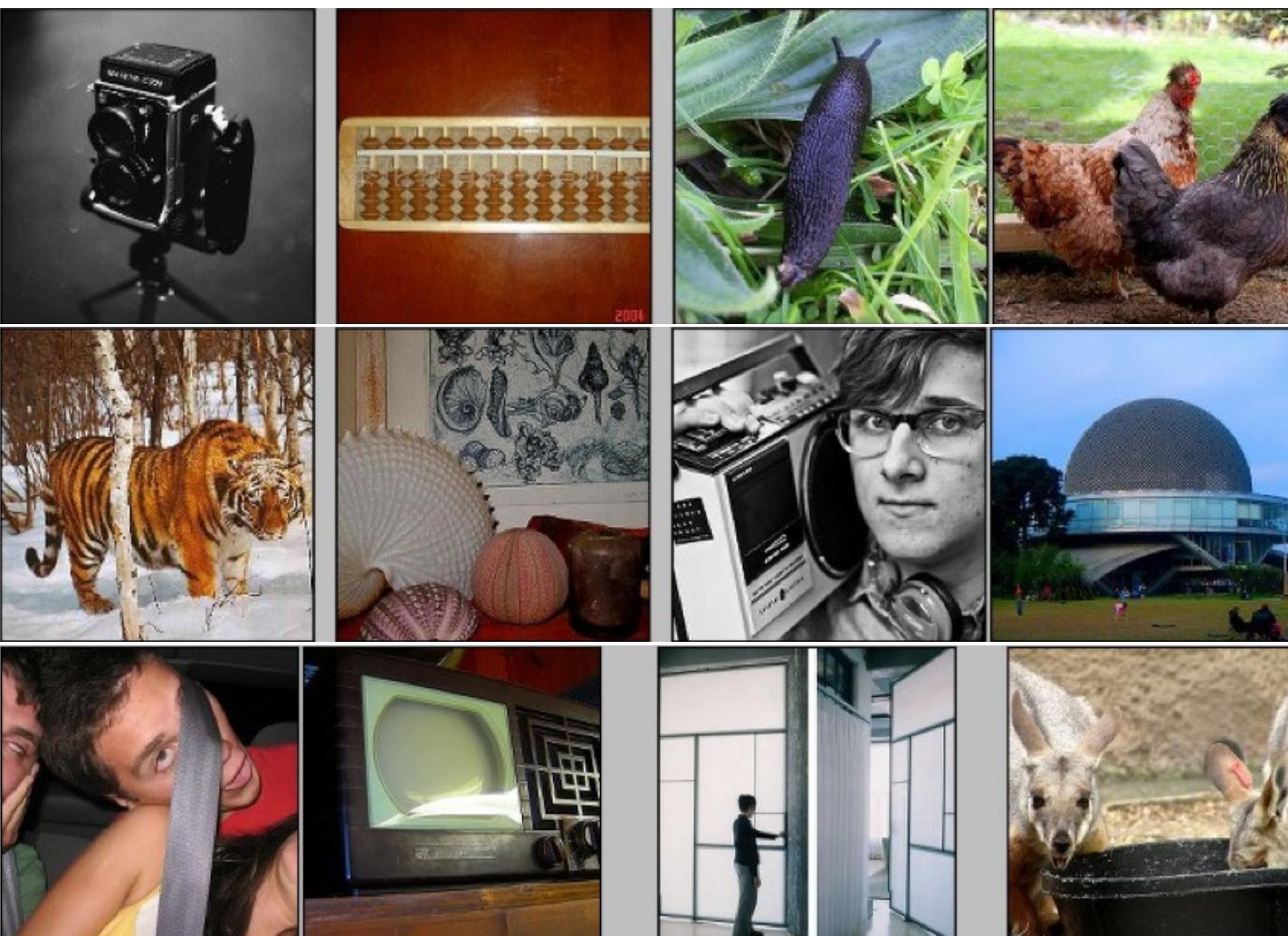
[Wu & He, “Group normalization”, ECCV 2018]

CNN Successes

- Handwritten text/digits
 - MNIST (0.17% error [Ciresan et al. 2011])
 - Arabic & Chinese [Ciresan et al. 2012]
- Simpler recognition benchmarks
 - CIFAR-10 (9.3% error [Wan et al. 2013])
 - Traffic sign recognition (0.56% error vs 1.16% for humans [Ciresan et al. 2011])
- But until recently, less good at more complex datasets
 - Caltech-101/256 (few training examples)



ImageNet Dataset

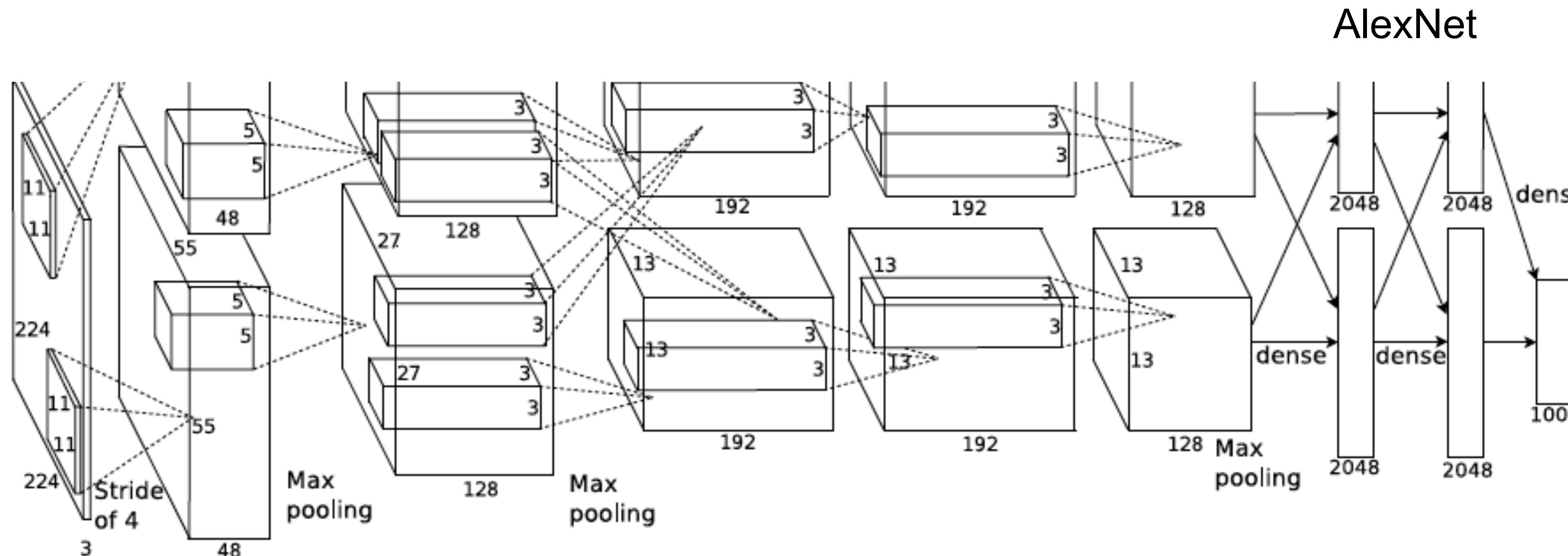


- ~14 million labeled images, 20k classes
- Challenge: 1.2 million training images, 1000 classes
- Images gathered from Internet
- Human labels via Amazon Mechanical Turk

[Deng et al. CVPR 2009]

ImageNet Challenge 2012 (ILSVRC)

- Similar framework to LeCun'98 but:
 - **Bigger model** (7 hidden layers, 60,000,000 params)
 - **More data** (10^6 vs. 10^3 images)
 - **GPU** implementation (50x speedup over CPU)
 - Trained on two GPUs for a week
 - Better regularization for training (DropOut)

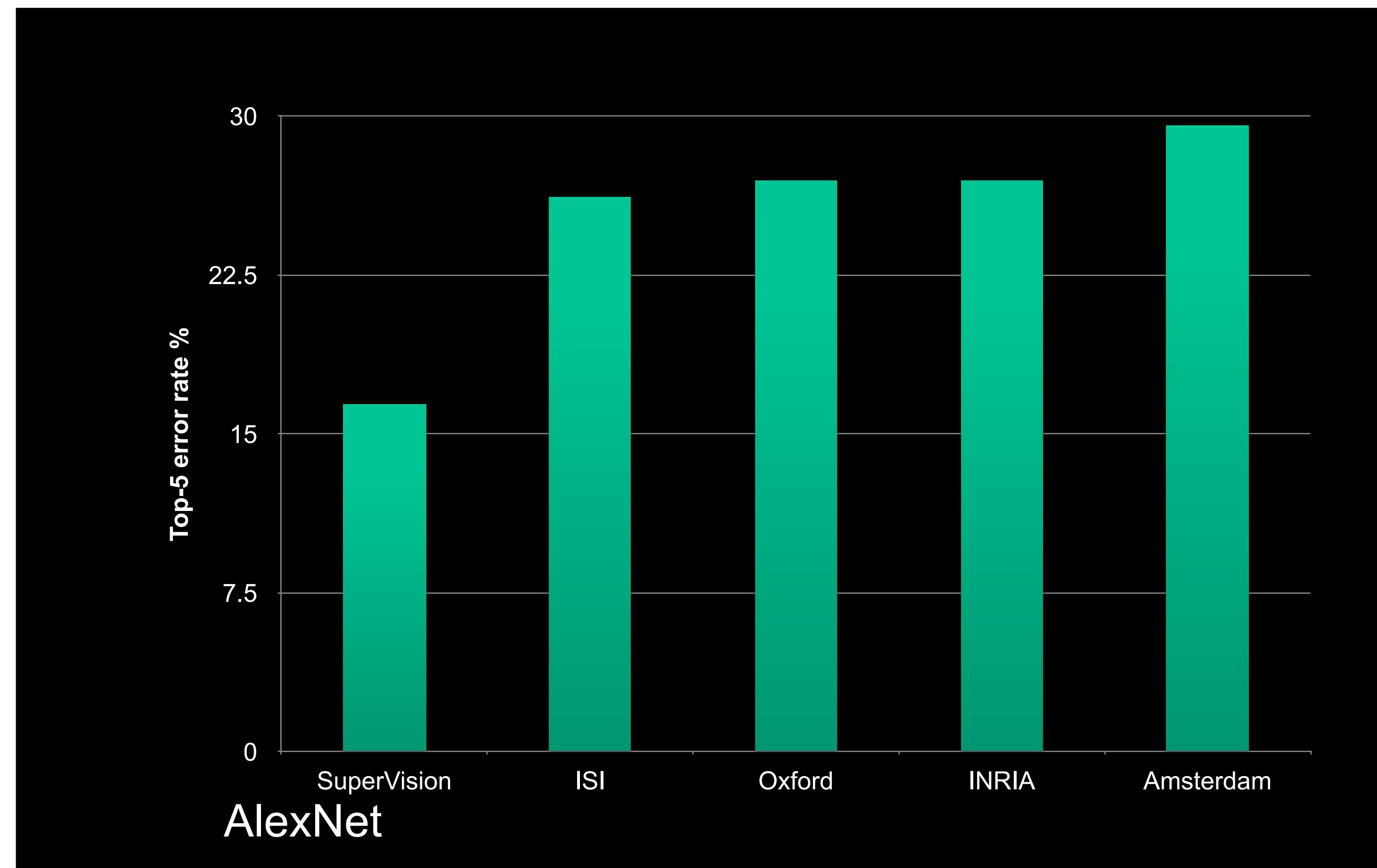


Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton,
[ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012](#)

ImageNet Challenge 2012 (ILSVRC)

AlexNet – 16.4% error (top-5)

Next best (non-convnet) – 26.2% error



Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - **Recap: Training NNs**
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Recap: Stochastic Gradient Descent (SGD)

The objective function is an average over all **N** training data points:

$$\theta_{t+1} \rightarrow \theta_t - \frac{\alpha_t}{N} \sum_i \frac{\partial \ell(\theta, X_i, Y_i)}{\partial \theta}$$

(gradient descent)

Key idea: approximate the gradient with **1** random datapoint:

$$\theta_{t+1} \rightarrow \theta_t - \alpha_t \frac{\partial \ell(\theta, X^{(i_t)}, Y^{(i_t)})}{\partial \theta}$$

(stochastic gradient descent)

Pick **K** random points instead of picking 1 (with $K \ll N$):

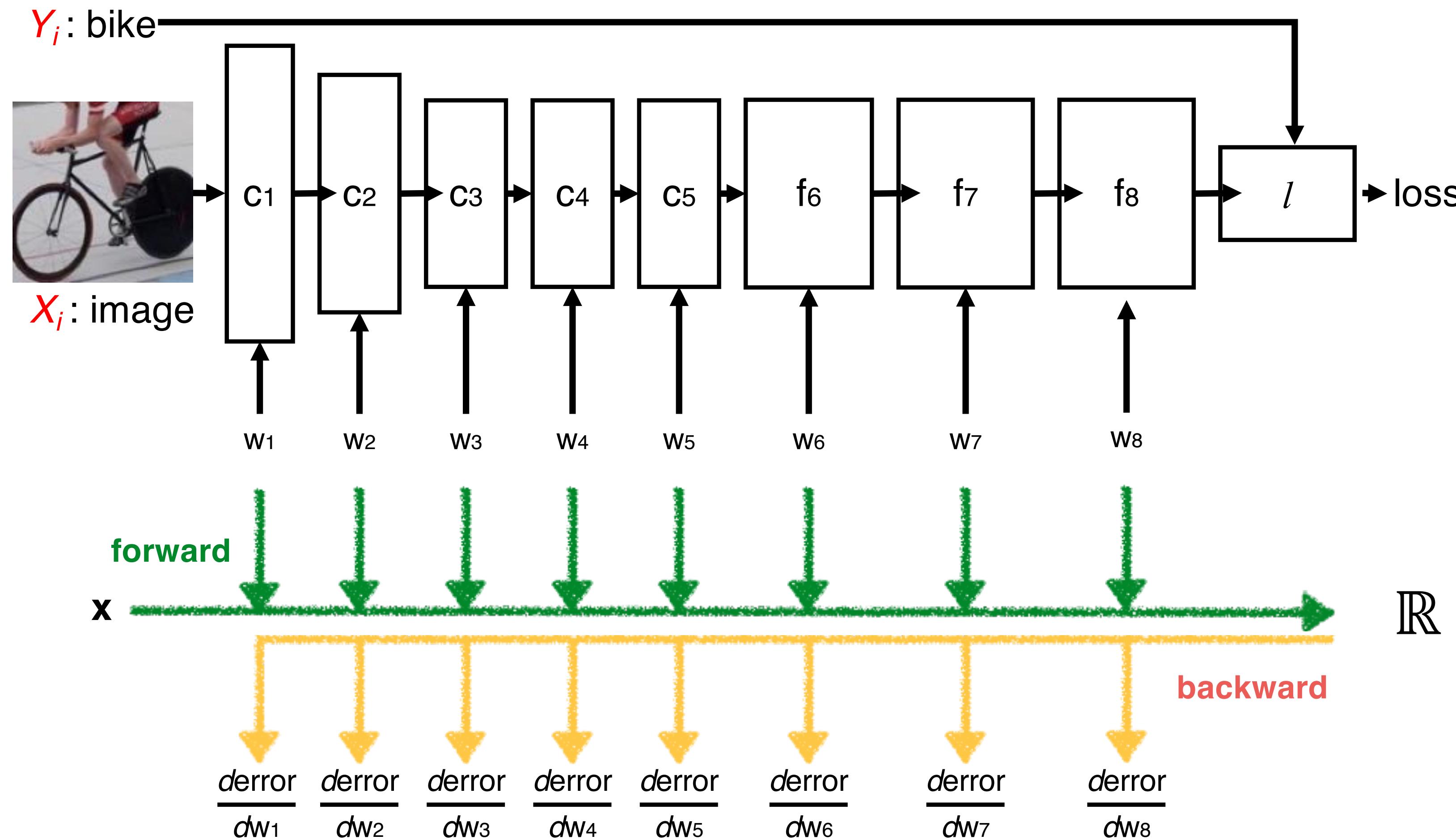
$$\theta_{t+1} \rightarrow \theta_t - \frac{\alpha_t}{K} \sum_{i=1}^K \frac{\partial \ell(\theta, X_i, Y_i)}{\partial \theta}$$

(stochastic gradient descent
with mini-batches)

=> commonly used

Recap: Backpropagation

Computing the gradients: While in theory, we just have the gradients of composite functions and for that apply **chain rule**, there is an **efficient** way to do it, called **backpropagation**.



Recap: Training a neural network

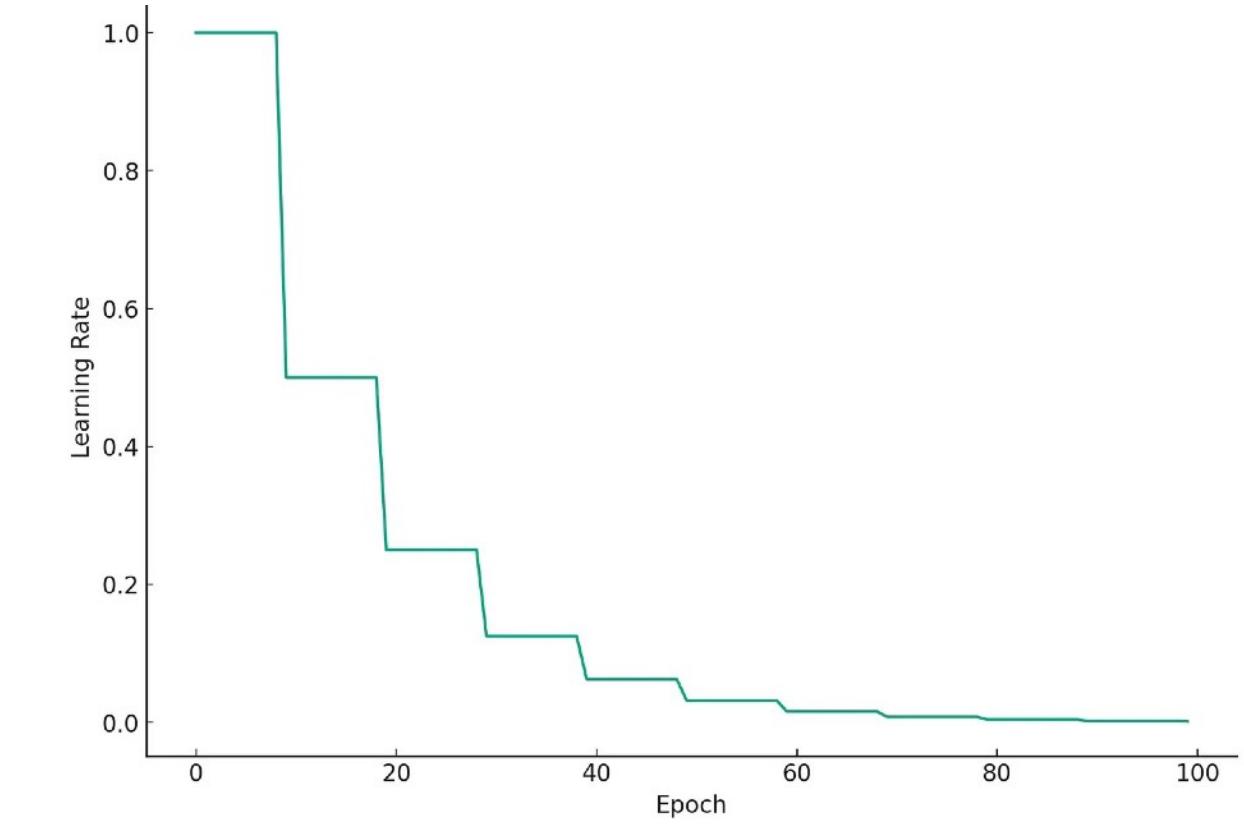
Given a (X, Y) pair:

- **Forward pass:** apply network to X to produce an output \hat{Y}
- **Evaluation:** Compute loss function, i.e., $\ell(\hat{Y}, Y)$
- **Backward pass:** compute the gradient with backpropagation
- **Update:** Take a step in the direction of the gradient

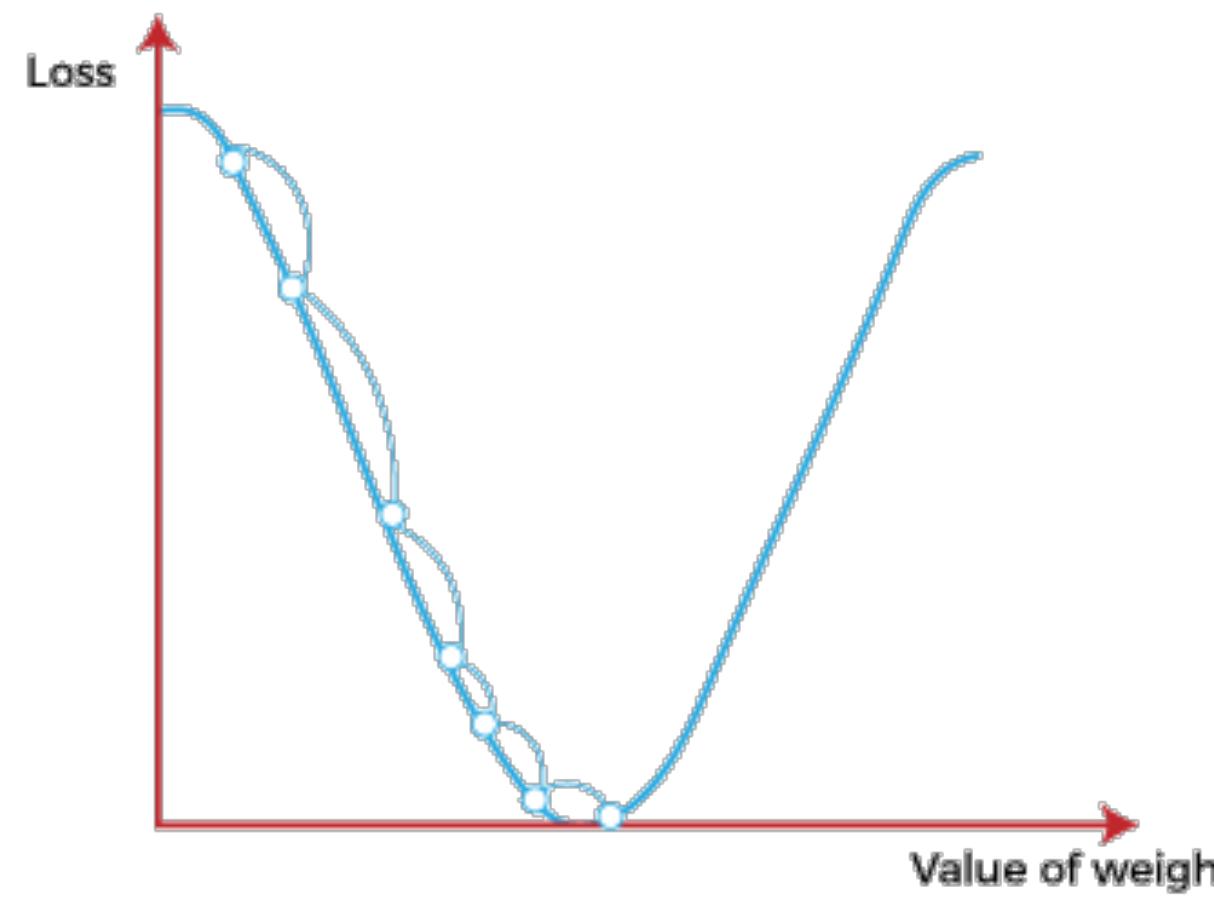
Recap: Learning rate policy

- There is no standard policy with NNs.
- The best is usually to pick a LR as high as possible without having divergence, keep it constant until convergence, then decrease it.
- Other more complex approaches are also widely used such as sinusoidal LR schedules, linear warmup etc.
- LR schedules even help with adaptive optimizers (e.g., AdaGrad)
- There are no guarantees with NNs, looking at training curves is crucial!

