

Introduction AI

Raoul Grouls, 28-02-2024

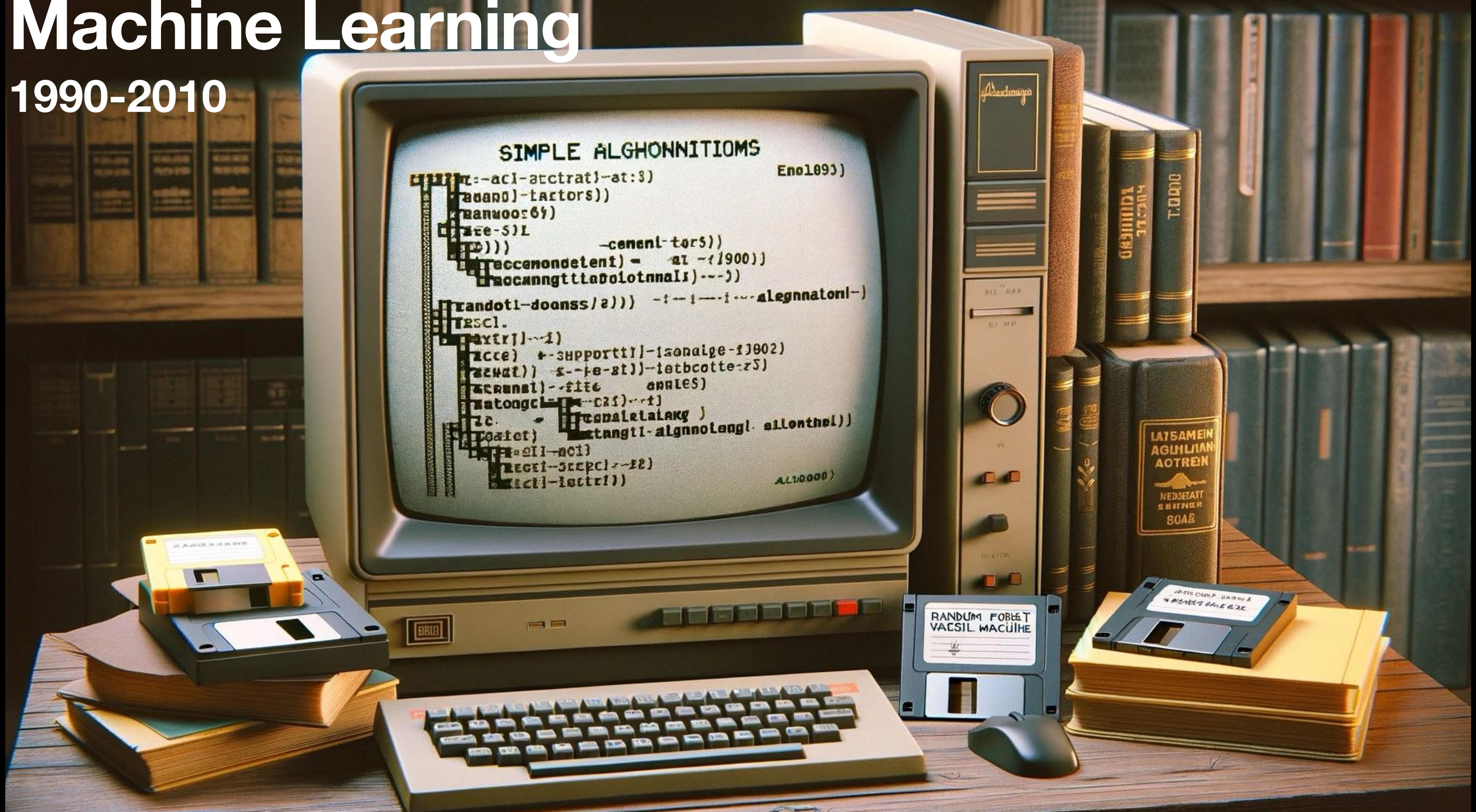
Recap

What would you use and why?

- You want to do traffic flow simulations
- There is not really a solution for your problem in literature; it is highly specific and fast changing; you have some vague ideas about different solutions.
- Your optimization problem would need insane amounts of memory, or time, to compute.

Machine Learning

1990-2010



Deep Learning



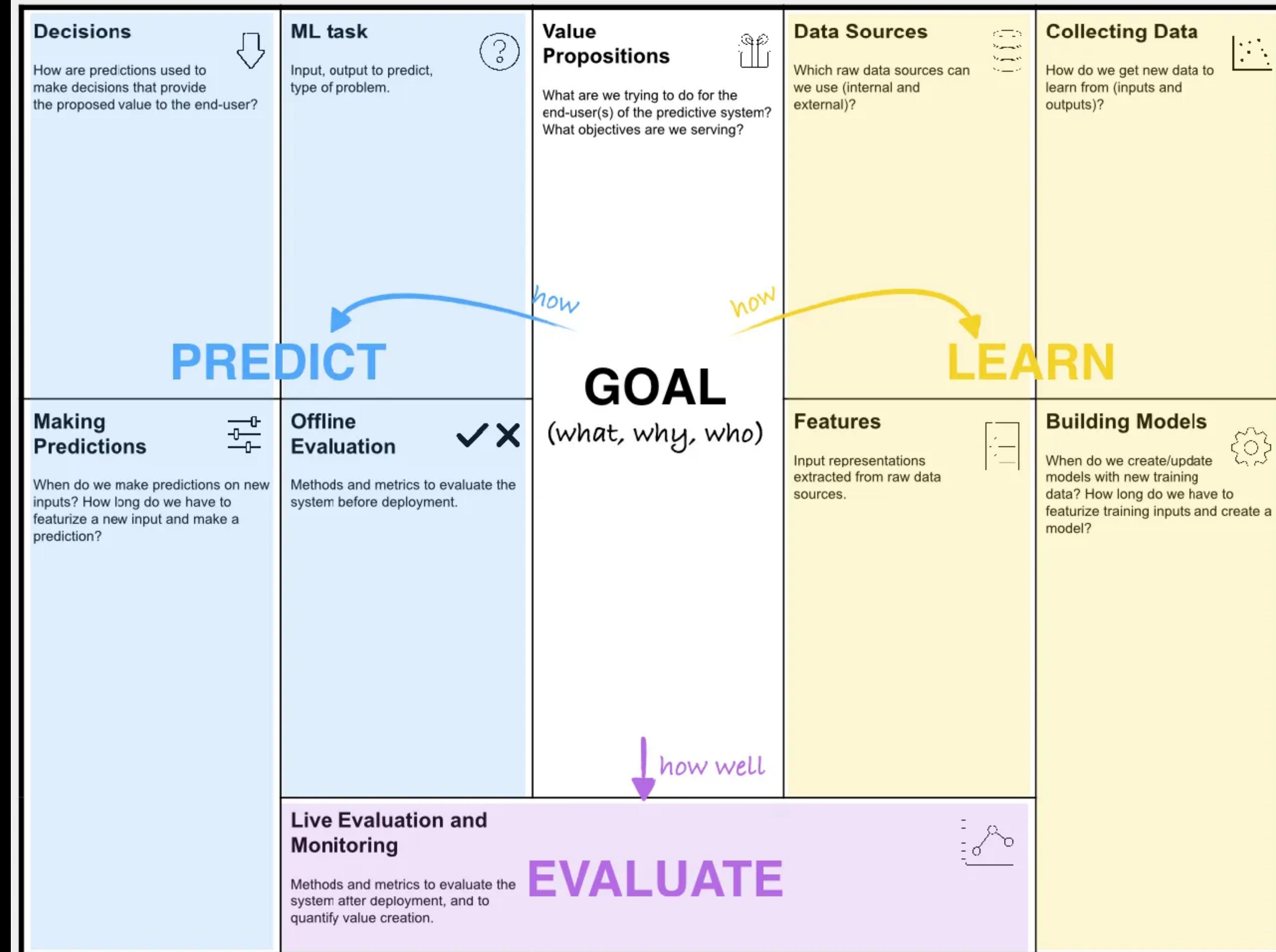
The Machine Learning Canvas (v0.4)

Designed for:

Designed by:

Date:

Iteration:



Predict

- Decisions: descriptive, predictive, prescriptive
- Task: classification, regression, clustering?
- Evaluation: what are good metrics?

Learn

- Data sources: what is good data?
- Collecting: how do we get new data?
- Features: what kind of information do you need?

