

Neural Networks and Learning Systems  
TBM126 / 732A55  
2021

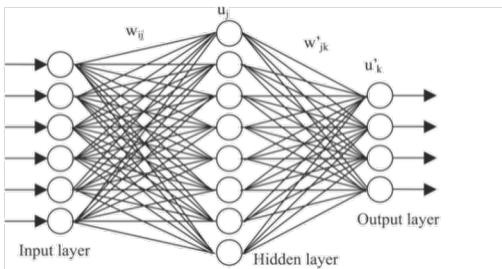
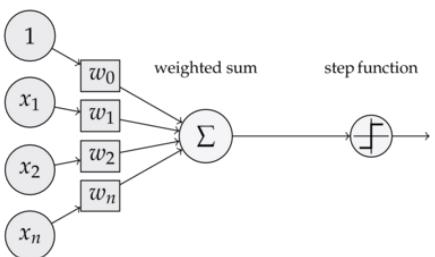
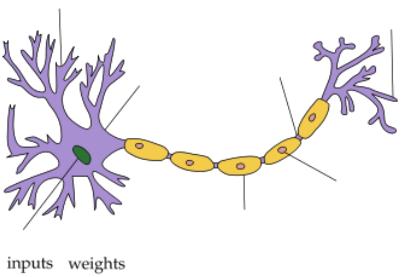
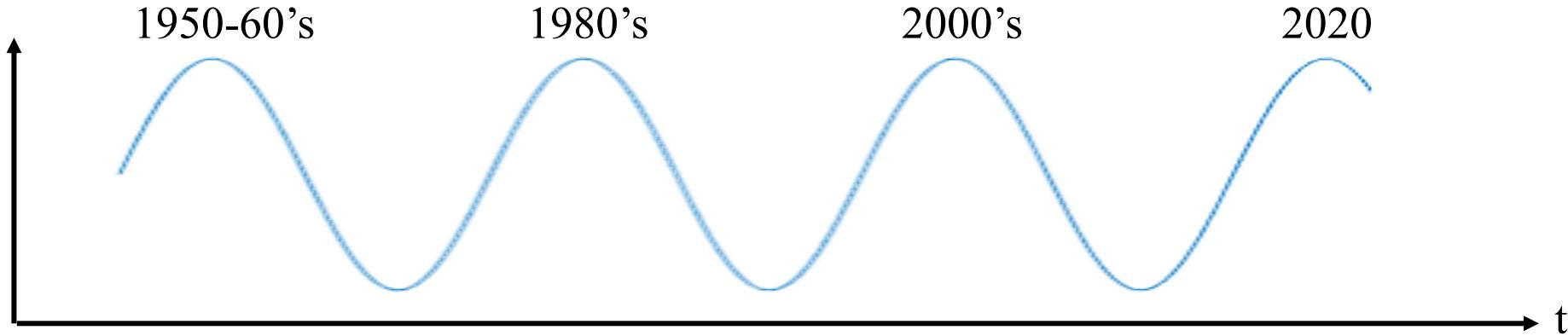
## Lecture 4

# Deep Neural Networks

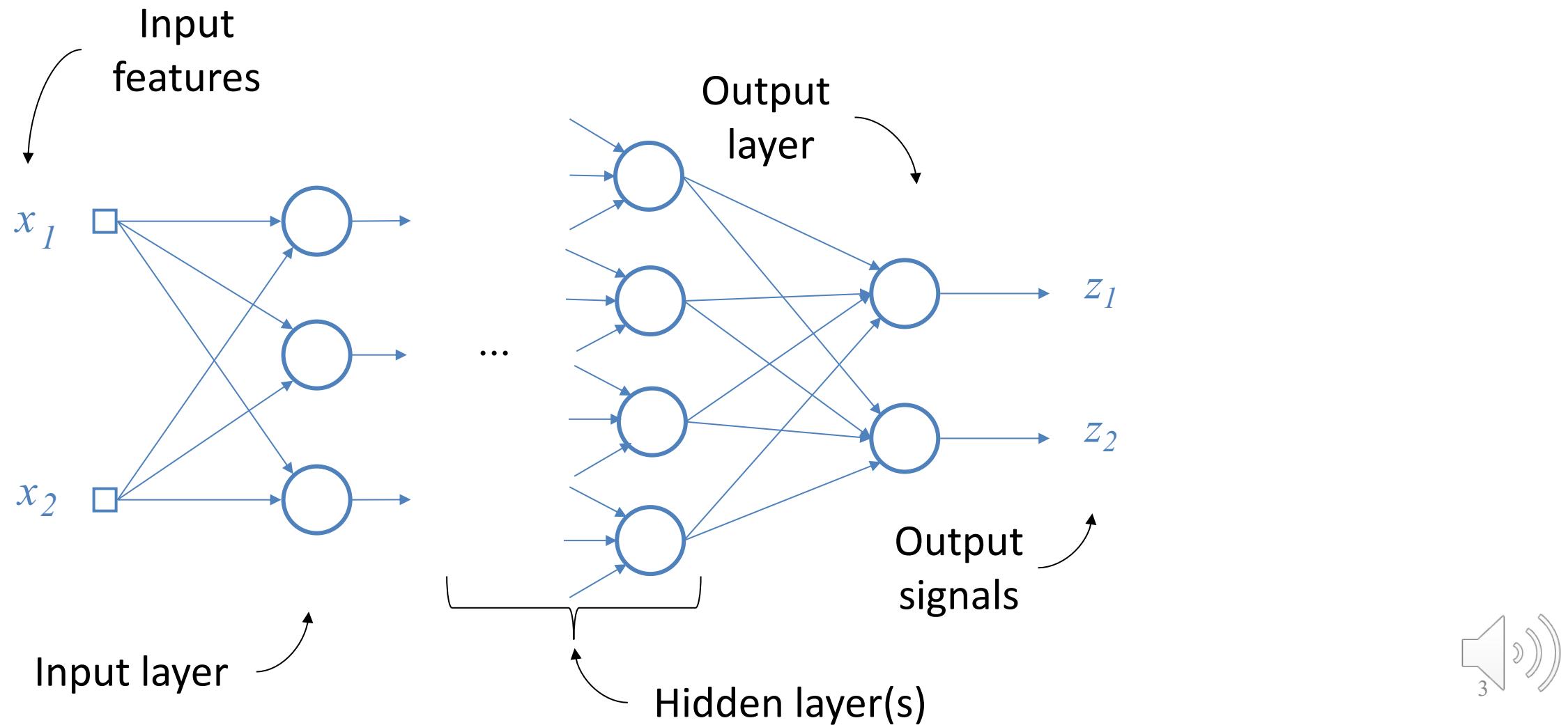
*Magnus Borga*  
*magnus.borga@liu.se*



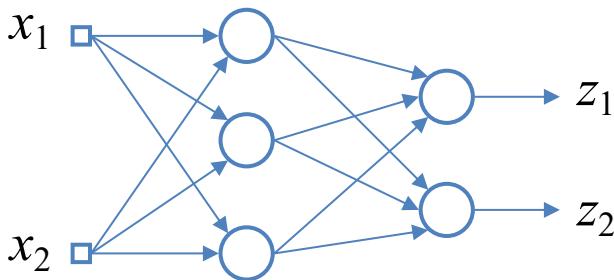
# ML – A history of hypes



# Recap: Neural Networks



# Recap: Error Back Propagation



## Loss function

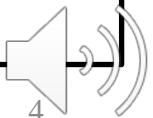
$$\varepsilon(\mathbf{w}) = \sum_{k=1}^K \sum_{m=1}^M [y_{mk} - z_{mk}(\mathbf{w})]^2$$

Annotations for the loss function equation:

- # training examples: Points to the summation over  $k$ .
- # output nodes: Points to the summation over  $m$ .
- all weights: Points to the term  $y_{mk}$ .
- desired output: Points to the term  $y_{mk}$ .
- actual output: Points to the term  $z_{mk}(\mathbf{w})$ .