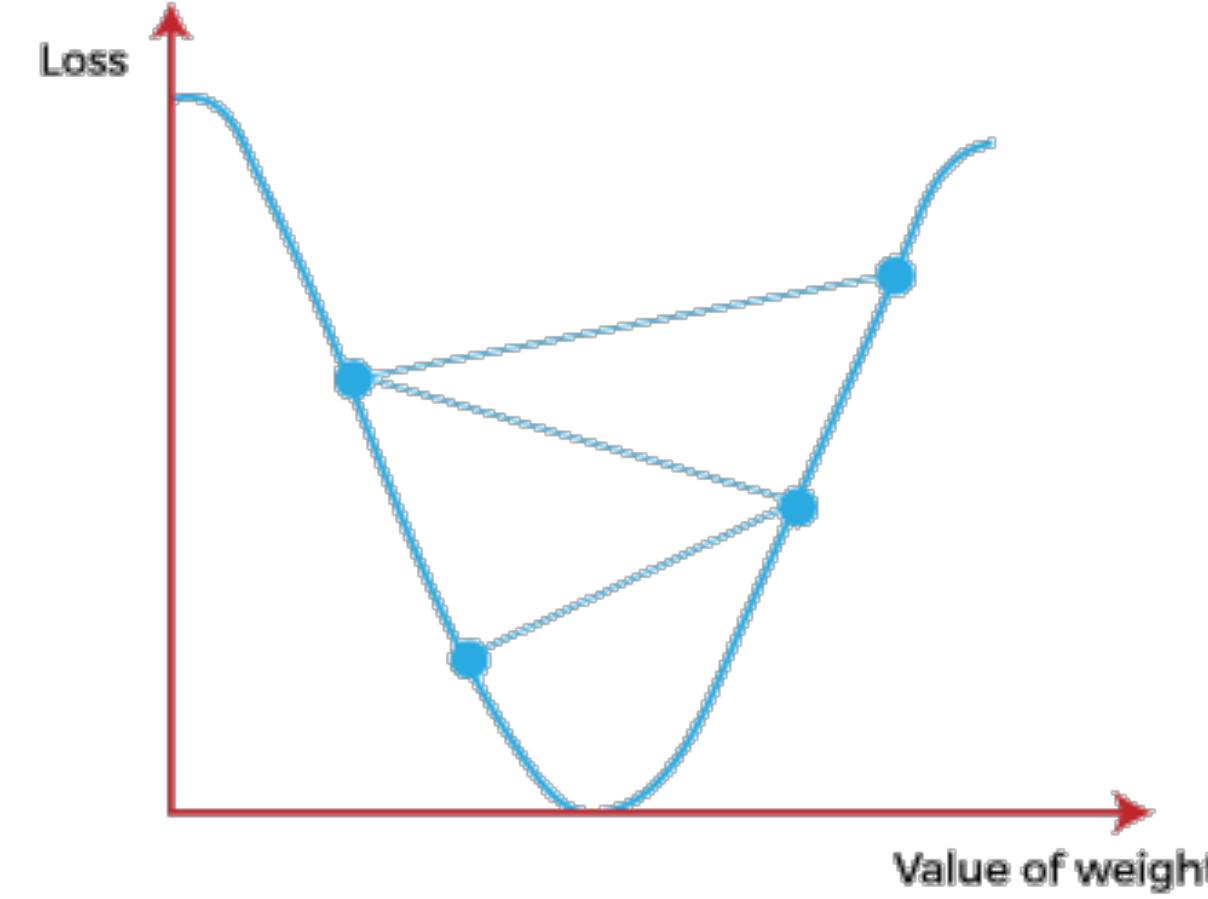
Step decay



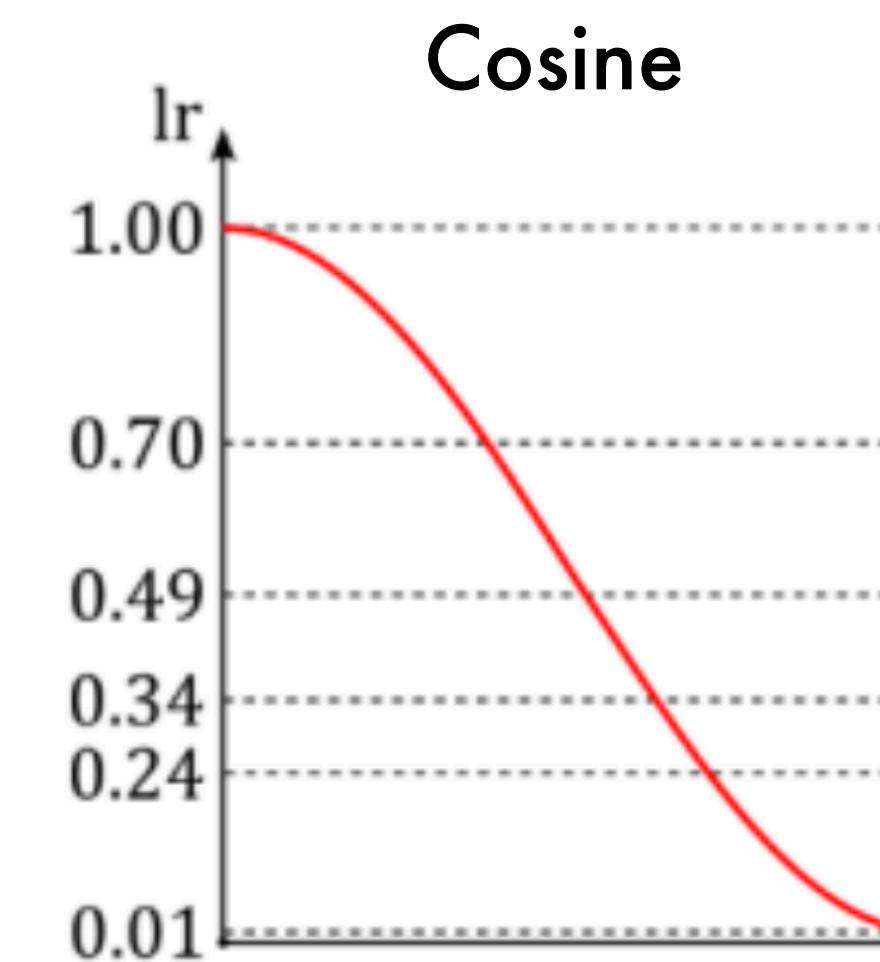
Small Learning Rate



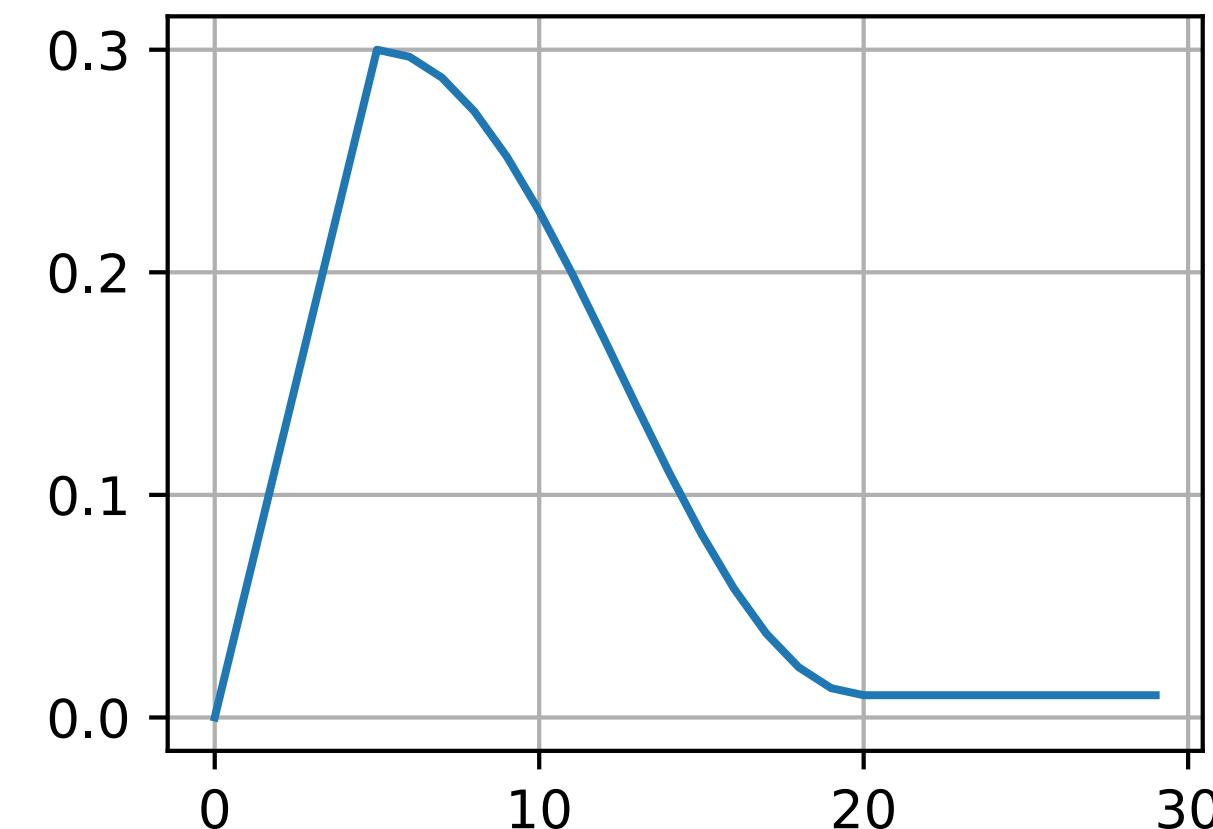
Large Learning Rate



Cosine

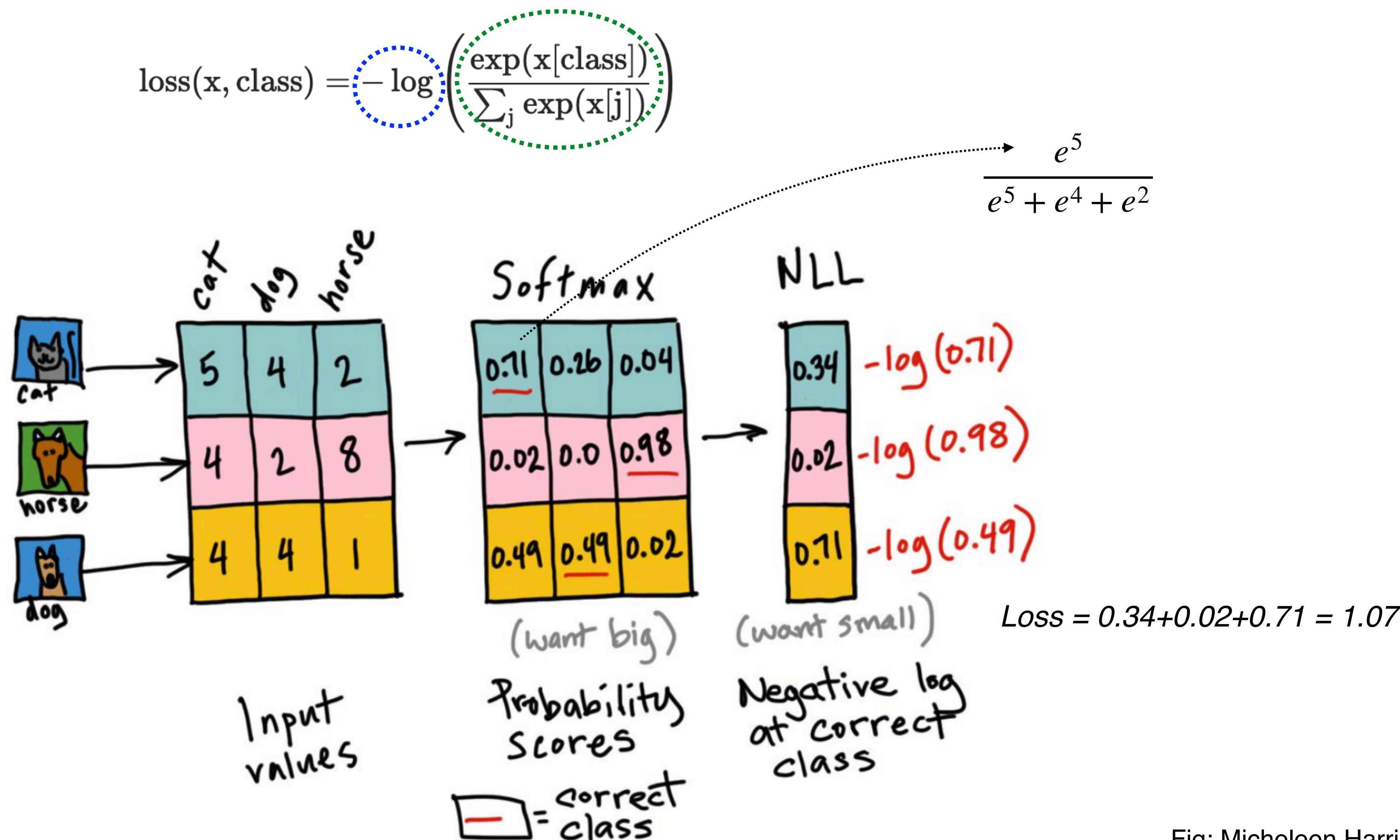


Linear warmup

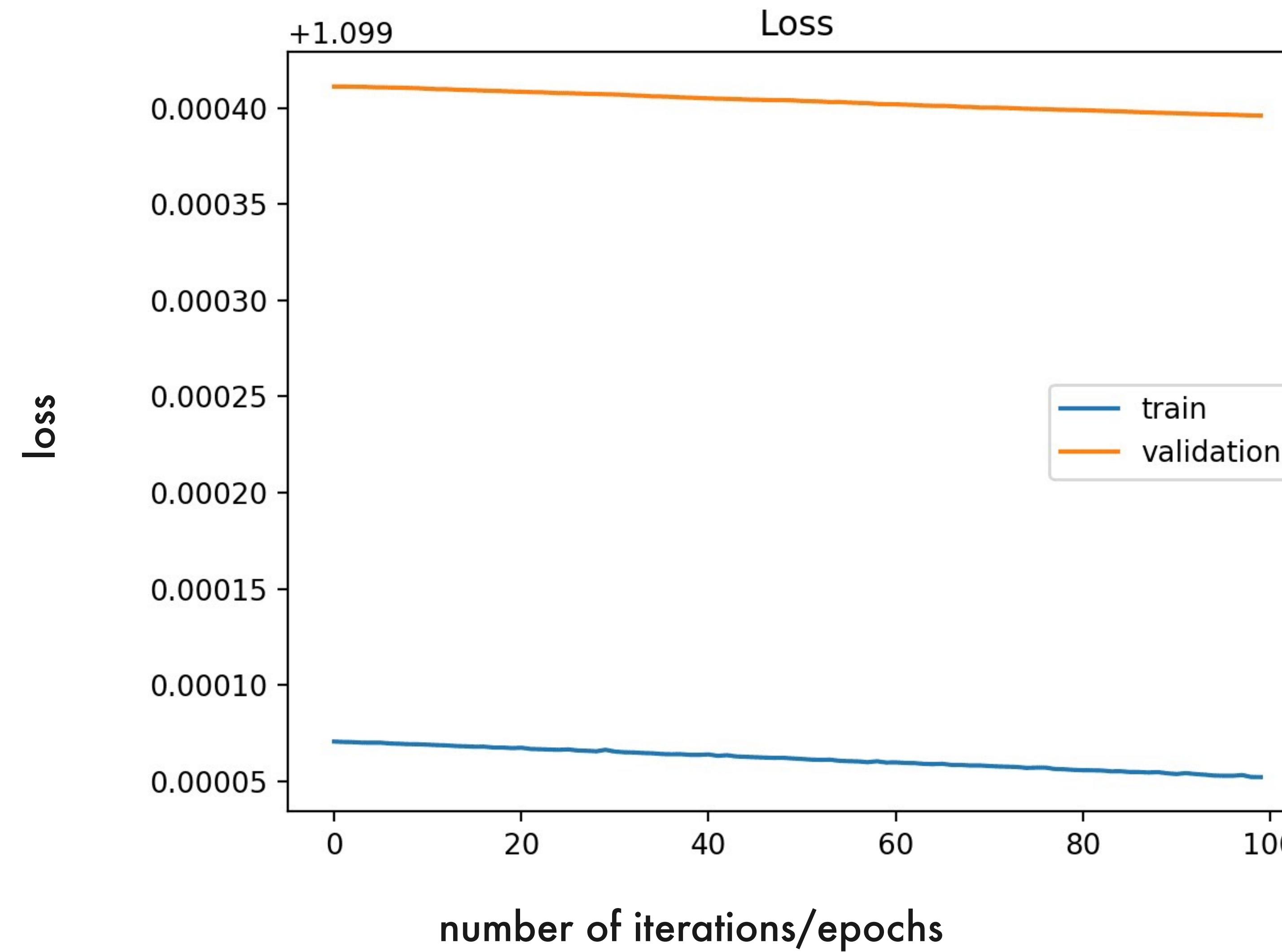


Recap: Classification loss

- Cross-entropy loss = softmax + negative log-likelihood



How can I resolve the problem with training?

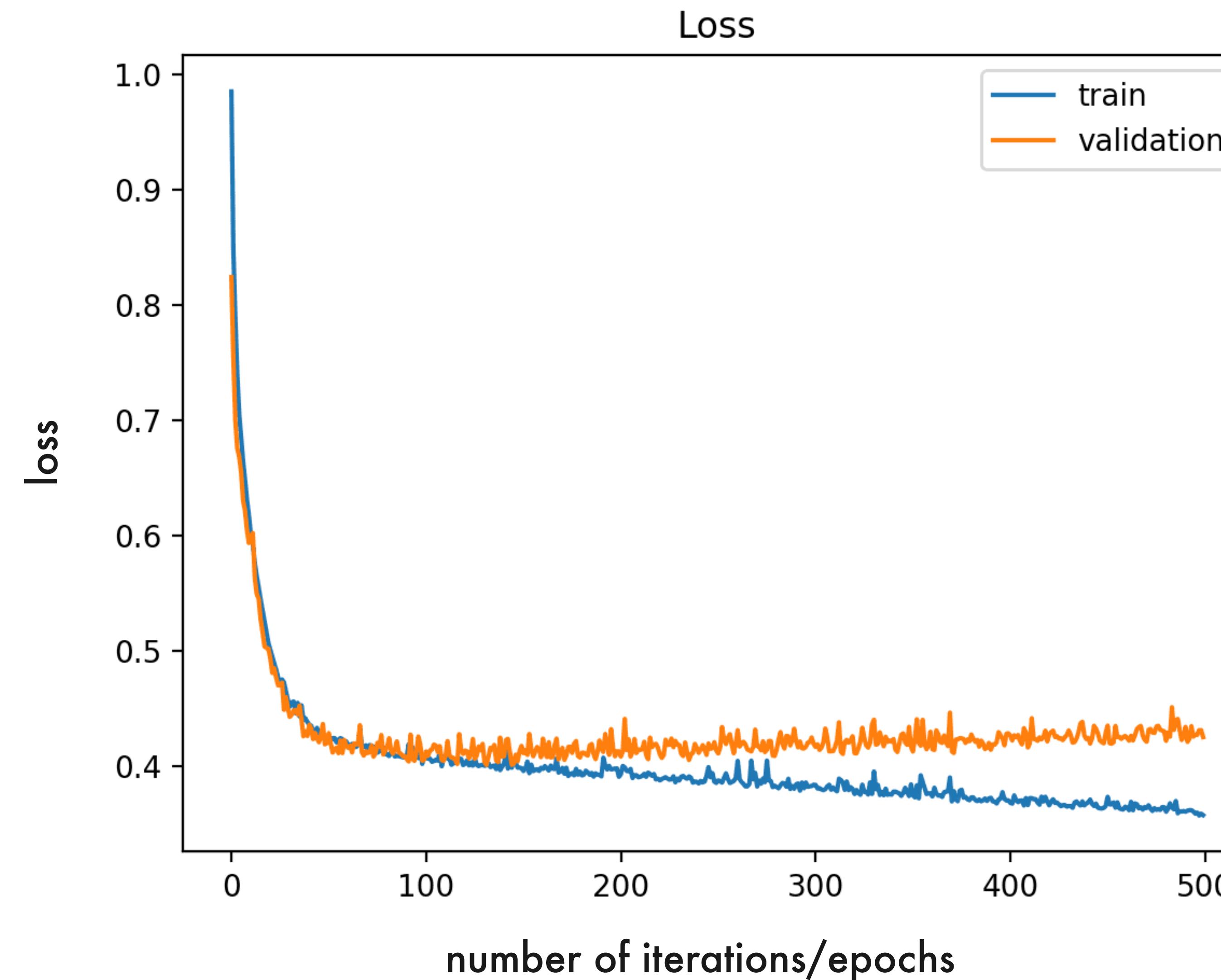


Credit: [Jason Brownlee](#)

[Image sources](#)

102

How can I resolve the problem with training?



Credit: [Jason Brownlee](#)

[Image sources](#)

How to avoid overfitting?

Deep networks have many parameters.

Some regularization techniques:

- Smaller network, i.e., less parameters
- Data augmentation
- Suboptimize, i.e., “early stopping”
- Force redundancy in hidden units, i.e., “dropout”
- Penalize parameter norms, i.e., “weight decay”

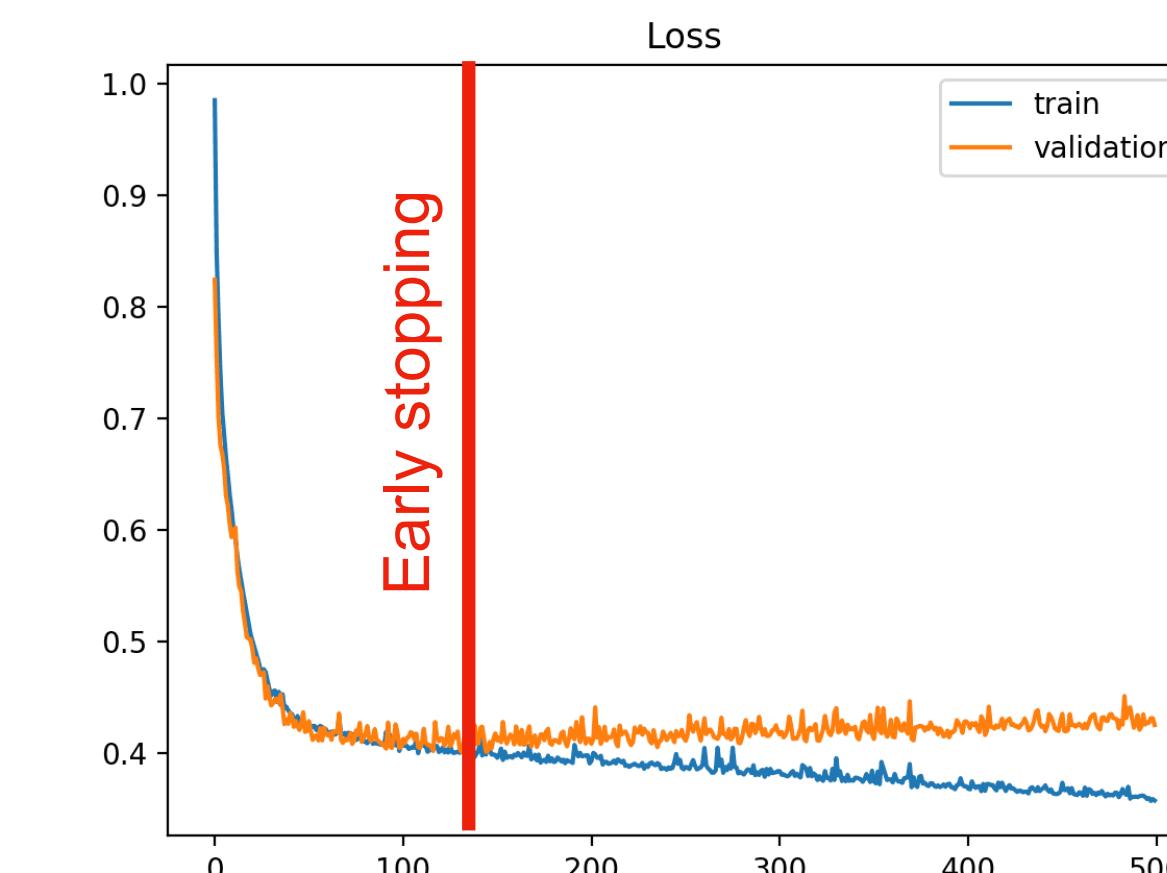
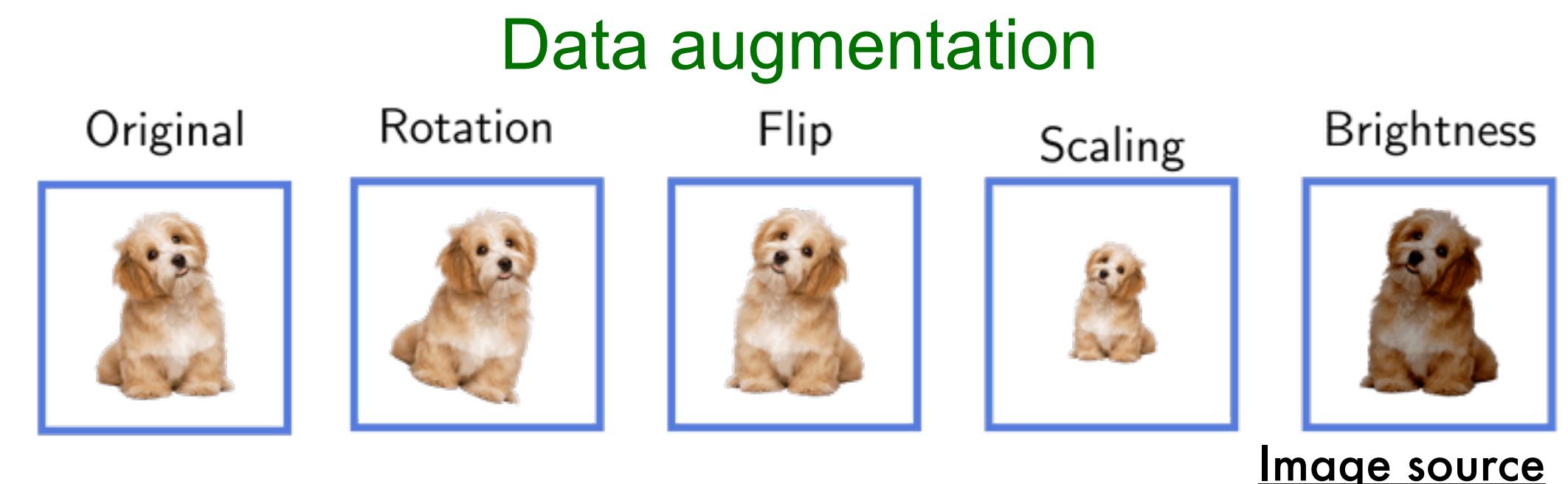
L2 penalty:

