

# Deep Learning II: Advanced Neural Network Architectures

8DM50 Machine Learning in Medical Imaging and Biology

Jelmer Wolterink

28-09-2020

# Deep Learning II

## Me

- Assistant professor @ University of Twente
- Deep learning for medical image analysis - cardiovascular image analysis (CT, MR)

## Today

1. **Advanced neural network architectures**
2. Interpretability and generative adversarial networks
3. Practical assignment in Keras

# In this lecture

## Convolutional neural networks

- Recap
- Advanced architectures

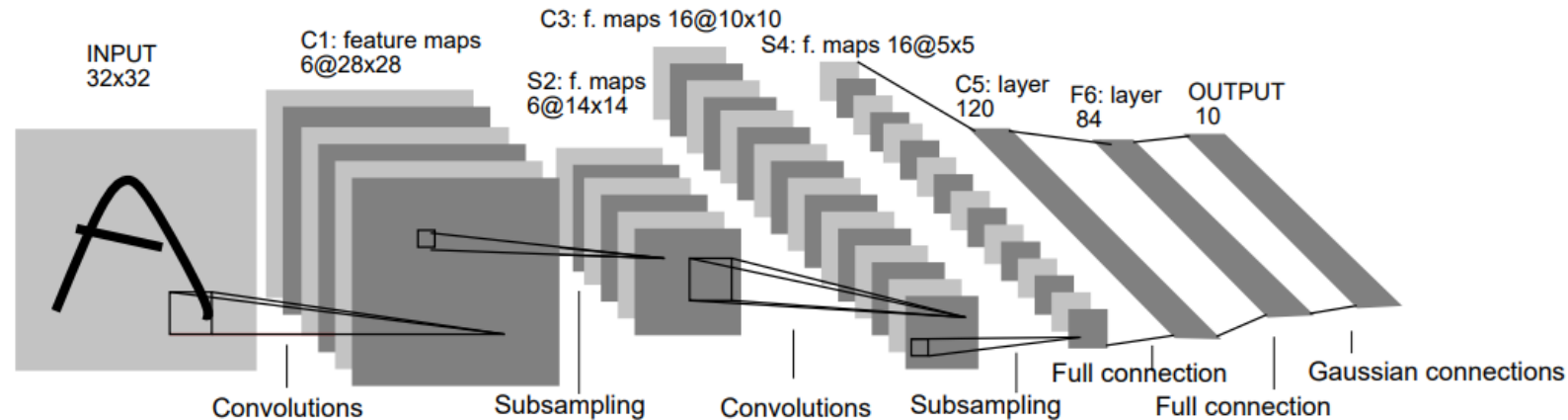
## Neural networks for sequential data

- Recurrent neural networks
- Long short term memory (LSTM) units

## CNNs for pixelwise prediction

- Patch-based segmentation
- Encoder-decoder architectures

# Recap: Convolutional neural networks



[Demo](#)

A standard convolutional network consists of

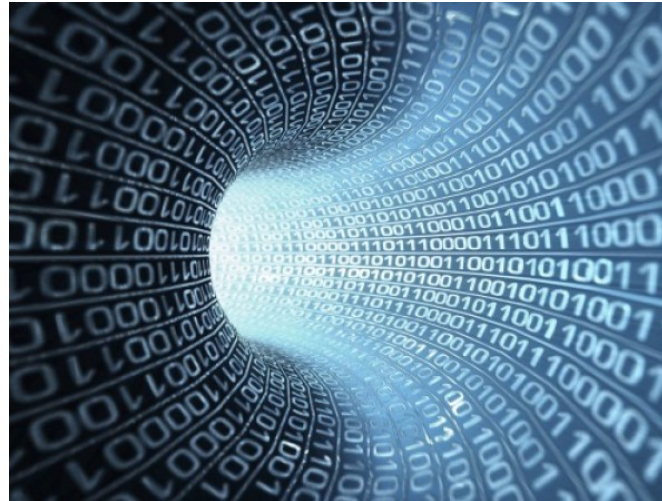
- **Convolutional** layers transform input into feature maps
- **Subsampling** operations reduce size of feature maps, e.g. max pooling
- **Fully-connected** layers perform classification (multi-layered perceptron)



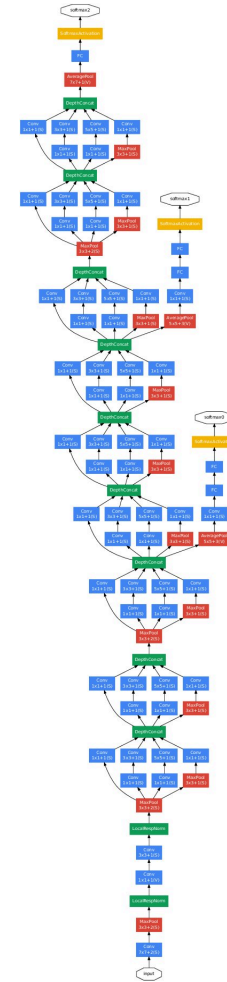
# What happened between 1998 and now?



Compute



Data



Algorithms

# Data: ImageNet challenge

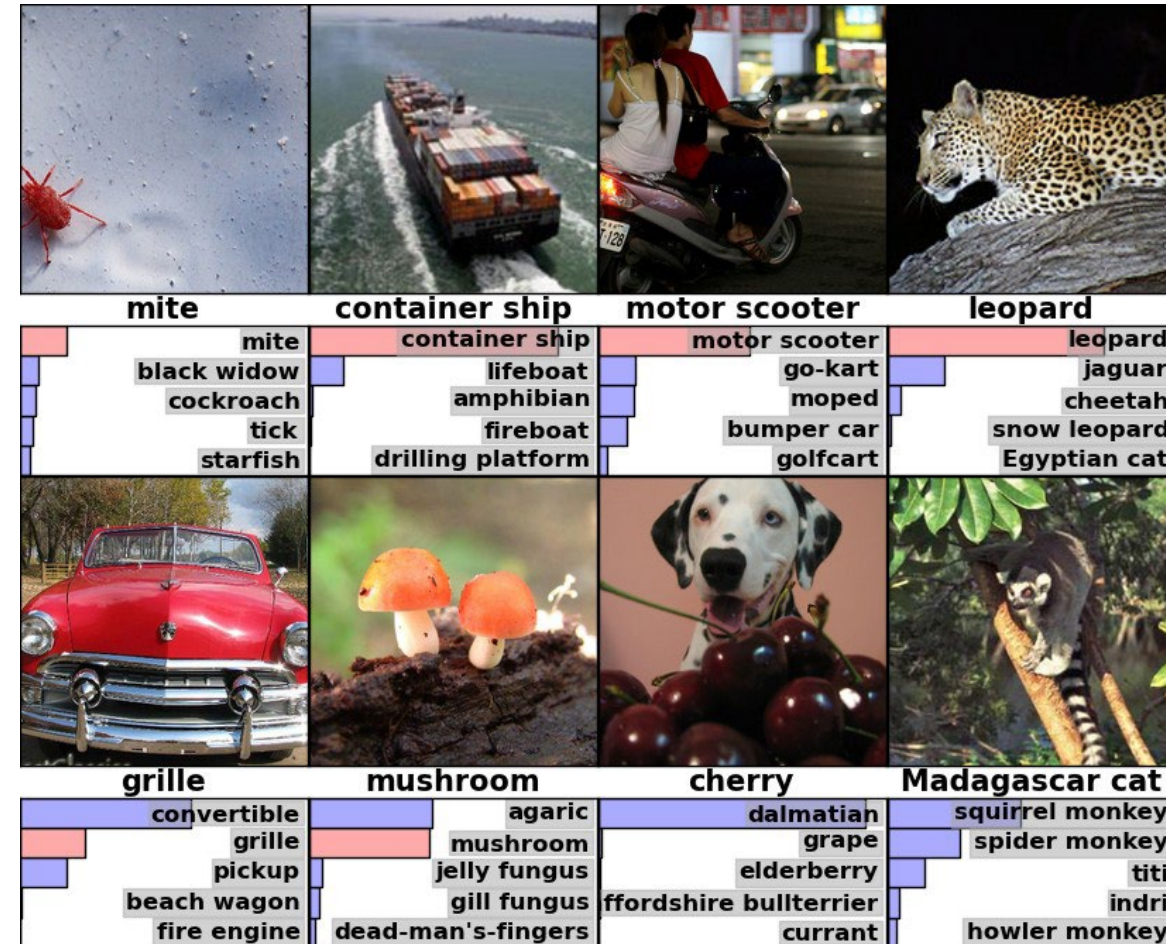
Benchmark for image classification/object detection