## Weight update

$$\Delta w_{ij} = -\eta \sum_{k=1}^K \frac{\partial \varepsilon(k)}{\partial w_{ij}}$$



# Why do we need deep networks?

- 1 hidden layer is enough to produce any classification boundary – in theory...
- In practice, shallow networks are limited in terms of how complex functions they can learn
- Complex boundaries can be more compactly represented with many layers - less nodes in total compared to a 2-layer solution.
- A deep network is able to also learn the basic feature extraction

# Challenges in training deep networks

- Huge number of parameters → Extreme degrees of freedom!  
This requires:
  - Huge amounts of training data
  - Powerful computational capacity
  - Ways of reducing the degrees of freedom
  - Regularization
- Basic back-propagation using sigmoid functions does not work in deep networks
  - “Vanishing gradient problem”
  - Instability during training – when parameters change in one layer, the data changes for the following layers

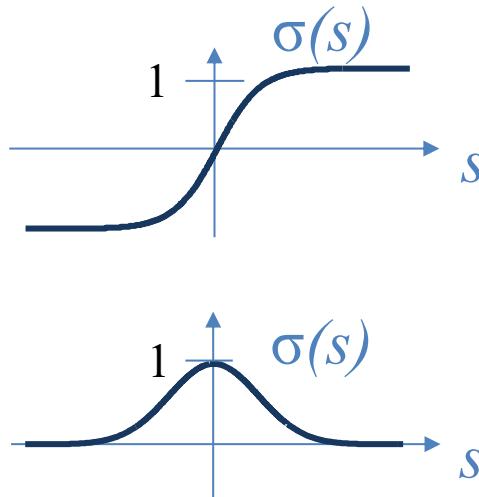


# Vanishing gradient problem

Common sigmoid activation functions such as tanh squeezes the output to a narrow range. This means that the gradient is very small for most input values.

$$\sigma(s) = \tanh(s)$$

$$\sigma' = 1 - \tanh^2(s) = 1 - \sigma^2$$



What happens if we apply the chain rule and multiply small error gradients many times when propagating back thru a deep network?

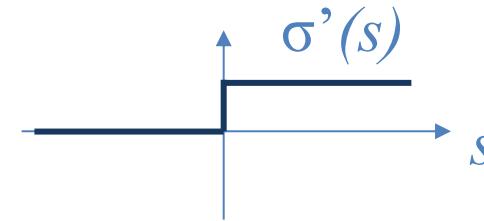
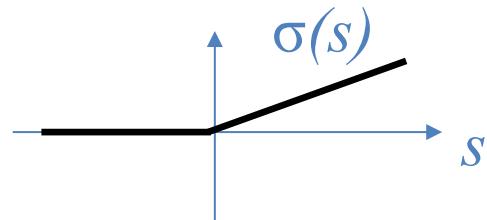


# Vanishing gradient problem

- Error gradients become very small for early layers in the network
- → weights in early layers are not updated!
- Solution: Use an activation function with a derivative that survives the back propagation



# ReLU – Rectifying Linear Unit



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

$$\sigma' = \begin{cases} 0, & s < 0 \\ 1, & s \geq 0 \end{cases}$$

- The ReLU activation function is often used in the hidden layers of deep neural networks to avoid the vanishing gradient problem
- The output layer often uses a Softmax function to be able to represent the probabilities for the different classes



# Instability during training

- Learning works best when the input is centred around zero and has a uniform range
- When parameters in one layer change, the input to the following layers may move away from the optimal range
- The optimization goal becomes a moving target.

# Batch normalization

- For each training batch, normalize the output from the hidden layers

$$\hat{x} = \frac{x - \bar{x}}{\sigma_x}$$

- The output is then scaled by

$$y = \gamma \hat{x} + \beta$$

- $\gamma$  and  $\beta$  are parameters learned by the network

- (It is still debated why batch normalization actually works)



# Drop out

- The large amount of parameters makes deep networks sensitive to overfitting
- One way of avoiding overfitting is to randomly turn off nodes in hidden layers
- This is called drop out
- Drop out prevents the network to rely on each parameter
- This is a kind of regularization of the solution
- Usually implemented as a dropout layer that randomly switch off output from nodes in the previous layer



# Convolutional Neural Networks

- Consider an image of  $512 \times 512$  pixels
- The naive approach is to consider each pixel as an independent input feature
- This gives an input feature vector with 262,144 dimensions
- But image features have some important properties
  - Locality
  - Shift invariance
- These properties are not utilized in this naive approach

# Locality

- Locally, images consist of simple features such as lines, edges and corners
- Objects are built of spatial relations between such local features
- In order to detect such local features, we only have to look in a very small neighbourhood of an image
- On a more global scale, we use more abstract features

# Recap (Le 1): Features – Iris dataset



Iris setosa



Iris versicolor



Iris virginica

Fisher's Iris Data

Sepal Length	Sepal Width	Petal Length	Petal Width	Species
5.7	4.4	1.5	0.4	<i>I. setosa</i>
5.8	2.6	4.0	1.2	<i>I. versicolor</i>
5.8	2.7	5.1	1.9	<i>I. virginica</i>