encourages the norm of
the parameters to be low

$$L(\theta) + \frac{\lambda}{2} \left\| \theta \right\|_2^2$$

$$\theta_{t+1} \rightarrow \theta_t - \alpha \frac{\partial L}{\partial \theta_t} - \alpha \lambda \theta_t$$

gradient



Agenda

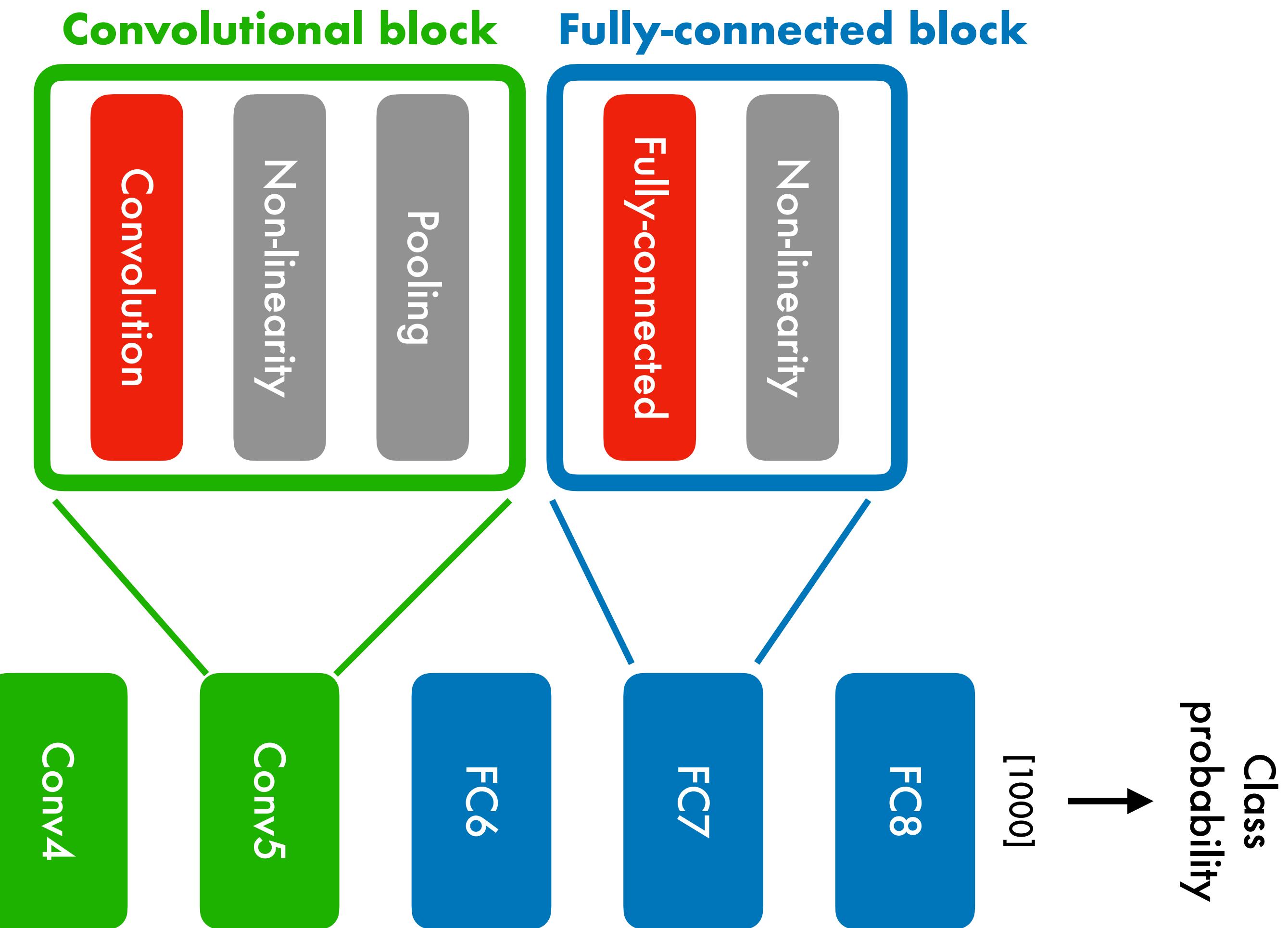
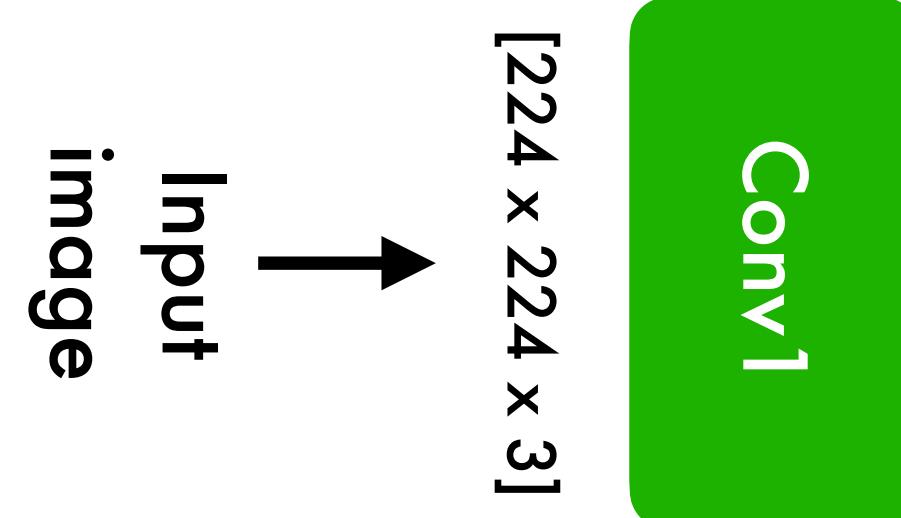
- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Visualizing CNNs

What does CNN learn once it is trained?

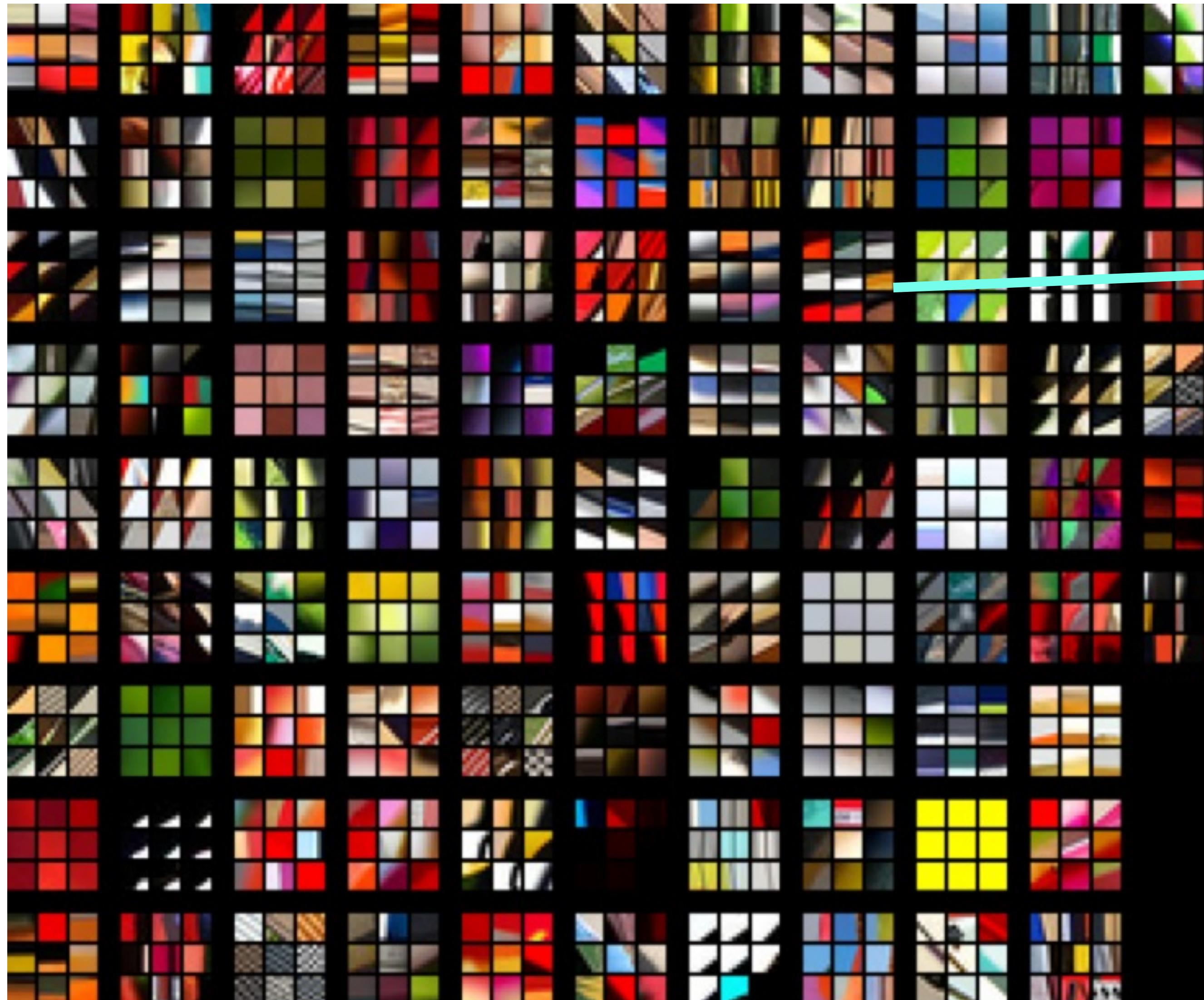
Recap: AlexNet

- **Fully-connected layer**
- **Convolution layer**
- **Pooling layer** (e.g., Max-pooling)
- **Non-linearity layer** (e.g., ReLU)
- ...

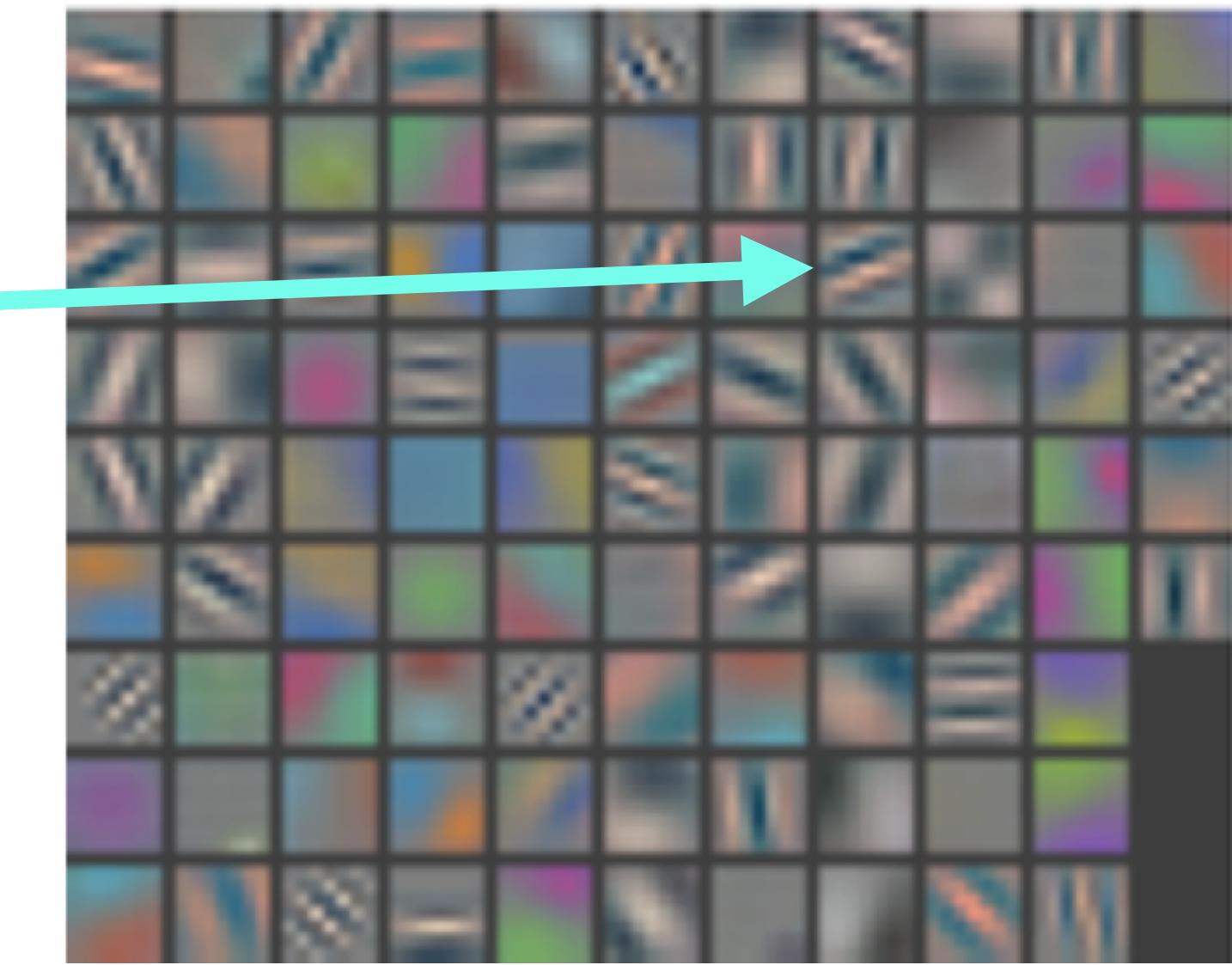


Layer 1: Top-9 patches

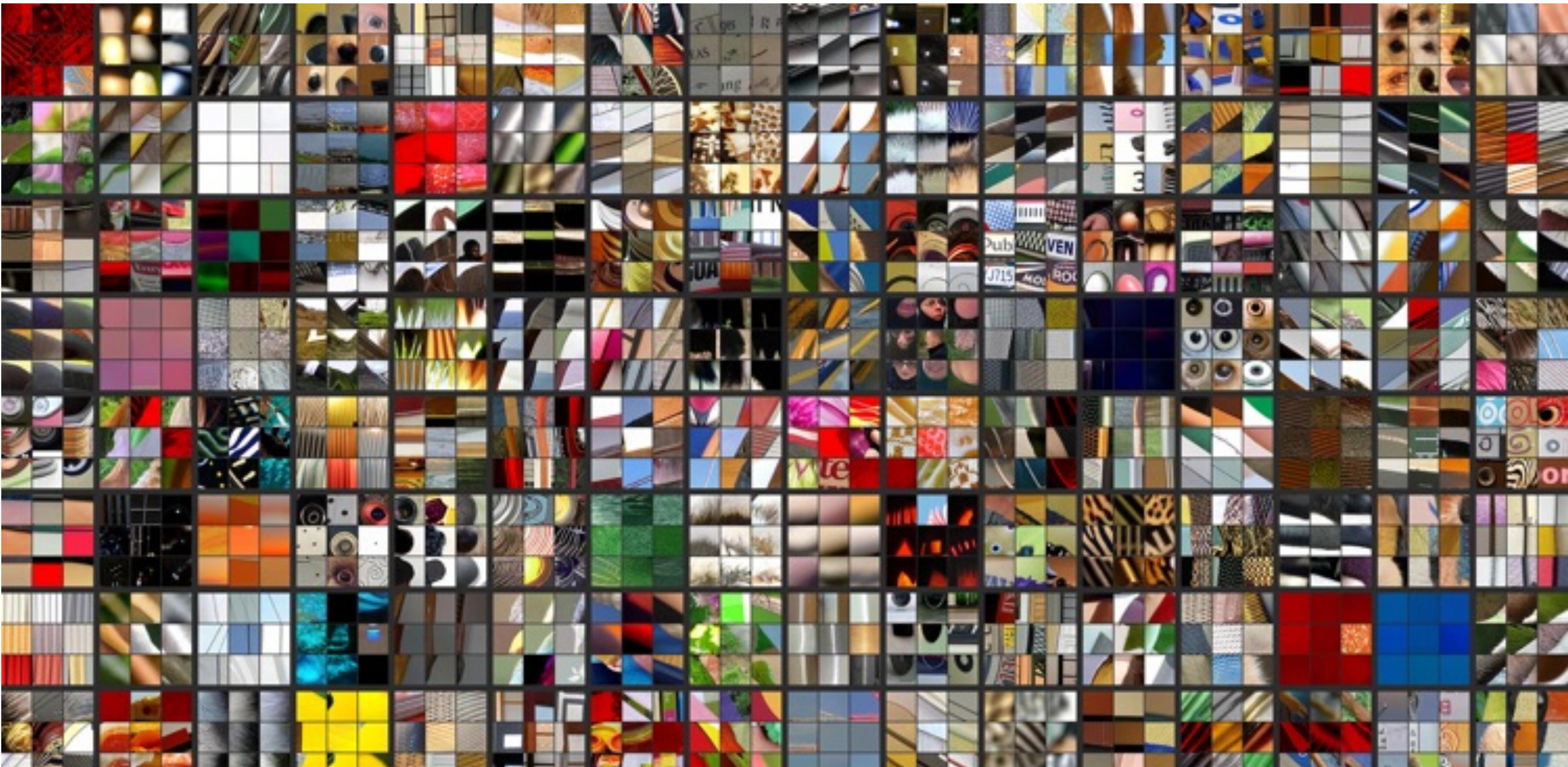
Patches from validation images that give maximal activation of a given feature map



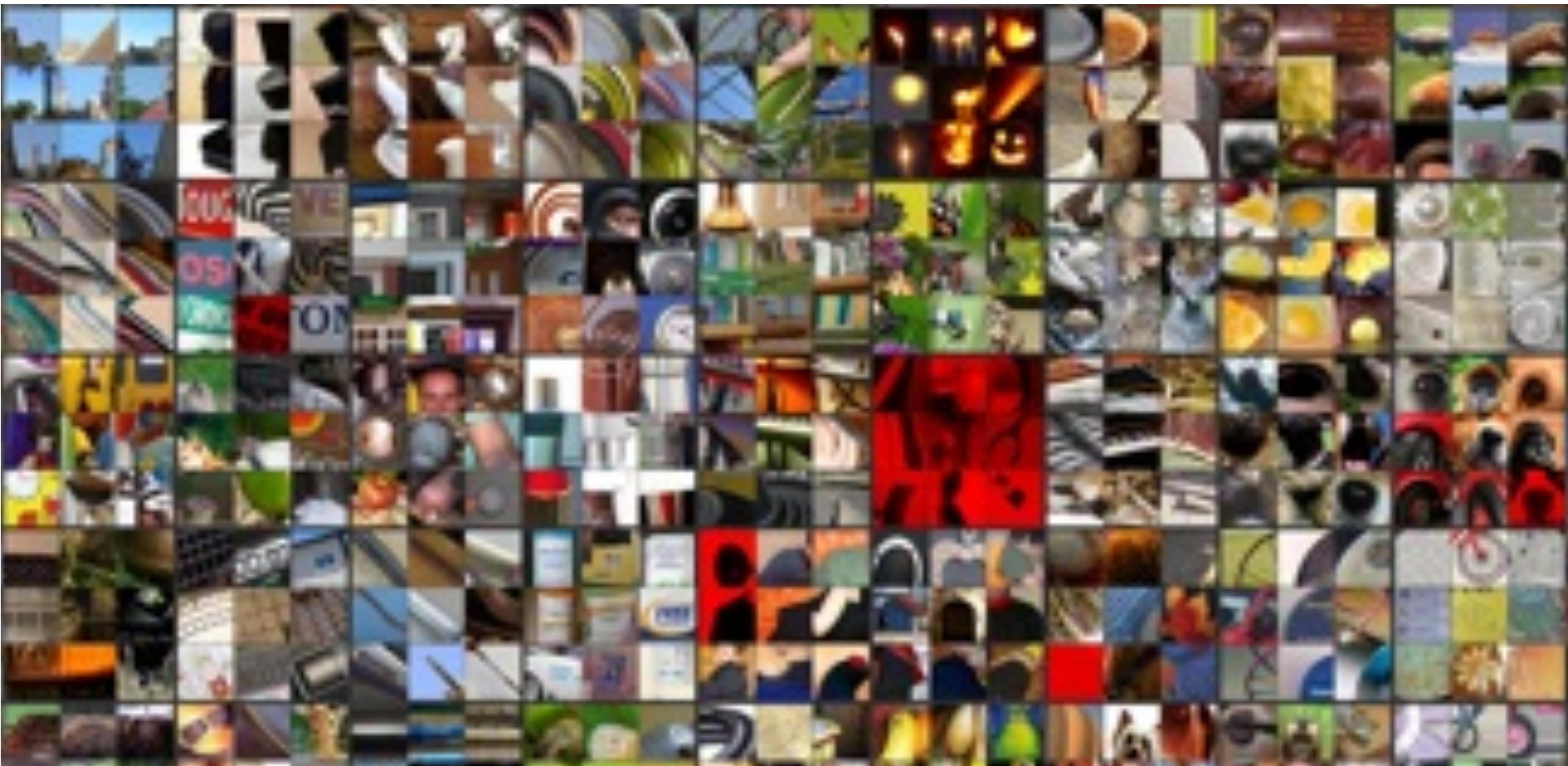
Learned filters



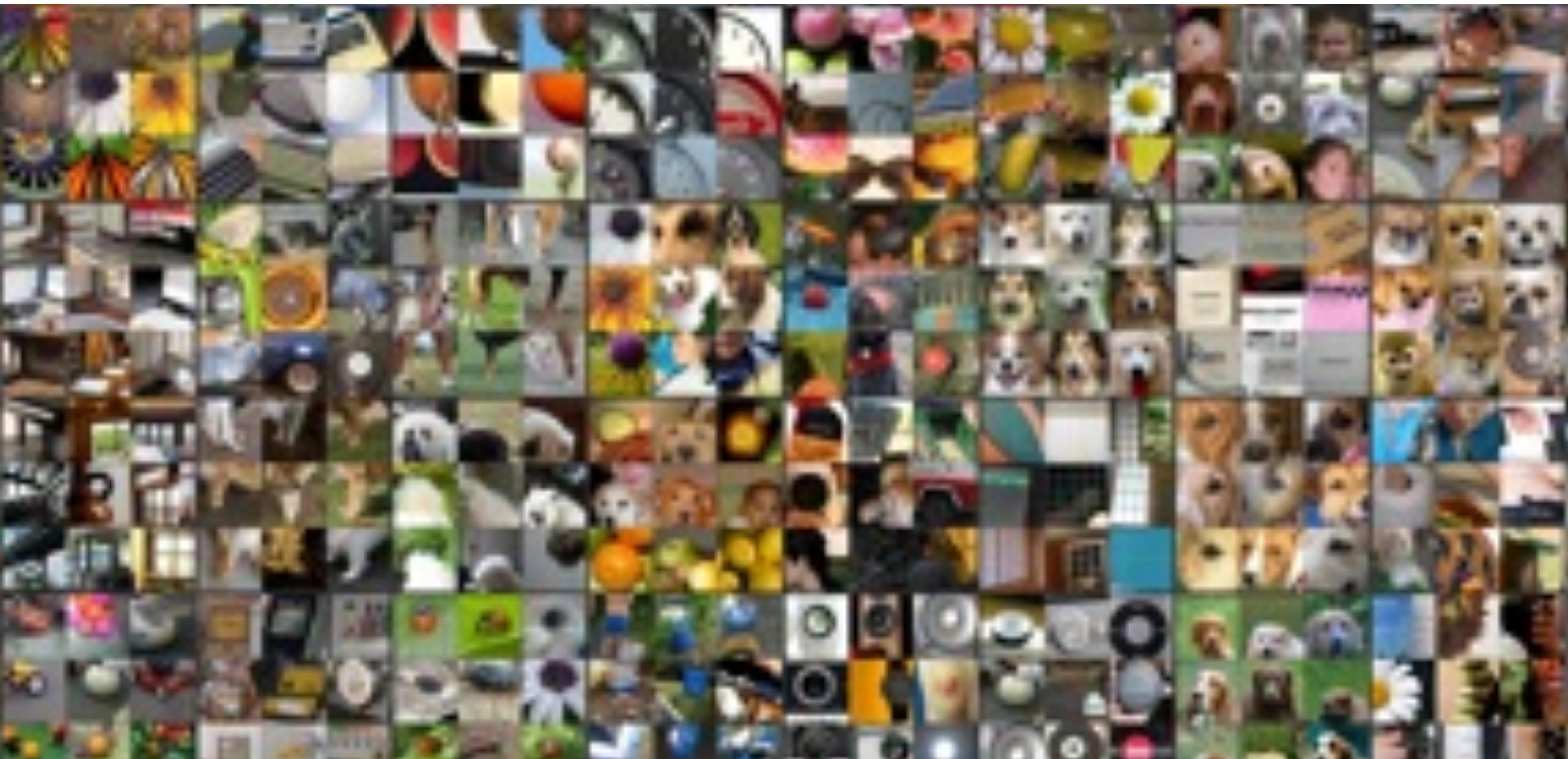
Layer 2: Top-9 patches



Layer 3: Top-9 patches



Layer 4: Top-9 patches



Layer 5: Top-9 patches



References: Visualizing and understanding NNs

Analysis tools

Visualizing higher-layer features of a deep network
Ethan et al. 2009
[intermediate features]

Deep inside convolutional networks
Simonyan et al. 2014
[deepest features, aka “deep dreams”]