## Data

- > 1,200,000 RGB images
- Images show one of **1000** classes

## Task

- Detect label of image
- Top-1\top-5 accuracy



- Substantially outperformed 'conventional' methods in 2012
- Convolutional + subsampling + fully connected layers
- Trained in parallel on two GPUs
- Training time
  - **2012**: 5 to 6 days (2 x GTX 580 3GB GPU)
  - **2017**: 24 minutes (supercomputer 32,000 cores)
- Large 11 x 11 convolution kernels

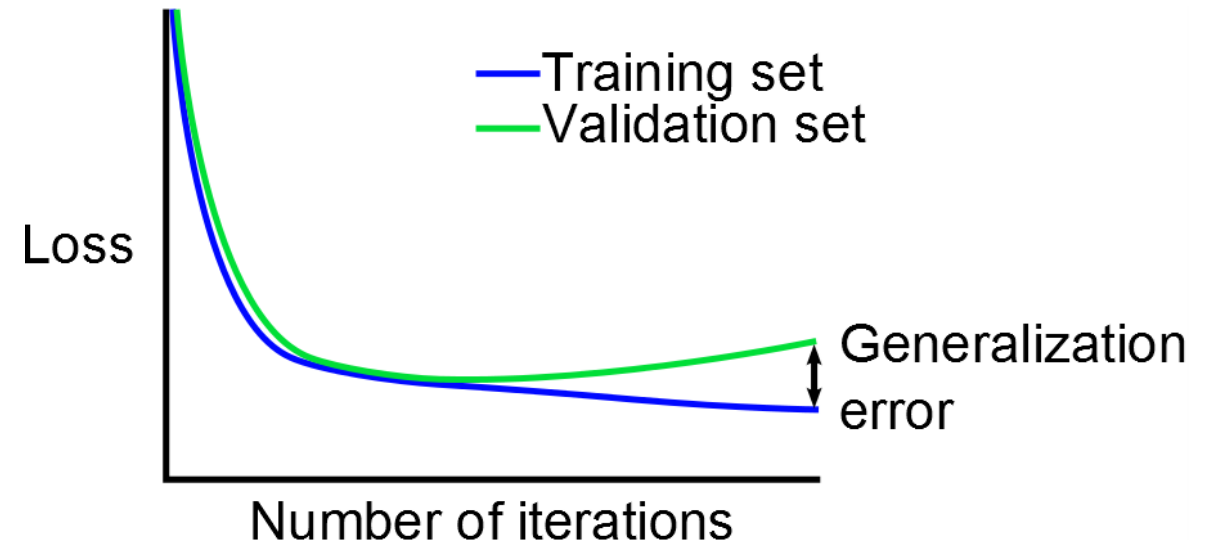
# Overfitting

## Reasons

- Too many parameters
- Not enough data

## Solution

- Reduce number of parameters



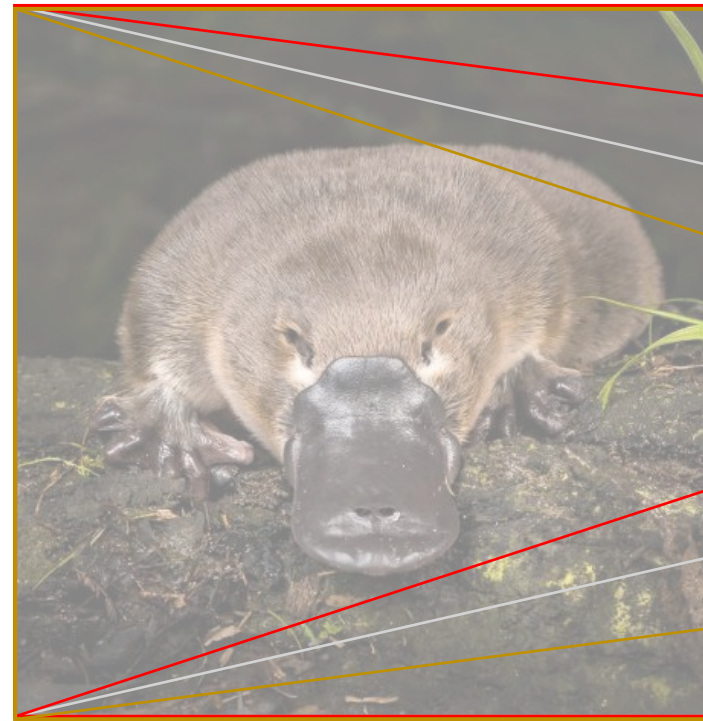


# Kernel size

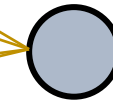
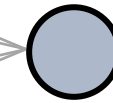
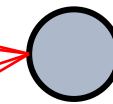
100

100

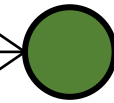
10,000 weights



Input



Hidden



Output

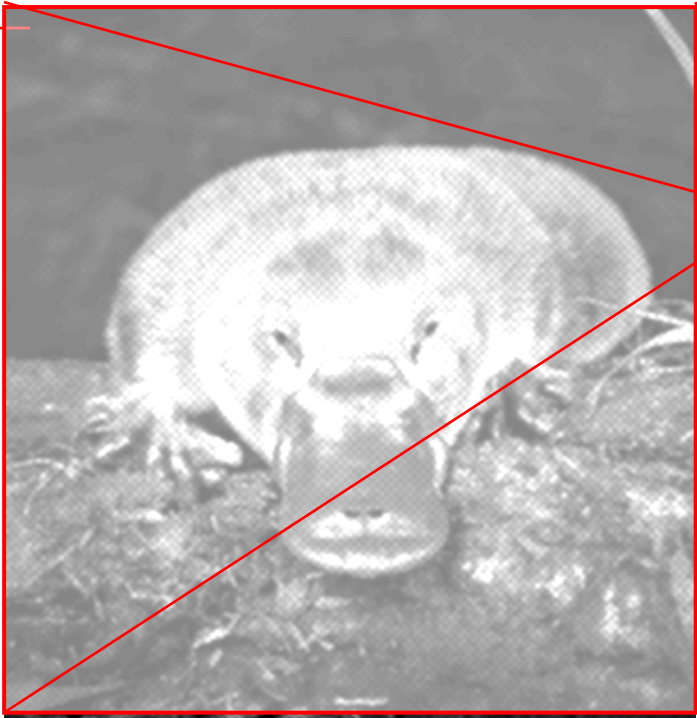
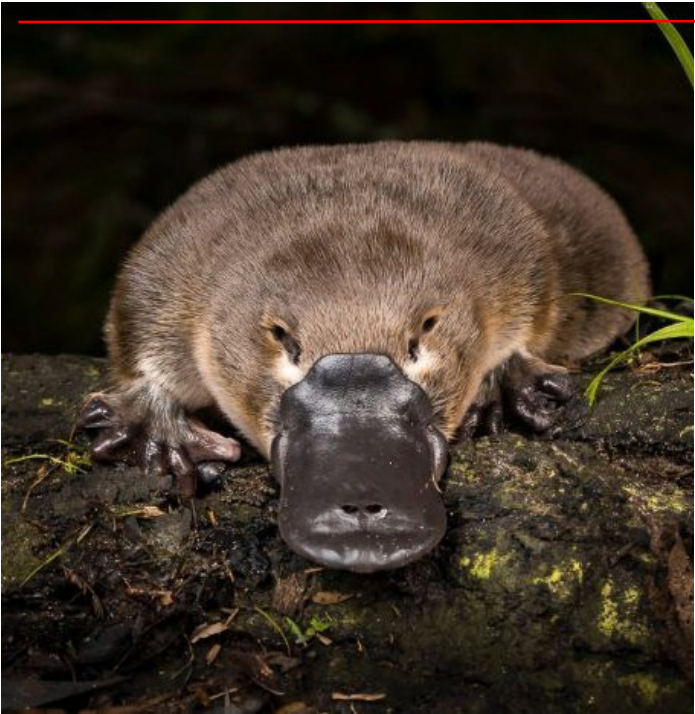
Platypus