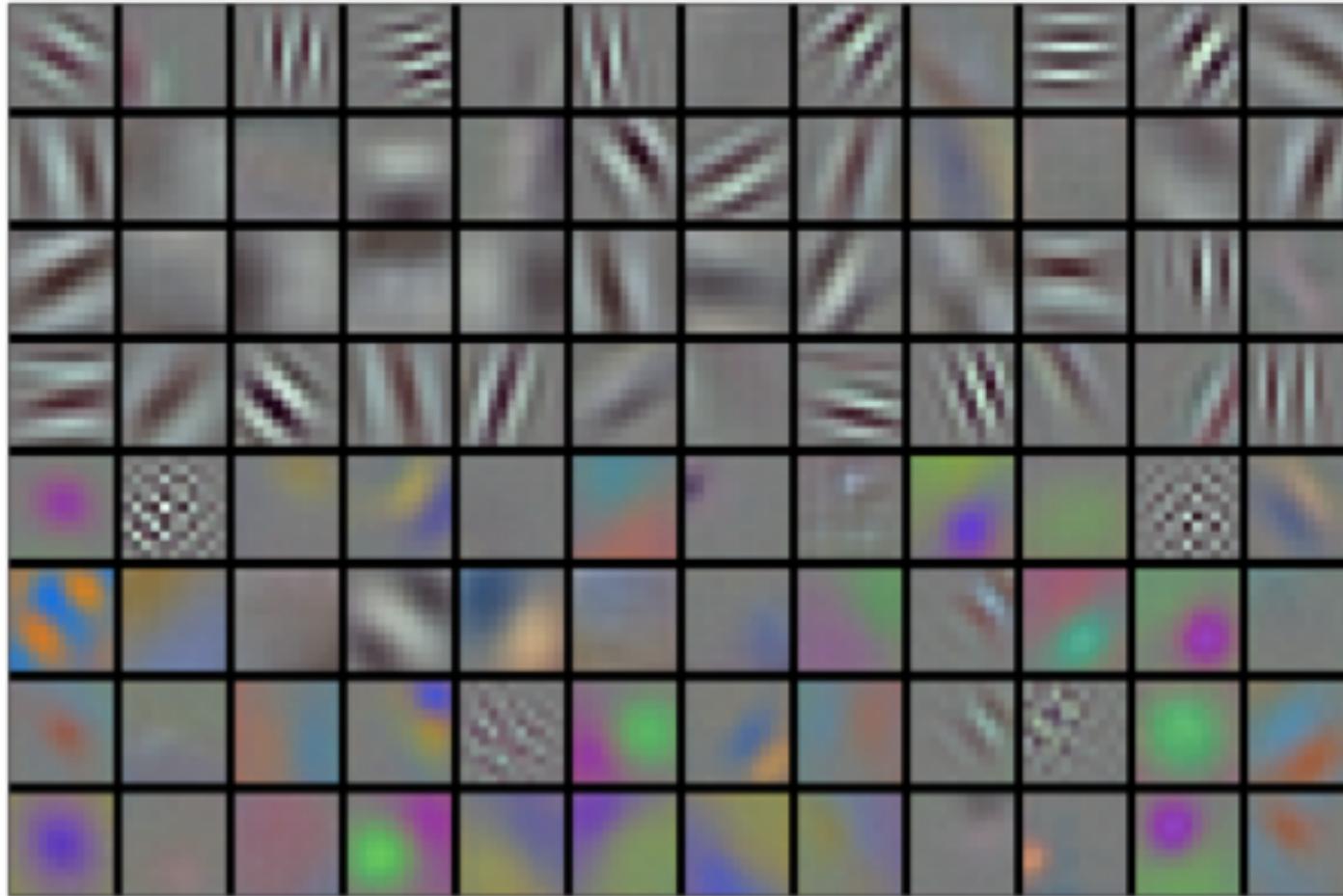
From Wikipedia

Feature vectors:

$$\mathbf{x}_1 = \begin{bmatrix} 5.7 \\ 4.4 \\ 1.5 \\ 0.4 \end{bmatrix} \quad \mathbf{x}_2 = \begin{bmatrix} 5.8 \\ 2.6 \\ 4.0 \\ 1.2 \end{bmatrix} \quad \mathbf{x}_3 = \begin{bmatrix} 5.8 \\ 2.7 \\ 5.1 \\ 1.9 \end{bmatrix}$$



# Primitive image features



# Shift invariance

- The description of an image feature is shift invariant
  - A corner, or a face, looks the same no matter where in the image it is located
  - While the naive approach can learn to detect an object in a certain position, it cannot recognize that object if it is moved in the image, because that would completely change the input vector
- Hence, we want the training of feature detectors to be independent on position, i.e. shift invariant



# Convolution

- In convolution a kernel is shifted over the image
- For each position, the scalar product between the kernel and the corresponding image patch is computed
- $(f * g)(x, y) = \sum_{i,j} f(x - i, y - j)g(i, j)$
- Remarks:
  - In neural networks, the kernel is not reversed, so it is actually a correlation, not a convolution
  - Note that the result is only well defined inside the image, where the complete kernel covers the image, i.e. the result is smaller than the image. This can be handled by padding the image before convolution



# Convolution – 1D Example

Input data

0	1	2	1	2	0	1	2	0
---	---	---	---	---	---	---	---	---

Convolution kernel

\*

-1	2	-1
----	---	----

Resulting data



--	--	--	--	--	--	--

# Convolution – Images

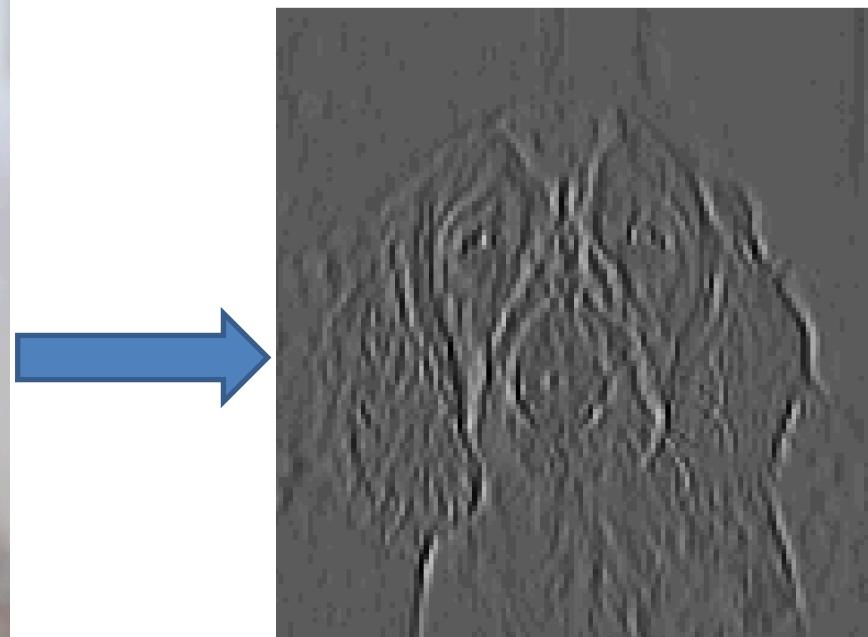
Convolution  
kernel

-1	2	-1
-1	2	-1
-1	2	-1

Image



Result (after down sampling)