DeConvNets
Zeiler et al. In ECCV, 2014
[intermediate features]

Understanding neural networks through deep visualisation
Yosinski et al. 2015
[intermediate features]

Artistic tools

Google’s “inceptionism”
Mordvintsev et al. 2015

Style synthesis and transfer
Gatys et al. 2015

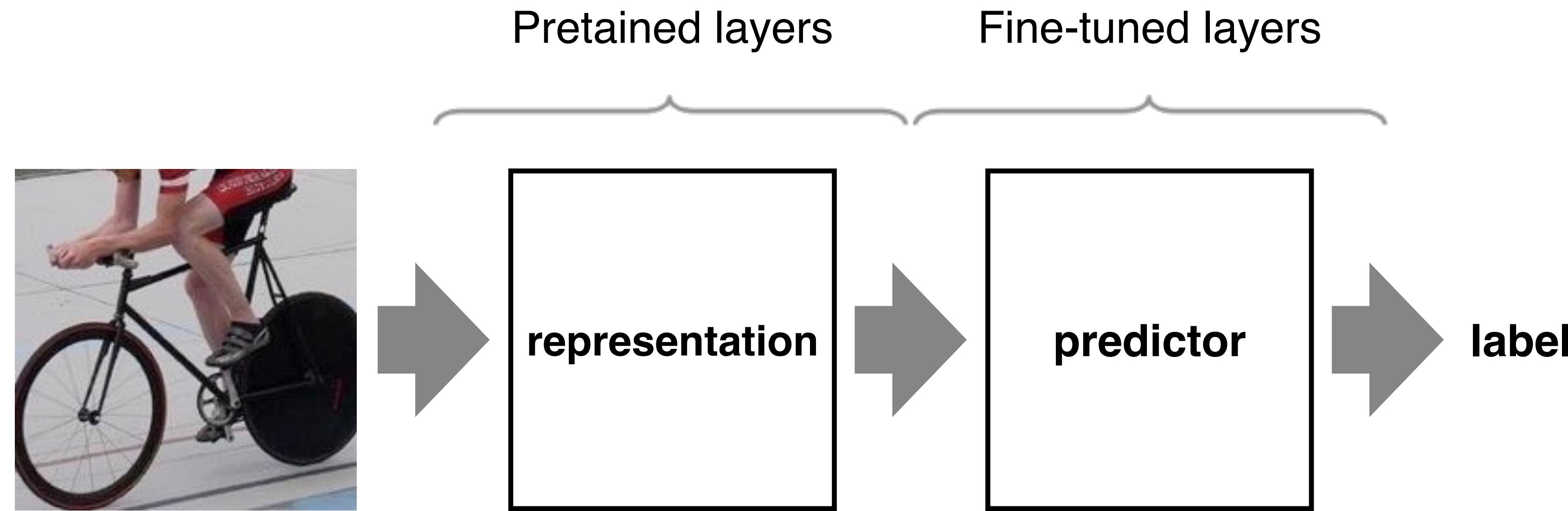
Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

Transferring learnt representations

“pretraining”

“Pre-training” and transfer learning



CNN as universal representations

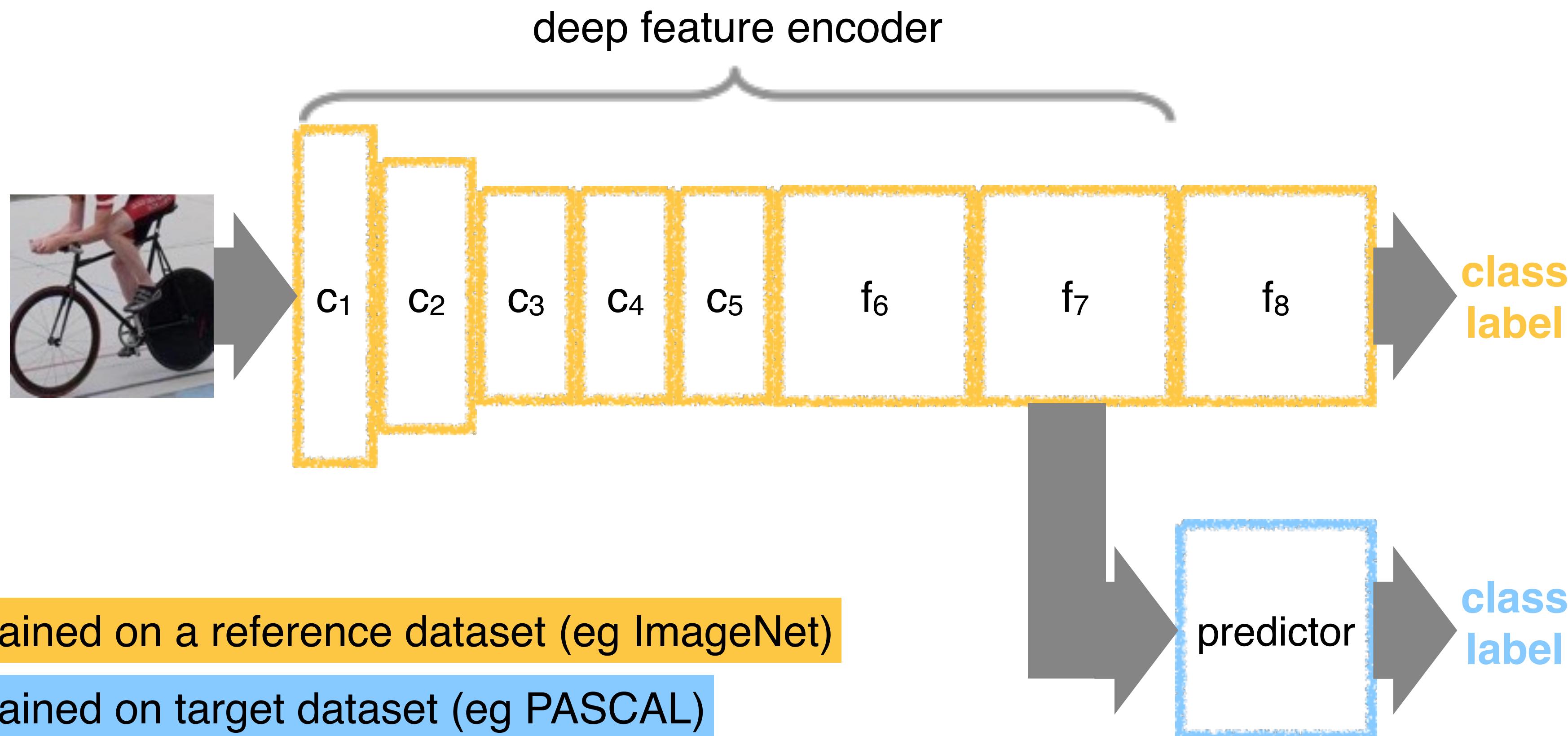
- First several layers in most CNNs are generic
- They can be reused when training data is comparatively scarce

Application

- Pre-train on ImageNet classification 1M images
- Cut at some deep conv or FC layer to get features

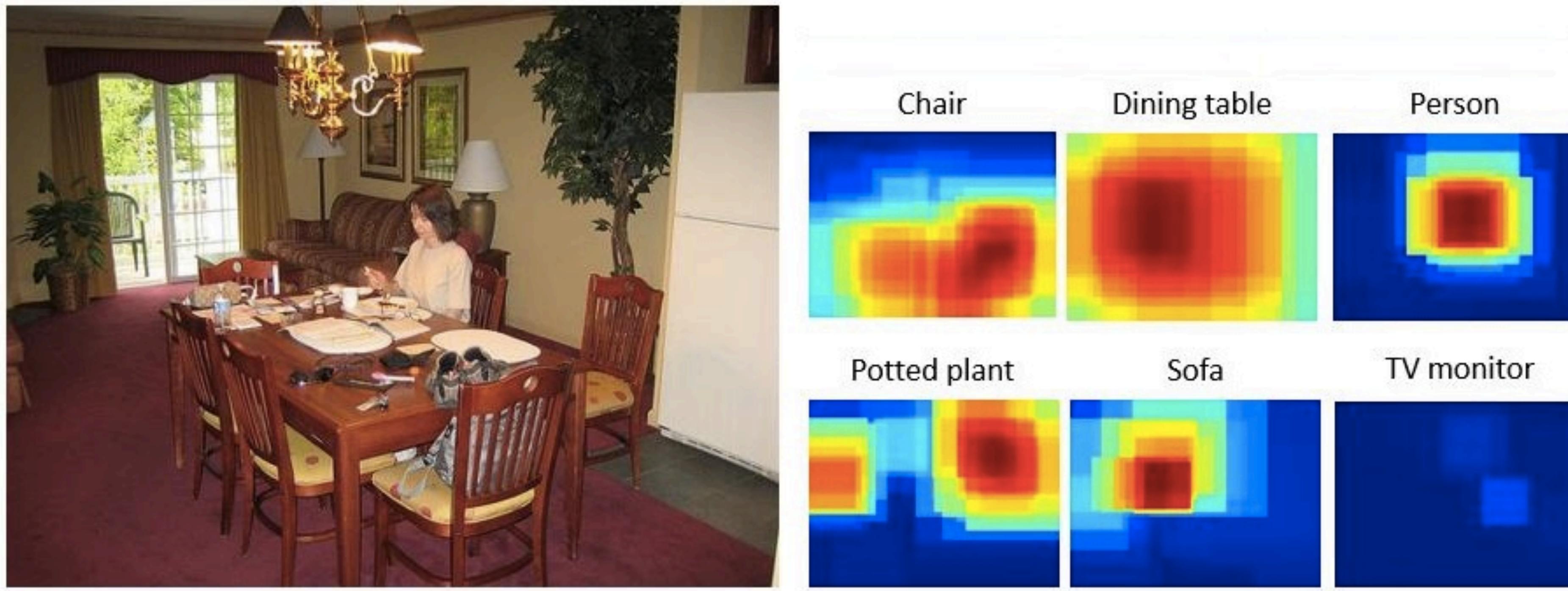
“Pre-training” and transfer learning

Deep representations are generic



A general purpose deep encoder is obtained by chopping off the last layers of a CNN trained on a large dataset.

Example



**Learning and Transferring Mid-Level Image
Representations using Convolutional Neural Networks**
M. Oquab, L. Bottou, I. Laptev, J. Sivic
In CVPR 2014

<http://www.di.ens.fr/willow/research/cnn/>

ImageNet classification challenge

ImageNet classification
challenge



Object centric
1000 classes
1.2M images

What about other recognition tasks and datasets?

ImageNet classification
challenge

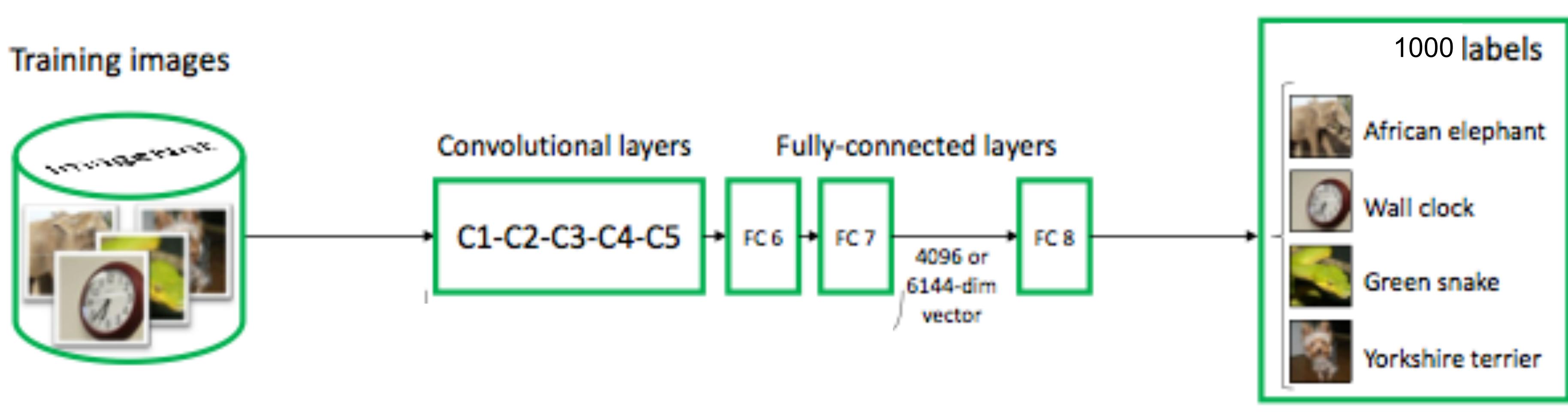


Object centric
1000 classes
1.2M images



Complex scenes
20 classes
10k images

Background – Convolutional neural network of [Krizhevsky et al. 2012]



Input: ~1M labelled images (1000 images / 1000 classes)

Number of parameters: ~60 million --- image representation

Training time: ~1 week on one GPU

Learn parameters using stochastic gradient descent on cross-entropy error function.

Can we transfer learnt parameters to other tasks with limited training data?

Challenge

The dataset statistics between the source task ([ImageNet](#)) and the target task ([Pascal VOC](#)) can be very different.

- Type of objects and labels
- Object size, object location, scene clutter
- Object viewpoints, imaging conditions

ImageNet



Maltese terrier

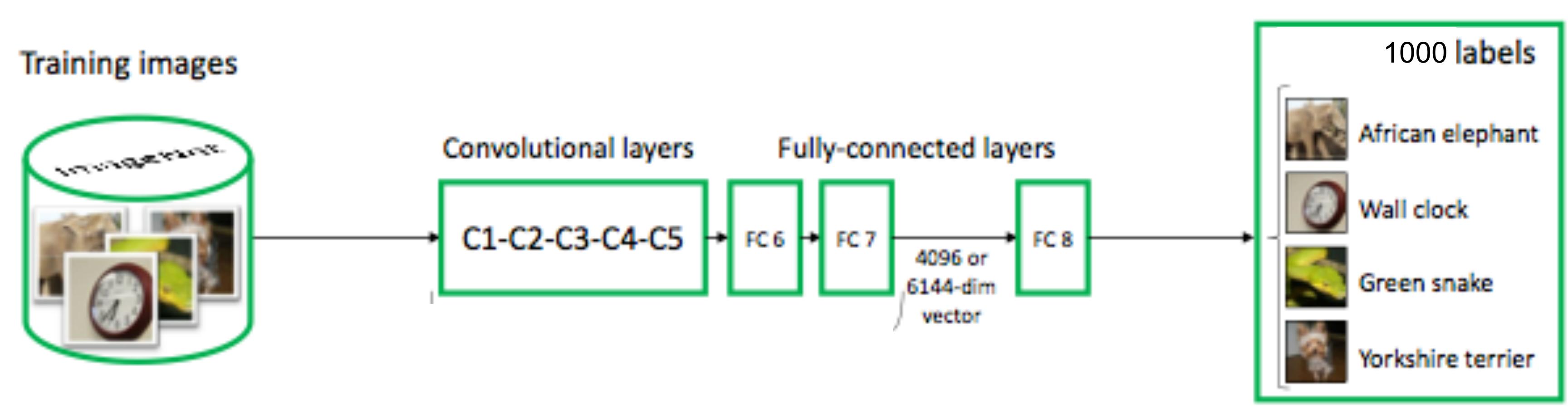
Pascal VOC



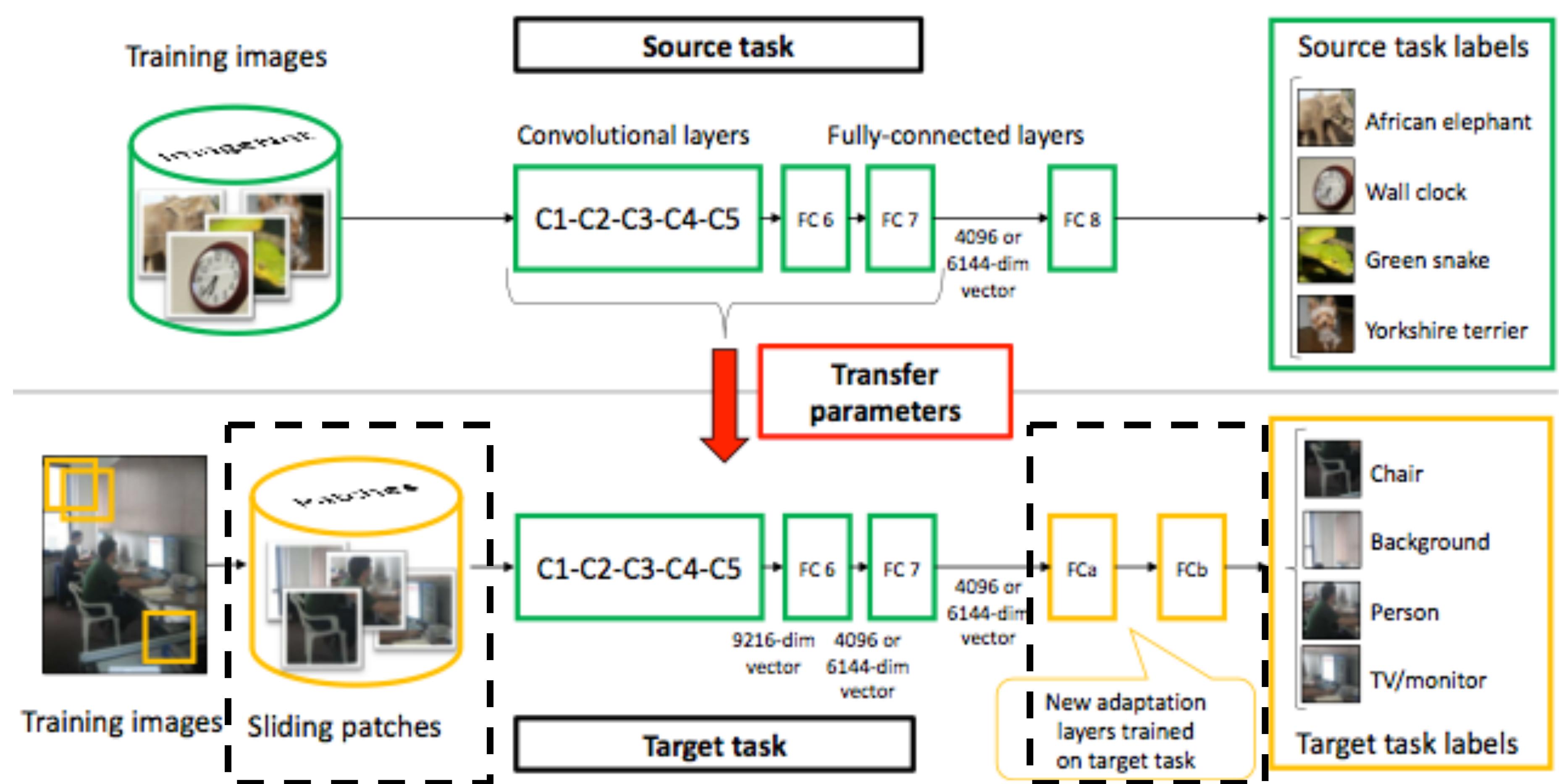
Dog



Approach



Approach [Oquab, Bottou, Laptev, Sivic, CVPR'14]



1. Design training/test procedure using sliding windows
2. Train adaptation layers to map labels

See also [Girshick et al.'13], [Donahue et al.'13], [Sermanet et al. '14], [Zeiler and Fergus '13]
Transfer learning workshop at ICCV'13, ImageNet workshop at ICCV'13

Pre-training helps

	plane	bike	bird	boat	btl	bus	car	cat	chair	cow	table	dog	horse	moto	pers	plant	sheep	sofa	train	tv	mAP
NO PRETRAIN	85.2	75.0	69.4	66.2	48.8	82.1	79.5	79.8	62.4	61.9	49.8	75.9	71.4	82.7	93.1	59.1	69.7	49.3	80.0	76.7	70.9
PRE-1000C	93.5	78.4	87.7	80.9	57.3	85.0	81.6	89.4	66.9	73.8	62.0	89.5	83.2	87.6	95.8	61.4	79.0	54.3	88.0	78.3	78.7
PRE-1000R	93.2	77.9	83.8	80.0	55.8	82.7	79.0	84.3	66.2	71.7	59.5	83.4	81.4	84.8	95.2	59.8	74.9	52.9	83.8	75.7	76.3
PRE-1512	94.6	82.9	88.2	84.1	60.3	89.0	84.4	90.7	72.1	86.8	69.0	92.1	93.4	88.6	96.1	64.3	86.6	62.3	91.1	79.8	82.8

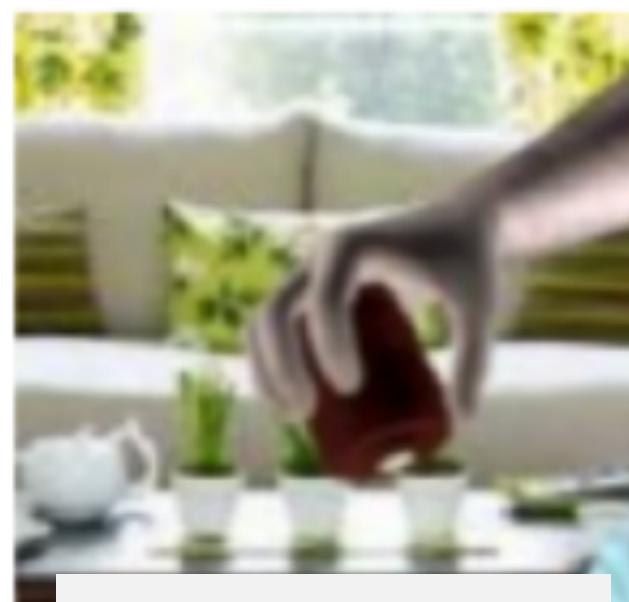
- Pascal VOC 2012 object classification

Action	jump	phon	instr	read	bike	horse	run	phot	comp	walk	mAP
NO PRETRAIN	43.2	30.6	50.2	25.0	76.8	80.7	75.2	22.2	37.9	55.6	49.7
PRE-1512	73.4	44.8	74.8	43.2	92.1	94.3	83.4	45.7	65.5	66.8	68.4
PRE-1512U	74.8	46.0	75.6	45.3	93.5	95.0	86.5	49.3	66.7	69.5	70.2