# Kernel size

100

1 weight

100



Platypus

Input

Hidden

Hidden

Output

# Kernel size

100

11 X 11 = 121 weights

100



Input



Hidden



Hidden

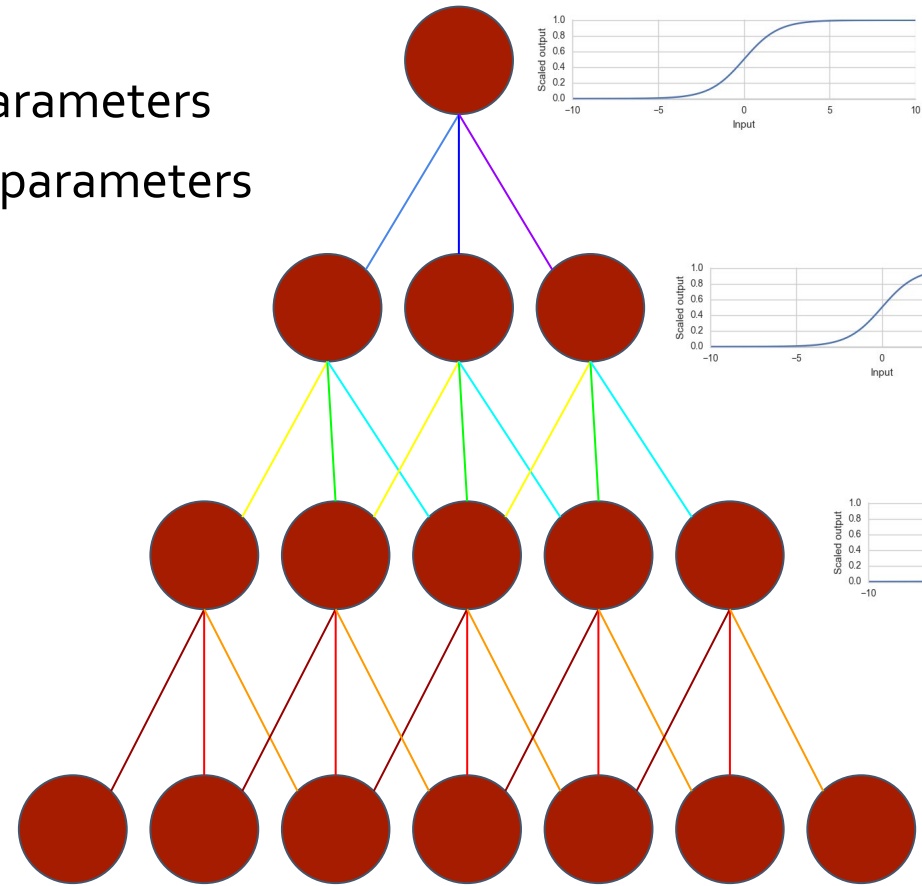
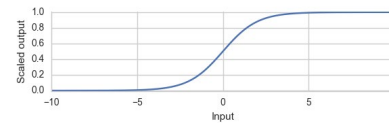
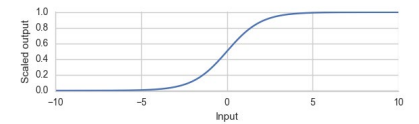
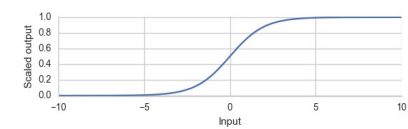
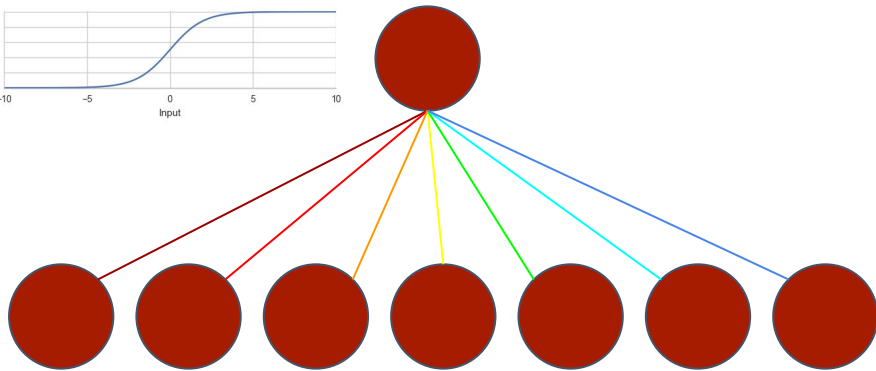
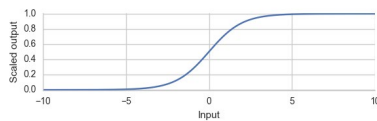


Platypus

Output

# Using smaller kernels

- Large kernels have many parameters:  $7 \times 7 = 49$  parameters
- Smaller kernels reduce parameters:  $3 \times (3 \times 3) = 27$  parameters
- More nonlinearities means more abstraction



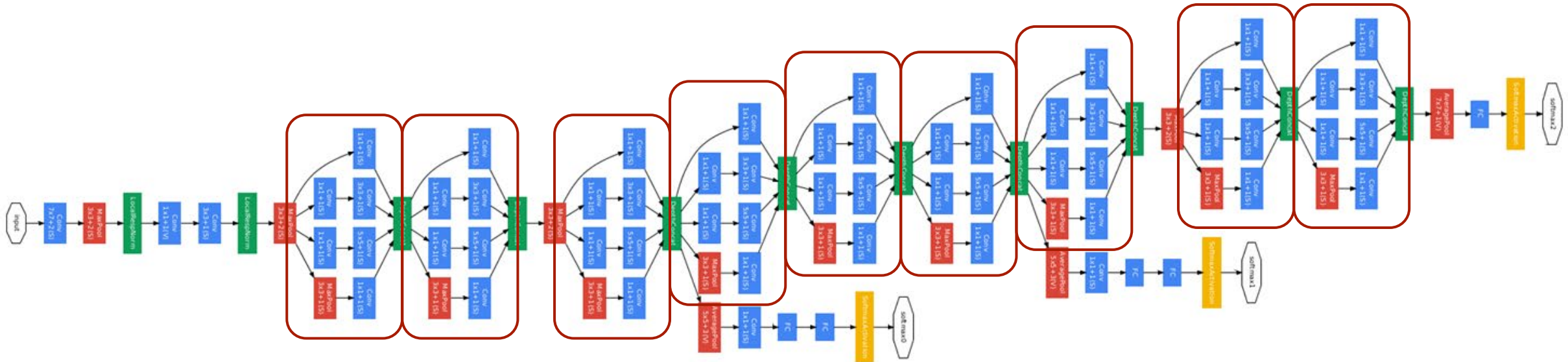


The diagram illustrates a deep convolutional neural network (CNN) architecture for image classification. It starts with an input image of size  $224 \times 224 \times 3$ . This is followed by a convolutional layer with  $224 \times 224 \times 64$  filters. The next layer is a max pooling layer, resulting in  $112 \times 112 \times 128$  feature maps. This is followed by another convolutional layer with  $56 \times 56 \times 256$  filters. The next layer is a max pooling layer, resulting in  $28 \times 28 \times 512$  feature maps. This is followed by a convolutional layer with  $14 \times 14 \times 512$  filters. The next layer is a max pooling layer, resulting in  $7 \times 7 \times 512$  feature maps. This is followed by a fully connected layer with  $1 \times 1 \times 4096$  nodes. The final layer is a softmax layer with  $1 \times 1 \times 1000$  nodes. A legend indicates the following operations: convolution+ReLU (white box), max pooling (red box), fully connected+ReLU (blue box), and softmax (orange box).