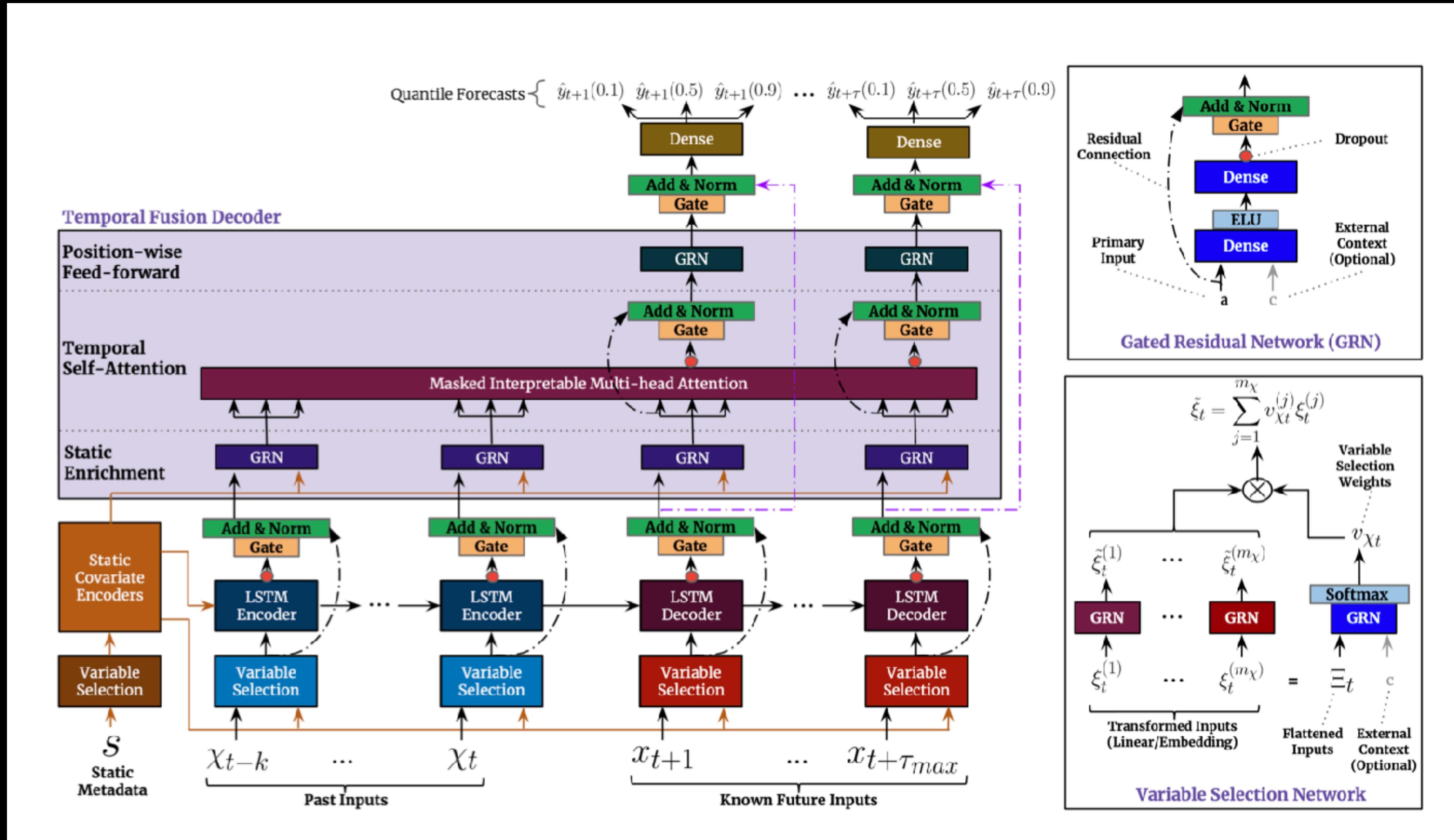
Neural Networks

Machine learning

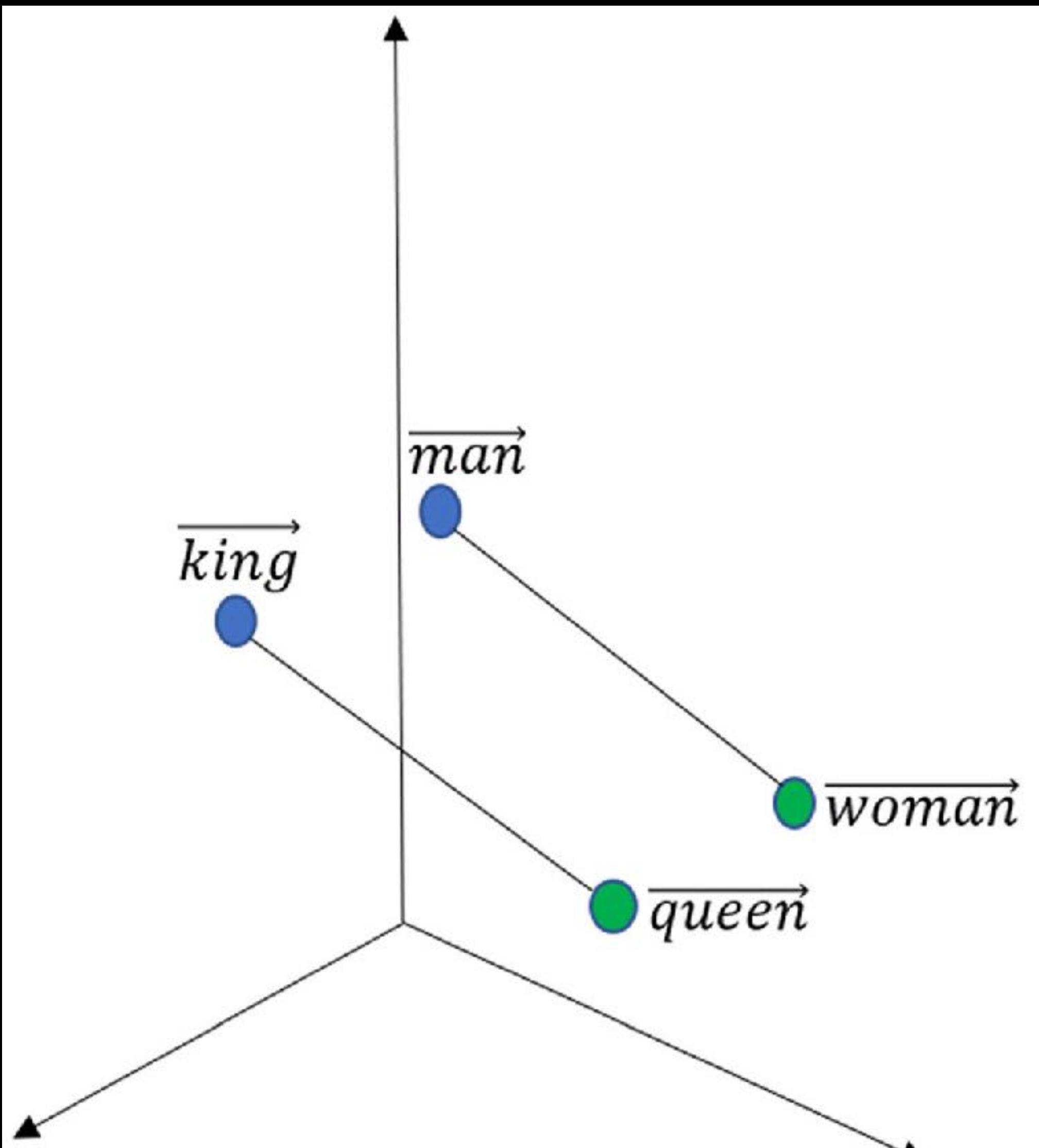
The summary

$$f: X \rightarrow y$$

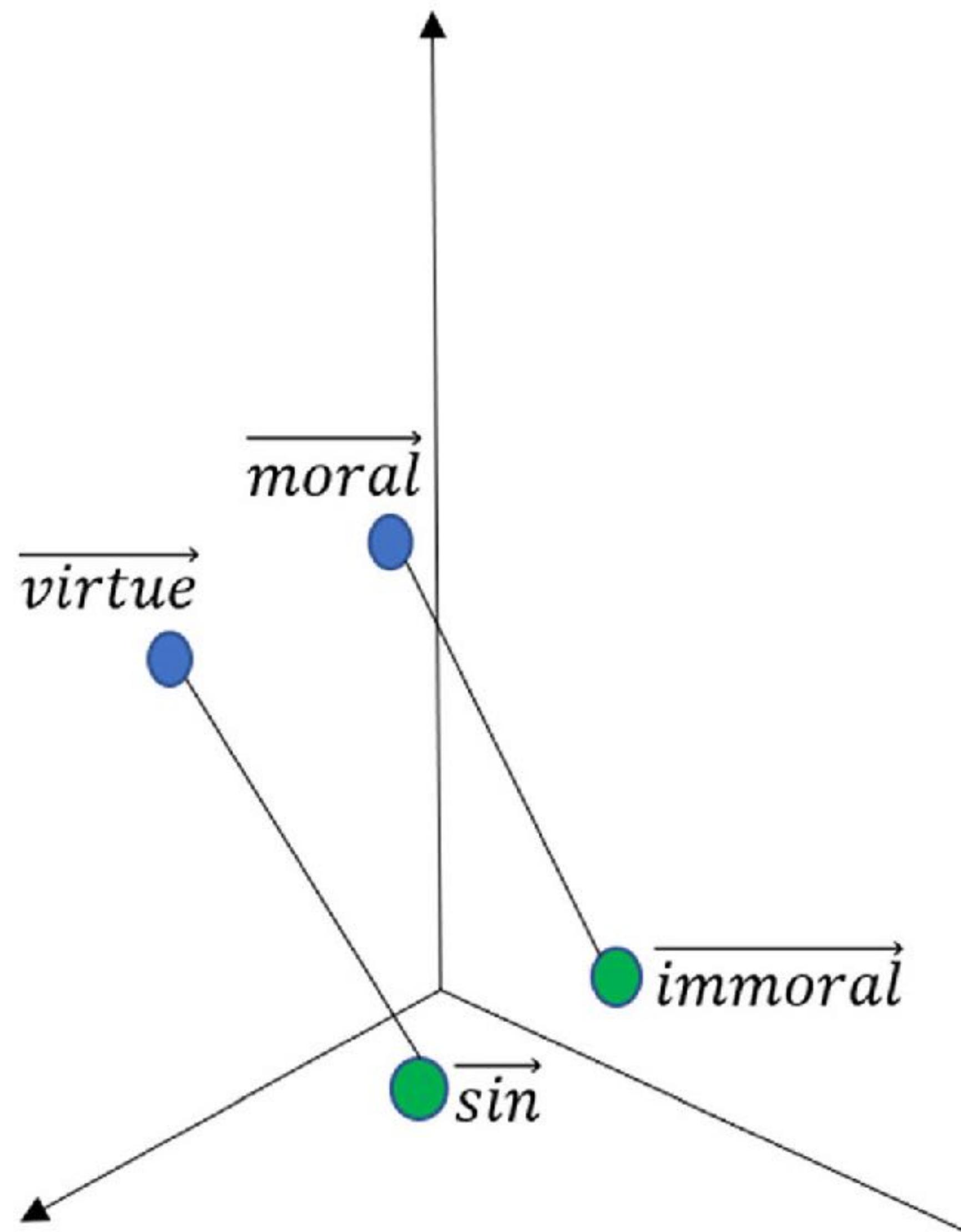
Complexity 2021 - TFT







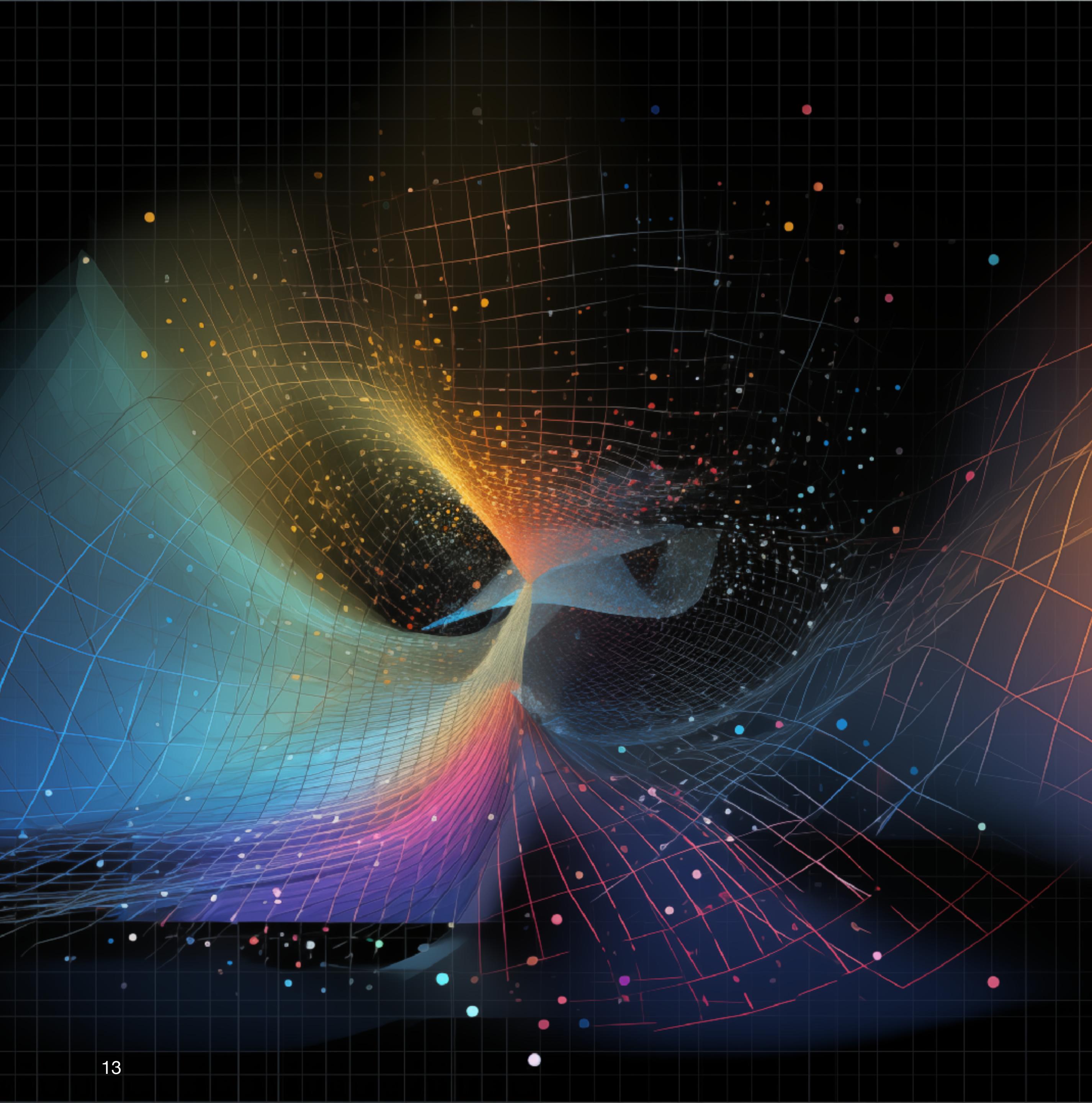
a “gender” dimension



a “morality” dimension

LLMs use a semantic vectorspace
with 768 dimensions : \mathbb{R}^{768}

Atoms in the visible universe : 10^{80}



Linear models

The basic mathematical shape of a linear model is:

$$Y = WX + b$$

- $Y = \{y_1, \dots, y_m\}$ are labels, so we have m labels.
- Every label Y_i has corresponding features $X_i = \{x_{i,1}, \dots, x_{i,n}\}$, so we have n features.
- We can store the observations in an (m, n) matrix X
- We can store n weights for every feature in a matrix W
- b is an additional bias weight



The matrix notation is concise. We could write this out in full as:

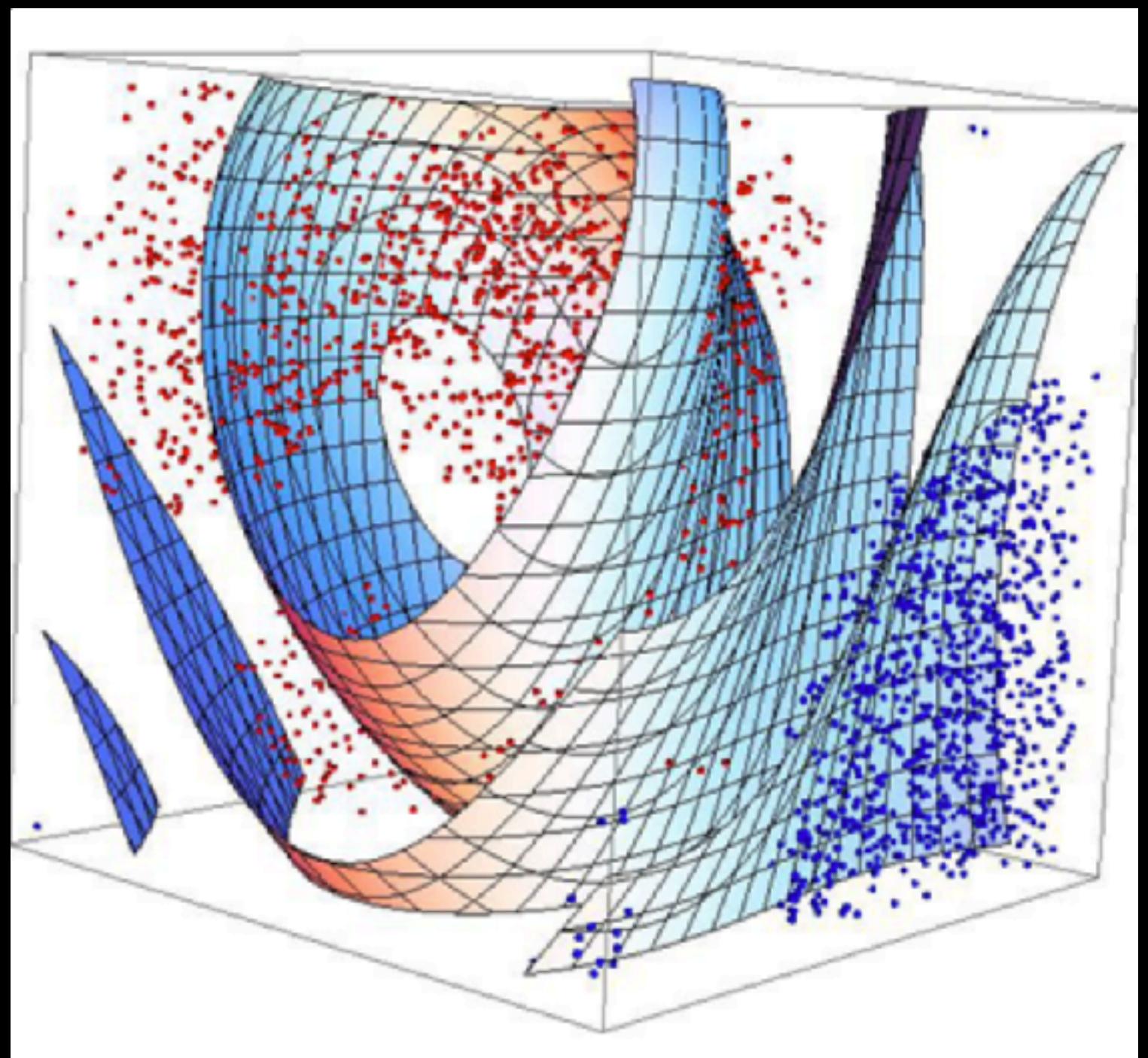
$$y_i = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Non-linear models

A lot of data is non-linear. This means we would need a curved hyperplane.