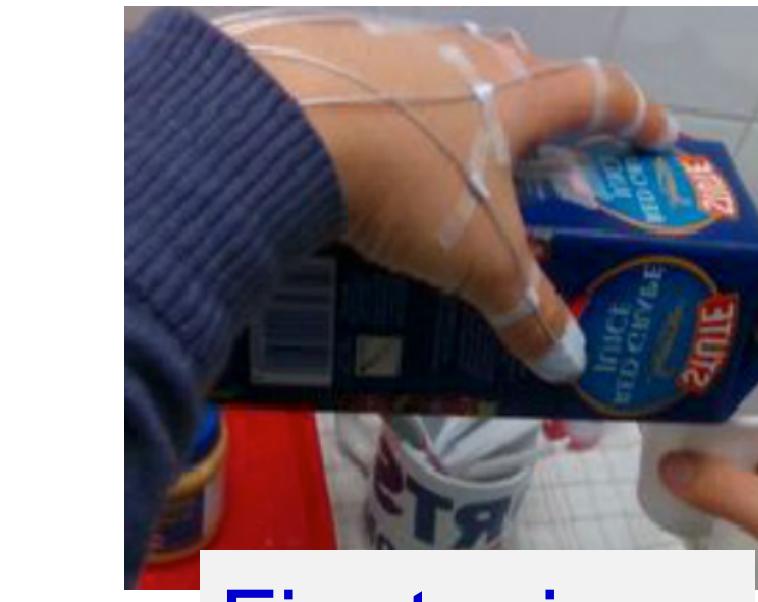
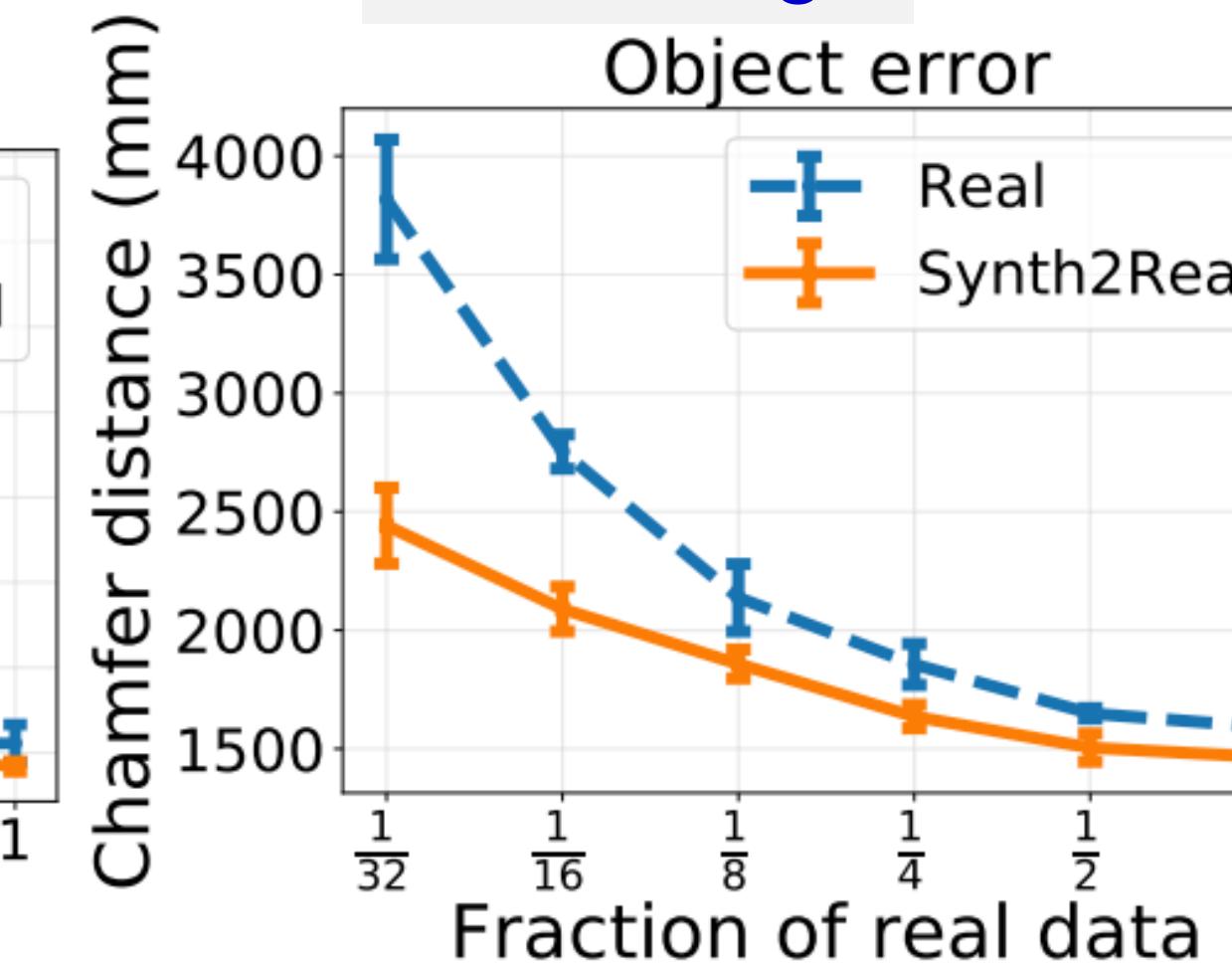
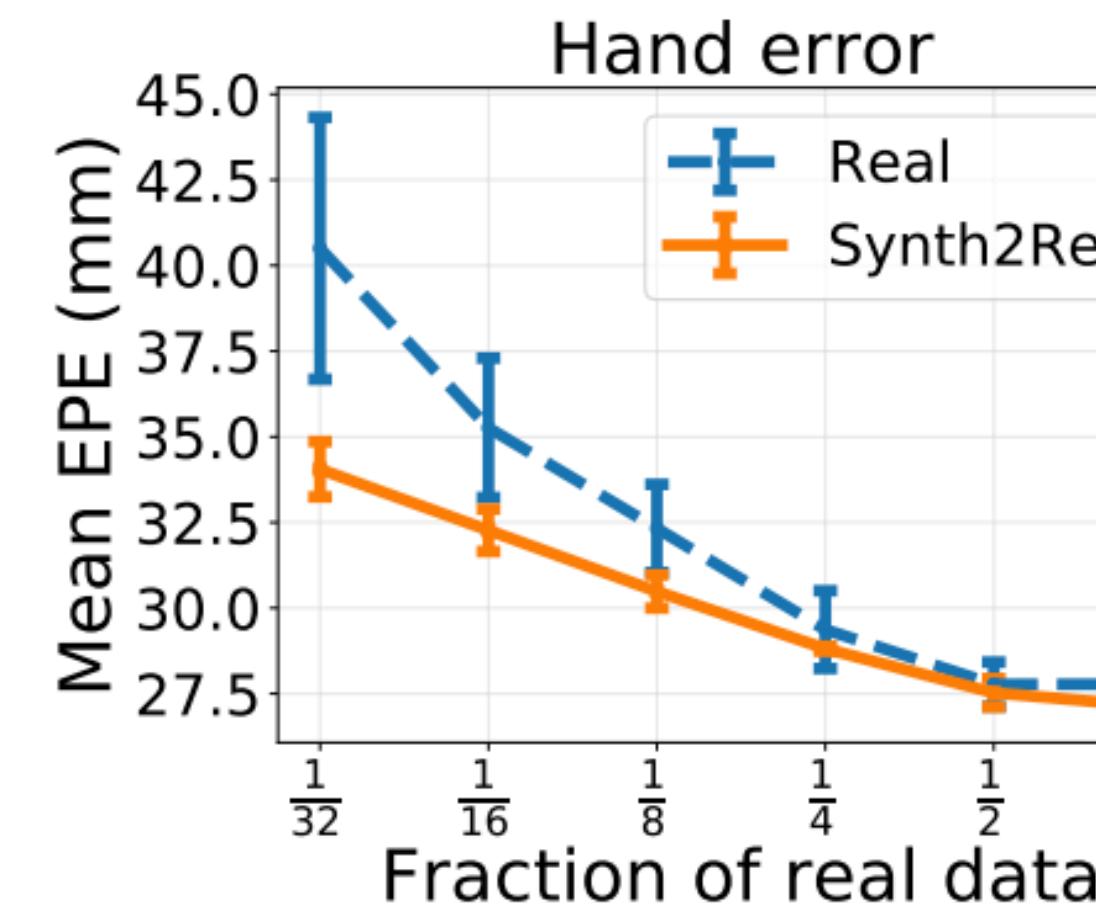
- Pascal VOC 2012 action classification

**Other "pre-training"
examples**

3D hand-object reconstruction



Pretraining



Finetuning

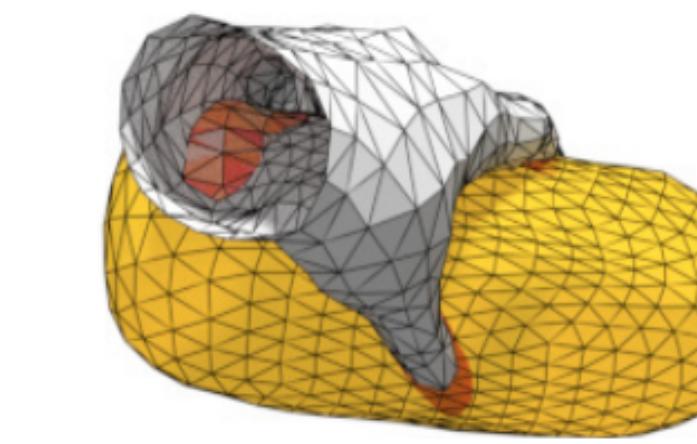


Figure 8: We compare training on FHB only (Real) and pre-training on synthetic, followed by fine-tuning on FHB (Synth2Real). As the amount of real data decreases, the benefit of pre-training increases. For both the object and the hand reconstruction, synthetic pre-training is critical in low-data regimes.

Text-to-Video Retrieval

Pretraining on millions of images & videos
Finetuning on MSRVTT with 9K training videos



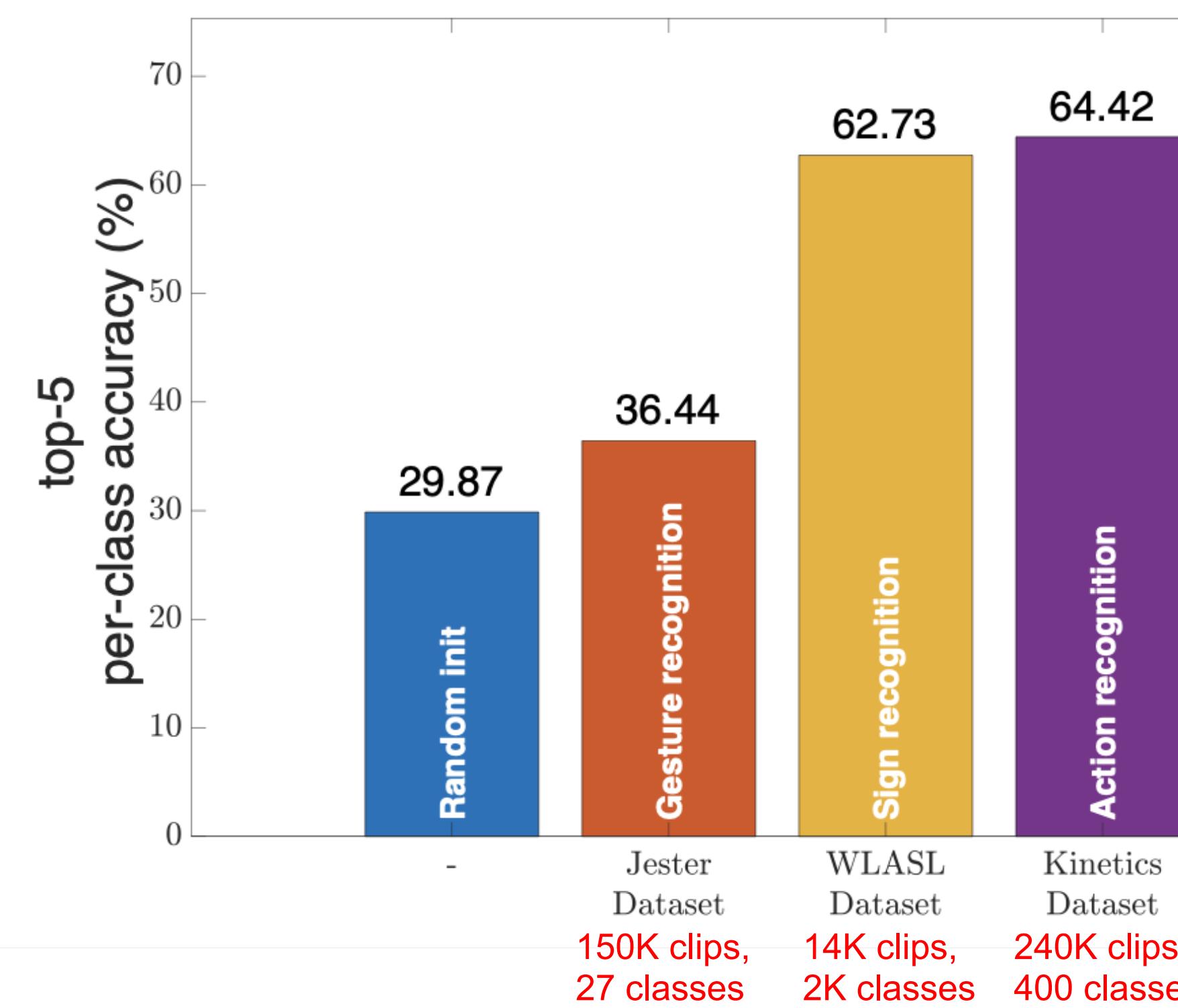
1. A man and a woman performing a musical.
2. A teenage couple perform in an amateur musical
3. Dancers are playing a routine.
4. People are dancing in a musical.
5. Some people are acting and singing for performance.

	Pre-training (for 1 epoch)	#pairs	$\uparrow\mathbf{R@1}$	$\uparrow\mathbf{R@10}$	$\downarrow\mathbf{MedR}$
	-	-	5.6	22.3	55
Noisy	ImageNet		15.2	54.4	9.0
	HowTo-17M subset	17.1M	24.1	63.9	5.0
	CC3M	3.0M	24.5	62.7	5.0
	WebVid2M	2.5M	26.0	64.9	5.0
	+ CC3M + WebVid2M	5.5M	27.3	68.1	4.0

Sign Language Recognition

Pretraining on various tasks on different datasets

Finetuning on 50K videos from BSL-1K sign language dataset



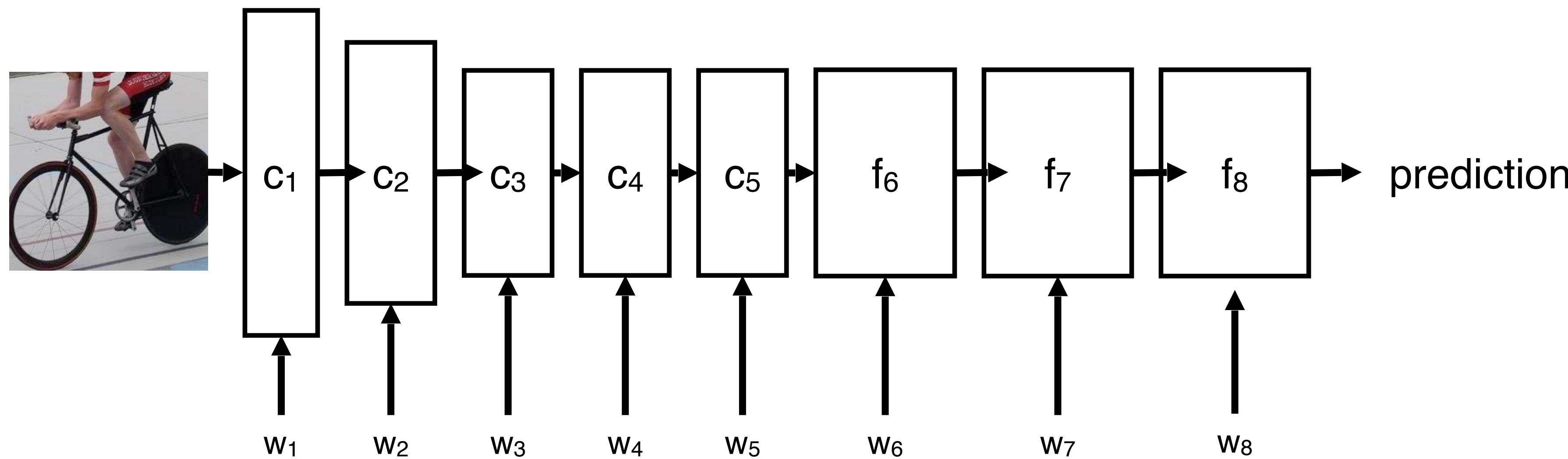
Pretraining Summary

- Common practice: Pretrain on large data, finetune on small data.
 - Remove the last class-specific layer (e.g. 1000 categories)
 - Add new layer(s) for the new task randomly initialized
 - Either “freeze” the pretrained parameters and train a simple classifier on top,
 - Or train “end-to-end” all parameters.
- Avoids overfitting
- Shortens training time
- Lots of pretrained models available online
- Task and domain-relevant pretraining is usually better

Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

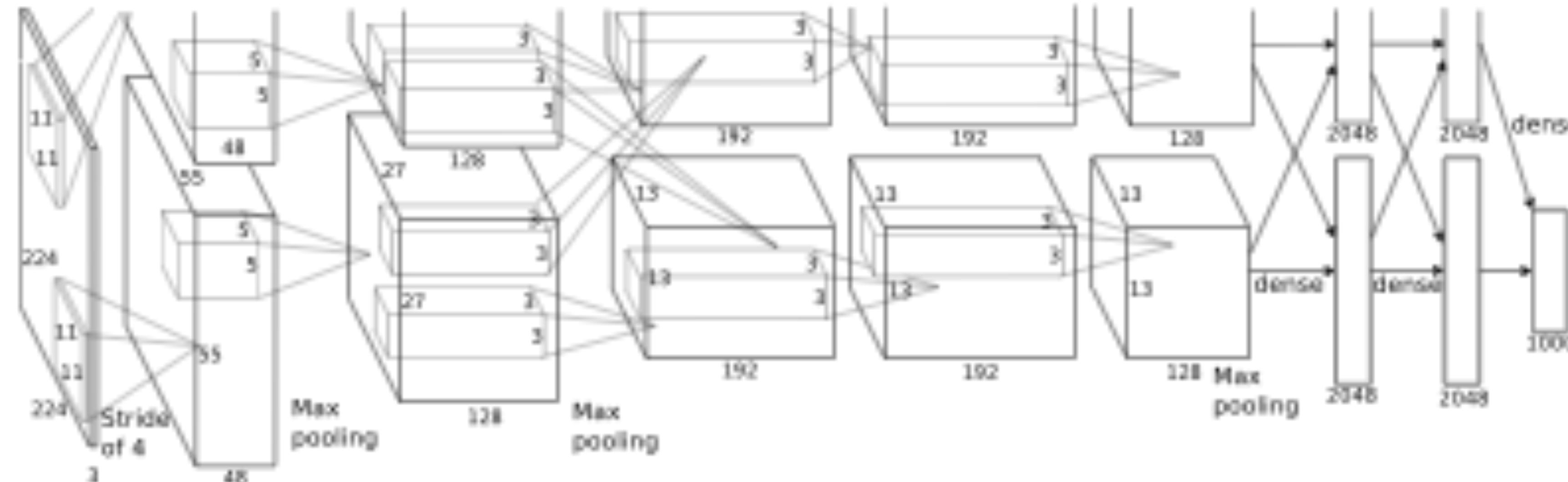
A CNN for image classification



Recall: the goal of this model is to map an input image to a class prediction.

Recall: The AlexNet model

A breakthrough in image understanding



[AlexNet by Krizhevsky et al. 2012]

Each large block represents a data tensor

Each smaller block represents a filter

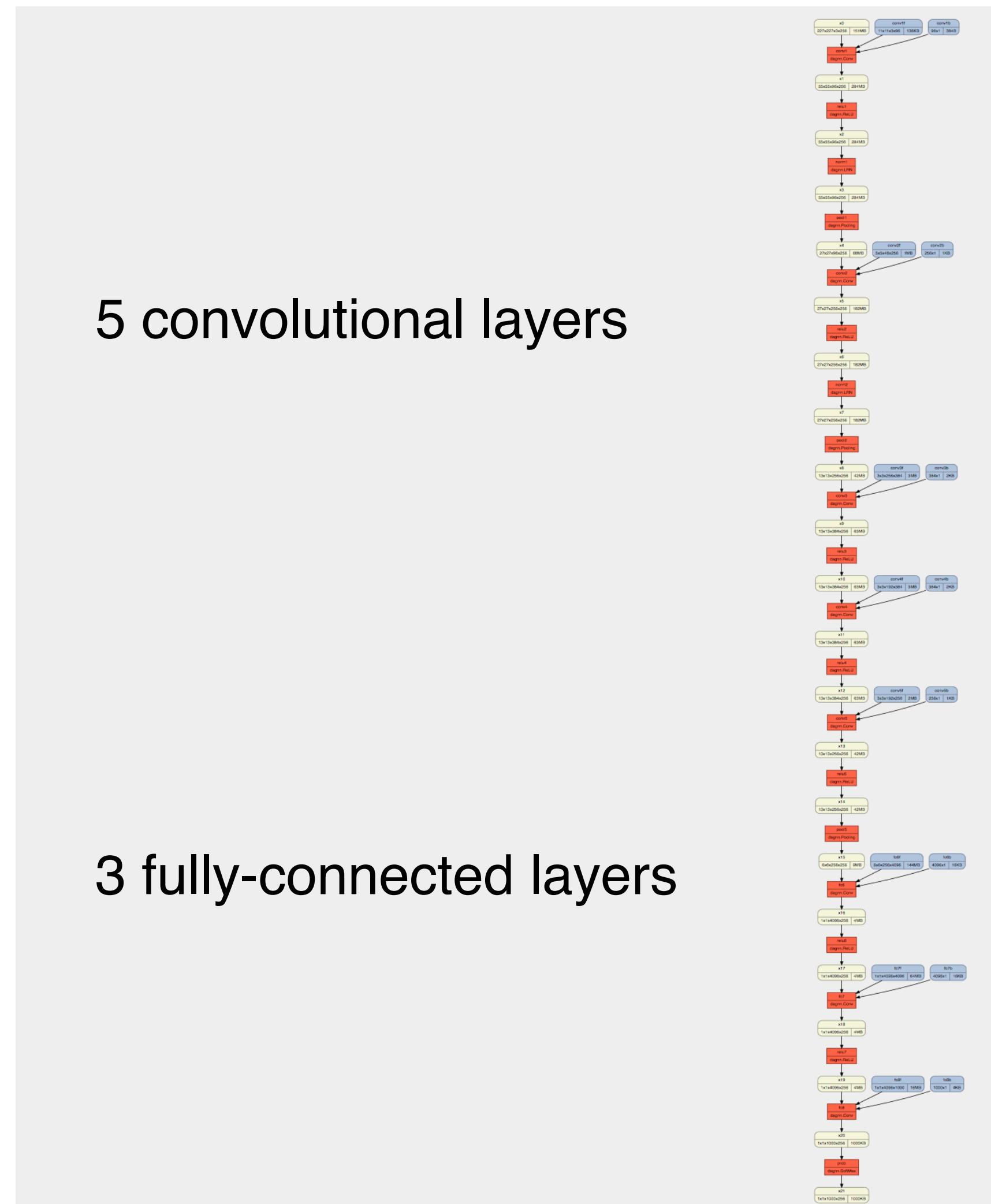
The filter size and stride are shown

The number of filters can be deduced from the number of feature channels

There are two parallel streams in this network (for efficiency reasons)

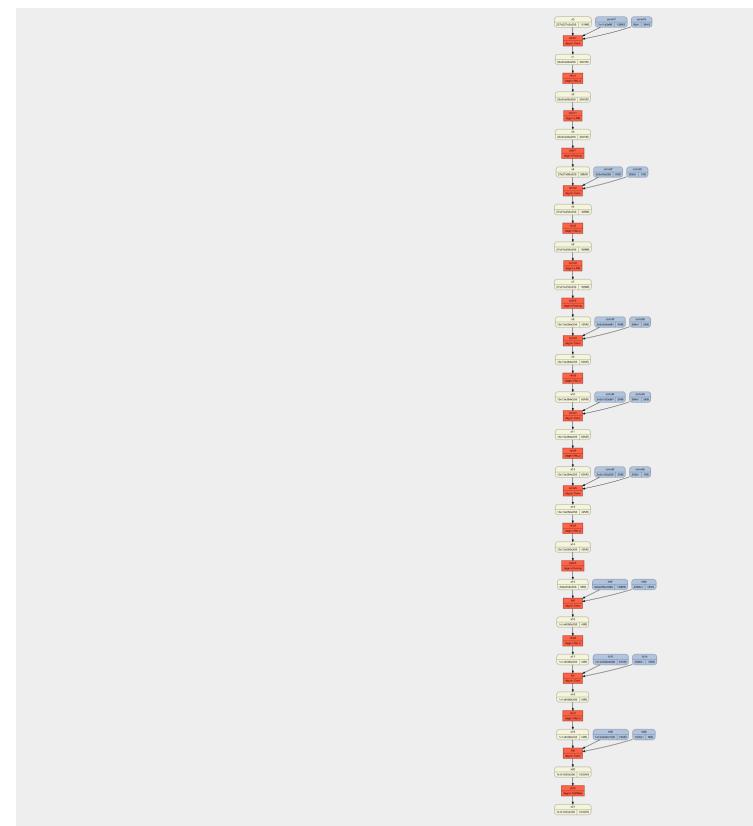
How deep is deep enough?

AlexNet (2012)

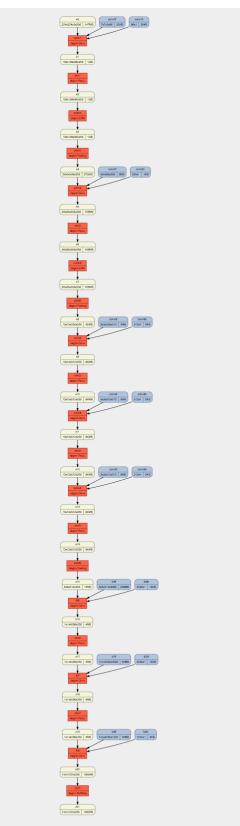


How deep is deep enough?