	D	
	16 weight layers	
ge)		
	conv3-64 conv3-64	
	conv3-128 conv3-128	
	conv3-256 conv3-256 <b>conv3-256</b>	
	conv3-512 conv3-512 <b>conv3-512</b>	
	conv3-512 conv3-512 <b>conv3-512</b>	

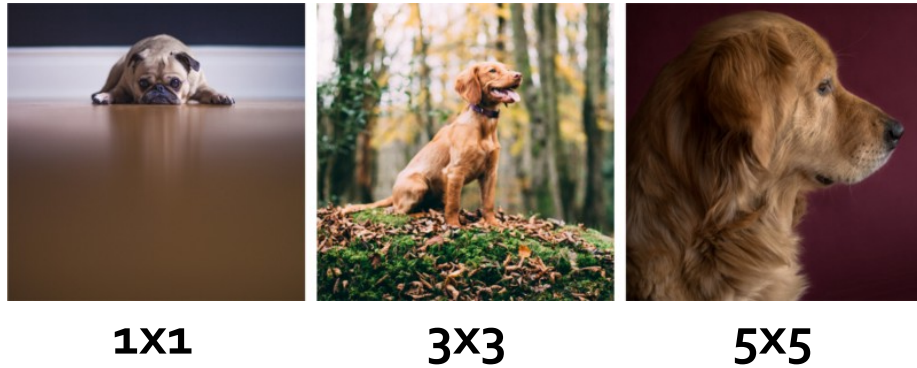
# GoogLeNet (Inception v1)

- 22 layer-network
- Very deep compared to LeNet/AlexNet
- SOTA on ImageNet (when published)

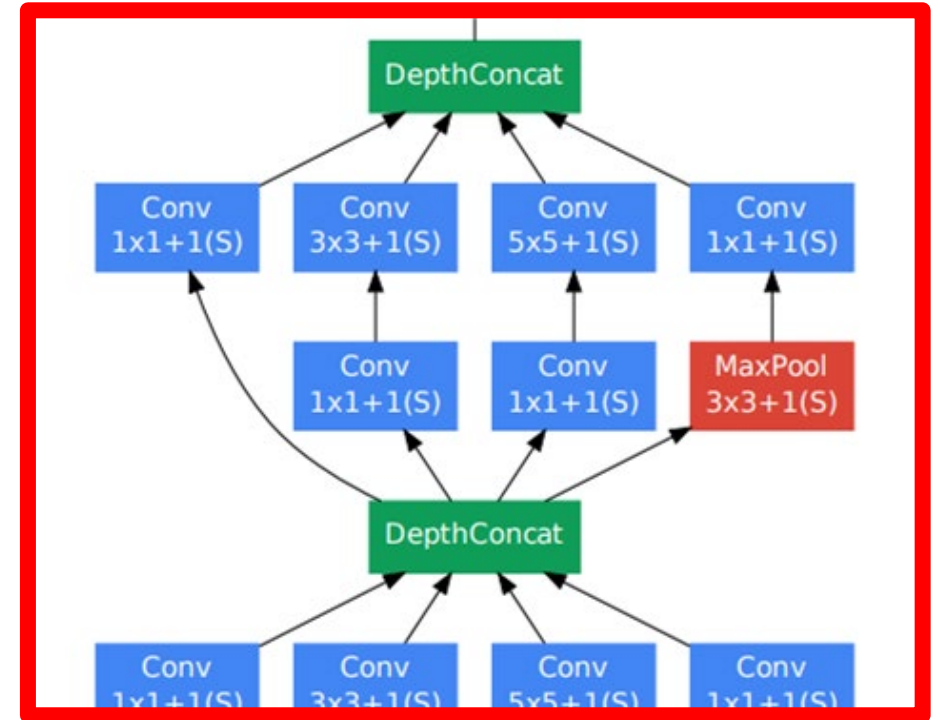


# Inception module

- Combine parallel multi-scale convolutions
- Let the model pick best filter size



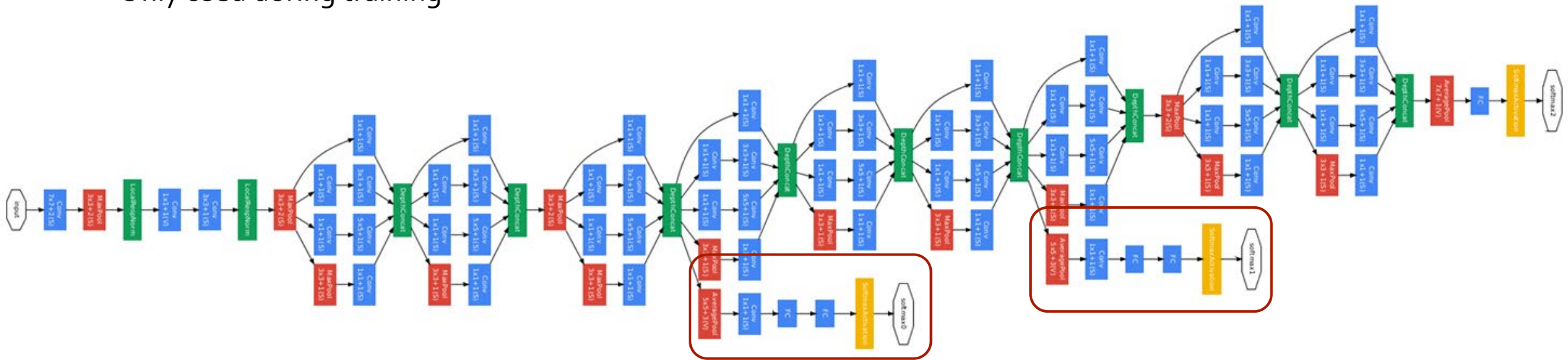
- Bottleneck layers: 1 x 1 convolutions
  - Aggregate feature maps
  - Prevent explosion in number of parameters



# Auxiliary classifiers

Auxiliary classifiers provide extra supervision

- Vanishing gradients
- Enforce useful features at intermediate layers
- Only used during training





# Residual connections

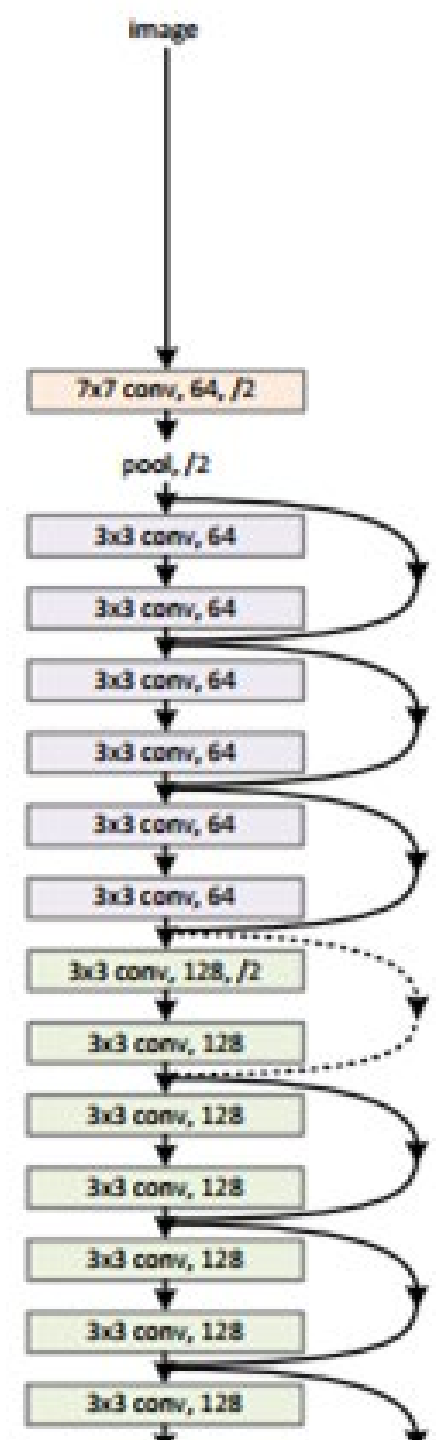
Very deep networks

- allow learning of better representations
- are difficult to optimize due to vanishing gradients