One trick to do this is the kernel-trick, which is commonly used with Support Vector Machines, which we touched upon in the short history.

Deep learning has found another trick



Deep Learning

Linear function

$$f(X) = WX + b$$

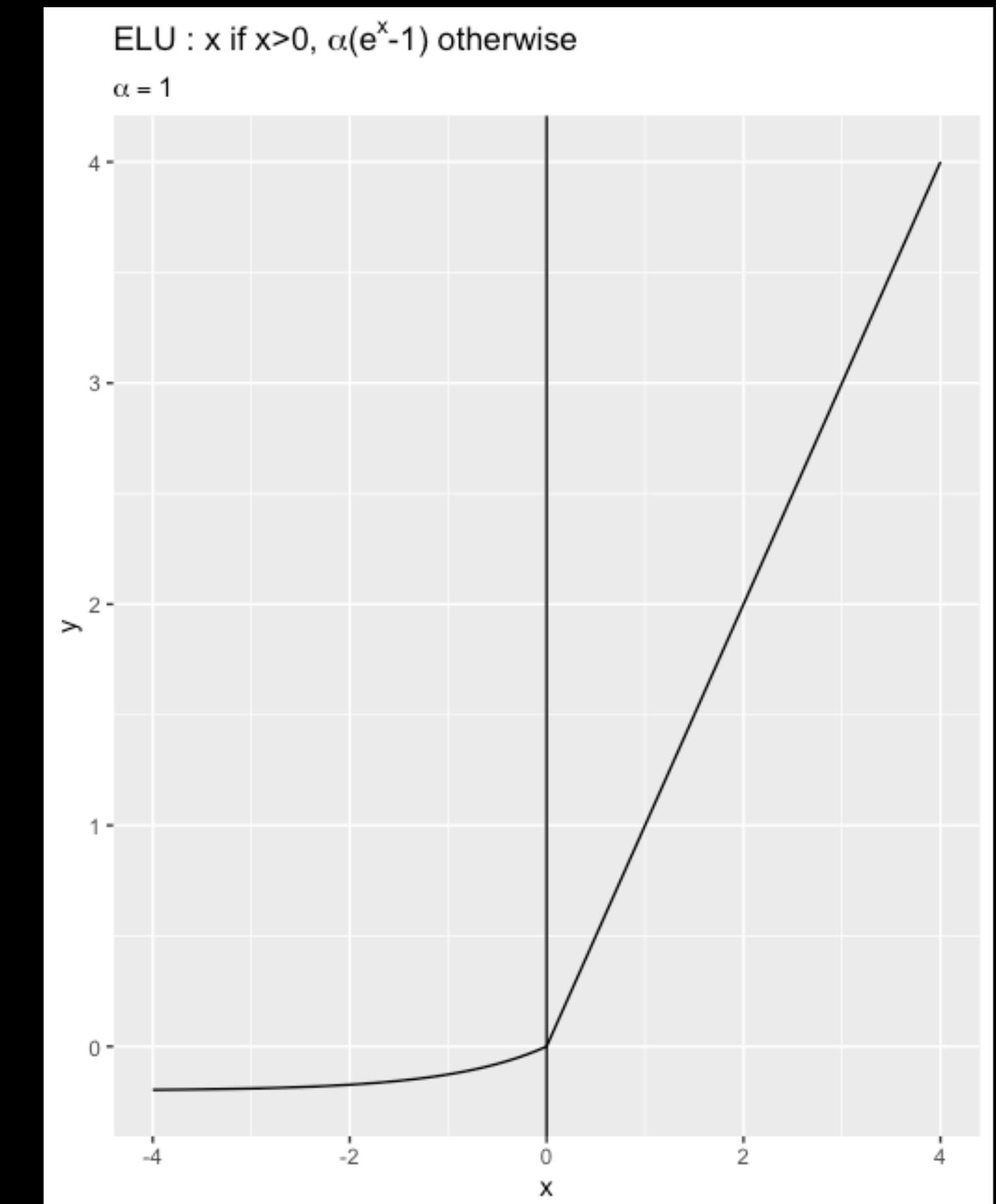
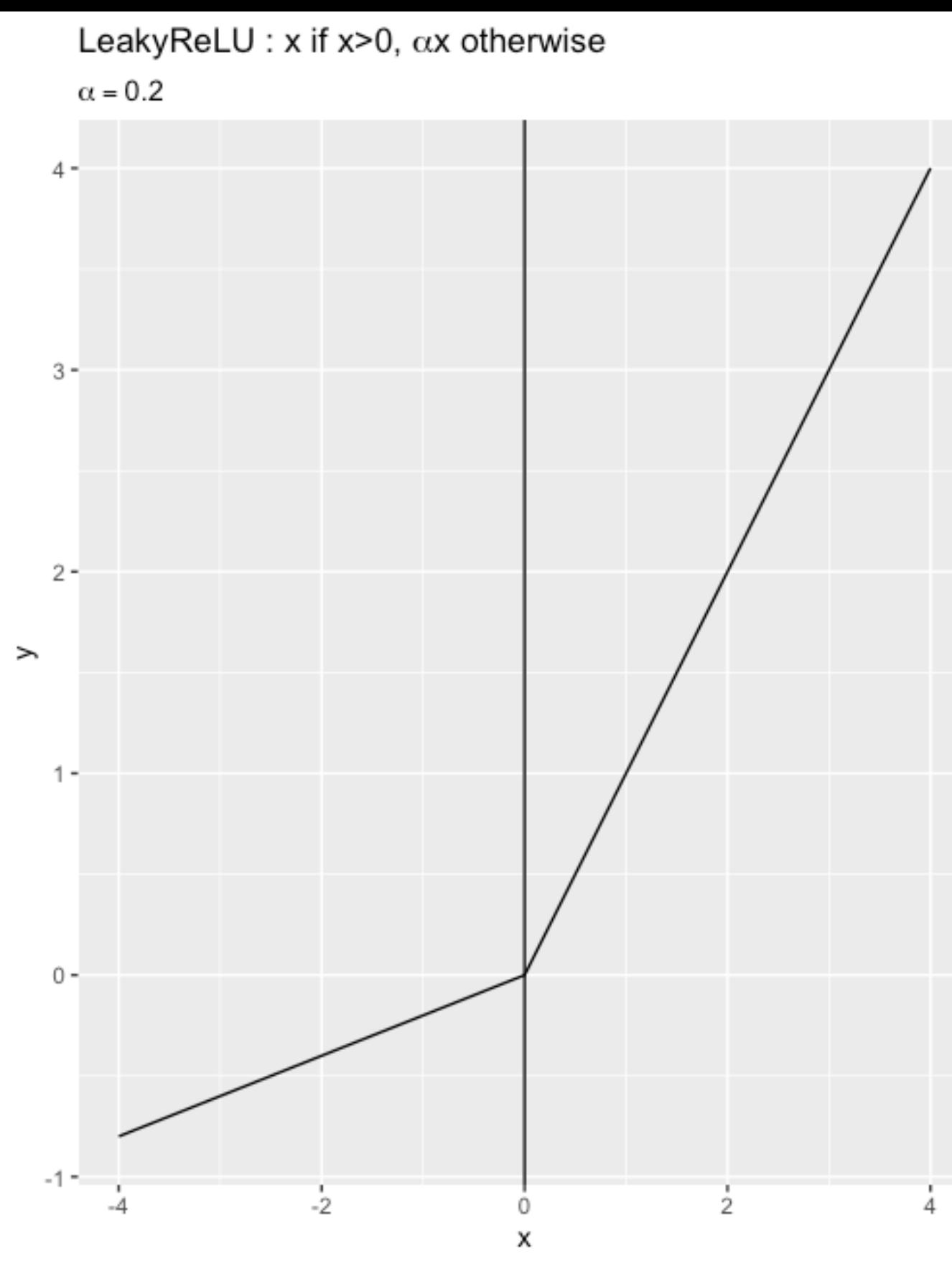
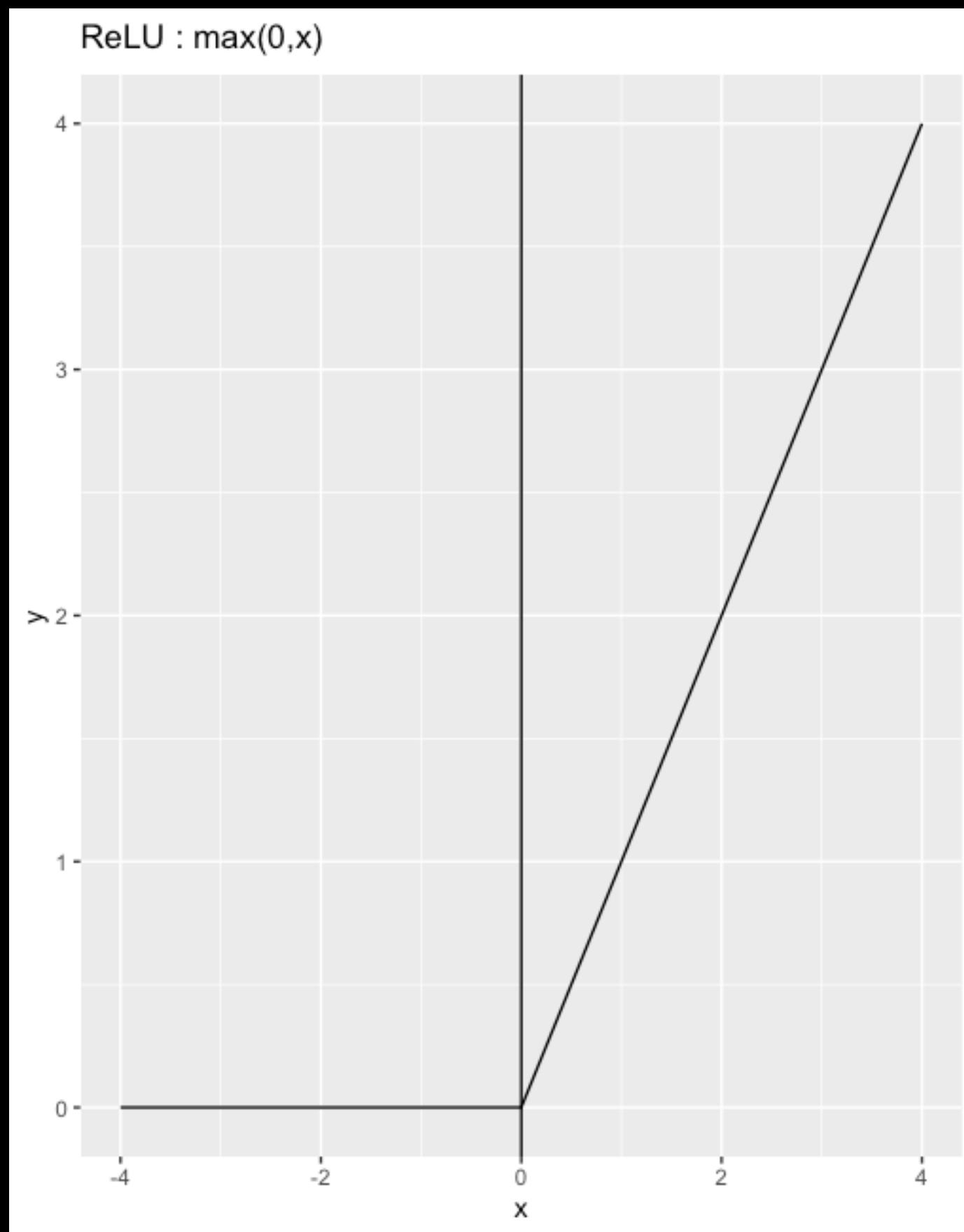
Nonlinear function