AlexNet (2012)



VGG-M (2013)

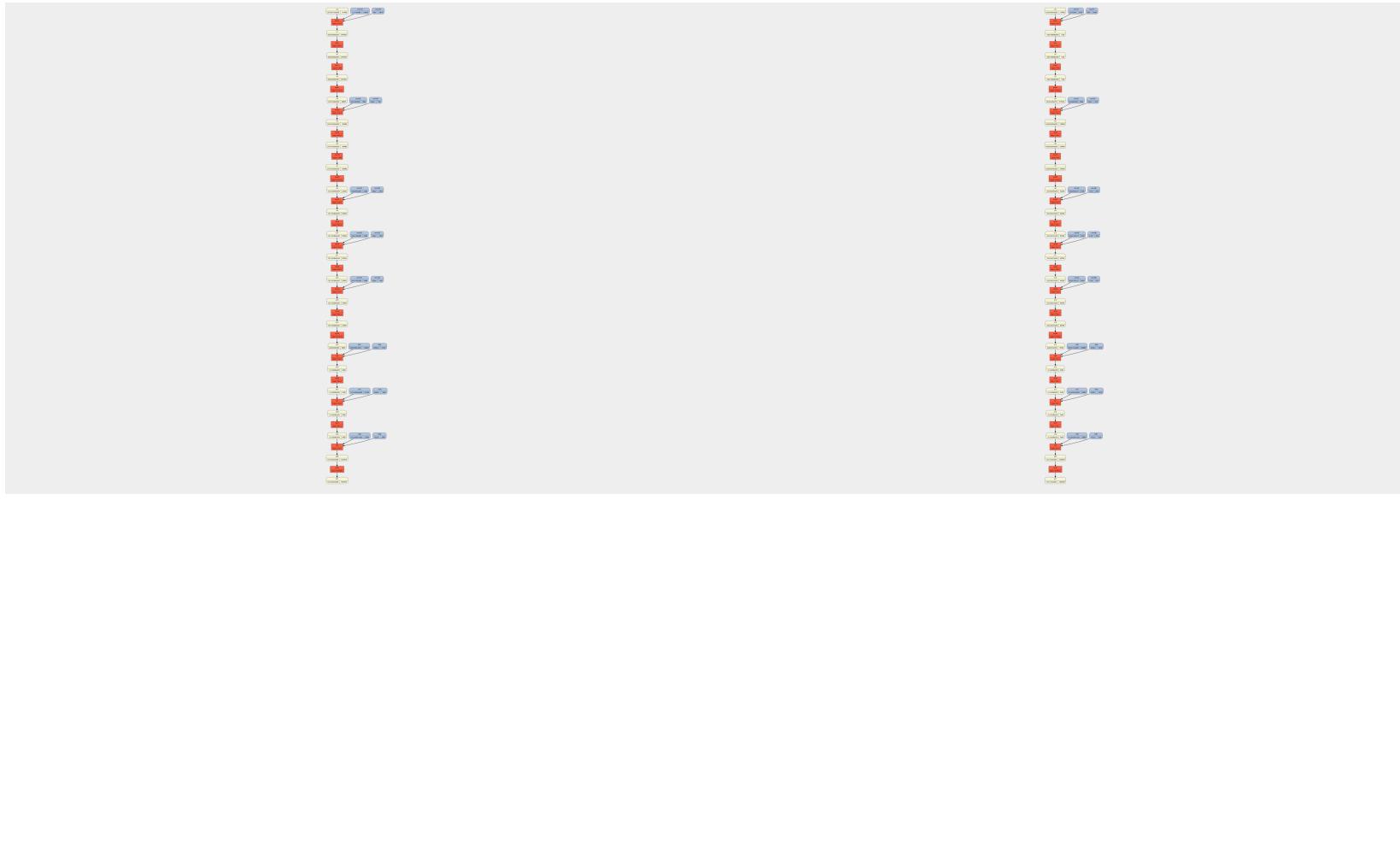


VGG-VD-16 (2014)

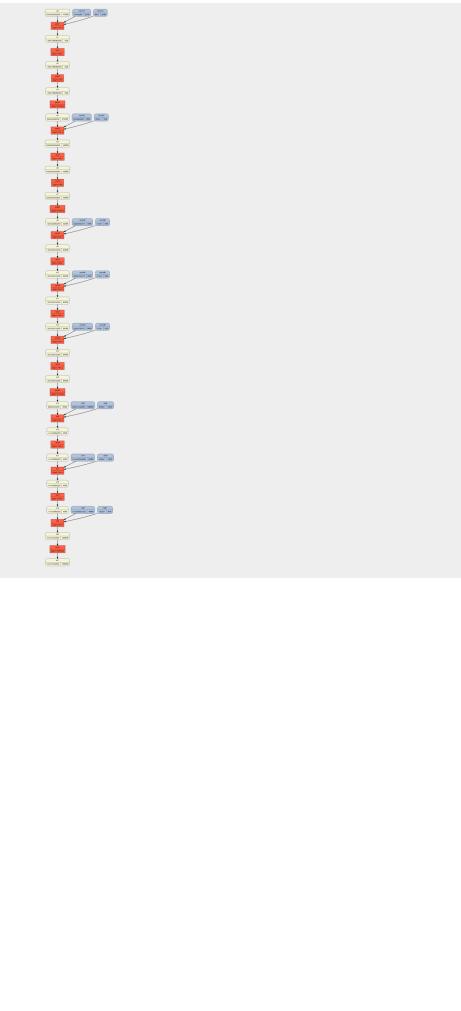


How deep is deep enough?

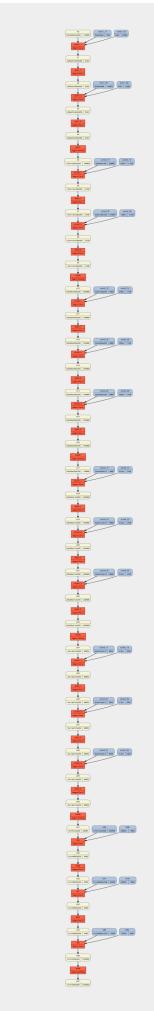
AlexNet (2012)



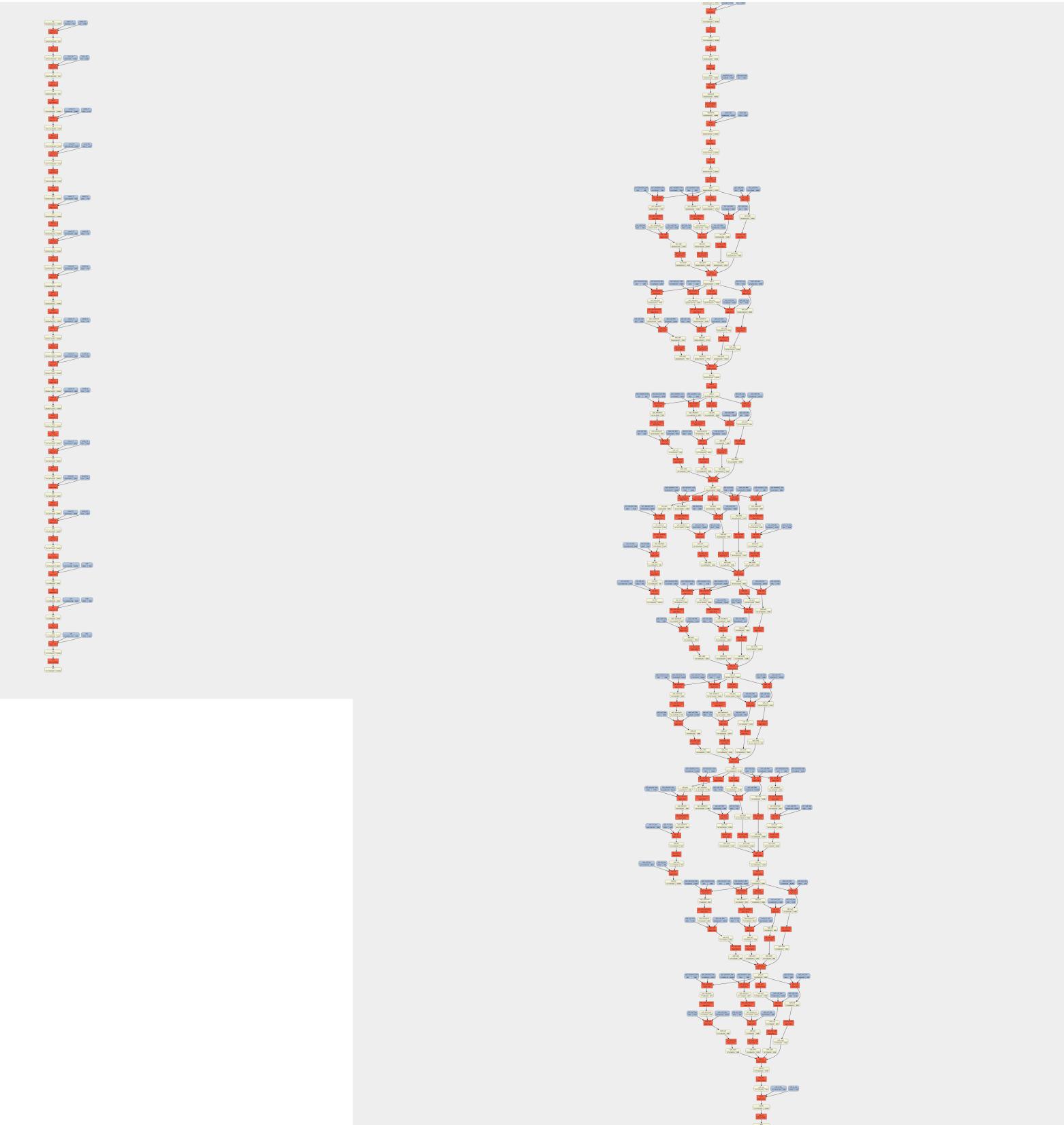
VGG-M (2013)



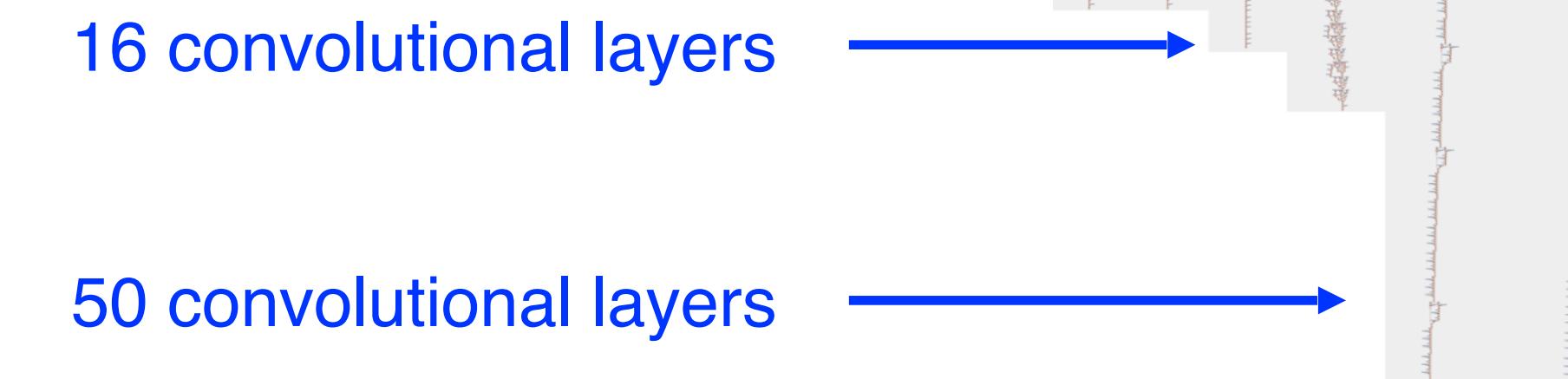
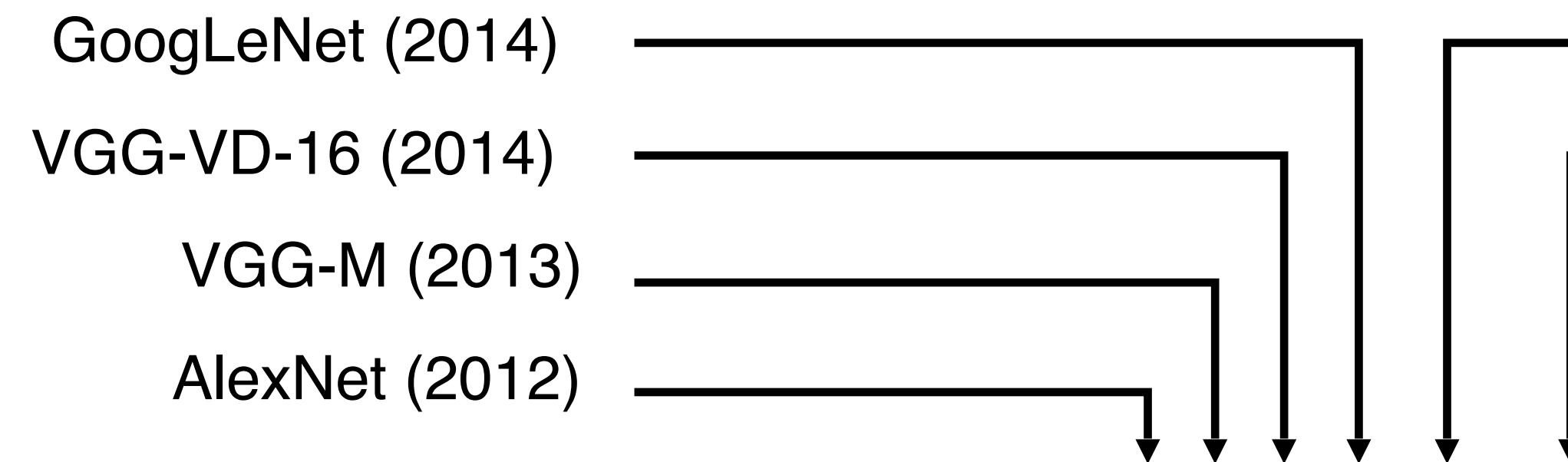
VGG-VD-16 (2014)



GoogLeNet (2014)



How deep is deep enough?



ResNet 50 (2015)

ResNet 152 (2015)

Krizhevsky, I. Sutskever, and G. E. Hinton.
ImageNet classification with deep convolutional neural networks. In Proc. NIPS, 2012.

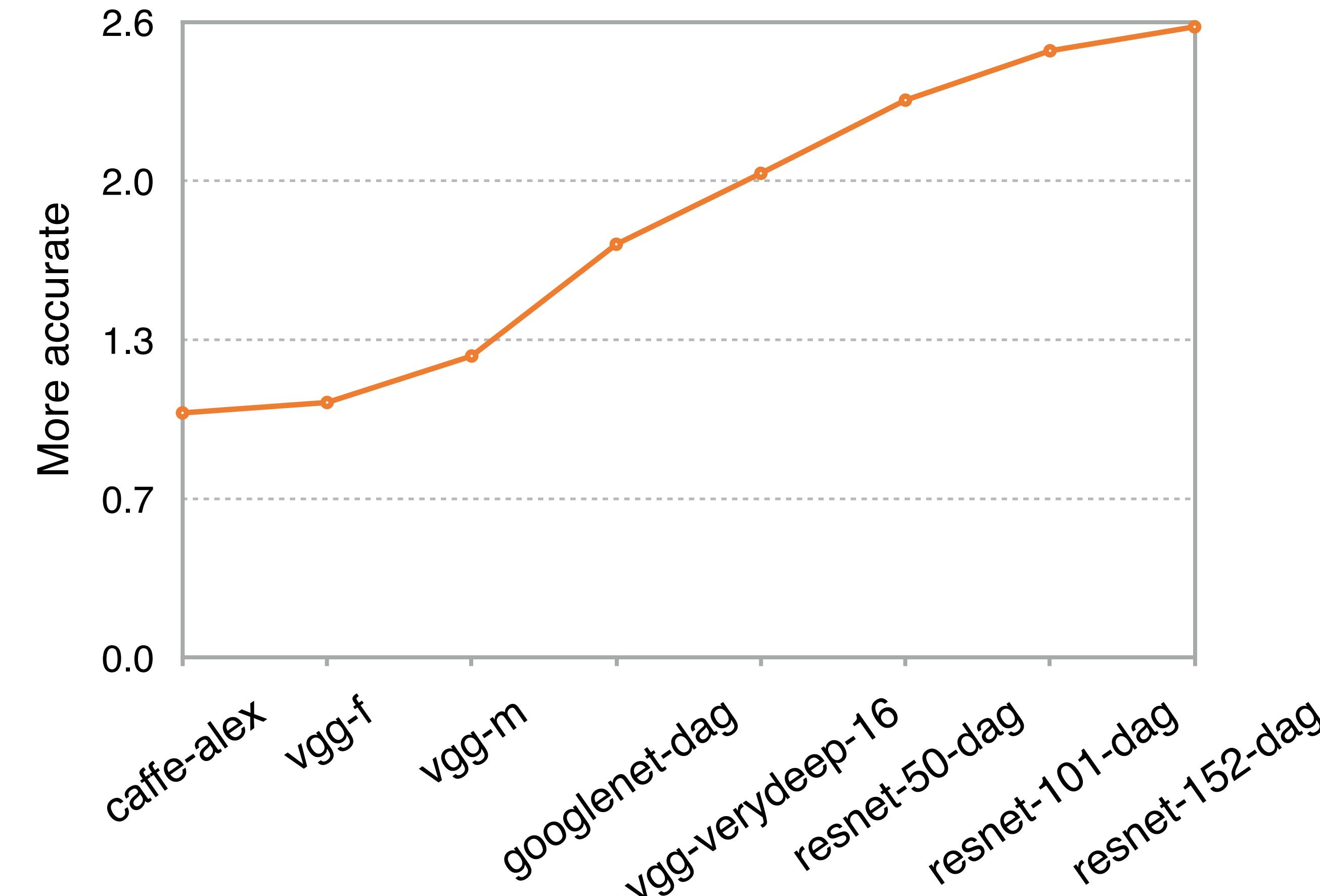
C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. *Going deeper with convolutions*. In Proc. CVPR, 2015.

K. Simonyan and A. Zisserman. *Very deep convolutional networks for large-scale image recognition*. In Proc. ICLR, 2015.

K. He, X. Zhang, S. Ren, and J. Sun. *Deep residual learning for image recognition*. In Proc. CVPR, 2016.

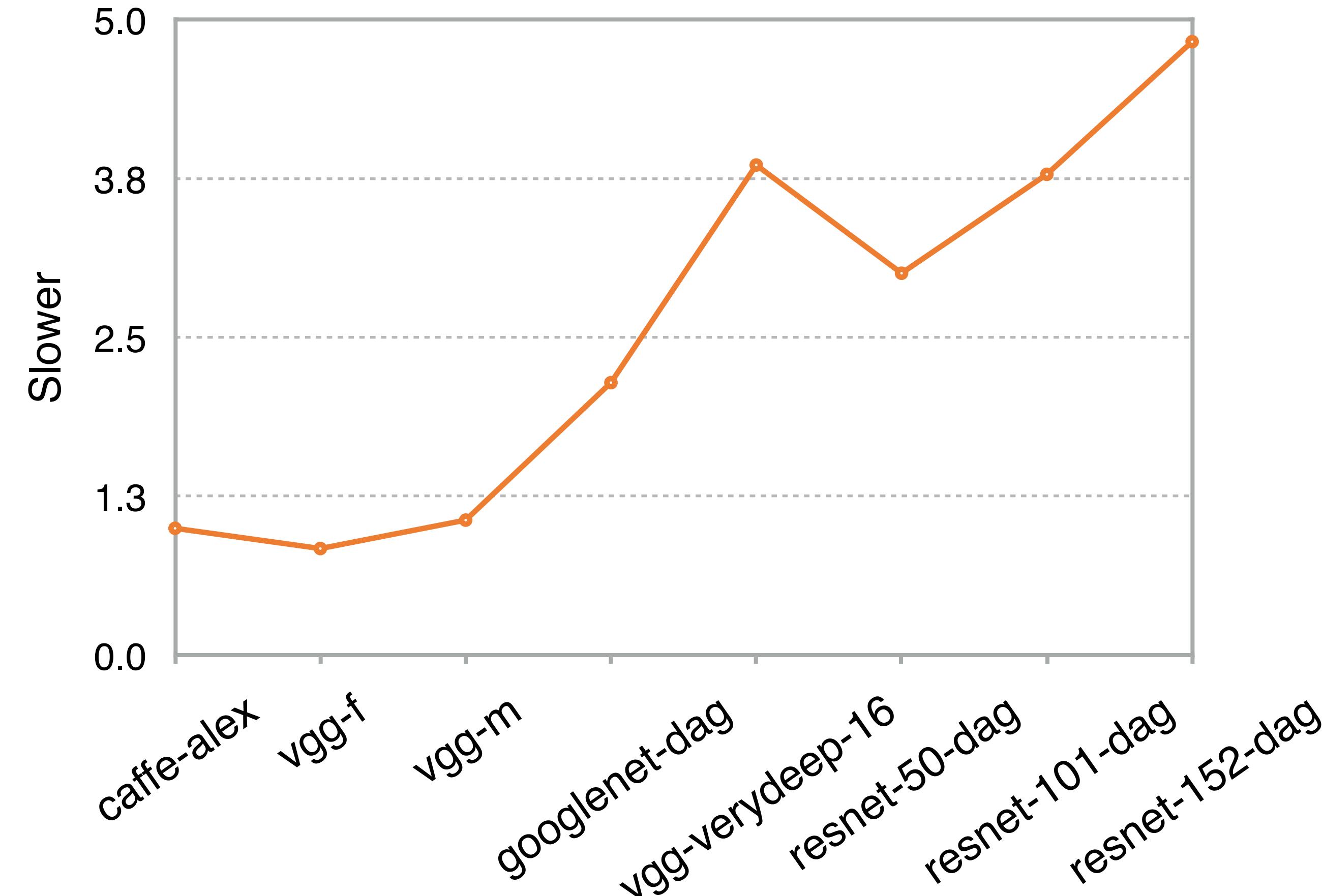
Accuracy

3 × more accurate in 3 years



Speed

5 × slower



Remark: 101 ResNet layers same size/speed as 16 VGG-VD layers

Reason: far fewer feature channels (quadratic speed/space gain)

Moral: optimize your architecture

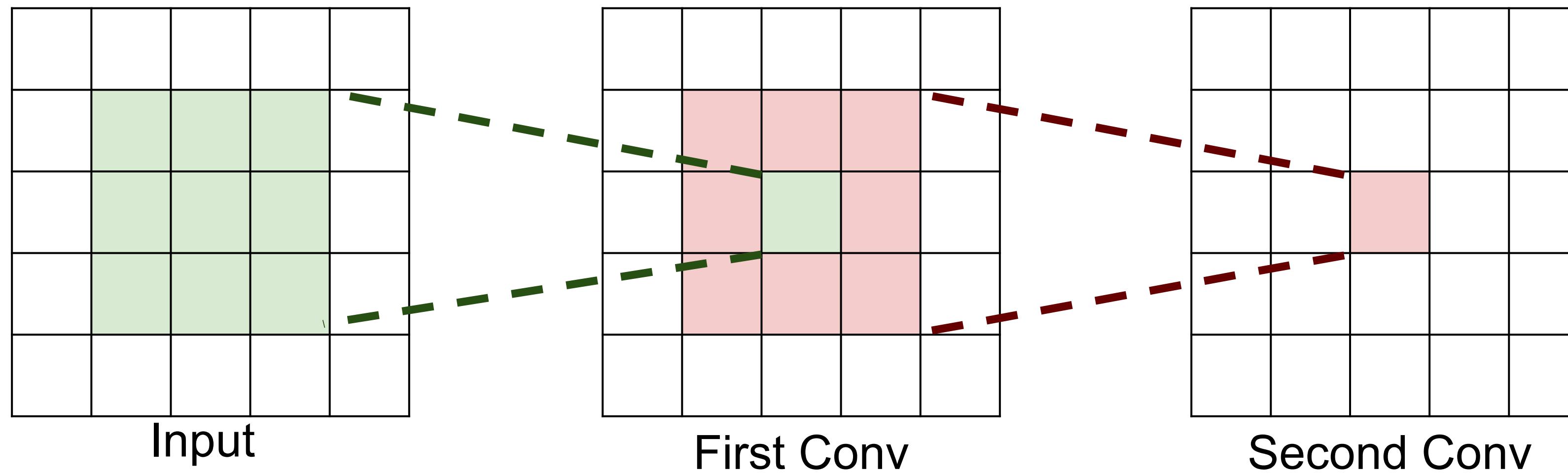
CNN architectures – notes and details

- Increased depth of recent architectures
- Number of parameters matter (How to count parameters?)
- Power of small filters, e.g. 3x3 convolutions
- ResNet architecture

The power of small filters

Suppose we stack two CONV layers with receptive field size 3×3

Q: What region of input does each neuron in 2nd CONV see?

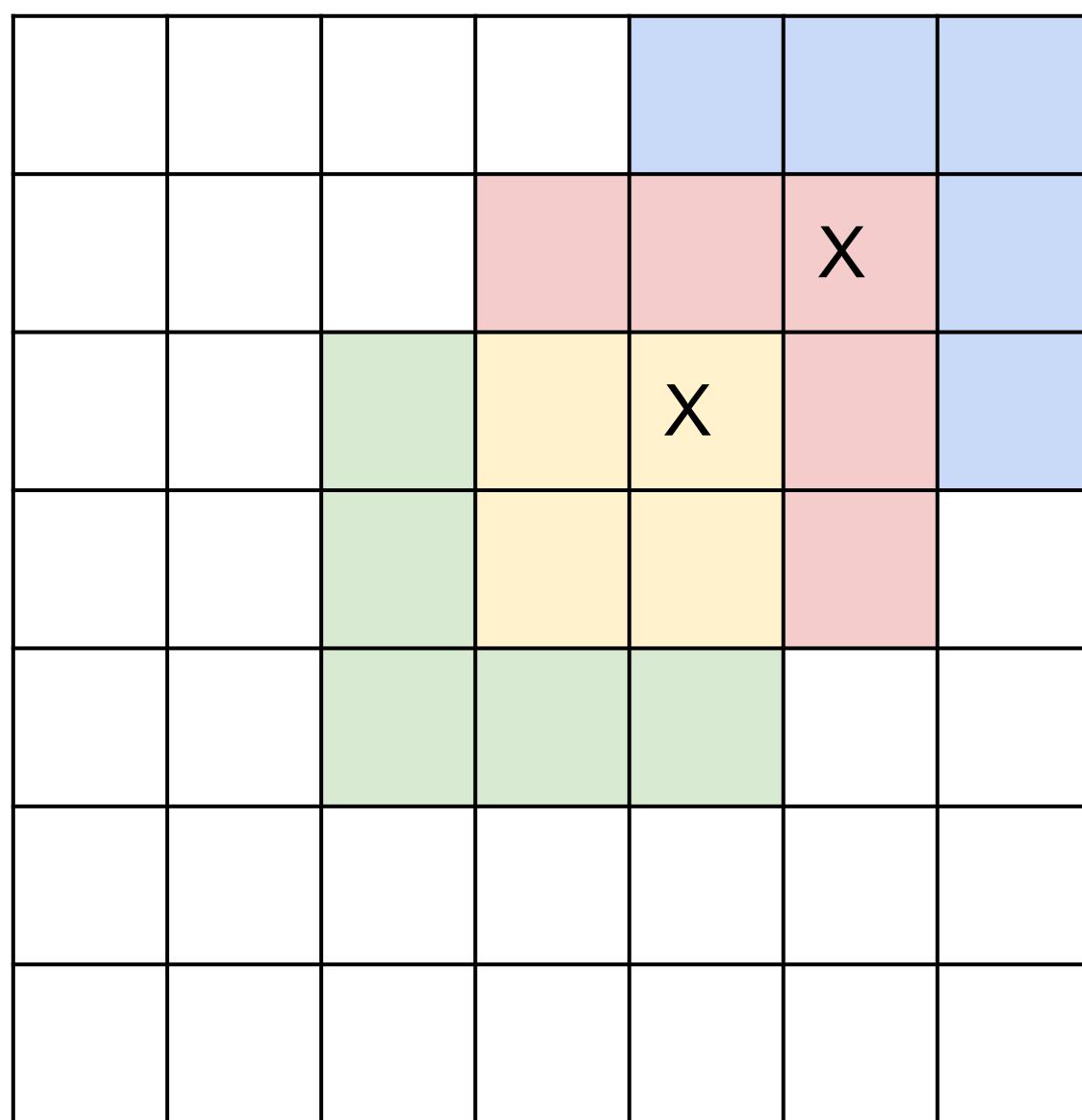


Answer: [5x5]

The power of small filters

Suppose we stack three CONV layers with receptive field size 3×3

Q: What region of input does each neuron in 3rd CONV see?