Residual connections can skip layers  $H(x) = x + F(x)$

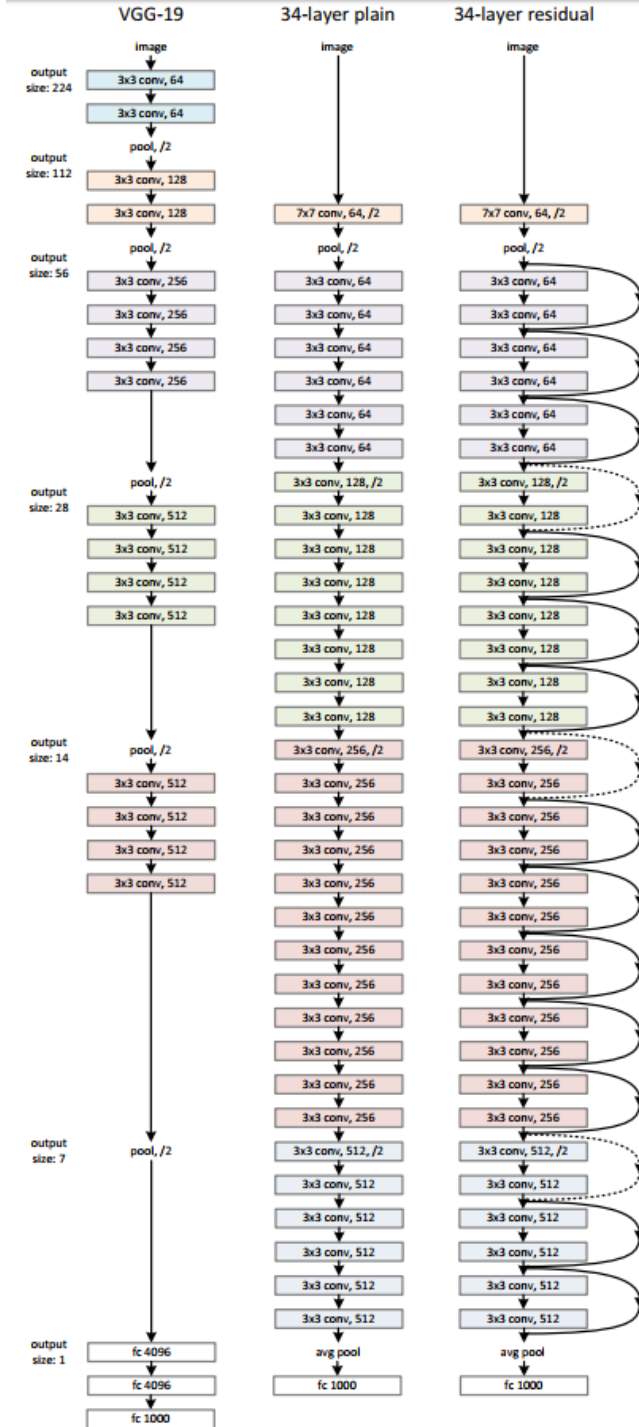
A deep network is at least as strong as it's shallower variant

- If adding layers doesn't help, just use the skip connection

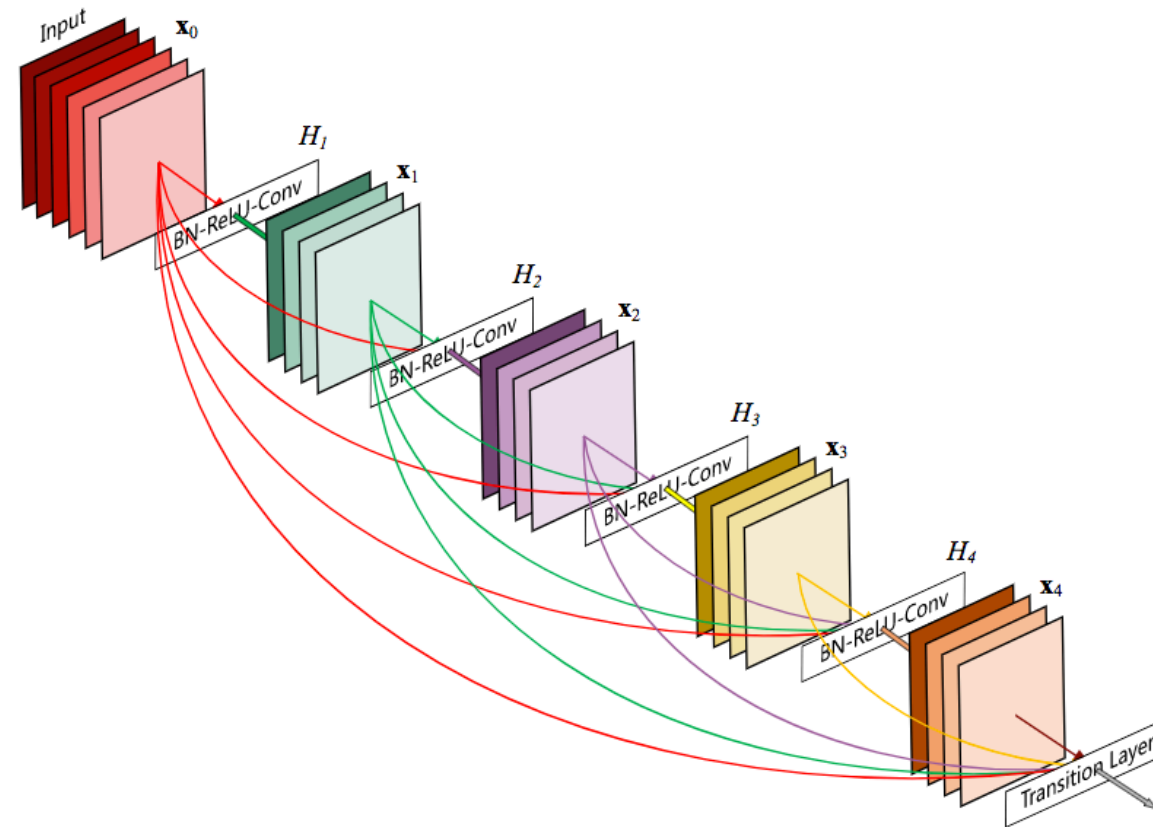


# Residual network (ResNet)

- Organize layers in blocks
- Use bottleneck layers
- Residual connections barely add computational complexity
- SOTA on ImageNet (when published)
- Inspired
  - Wide residual nets (50-layer wide ResNet > 152-layer ResNet)
  - DenseNets: get identity mapping from all previous layers