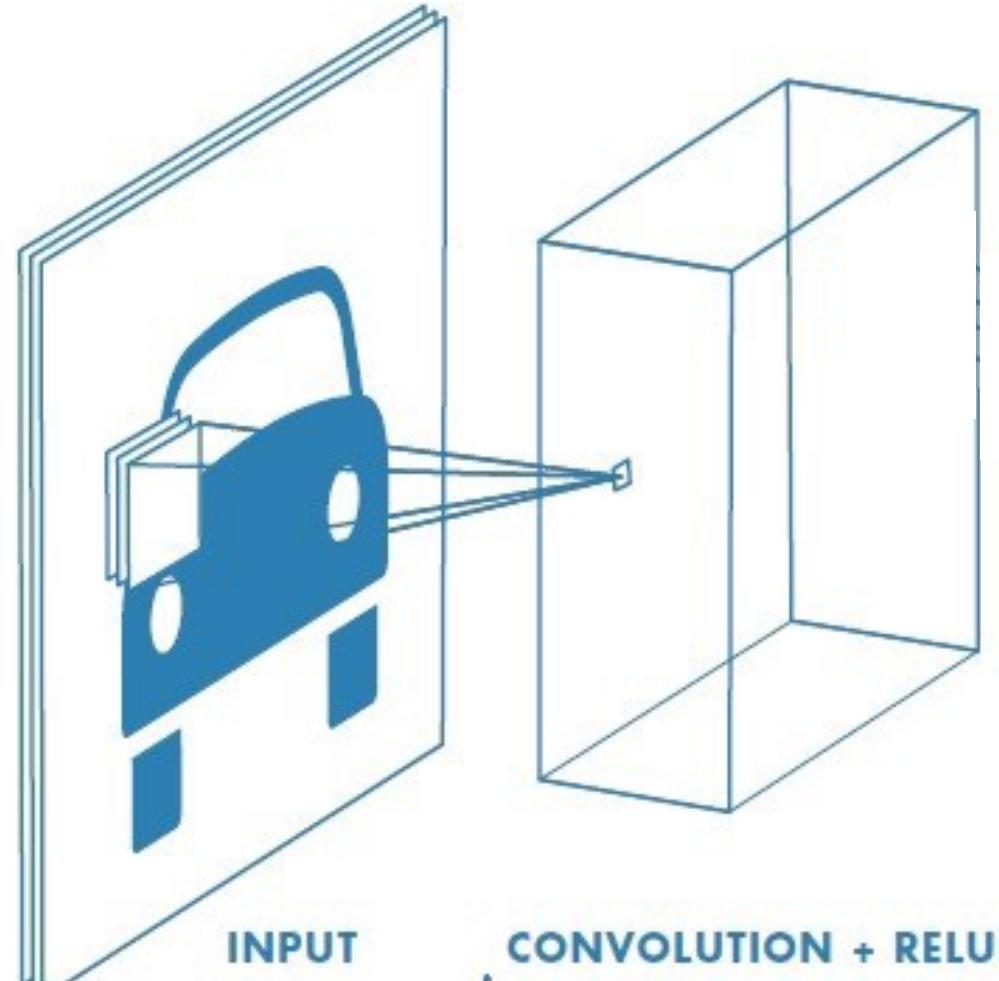


# Convolutional networks

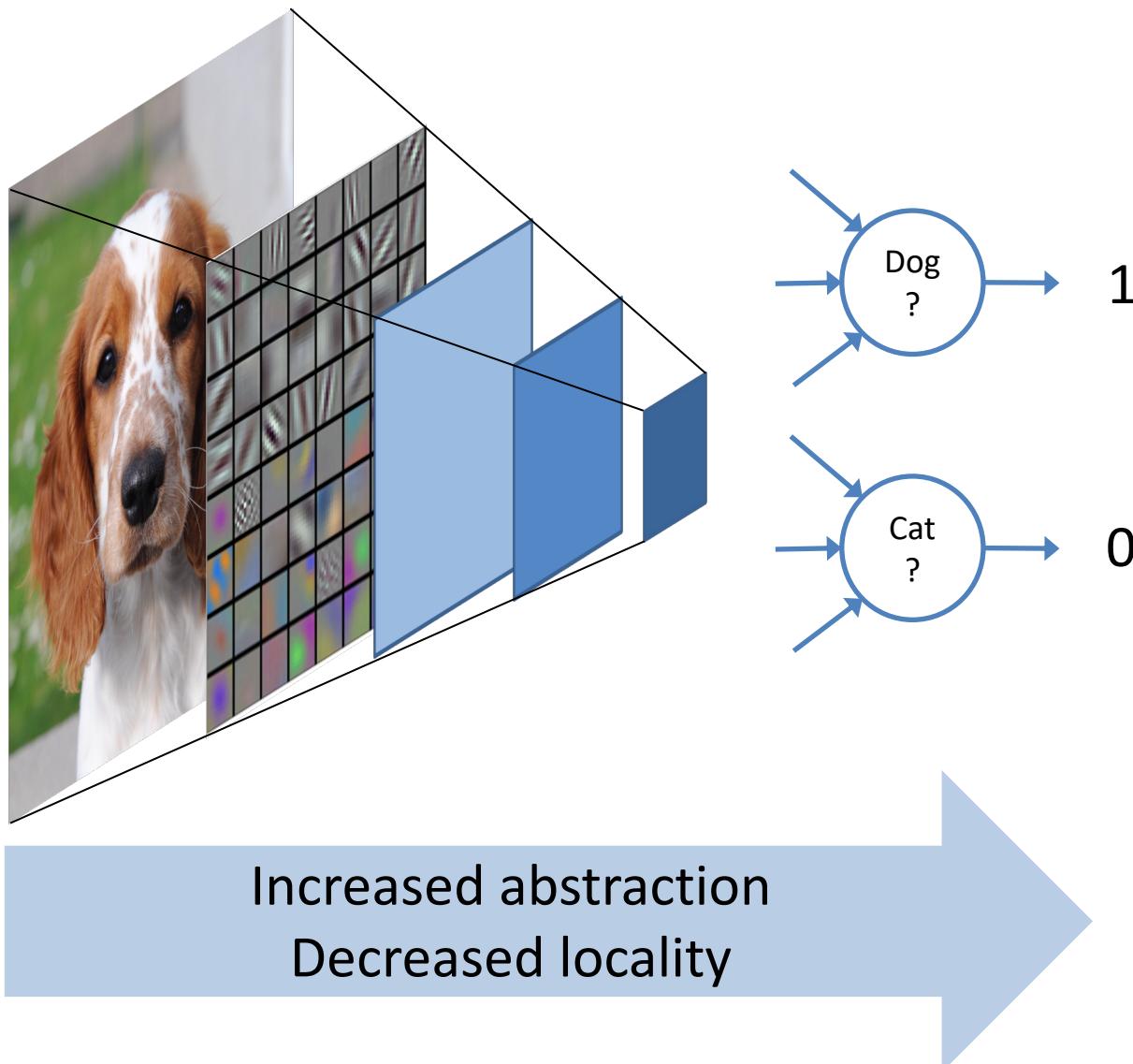
- In a convolutional layer, only the coefficients in the kernels need to be learned by the network
- The kernels are much smaller than the image
  - This reduces the number of parameters to learn and thereby the risk of overfitting
- The kernel is
  - spatially local
  - shift invariant

# Convolutional networks

- Each layer consists of multiple kernels each generating a separate feature image
- Each kernel has an additional dimension to cover all channels, e.g. colour channels
- The output of each layer is a stack of such feature images