$$\sigma(X) = \max(0, X)$$

Neural network

$$\hat{y} = f_n \circ \sigma \circ f_{n-1} \circ \dots \circ \sigma \circ f_1$$

Activations non-linearities



Deep Learning

$$X = \{\vec{x}_1, \dots, \vec{x}_n \mid \vec{x} \in \mathbb{R}^d\}$$

Data

$$y = \{y_1, \dots, y_n \mid y \in \{0,1\}\}$$

Trainable Weights W, b

$$(Non)linearity \quad f(X) = WX + b \quad \sigma(X) = max(0, X)$$

Predict

$$\hat{y} = f_n \circ \sigma \circ f_{n-1} \circ \dots \circ \sigma \circ f_1$$

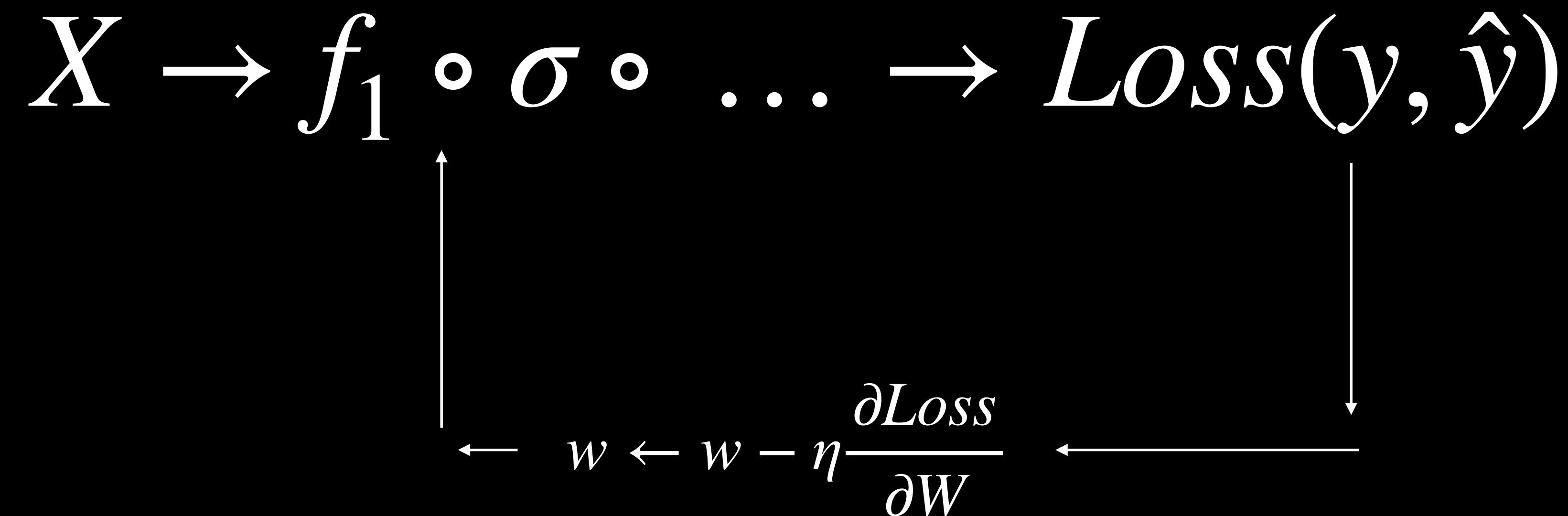
Loss

$$Loss(y, \hat{y})$$

Optimize

$$w \leftarrow w - \eta \frac{\partial Loss}{\partial W}$$

Training



Convolutional Neural Networks

Convolutions - the motivation

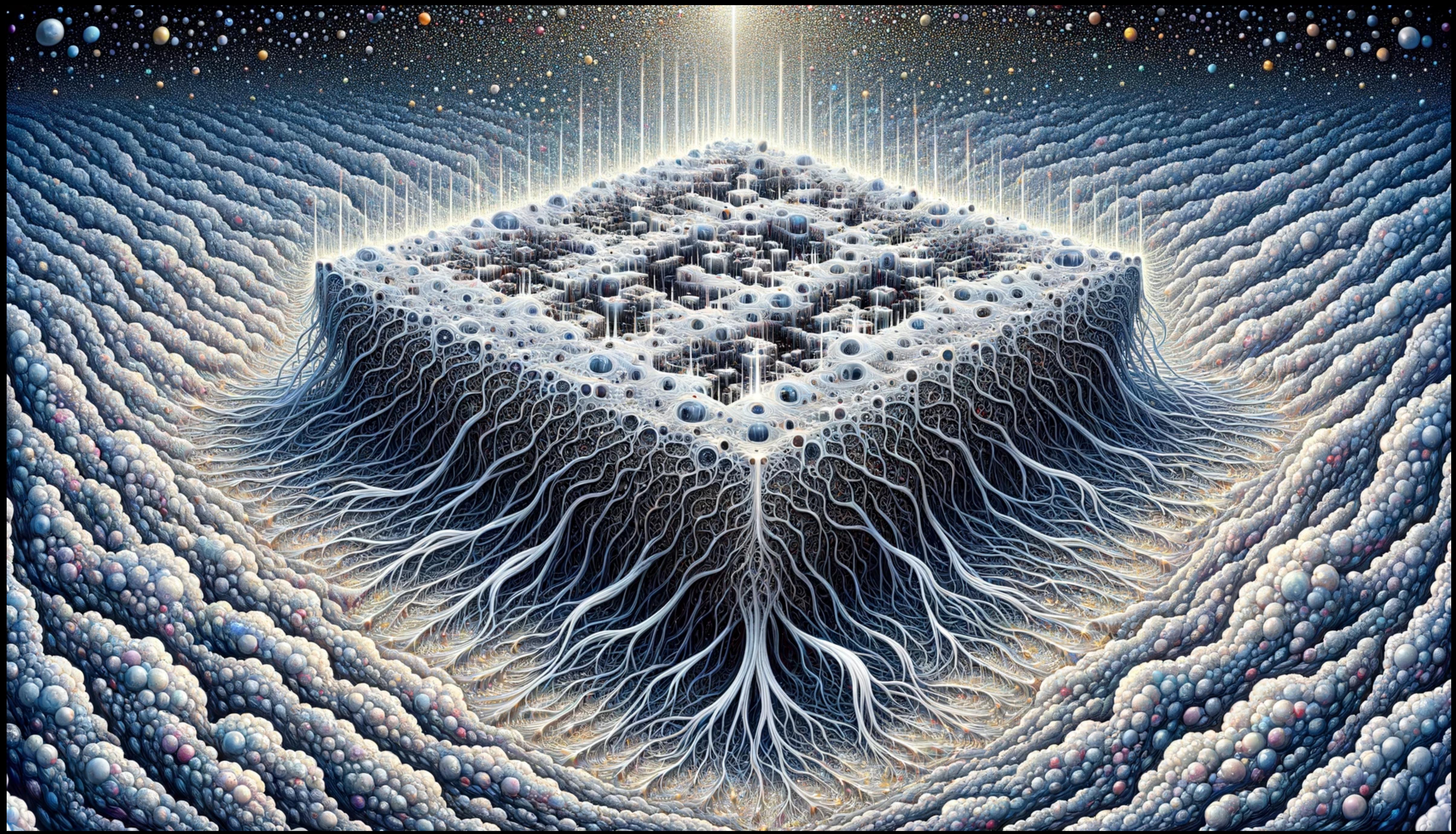
An image of size 28x28 has 784 pixels

With 256x256 you are already at 65536 pixels

Treating every pixels as a feature will blow up the amount of parameters for your model:

Assuming you halve the dimension, a batch with dimensions (32, 65536) will need a (65536, 32000) weight matrix.

That are over 2×10^9 parameters, just for the first layer...