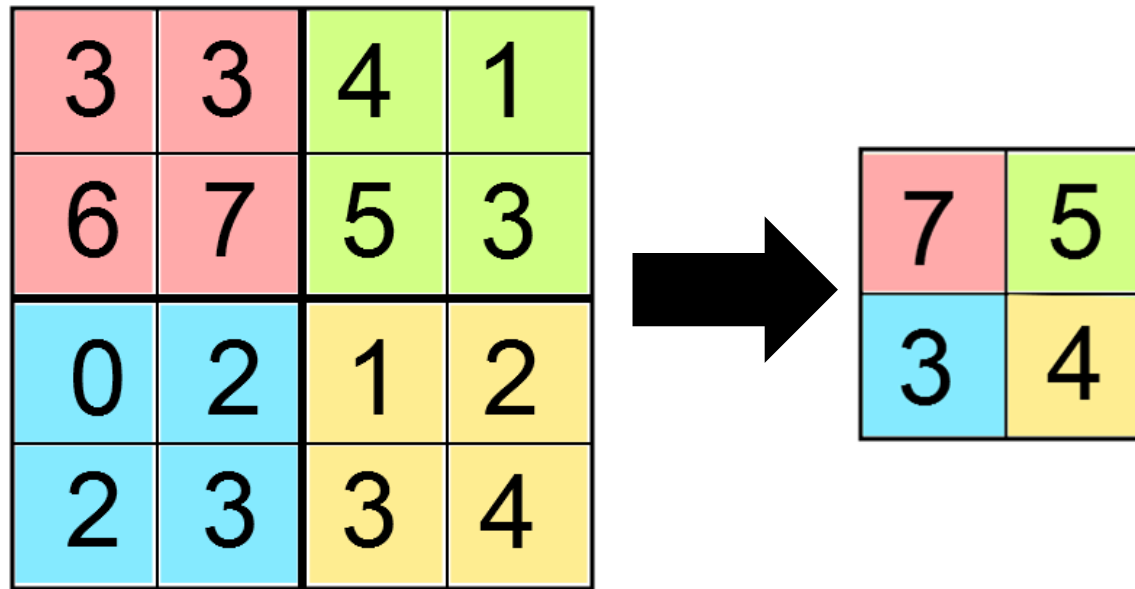


# Dense networks

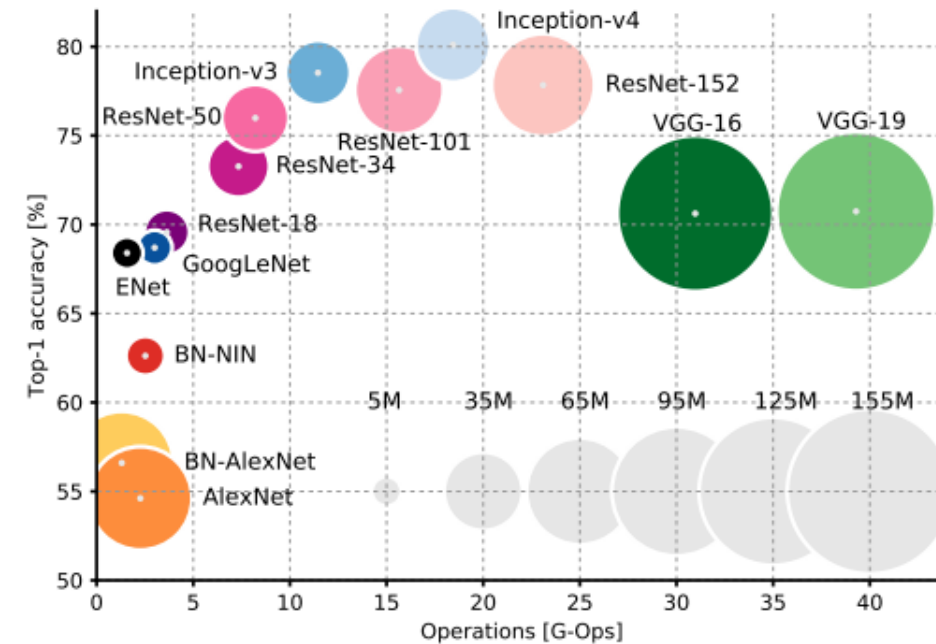
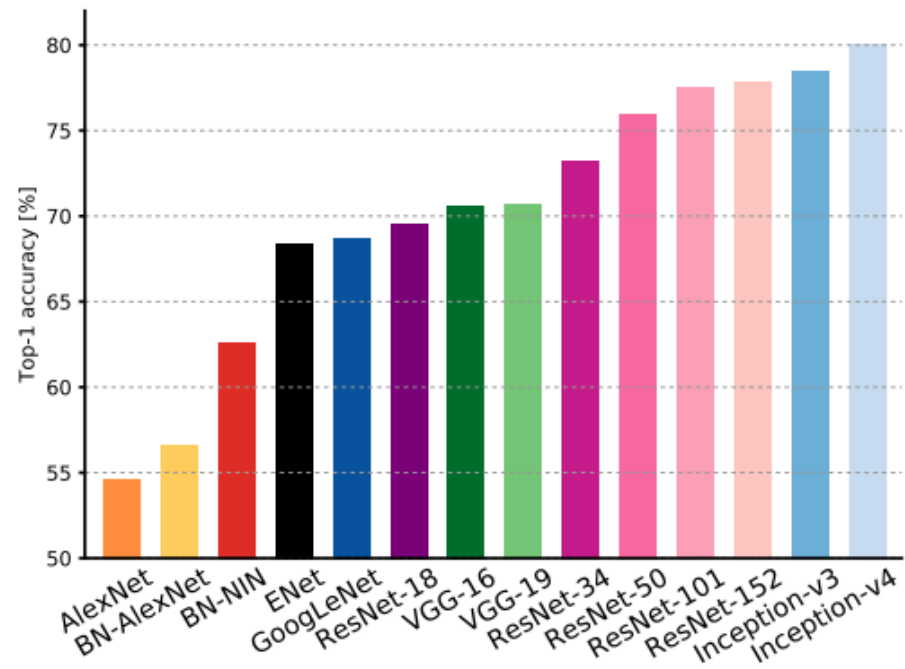


# Downsampling

- We often want to go from a large image to a single prediction
- Use downsampling operations like pooling
- Pooling is not trainable

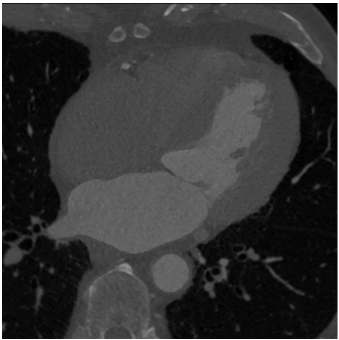


# Complexity vs. accuracy

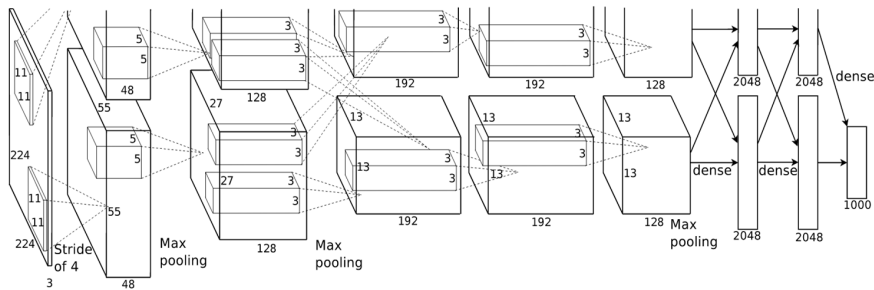


# Example: Organ localization in CT

2D image



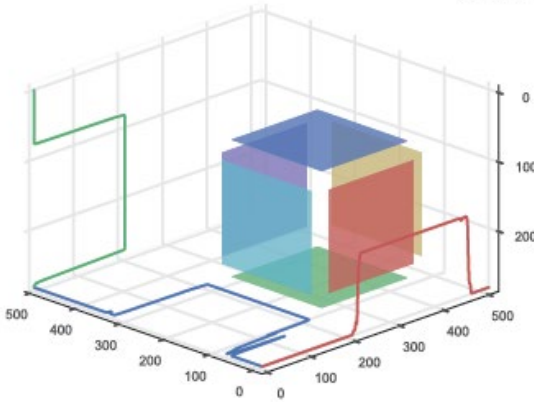
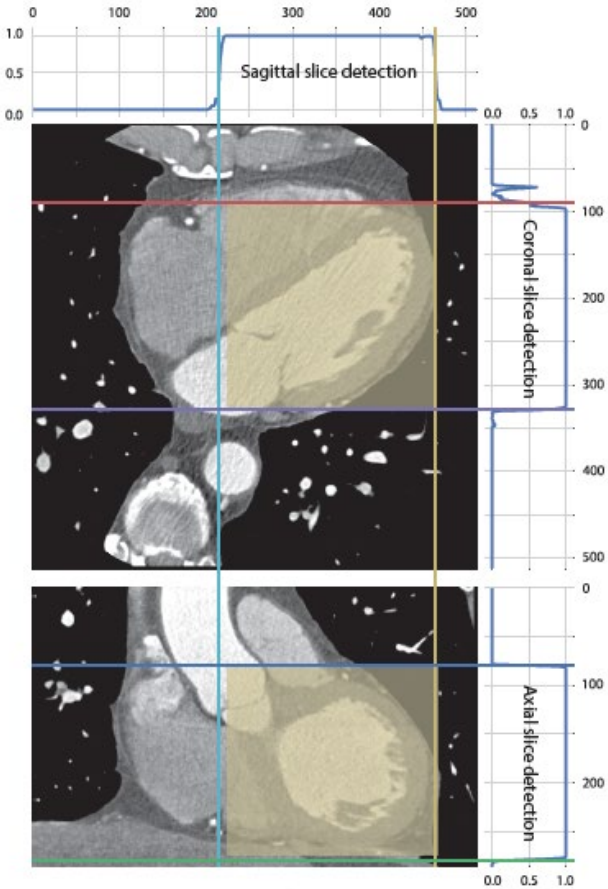
AlexNet



Ventricle?