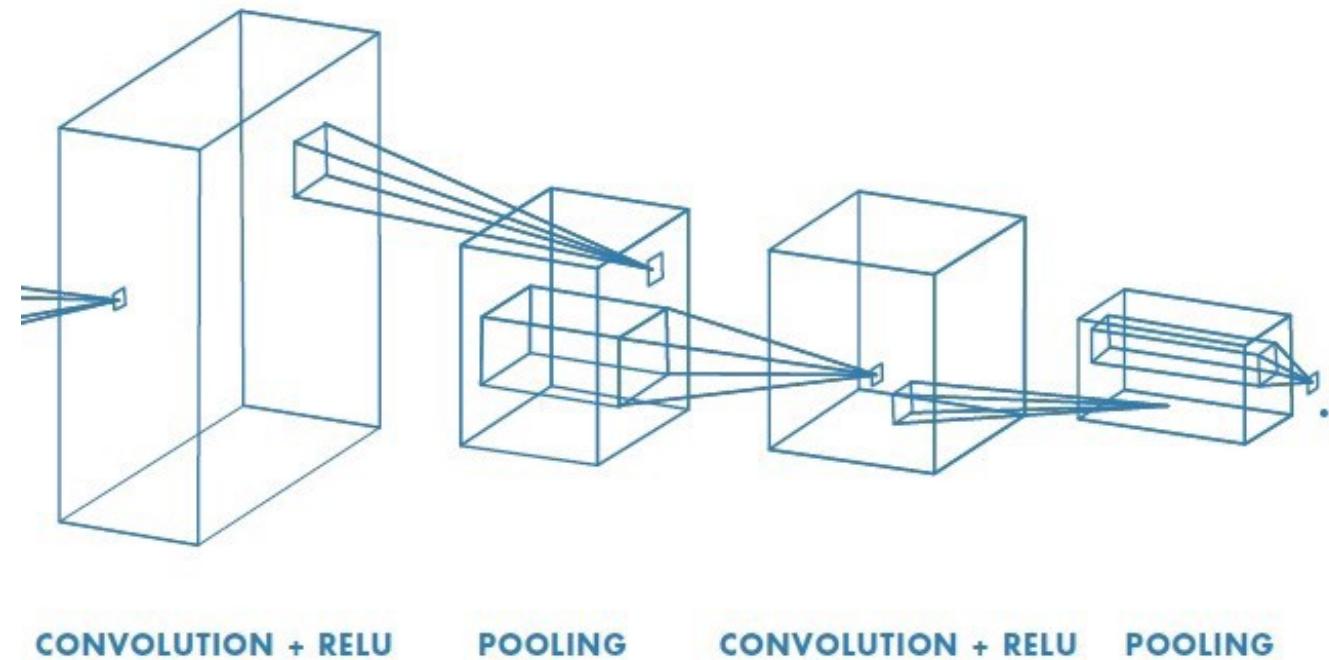


# Abstraction hierarchy



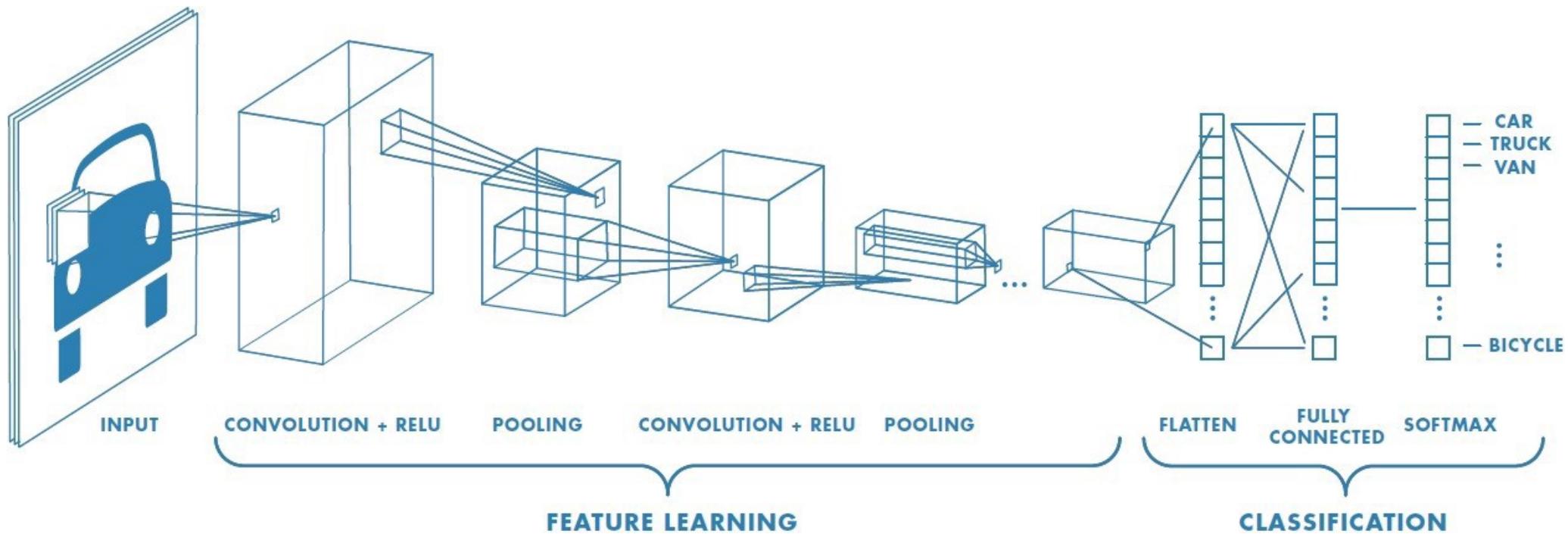
# Pooling

- In order to decrease locality and increase translation invariance as information propagates forward thru the network, the output from previous layer can be pooled together in a pooling layer
- Pooling combines output from nearby units
  - Mean pooling – average of the units
  - Max pooling – max value of the units
- Note – Pooling is done separately for each feature image

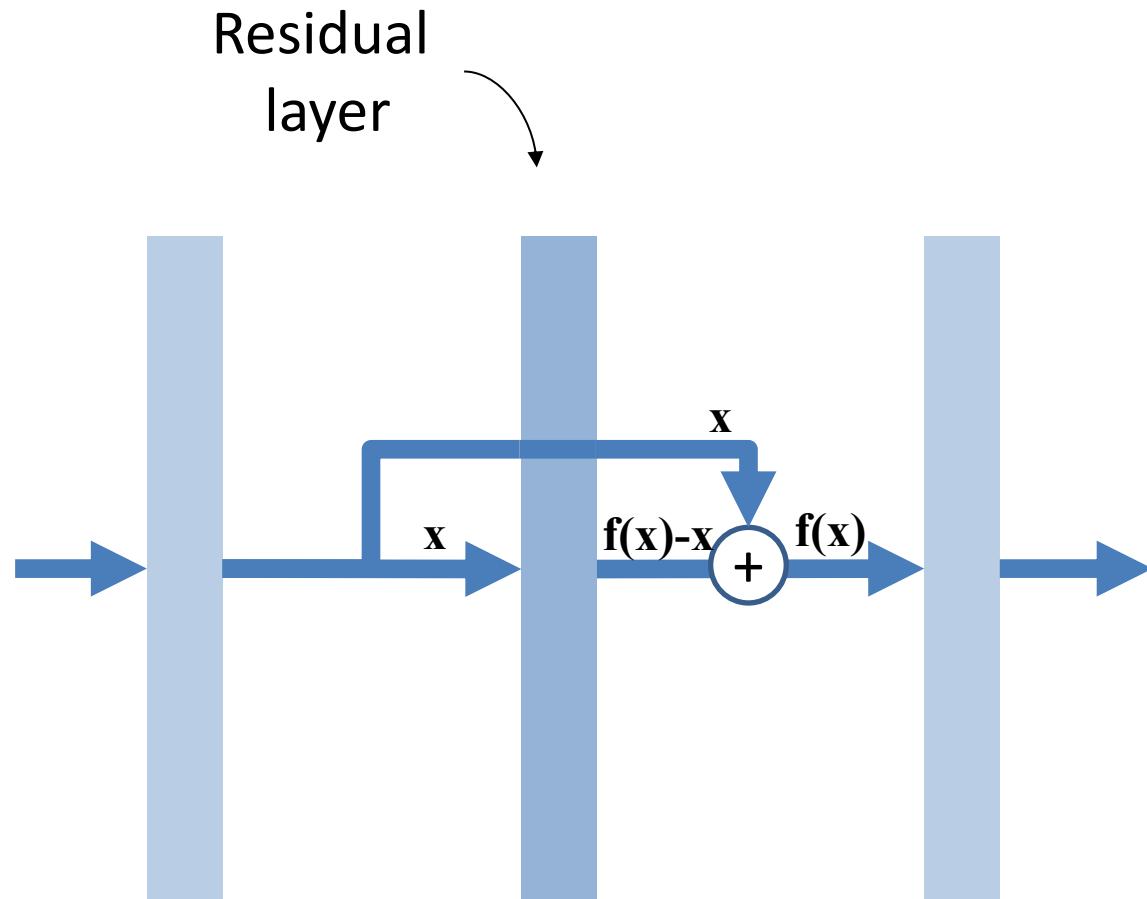


# Convolutional networks

- Early layers learn useful features
- Later layers learn how to combine features for classification