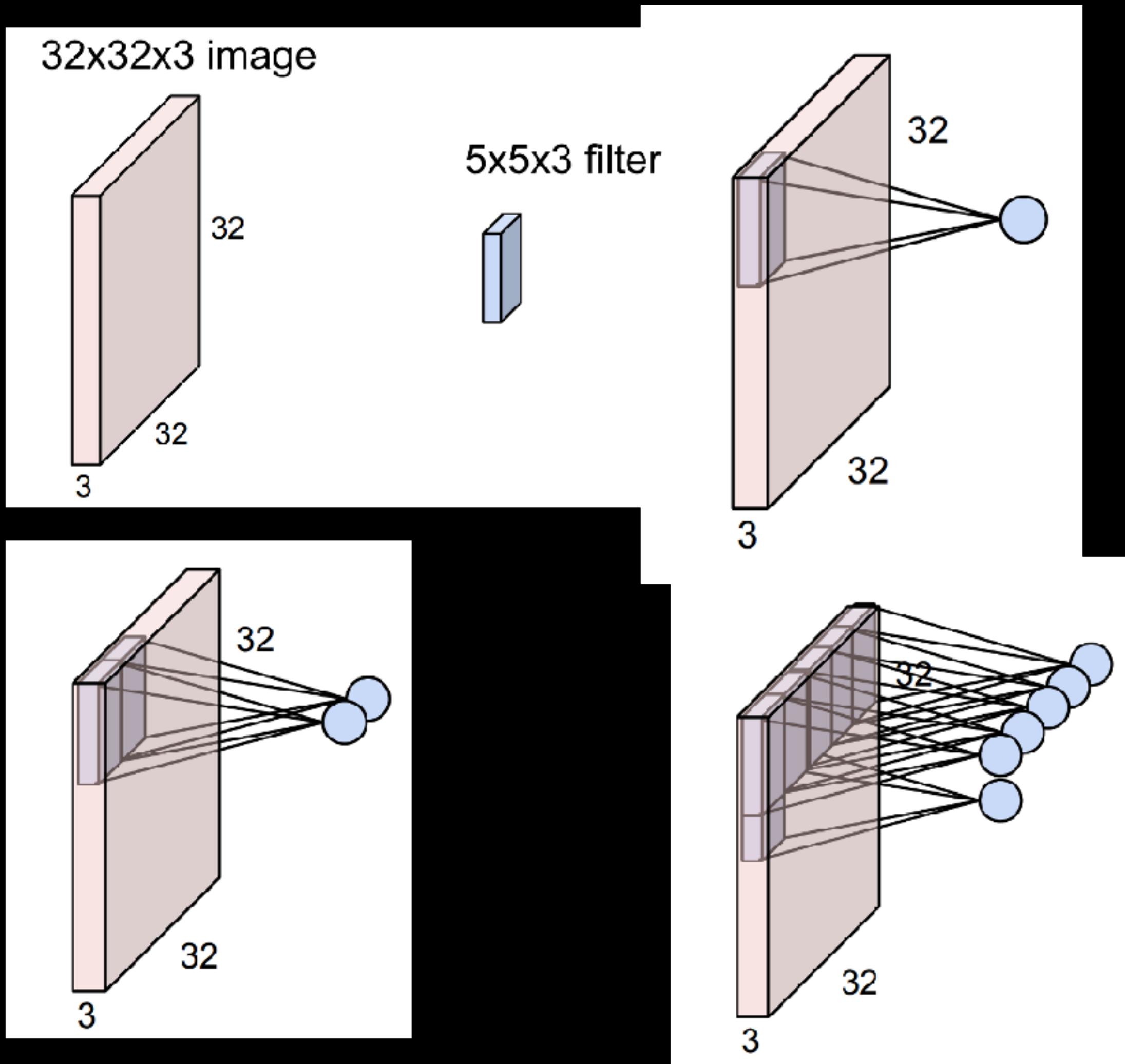




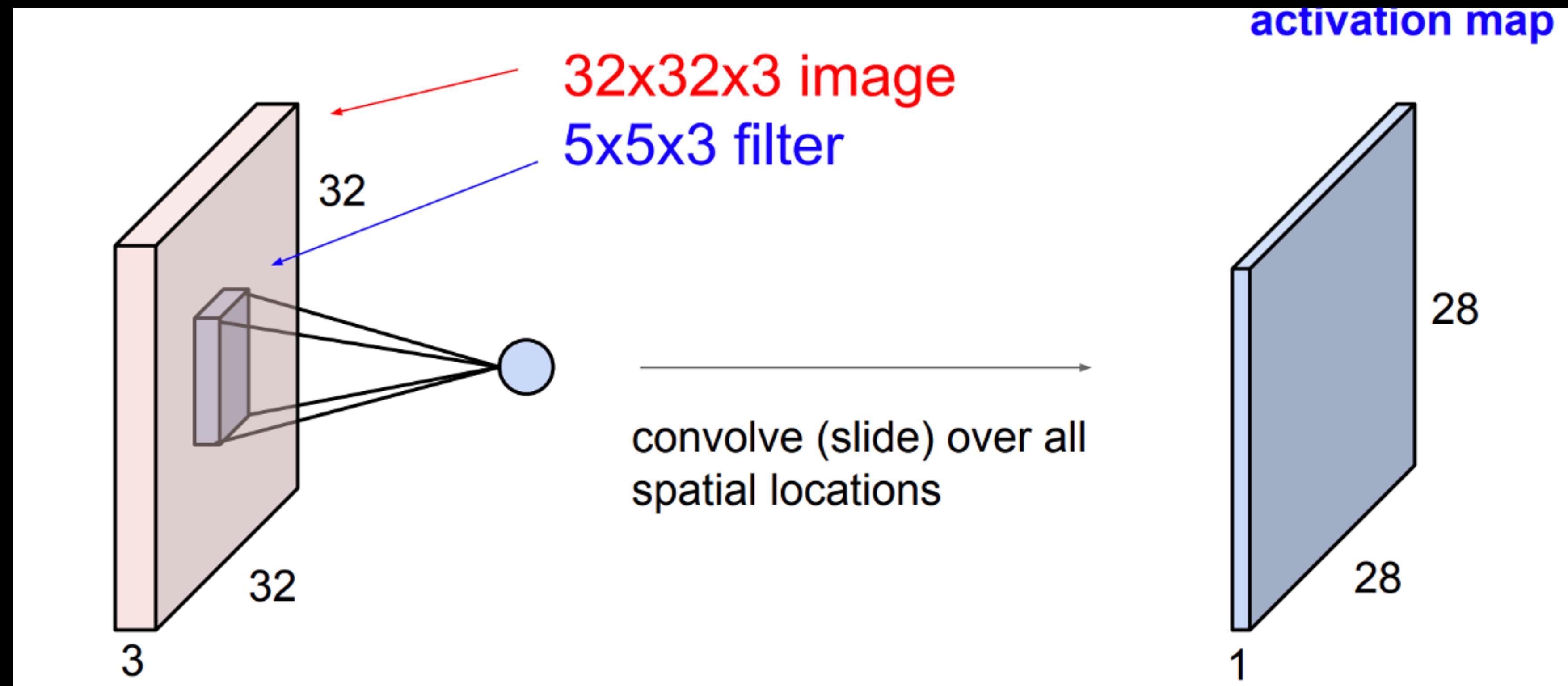
Convolutions

- Calculate the dot product between the filter and the image
- E.g.,

$$\begin{bmatrix} 1 & 10 \\ 2 & 20 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 2 & 0 \end{bmatrix} =$$
$$(1 * 1) + (10 * 2) + (2 * 2) + 0 = 25$$



Convolutions

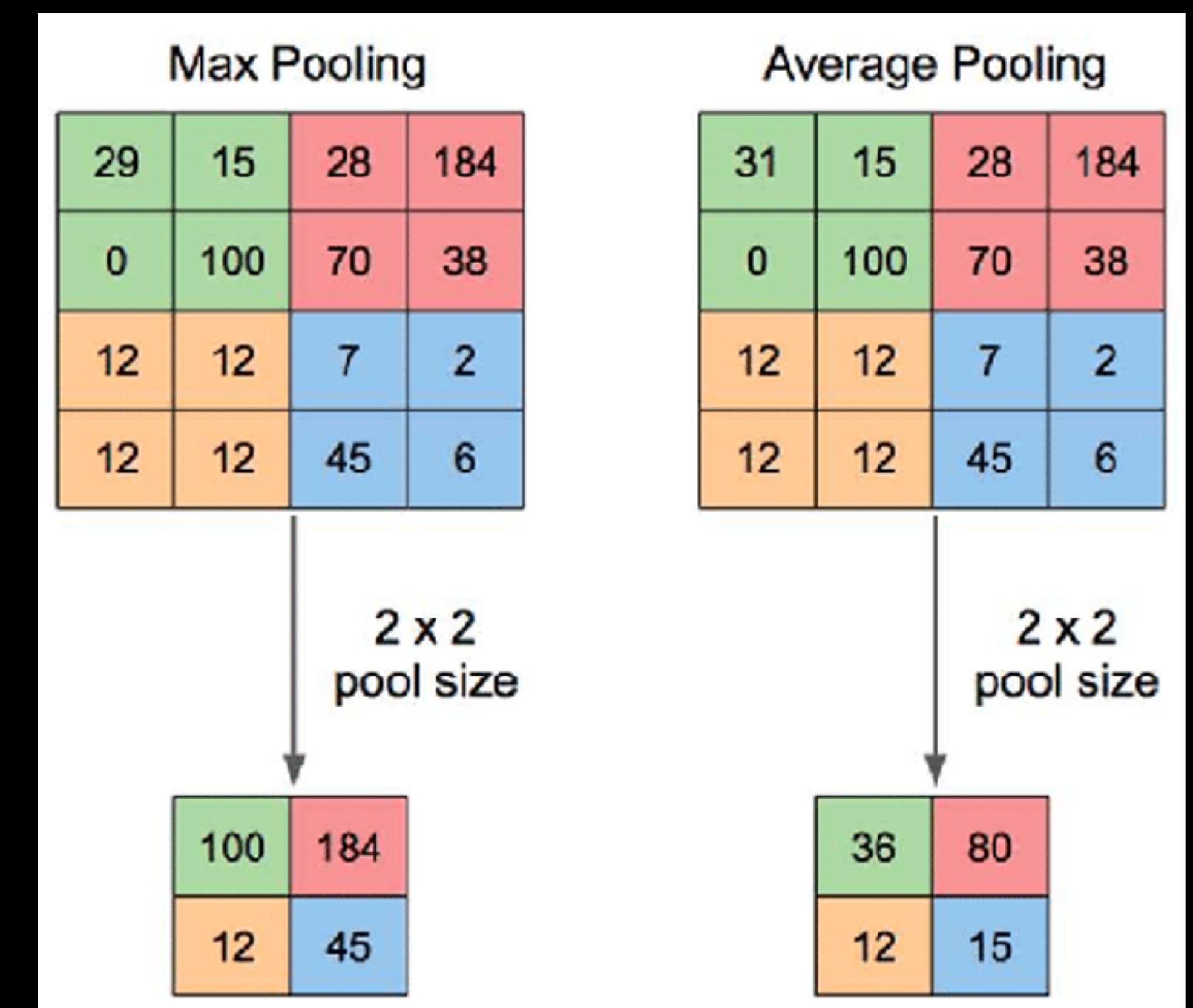


Convolutions

- In deep learning, we won't construct the filters for ourselves!
- We only define the parameters (eg size of the kernel, how many) and let the model learn the convolutions, based on the data and labels.

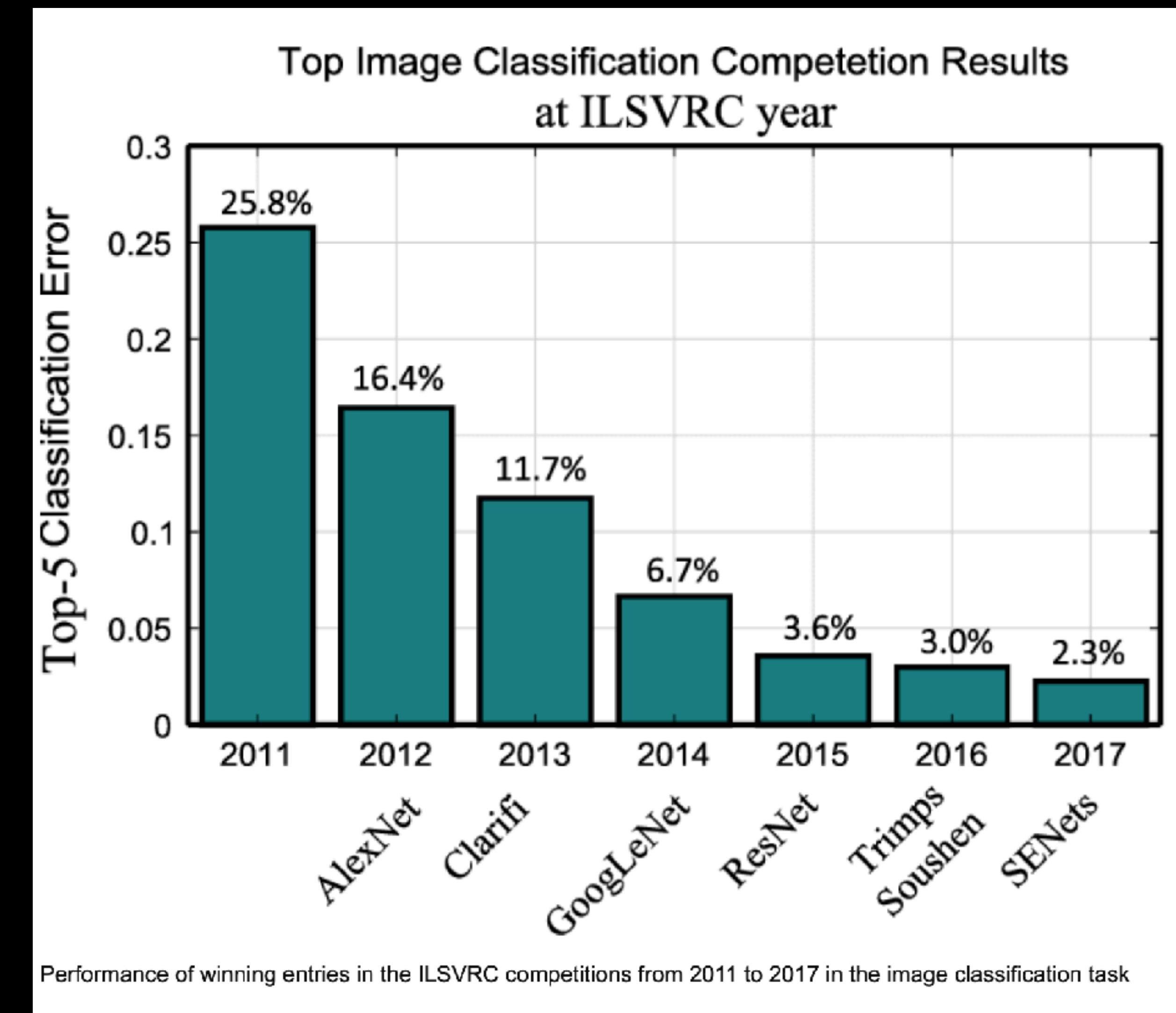
Pooling

- A way to downsample the input
- max pooling: focus on the highest value
- average pooling: smooth the pixels