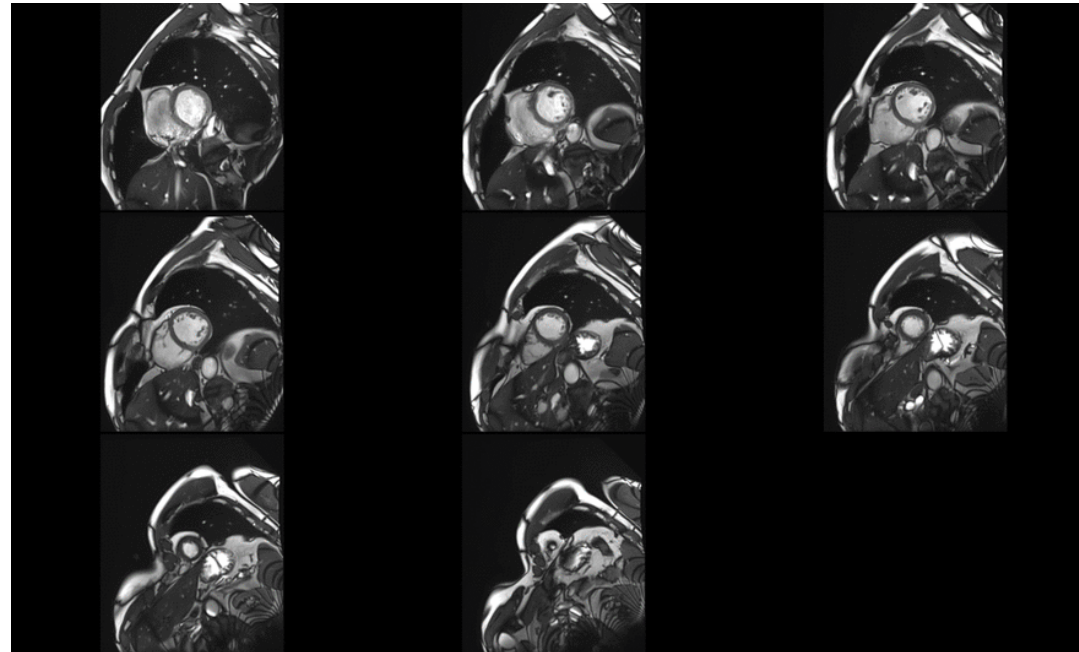
- Yes
- No

For each image slice



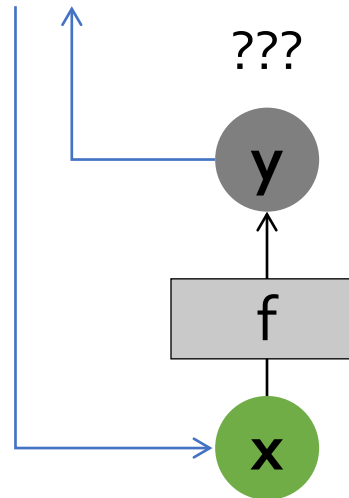
# Sequences

- A lot of data is sequential
- E.g. videos, audio, text, ECG, medical images, ...
- Can we use this in our neural network?



# Recurrent neural networks (RNNs)

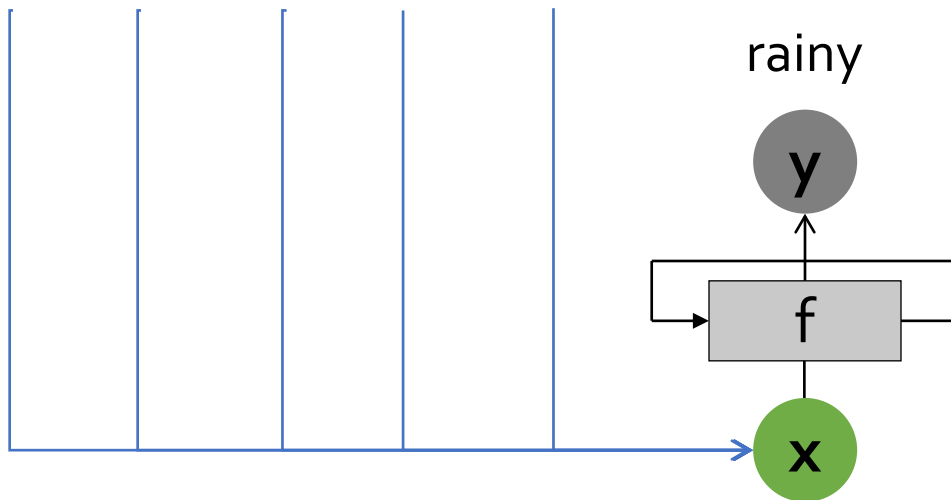
- It would be good to use information that came before
- A feedforward neural network has no 'memory'
- Consider training a neural network to predict the next word
  - "It's September, the weather is ..."



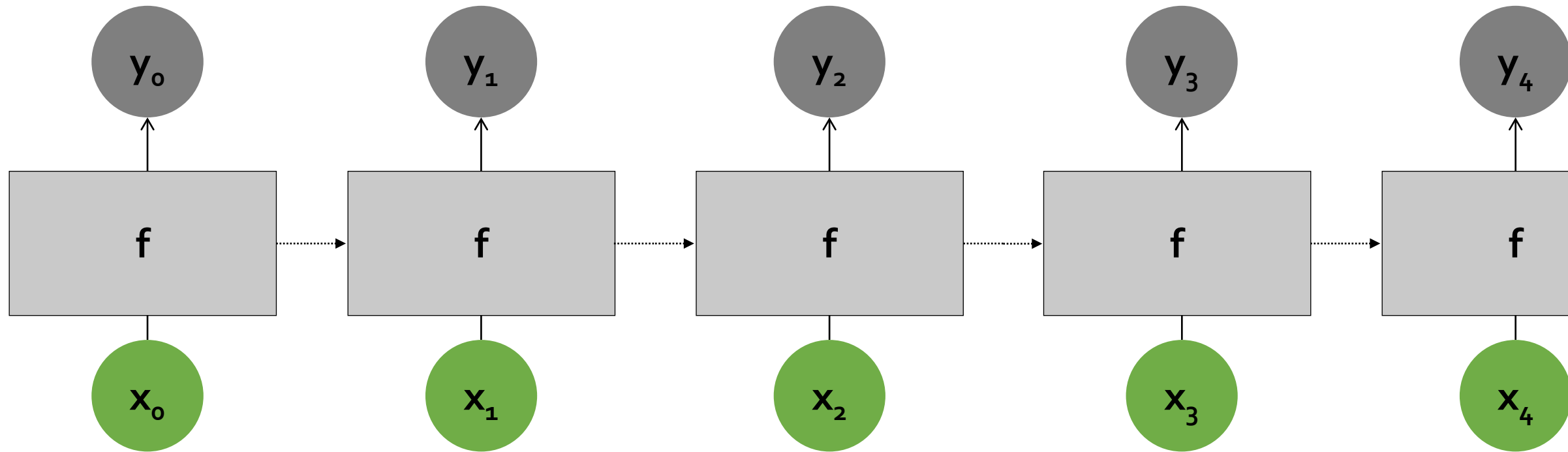


# Recurrent neural networks (RNNs)

- It would be good to use information that came before
- A feedforward neural network has no 'memory'
- Consider training a neural network to predict the next word
  - "It's September, the weather is ..."