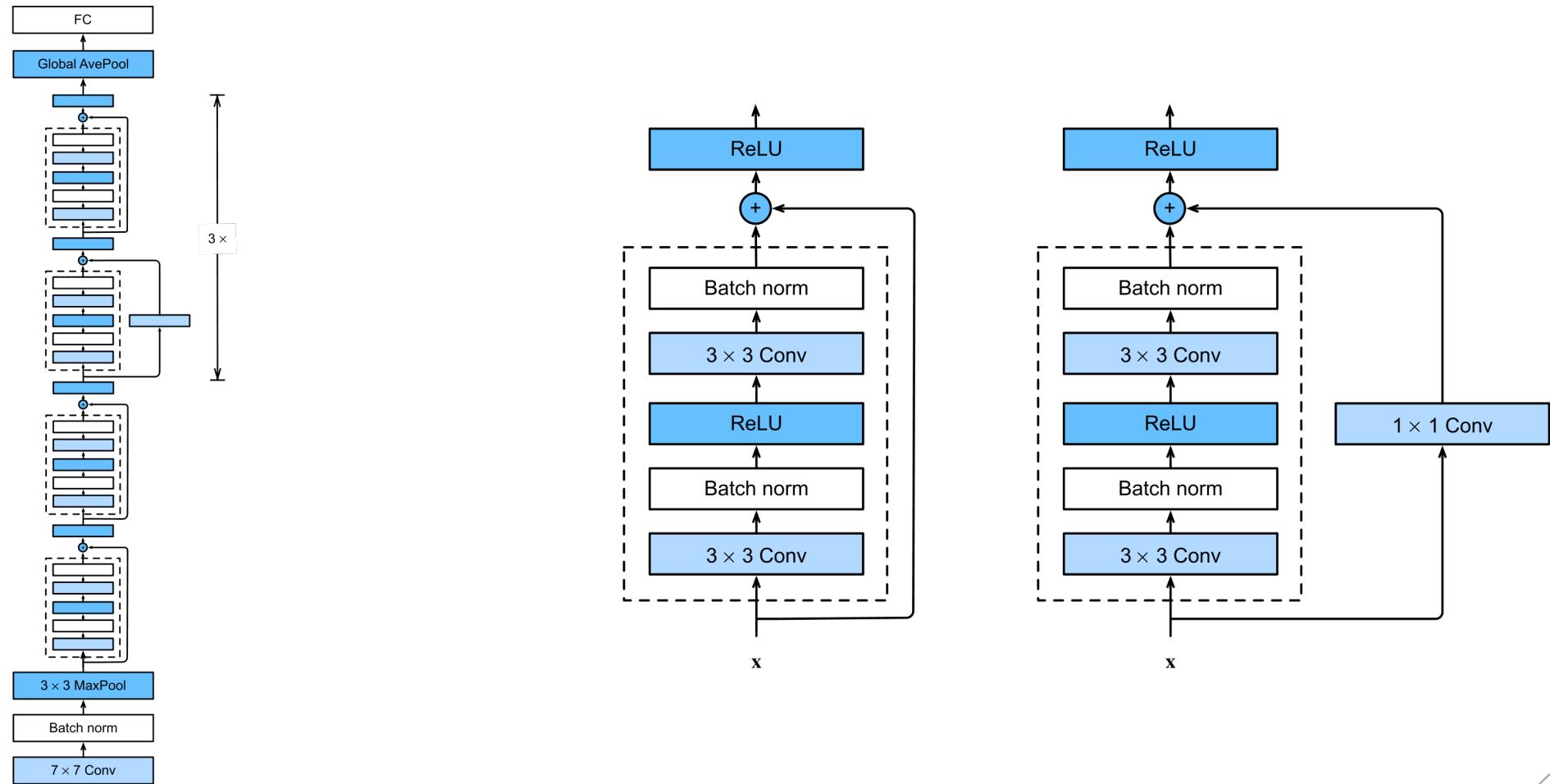
# Residual Networks



- Intuition: The identity function should be the simplest to learn
- The residual layer only needs to learn the residual  $f(x) - x$
- Early in the training, the short cuts make the network act as a more shallow network
- Easier to back propagate the error in very deep networks



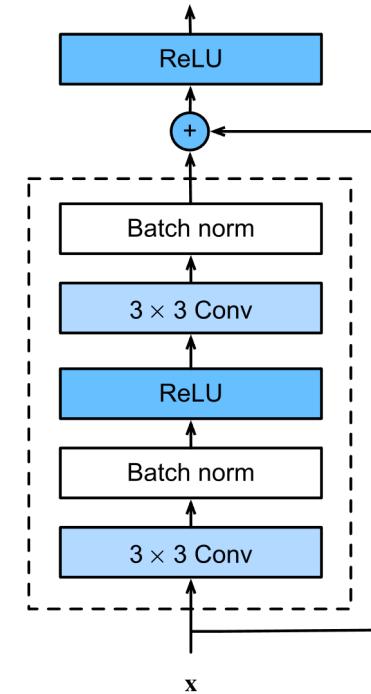
# ResNet



Images from D2l.ai

# Blocks

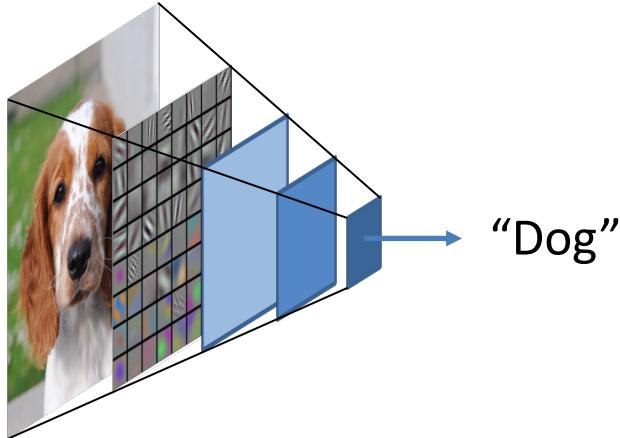
- To simplify the design of complex deep neural networks, they are based on several blocks
- Each block can consist of one or several layers
- See e.g. ResNet
- This enables design of networks with 100's of layers