Imagenet Challenge

- Human performance is about 5%



AlexNet

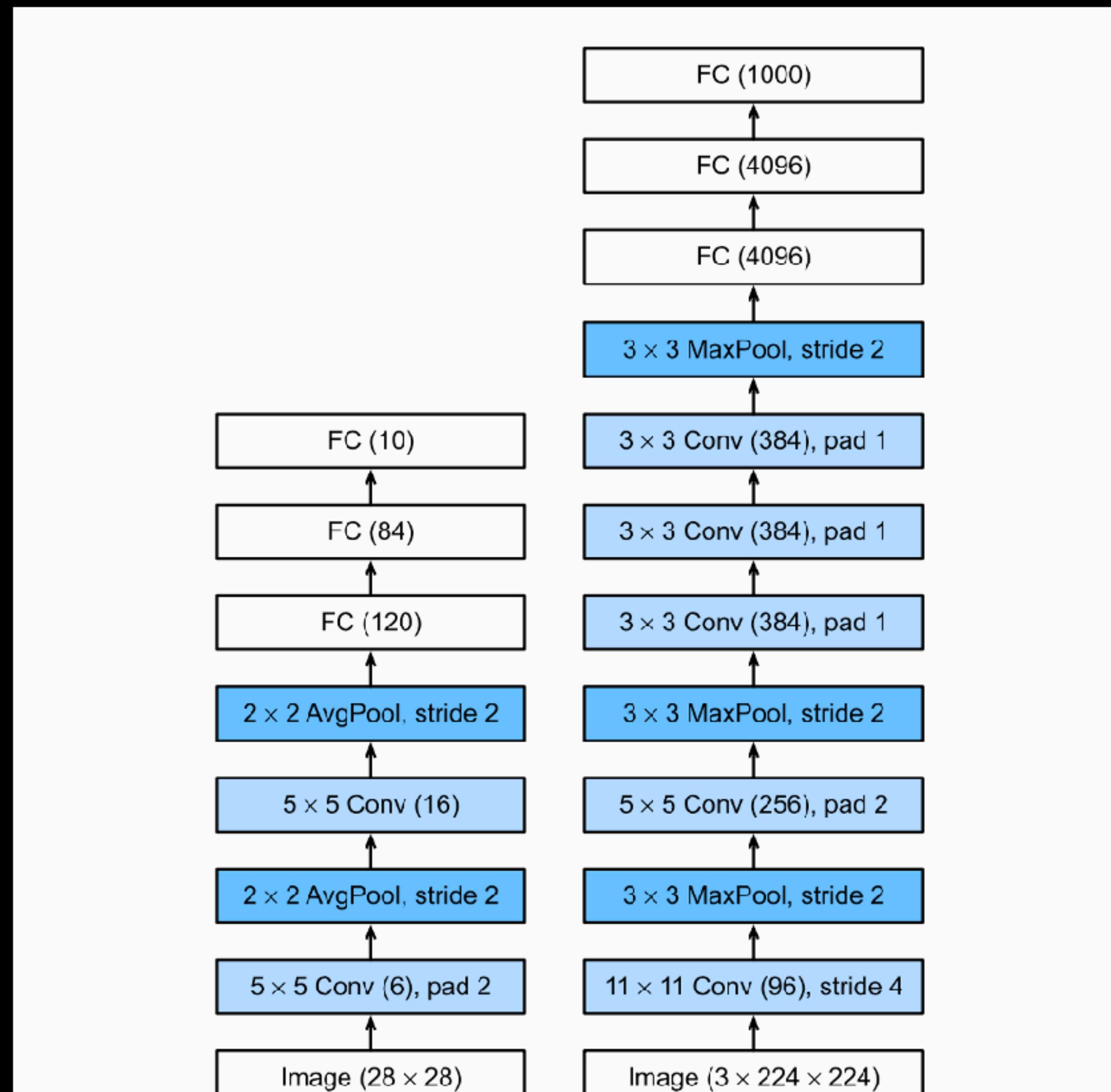
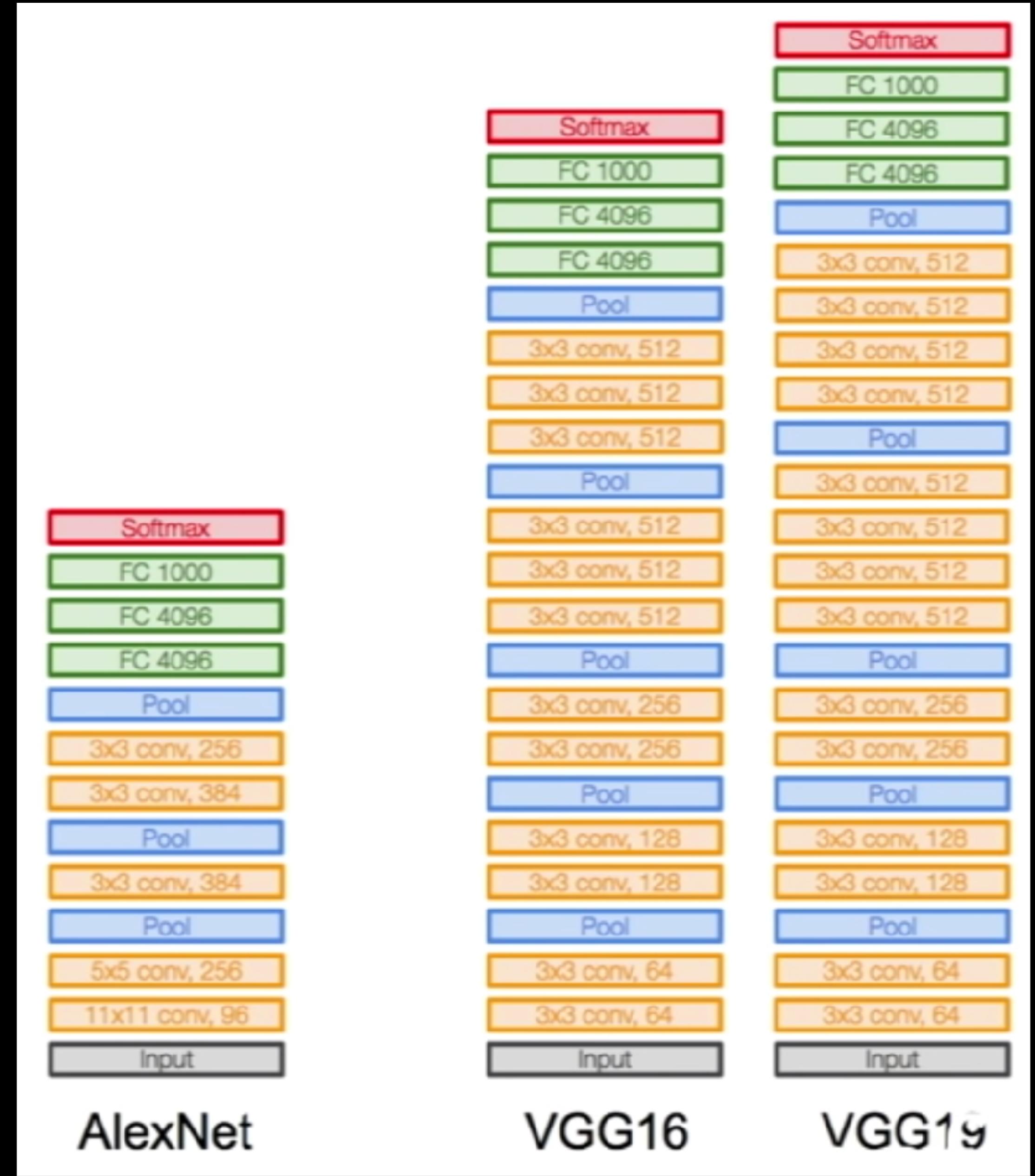


Fig. 7.1.2 From LeNet (left) to **AlexNet** (right).

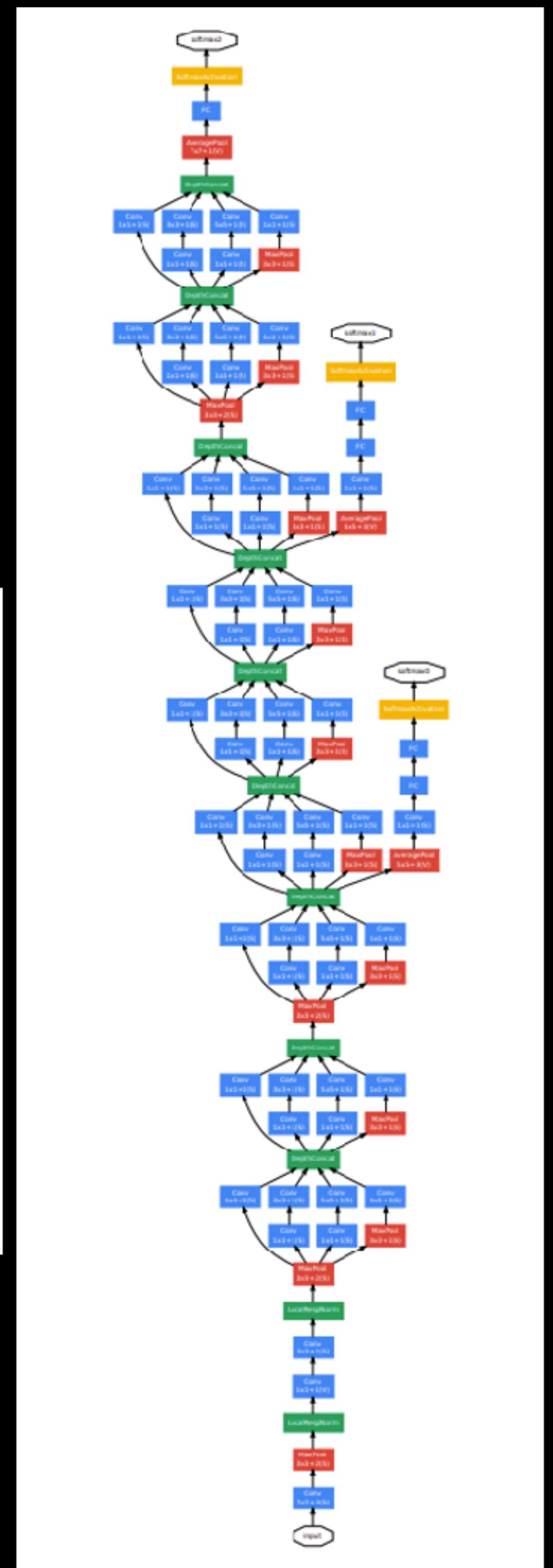
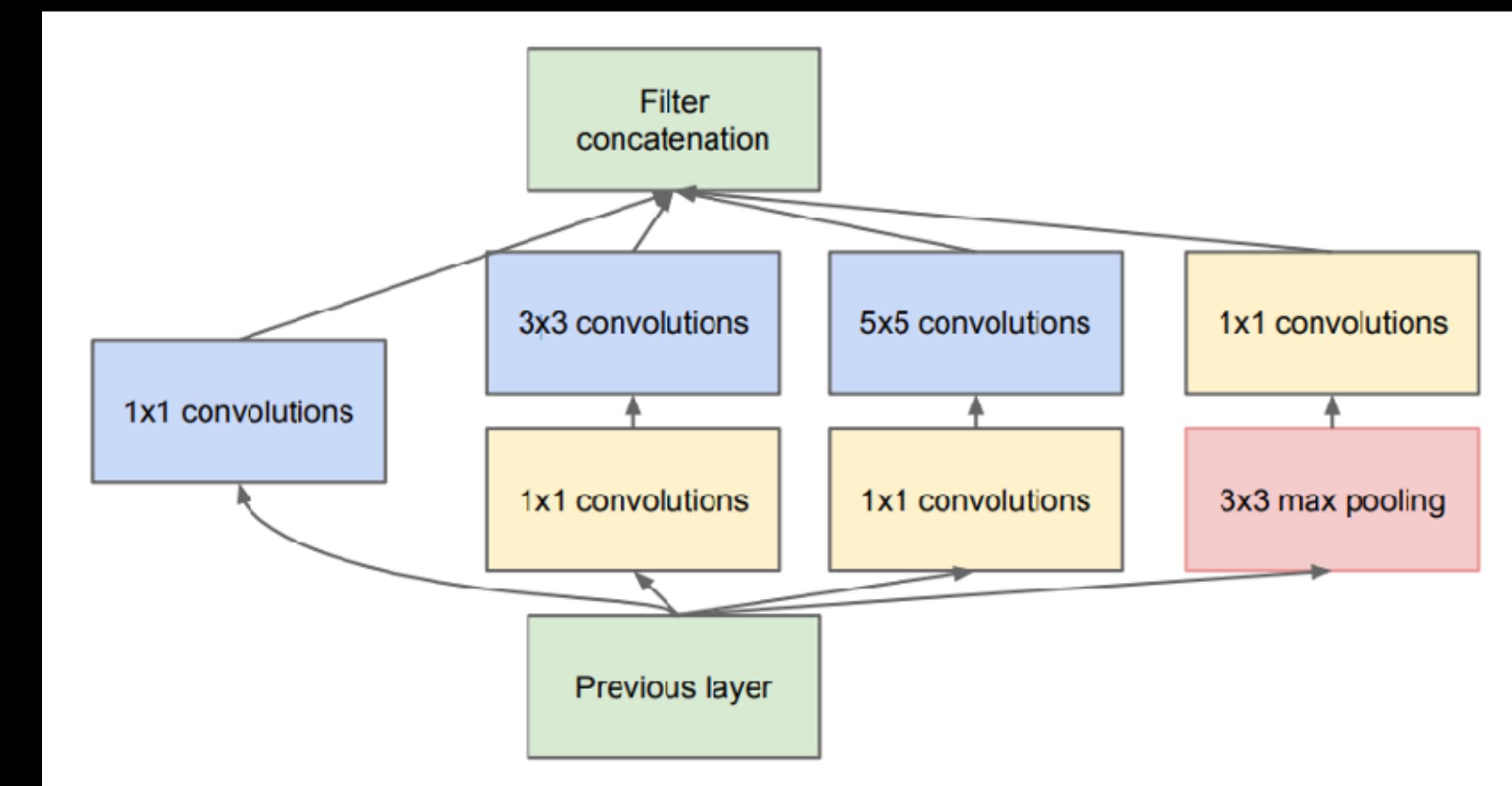
VGG