Answer: [7x7]

The power of small filters

Suppose input has depth C & we want output depth C as well.

1x CONV with 7x7 filters

Number of weights:

$$\begin{aligned} & C * (7 * 7 * C) \\ & = 49 C^2 \end{aligned}$$

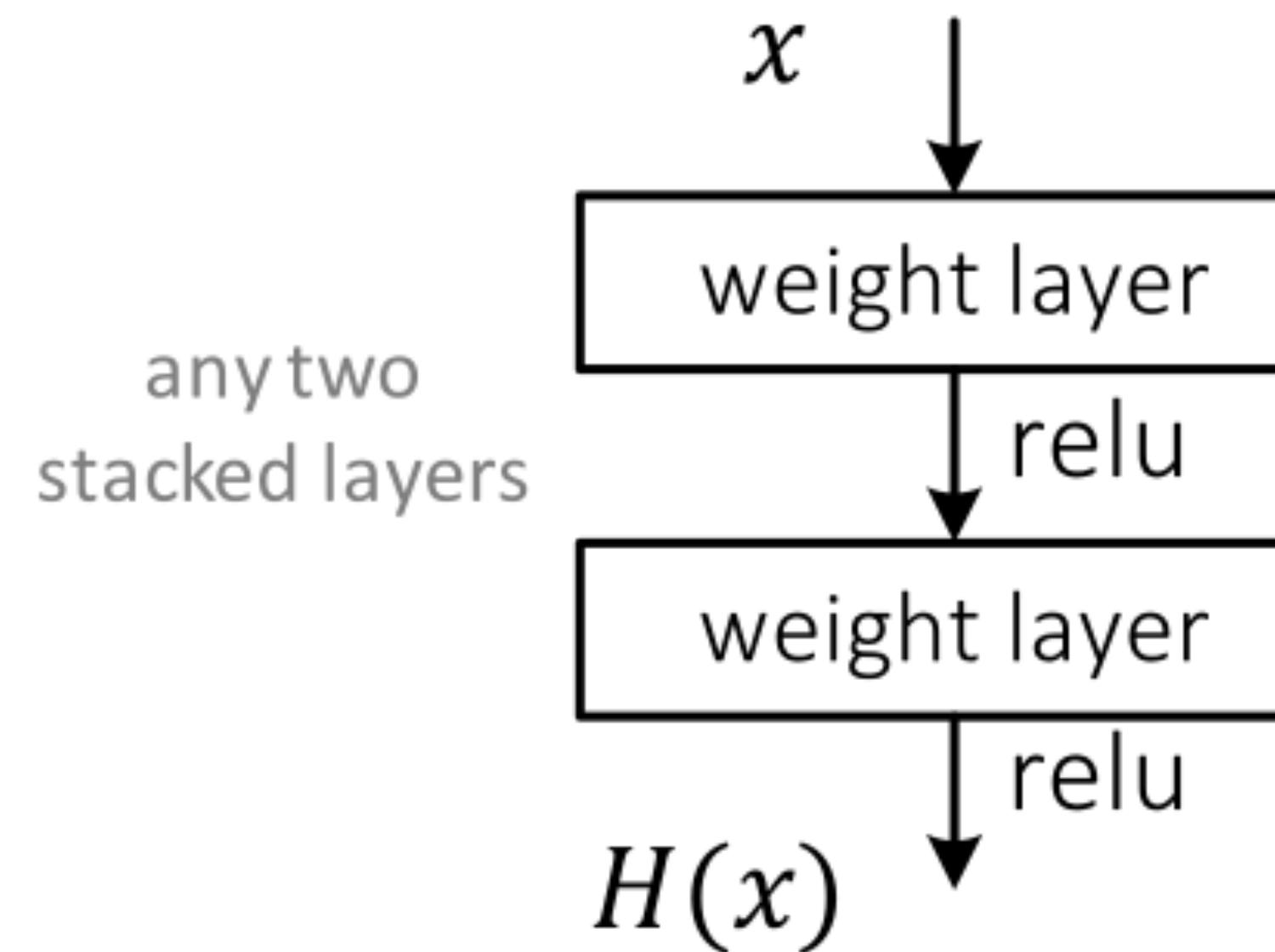
3x CONV with 3x3 filters

Number of weights:

$$\begin{aligned} & C * (3 * 3 * C) + C * (3 * 3 * C) + C * (3 * 3 * C) \\ & = 3 * 9 * C^2 \\ & = 27 C^2 \end{aligned}$$

Residual networks [ResNets]

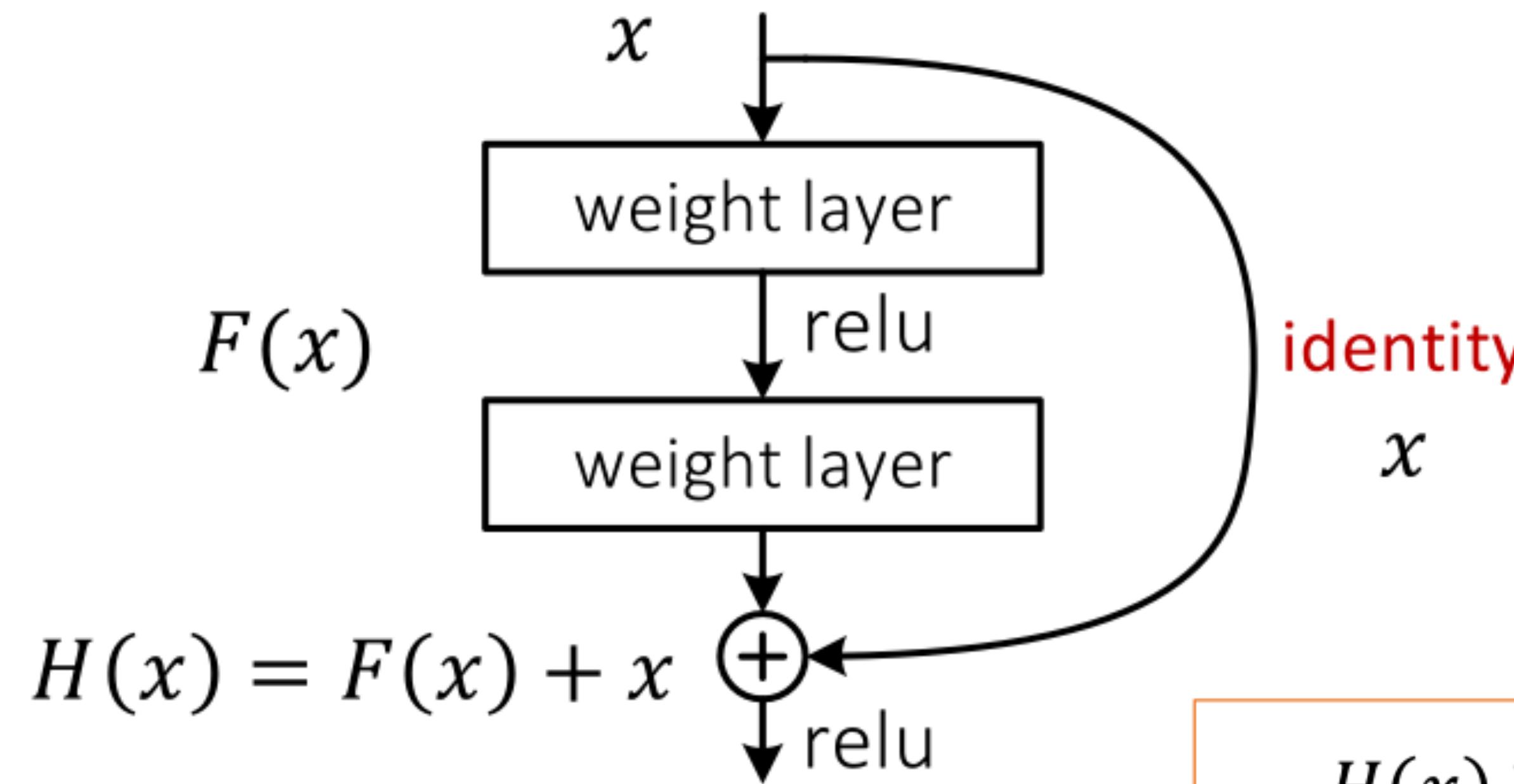
Plain net



$H(x)$ is any desired mapping,
hope the 2 weight layers fit $H(x)$

Residual networks [ResNets]

Residual net

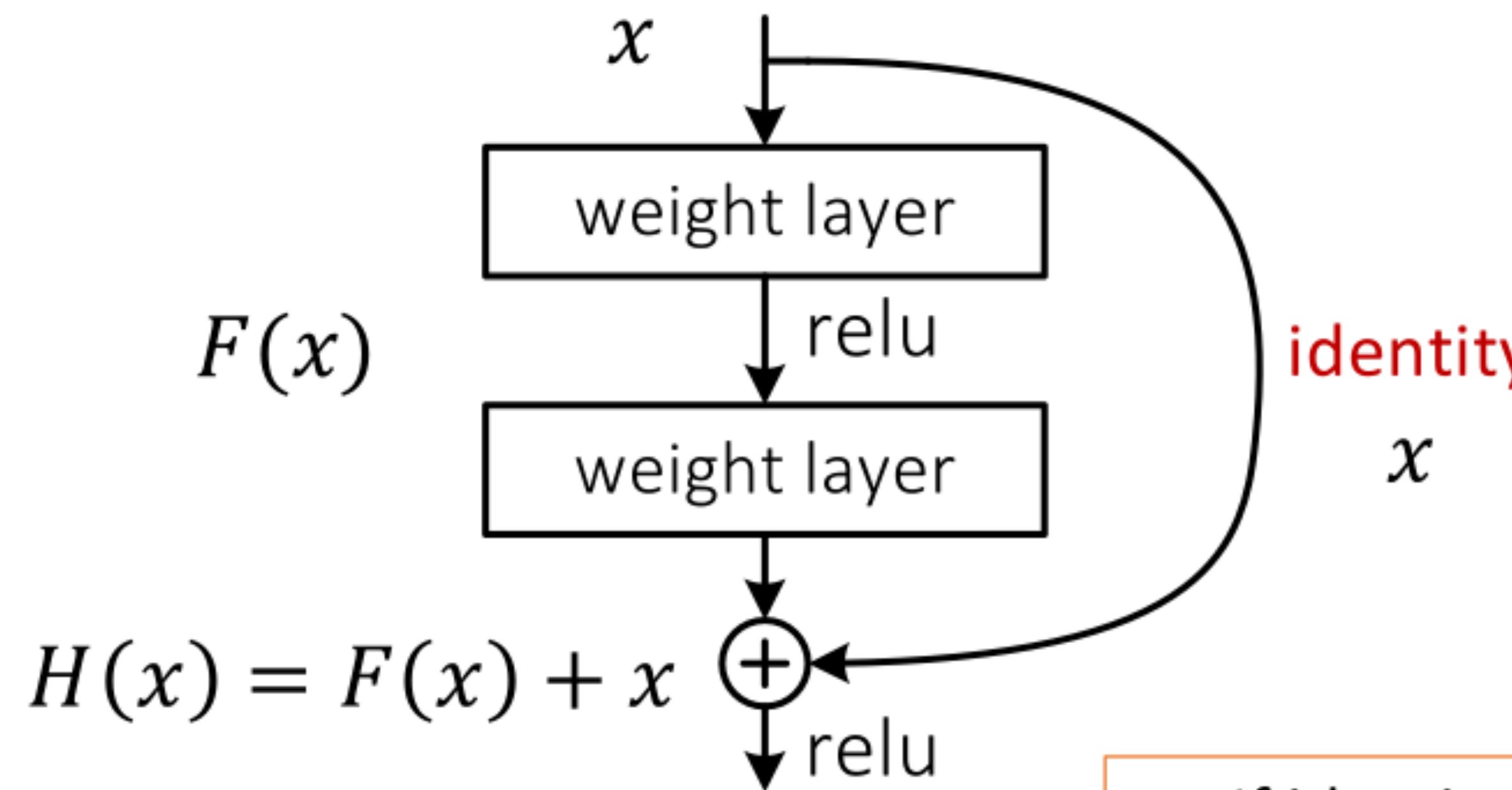


identity
 x

$H(x)$ is any desired mapping,
hope the 2 weight layers fit $H(x)$
hope the 2 weight layers fit $F(x)$
let $H(x) = F(x) + x$

Residual networks [ResNets]

$F(x)$ is a **residual** mapping w.r.t. **identity**

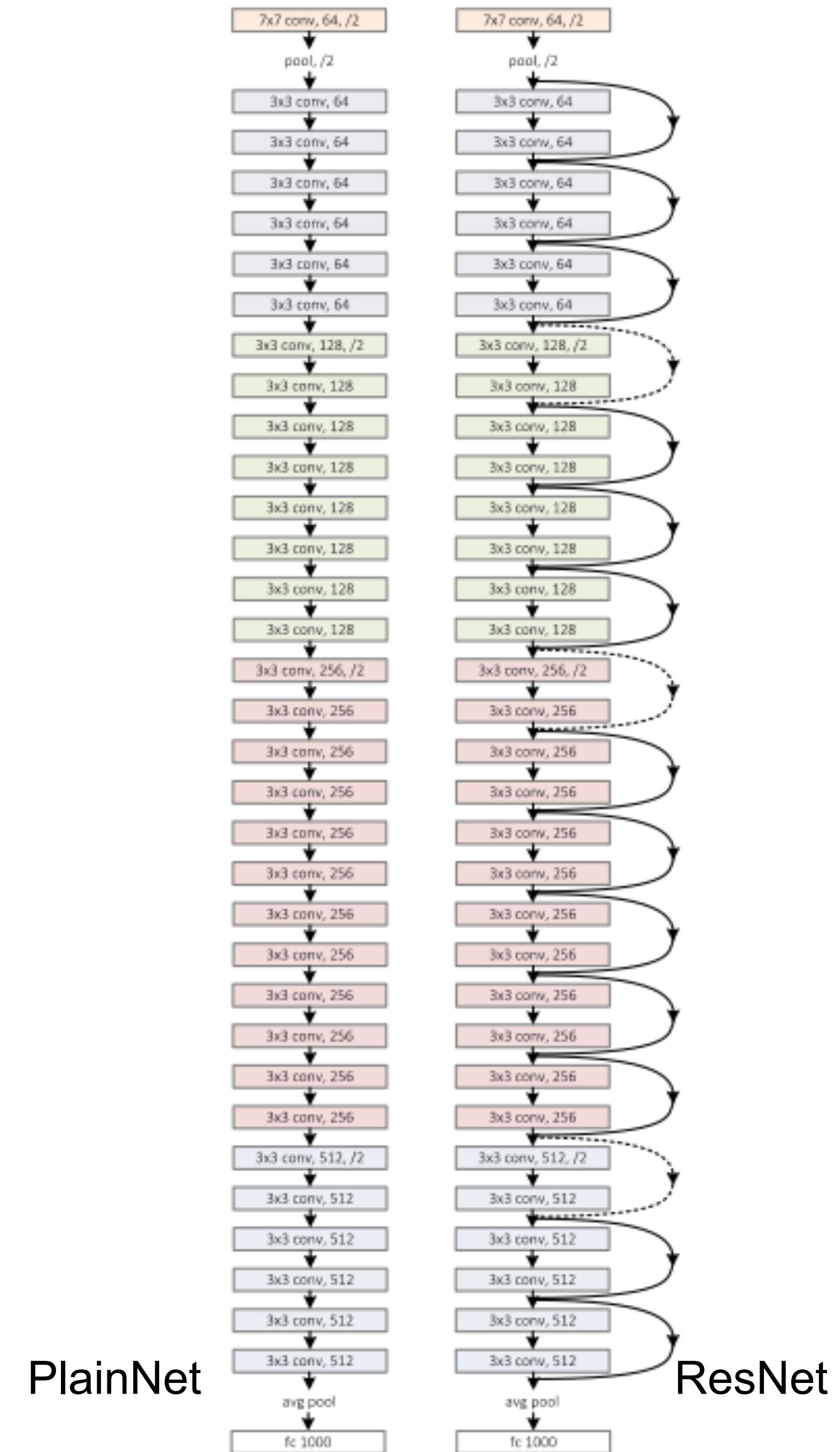


- If identity were optimal, easy to set weights as 0
- If optimal mapping is closer to identity, easier to find small fluctuations

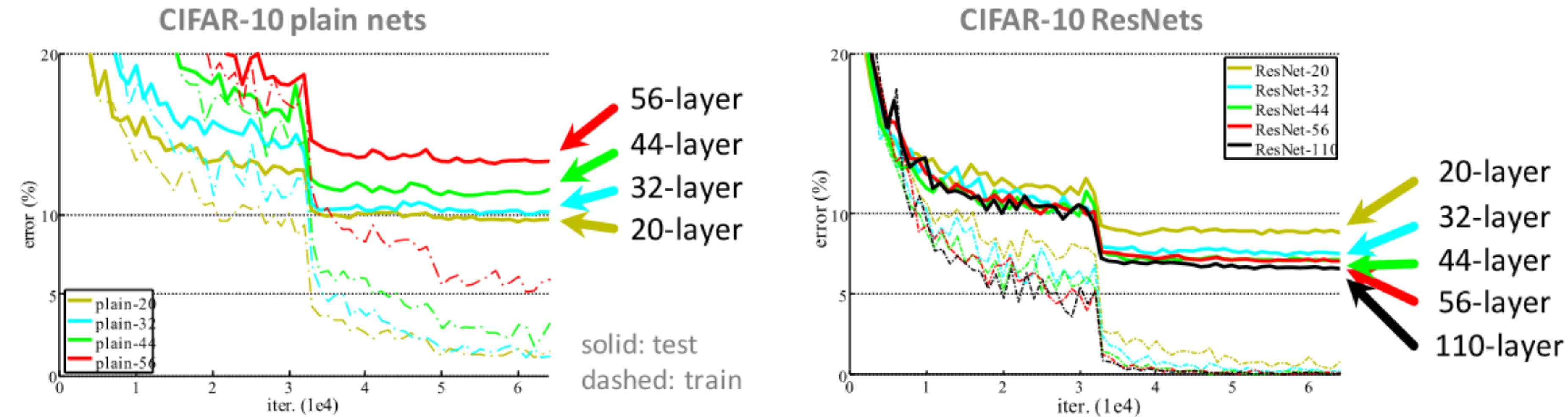
Network design

Basic design (vgg-style)

- almost all 3x3 conv
- Spatial size /2 => # filters x2
- Simple design, just deep
- No fully connected layers
- No dropout