A block in ResNet

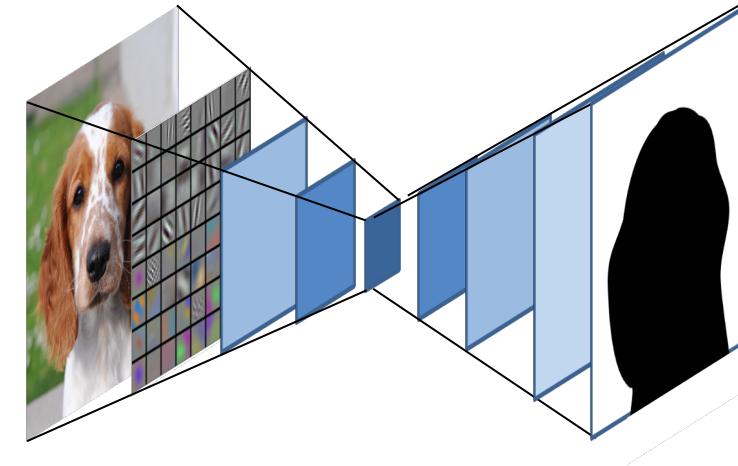
Image from D2l.ai

# Image segmentation



## Image classification

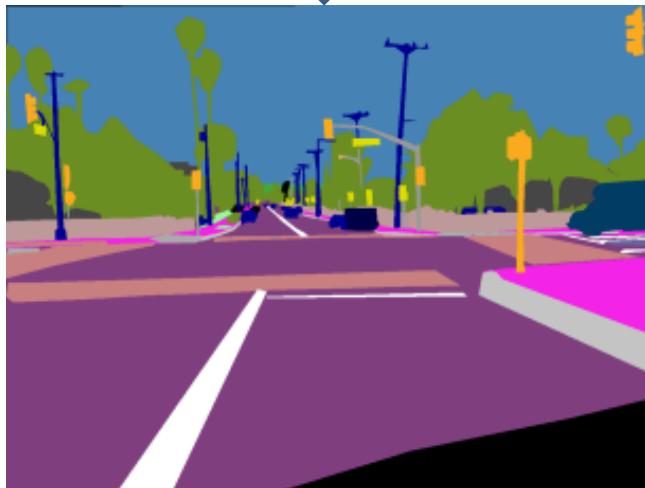
- "What does this image depict?"
- High-dimensional input
- Low-dimensional output



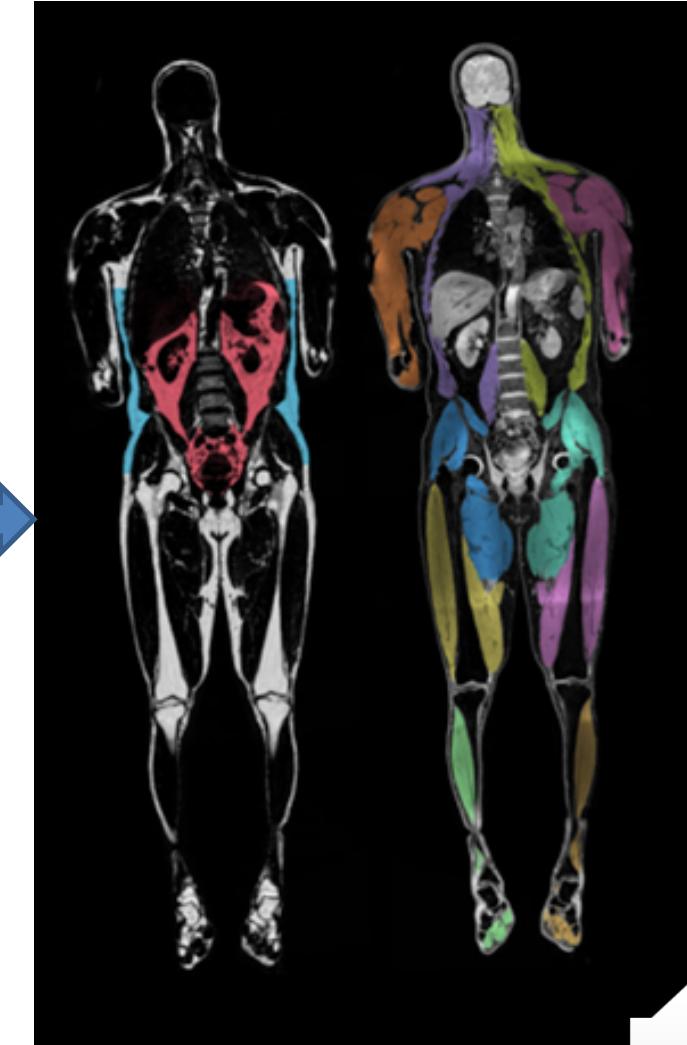
## Image segmentation

- Where is the dog?
- Classification of each pixel
- Foreground / Background
- Output dimensionality same as input