Changes in:
Filter size
channel numbers
depth
138 million parameters



GoogleNet

Parallel filters
Bottleneck layers
6.7 million parameters

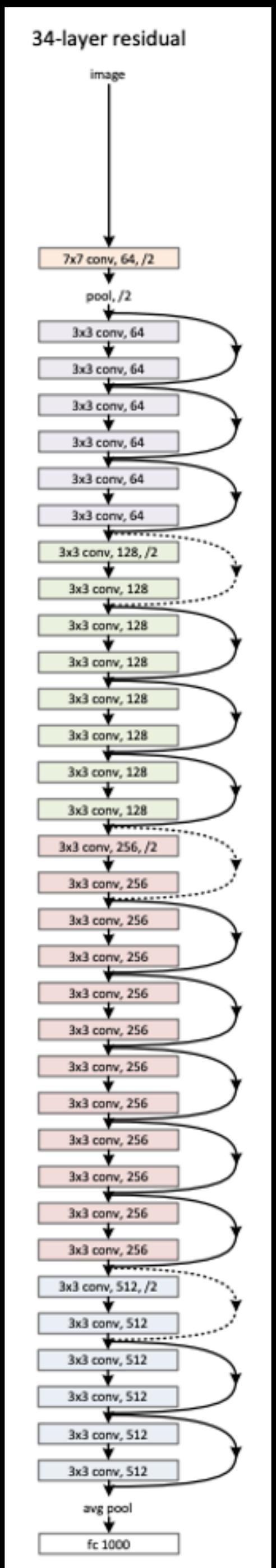


ResNet

batch normalisation

Residual blocks

11 million parameters for ResNet18

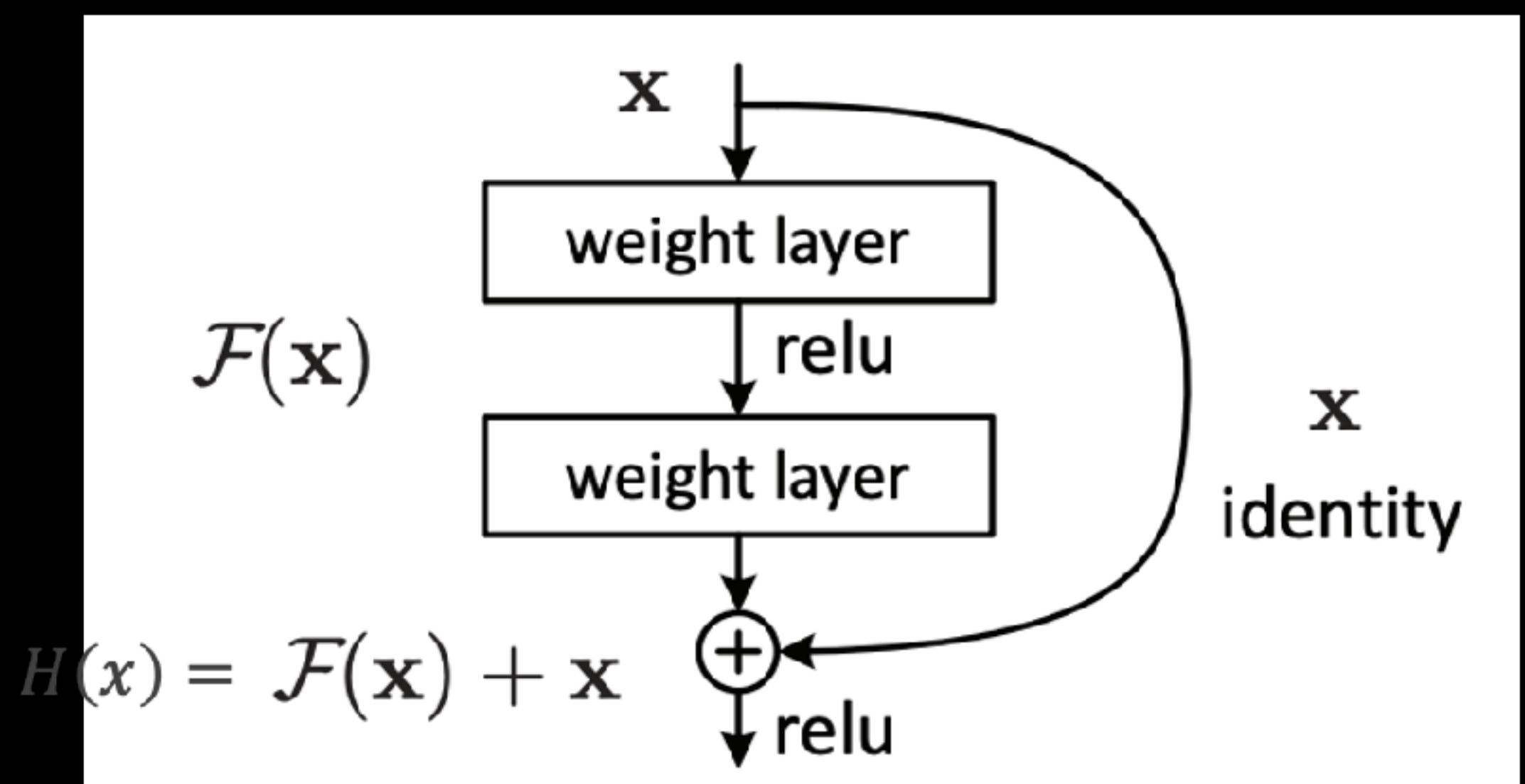


ResNet

Learn the residual, instead of the full mapping

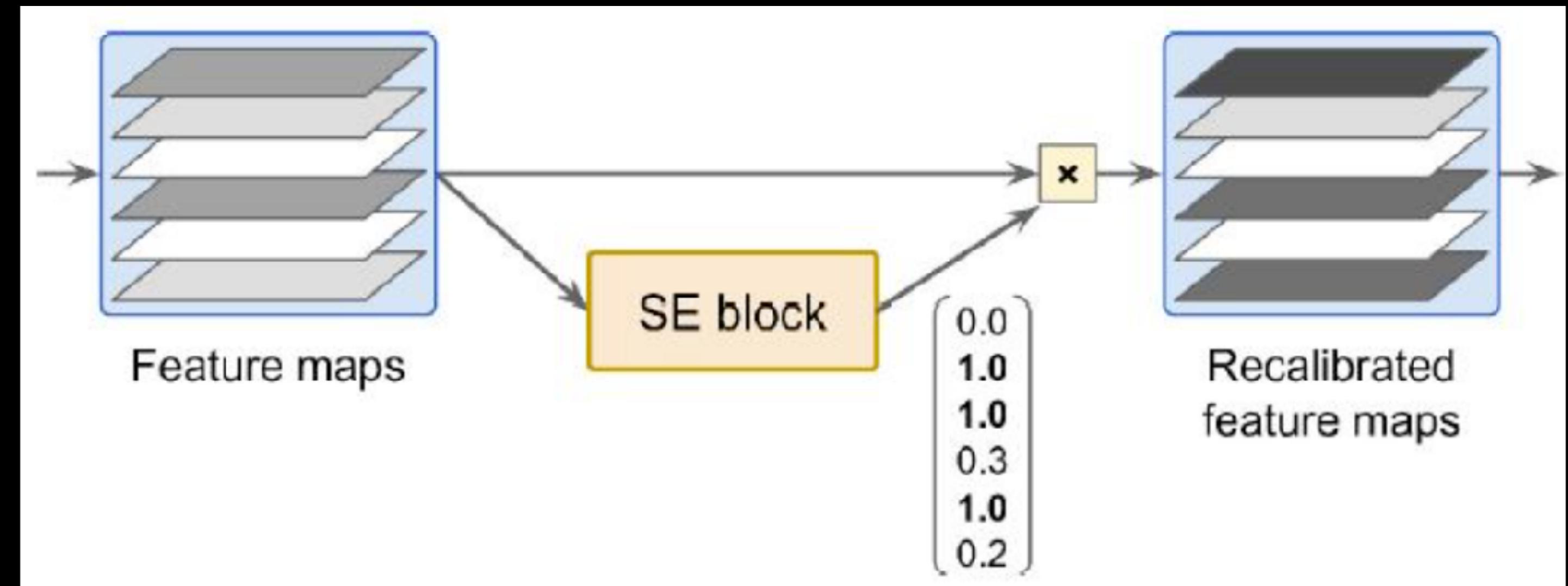
E.g., “use the same settings, but just change x ” is easier than starting from scratch.

It is easier for the gradient to flow back



SEnet

- Which features are most likely to be activated together?
- E.g., if you see a left eye, you expect a right eye. Even if it is hidden in the image.

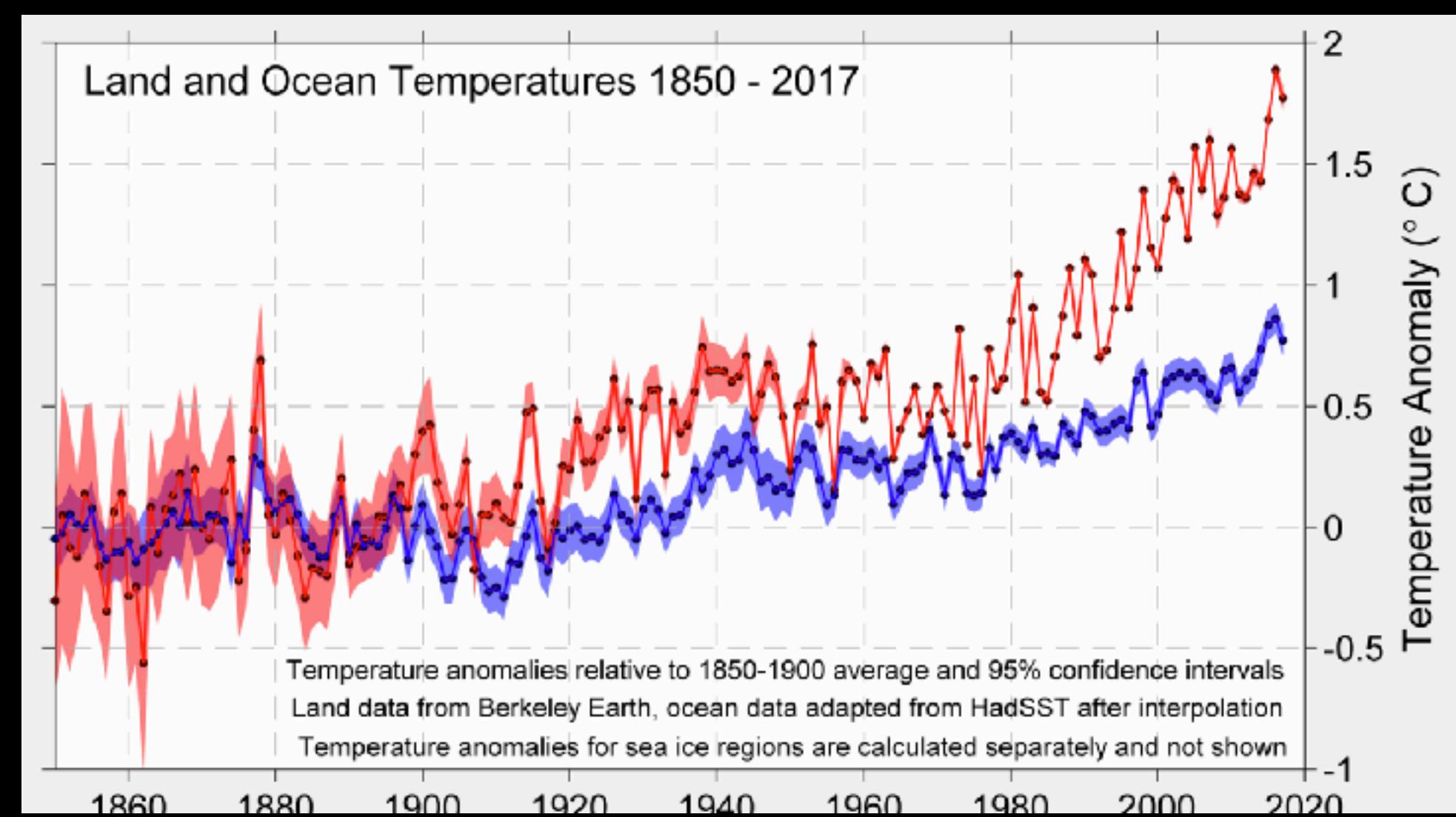
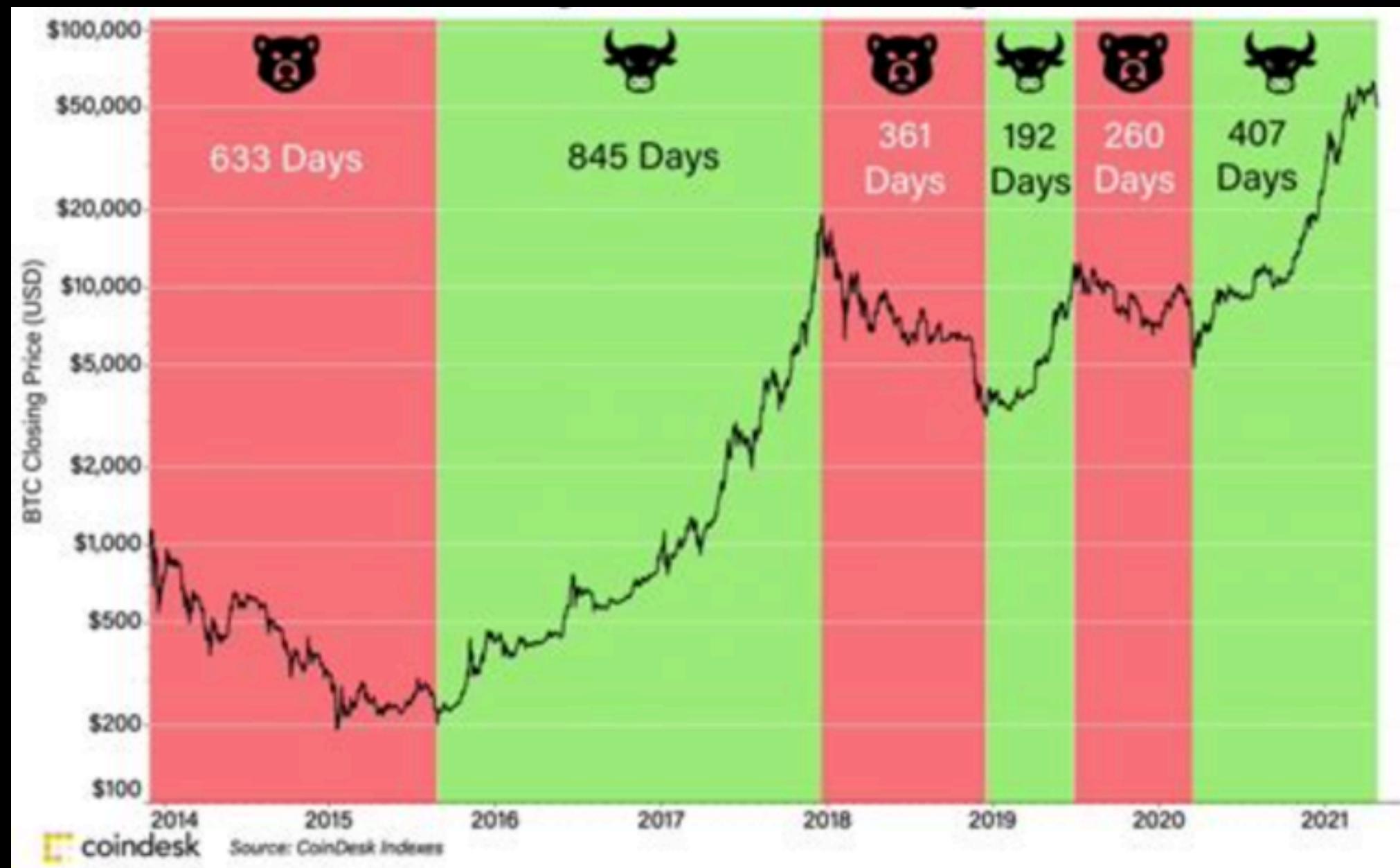


Motivation for RNNs

With sequences, the past offers context:

- Ik krijg geld van de **bank**
- Ik wil een nieuwe **bank** aanschaffen

We need the past to make sense of the future.