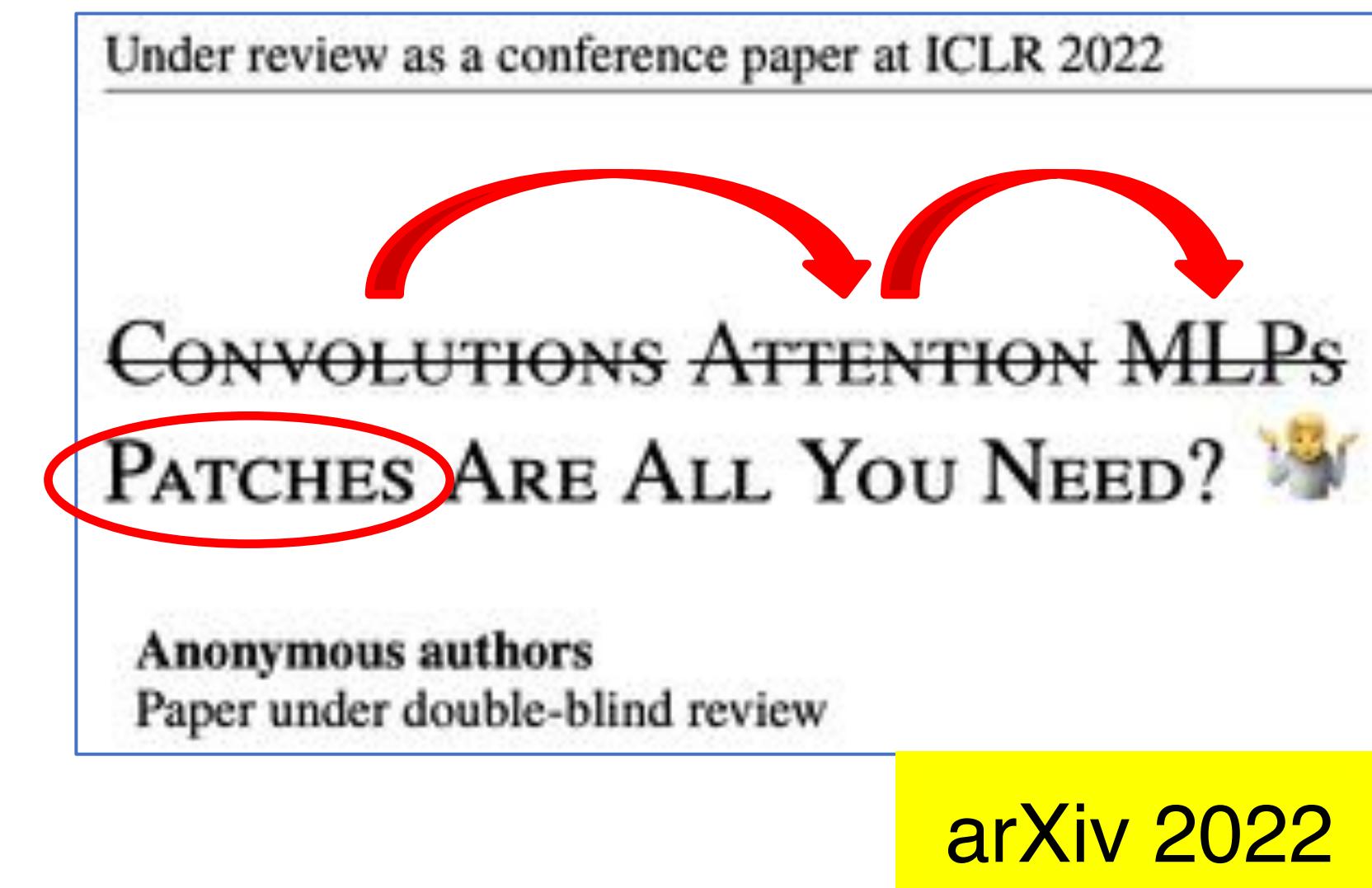
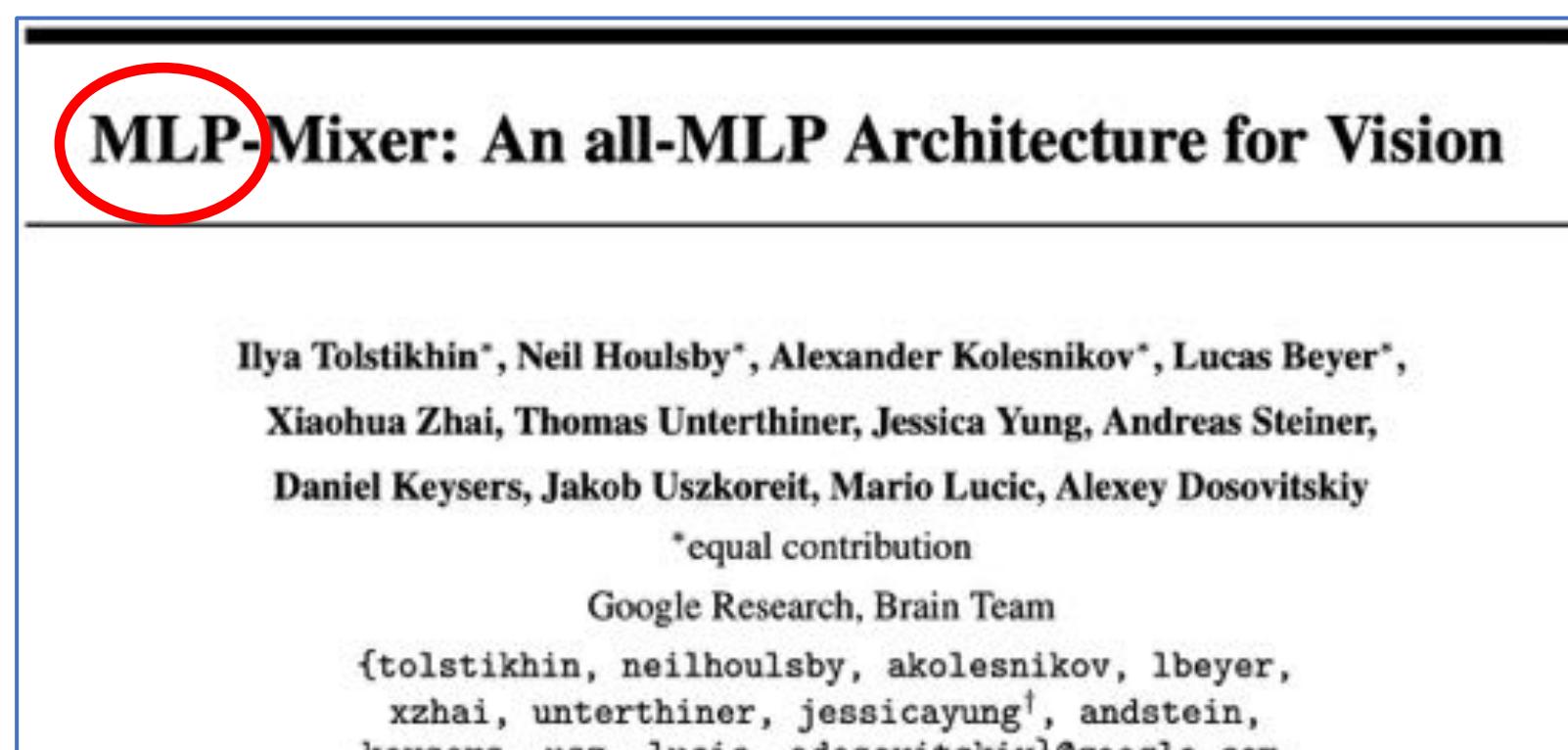
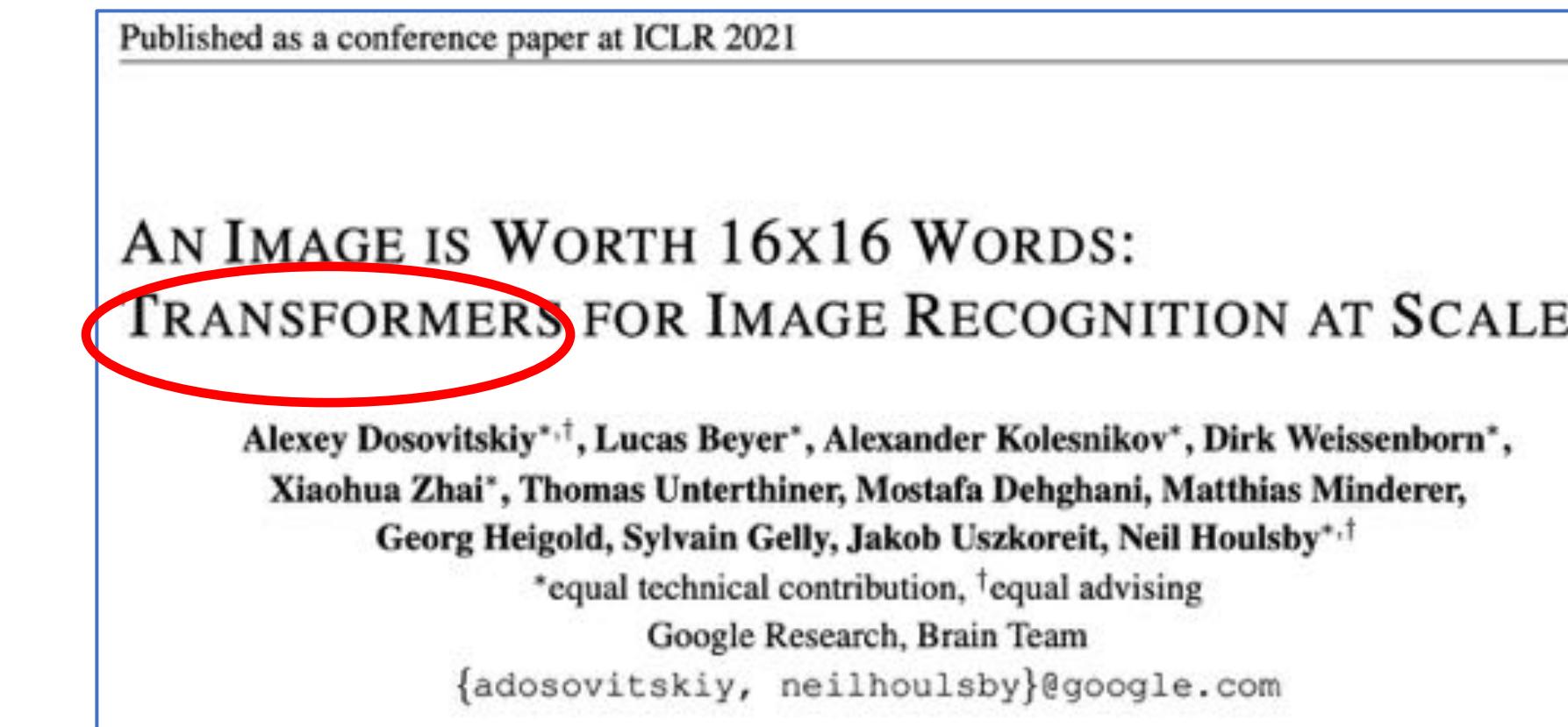
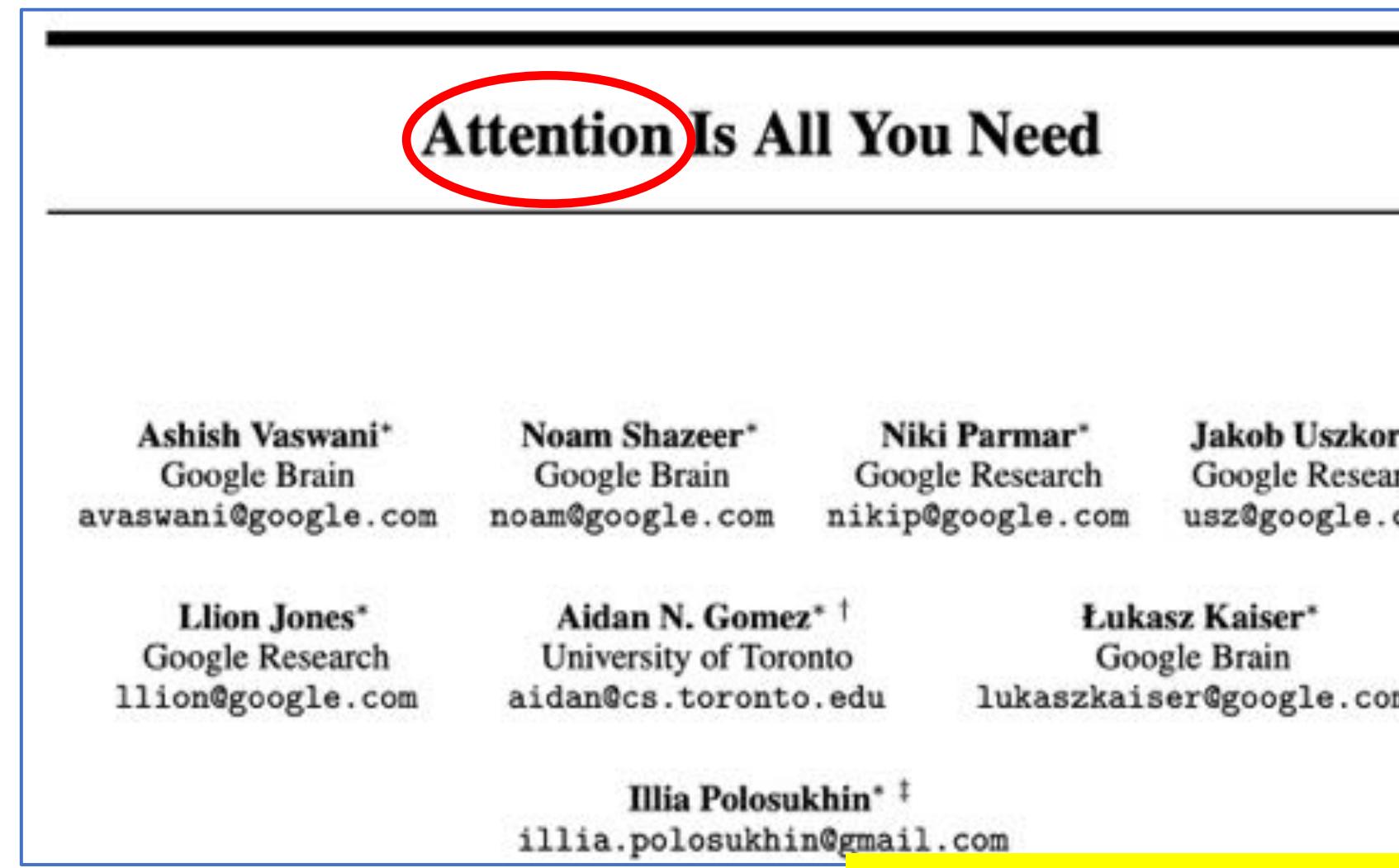
CIFAR-10 experiments



- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error** and **lower test error**

Do we need convolutions?

[next week]
Beyond CNNs



Agenda

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**

The field makes progress

Beyond

Classification

Computer vision tasks



Extracting meaning from
visual signals *

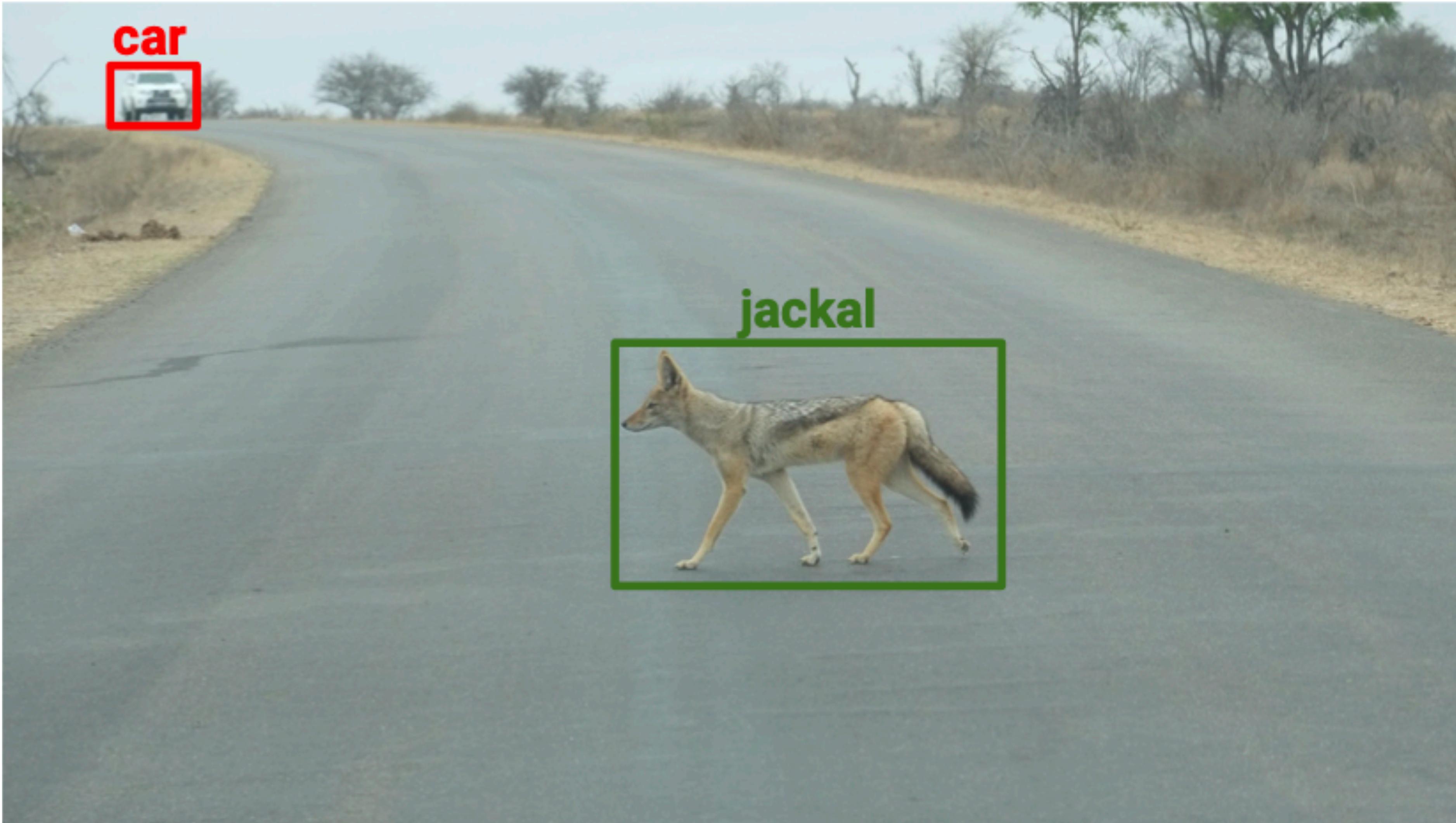
Object recognition,
Object detection,
Pixel-level segmentation,
3D localization,
etc.

*Visual signal: Image, video, depth, 3D point cloud, MRI, scans, ...

Example tasks



Object recognition and localization (detection)



Visual question answering



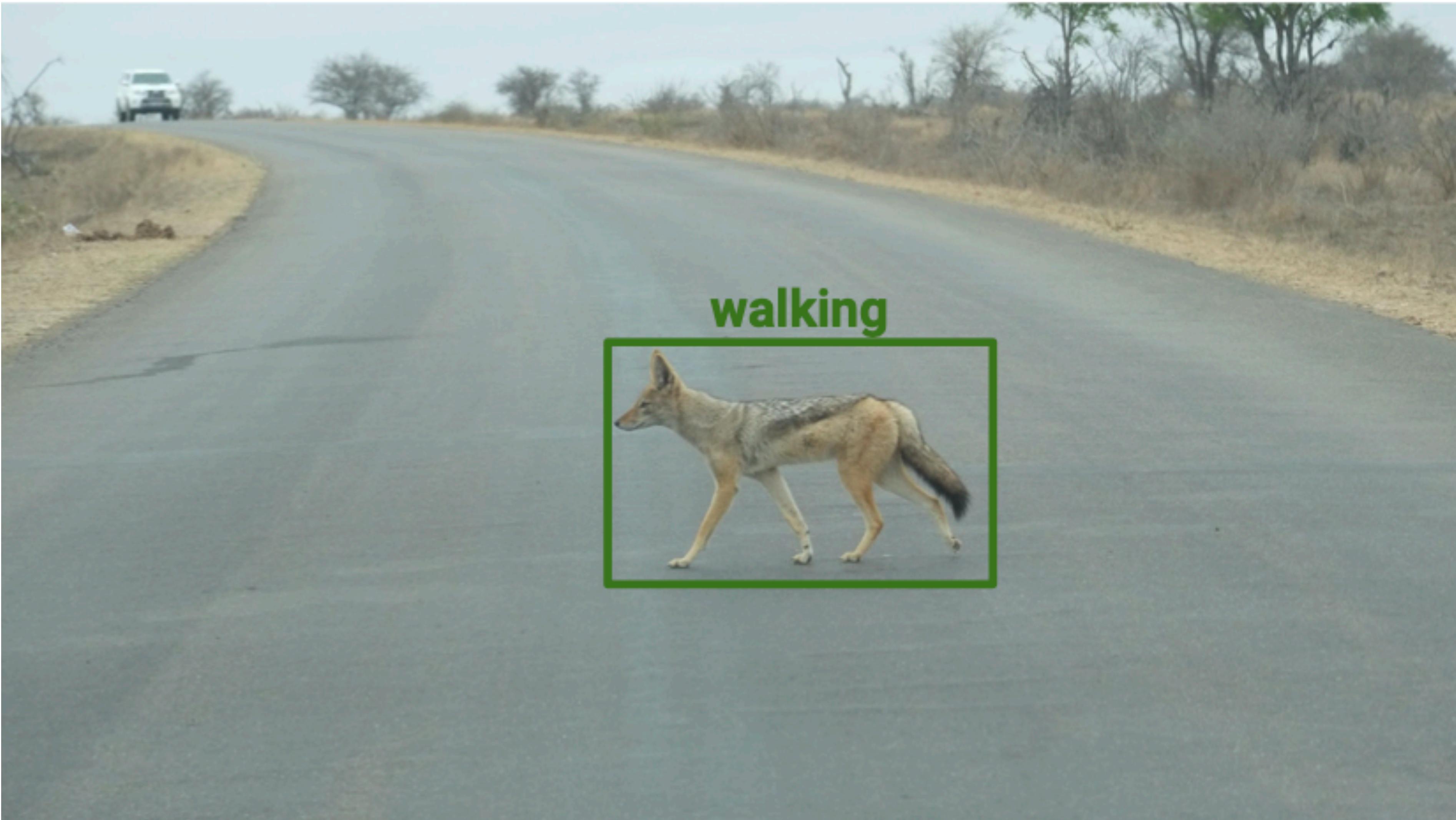
Q: Is this an outdoor scene?

A: Yes

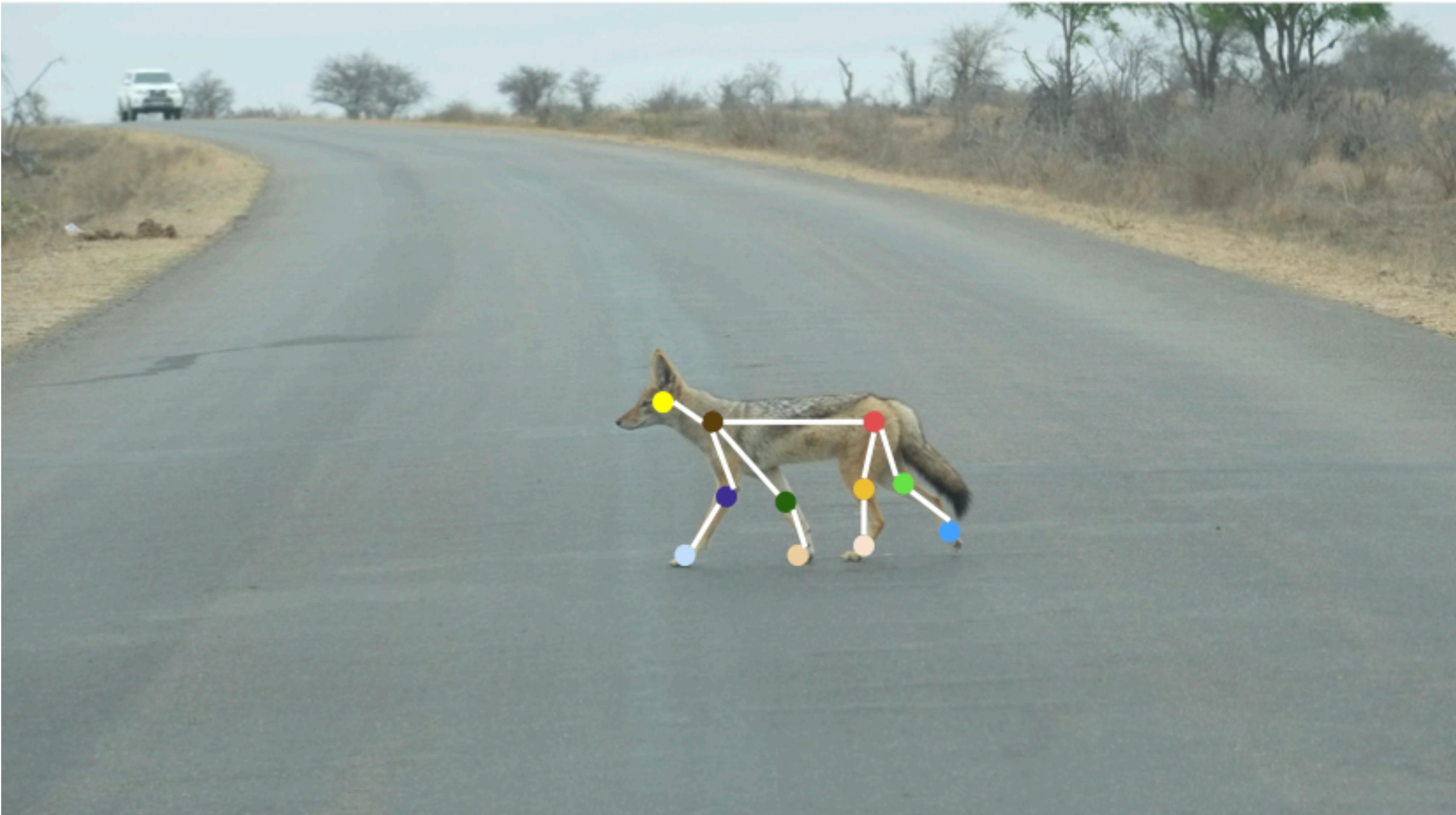
Q: What is the weather like?

A: Cloudy but dry

Activity recognition



Pose estimation



Captioning

A jackal walking across a rural asphalt road



Semantic segmentation



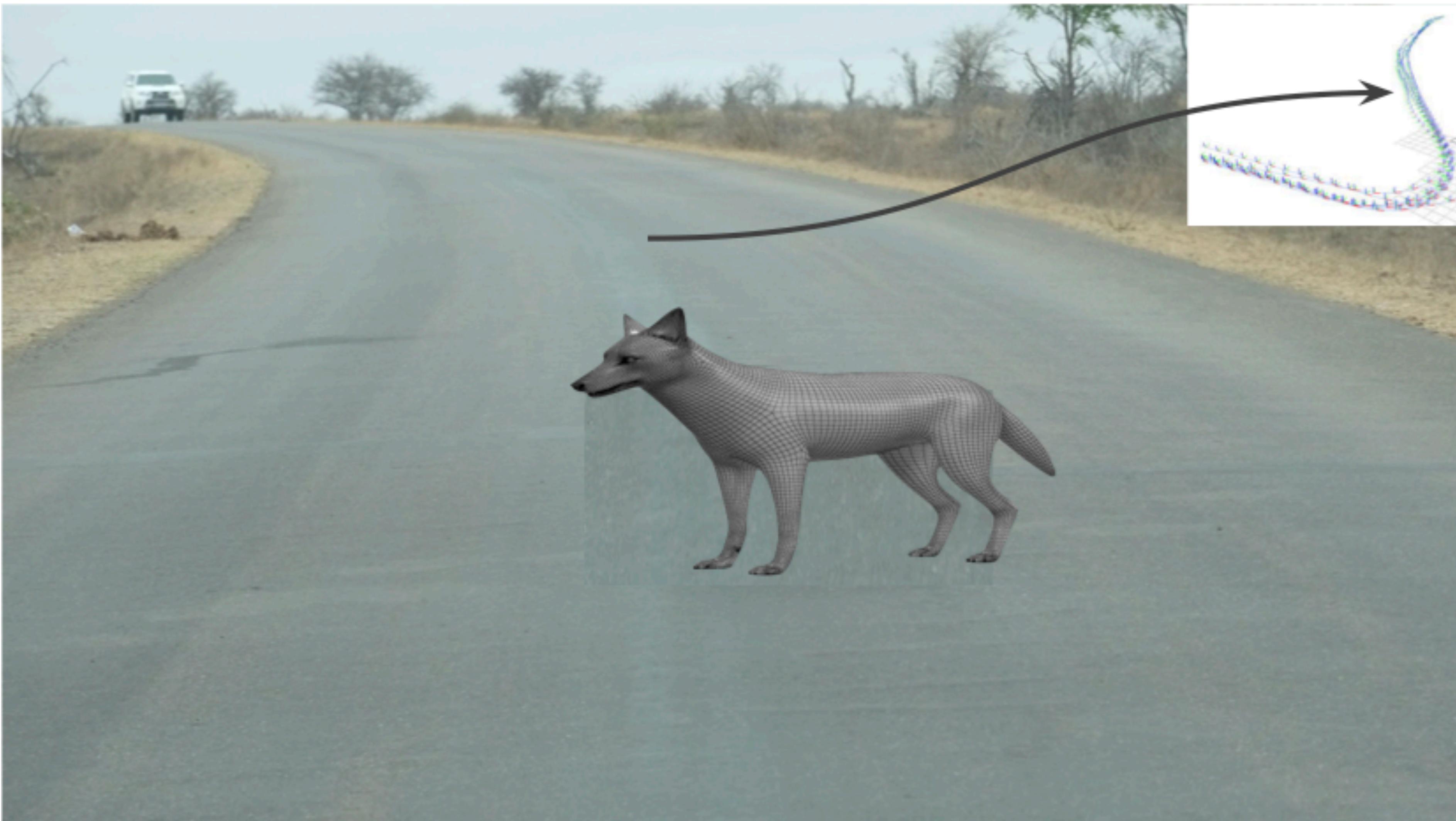
Depth estimation



3D shape estimation



Visual localization



Summary of today

- **1. Recap: Bag of Visual Words, Analogy with NNs**
- **2. Neural networks (NNs) for computer vision:**
 - Applications
 - A brief history: from perceptron to MLPs to CNNs
- **3. Convolutional neural networks (CNNs)**
 - Standard layers
 - Recap: Training NNs
 - Visualizing CNNs
 - Pretraining & finetuning NNs
 - Typical CNN architectures
- **4. Beyond classification - preview**