# Image segmentation – Examples

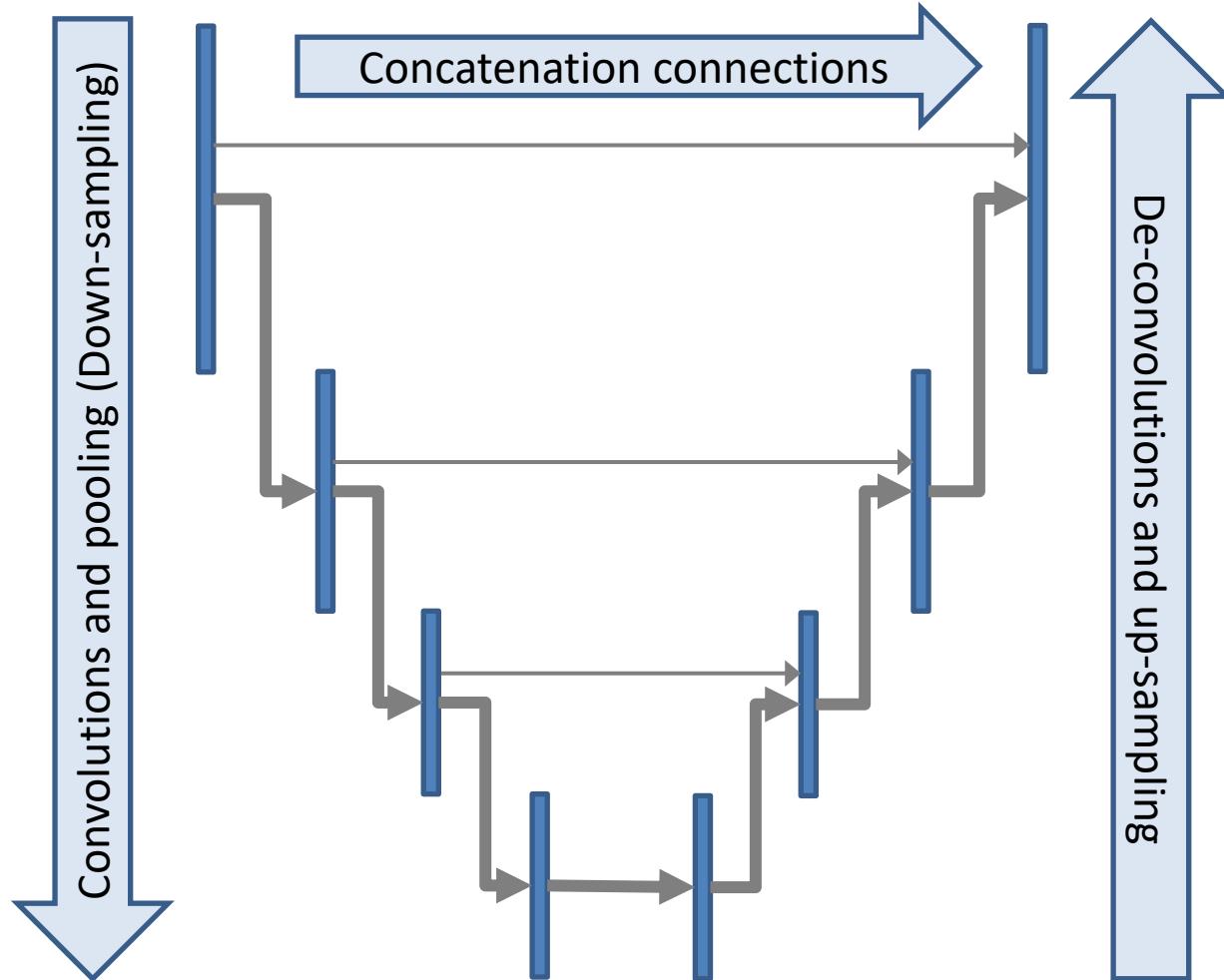


groundai.com



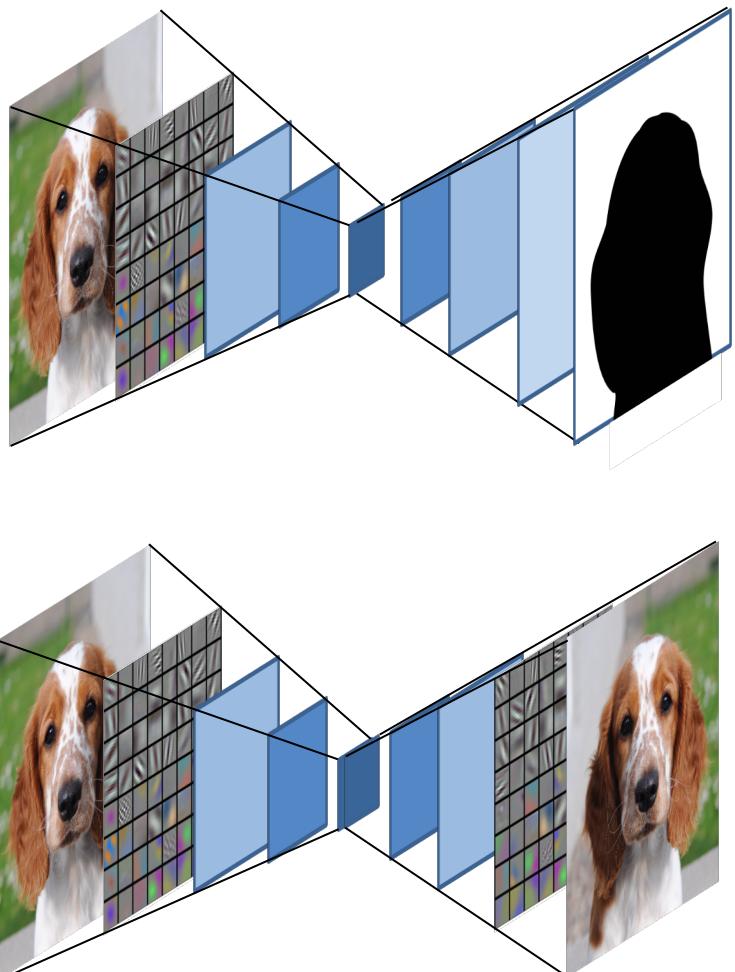
AMRA Medical

# U-Net



- Convolutional network for image segmentation
- The first part is a normal contracting network with convolutions and pooling
- The second, expanding part performs up-sampling and de-convolution
- Features from the contracting part are concatenated with up-sampling results in the expanding part

# Auto Encoders



## Image segmentation

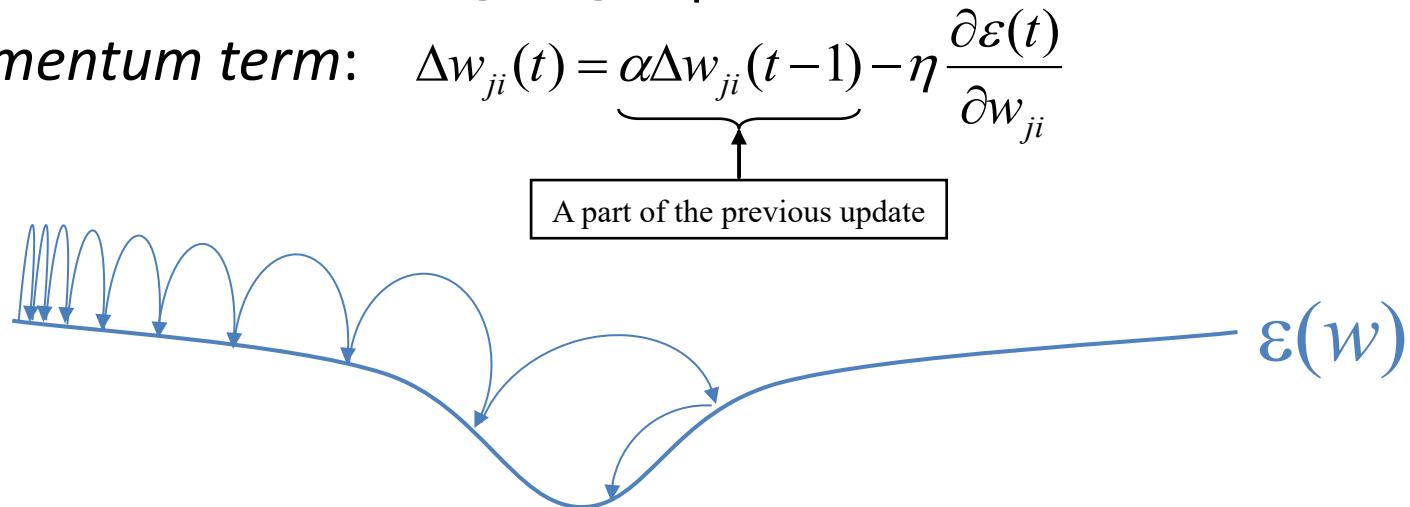
- Training data  $\{x_i, y_i\}$
- Binary output (foreground / background)

## Auto encoder

- Training data  $\{x_i, x_i\}$
- Does not require labelled training data – Unsupervised learning!
- The smallest layer must be able to represent all training data
- The bottle neck forces an efficient coding of the original data (images)

# Faster convergence

- Normalize input features, e.g. batch normalization
- Change to another activation function, e.g. ReLu
- Residual layers – skip connections
- Separate and adaptive step length  $\eta$  for each weight:
  - If the derivative has the same sign in several consecutive steps,  $\eta$  should increase. If the derivative change sign,  $\eta$  should decrease.
- Introduce a *momentum term*: 
$$\Delta w_{ji}(t) = \underbrace{\alpha \Delta w_{ji}(t-1)}_{\text{A part of the previous update}} - \eta \frac{\partial \varepsilon(t)}{\partial w_{ji}}$$



# Deep Neural Networks – Summary

- Based on the multi-layer perceptron and back propagation
- Additional smart tricks
  - ReLU activation function
  - Batch normalization
  - Weight sharing e.g. convolutional networks
  - Skipping layers – Residual Networks
  - Larger building blocks
  - ...
- Requires huge amounts of training data and lots of computational power!