Simple RNN

To incorporate the hidden state, we simply add it:

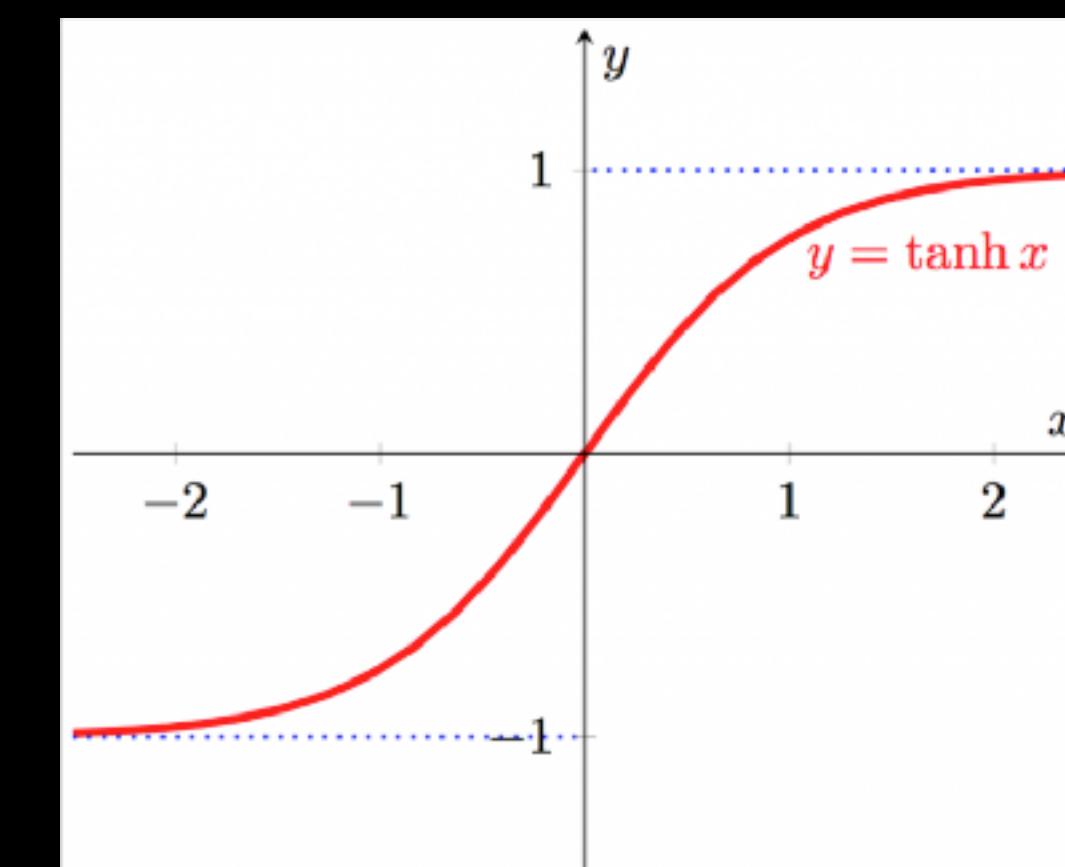
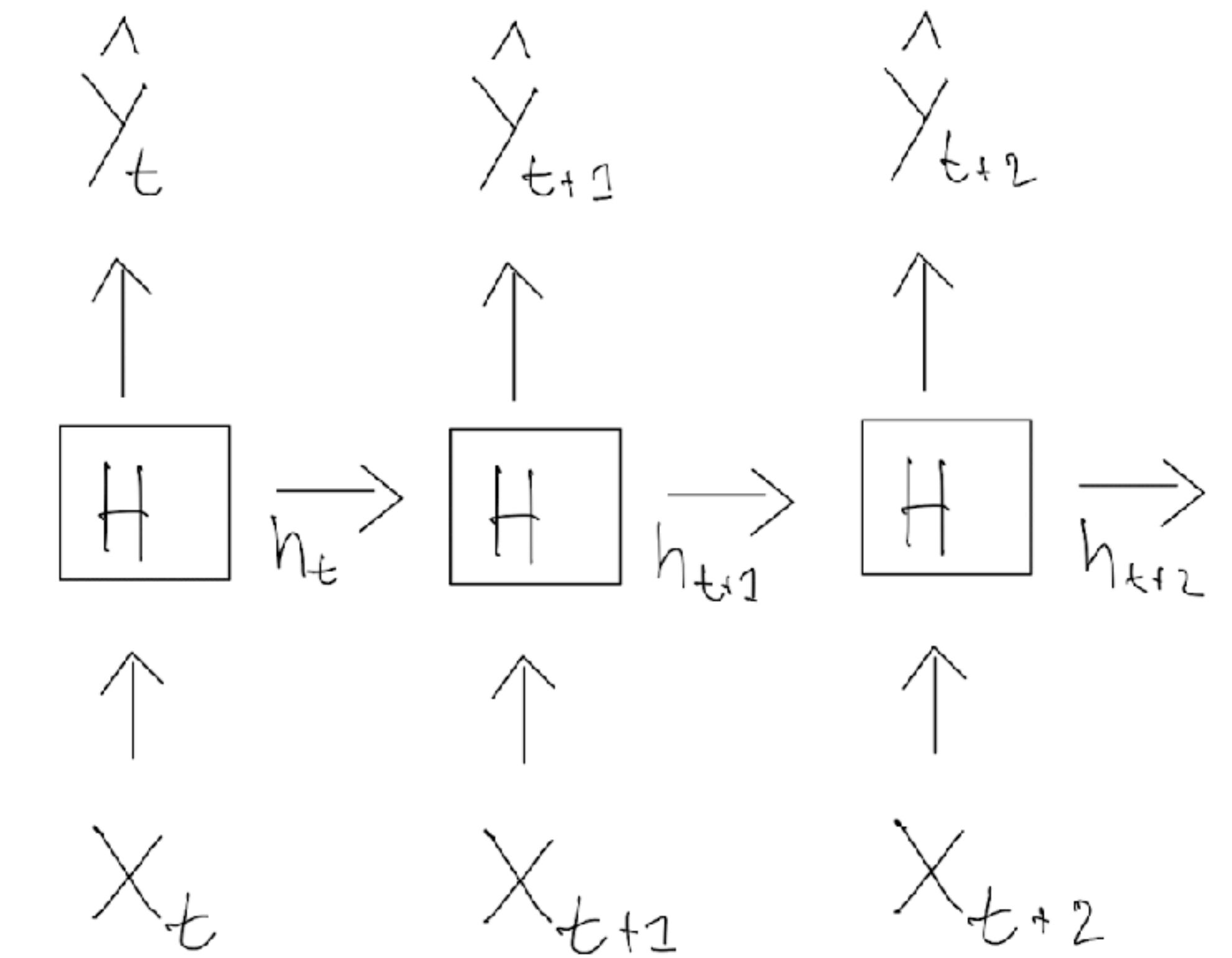
$$h_t = \sigma(W_x X_t + W_h h_{t-1} + b)$$

This is equivalent to

$$h_t = \sigma(W[X_t, h_{t-1}] + b)$$

where $[X, h]$ means concatenate

σ is an activation function, typically $tanh$

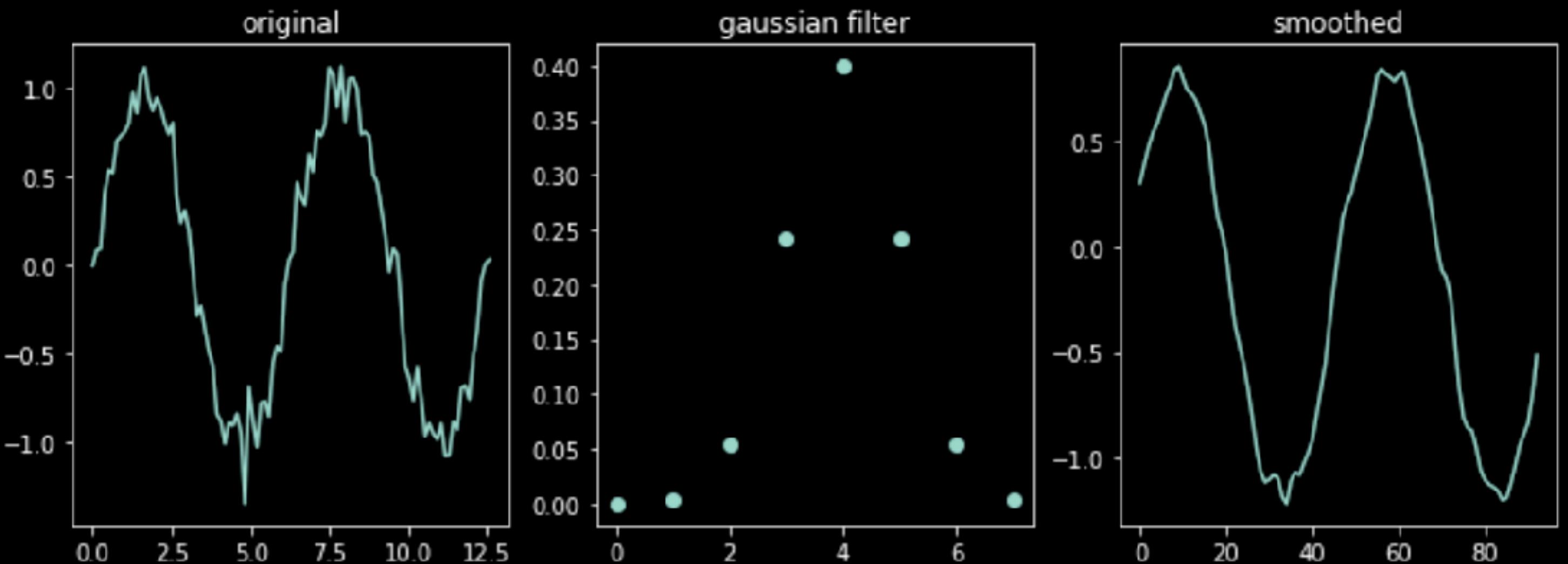


Transformers

Motivation Transformers

local weighing is not enough

- We have been using convolutions
- These look at local context
- But what about more distant relations?



Transformers: mixing in context

Let ξ_i be the vector that describes the semantic meaning of word i .

We can imagine we want to “mix in” other words, some more than others. A simple way would be with multiplication.

This way, we obtain a vector of numbers, one for every word in the sequence.

$$\begin{array}{cccc} \xi_1 & \xi_1\xi_1 & \xi_1\xi_2 & \xi_1\xi_3 \\ \xi_2 & \xi_2\xi_1 & \xi_2\xi_2 & \xi_2\xi_3 \\ \xi_3 & \xi_3\xi_1 & \xi_3\xi_2 & \xi_3\xi_3 \end{array}$$

Reweighting relevance

For the meaning of ξ_1 , let ξ_2 be irrelevant, and ξ_3 be highly relevant.

We would want $\xi_1 \xi_2$ to be closer to zero, and $\xi_1 \xi_3$ to be closer to one.

If we normalize all weights with a softmax, we make sure everything adds to 1

We can now reweigh the original ξ_1 to obtain $\tilde{\xi}_1$

And it will contain a little bit of ξ_2 context but much more of ξ_3 context

$$W = \text{softmax}([\xi_1 \xi_1 \quad \xi_1 \xi_2 \quad \xi_1 \xi_3])$$

$$W = [0.4 \quad 0.1 \quad 0.4]$$

$$\tilde{\xi}_1 = w_1 \xi_1 + w_2 \xi_2 + w_3 \xi_3$$