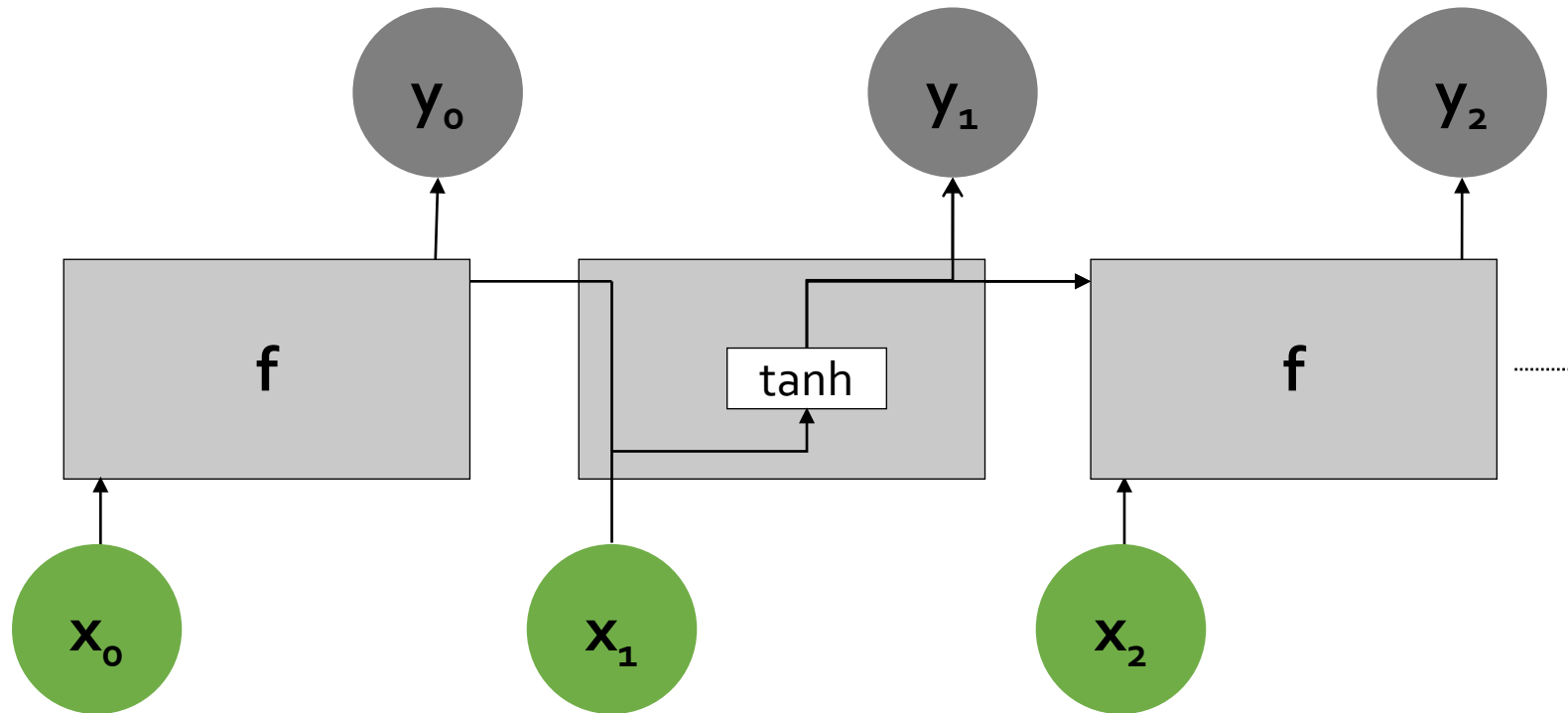


# Unrolling a recurrent neural network



# Unrolling a recurrent neural network

- Traditional RNNs have poor memory
- Previous outputs will get overwritten

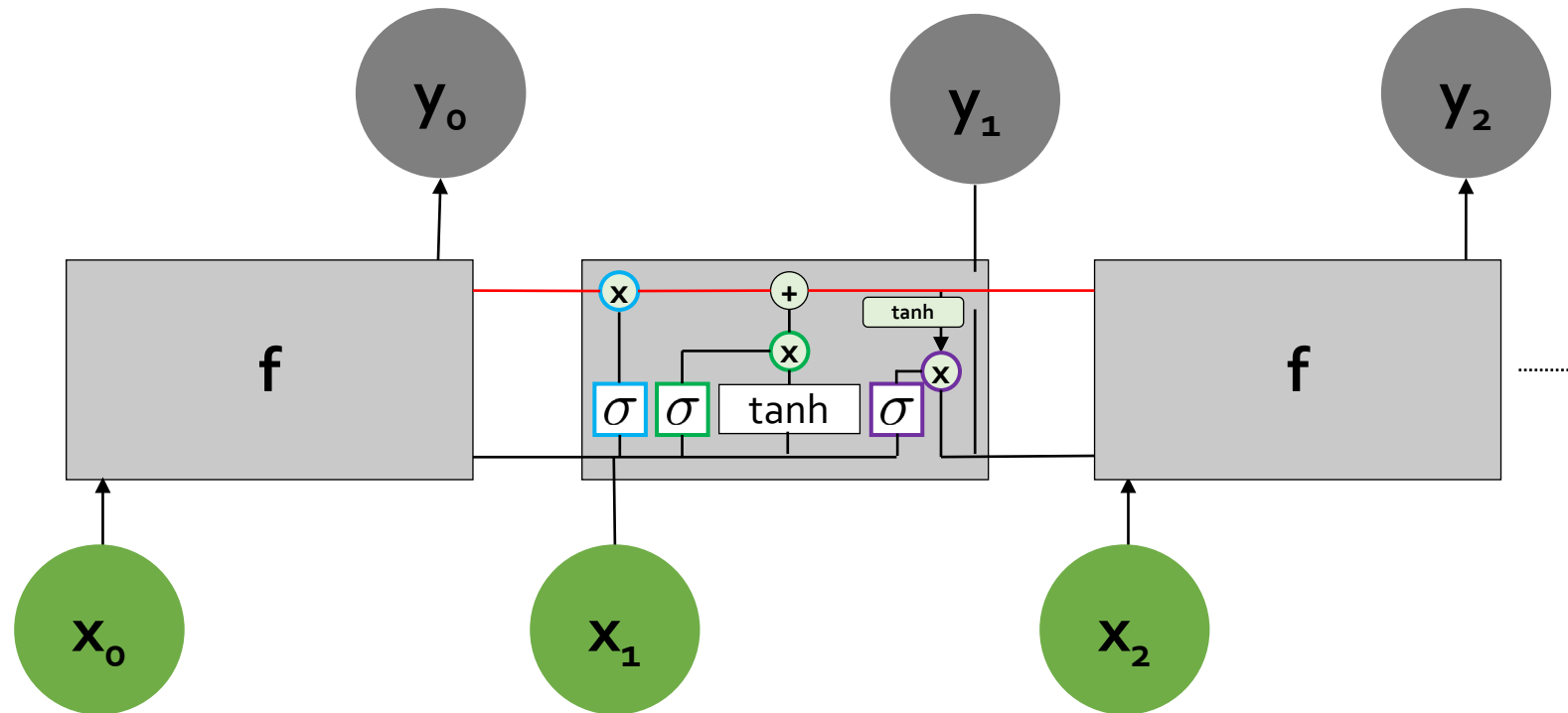


# Long short-term memory (LSTM)

## 1. Cell state

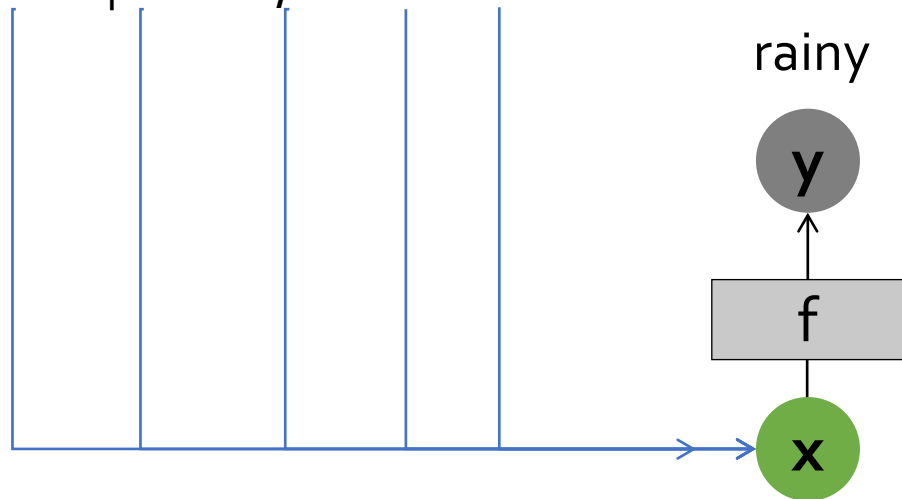
## 2. Gates

- Forget gate
- Input gate
- Output gate



# Recurrent vs. feedforward

- Recurrent networks are intuitively appealing, but
  - feedforward networks are faster (parallel), simpler and they often very competitive
  - “It’s September, the weather is ...”

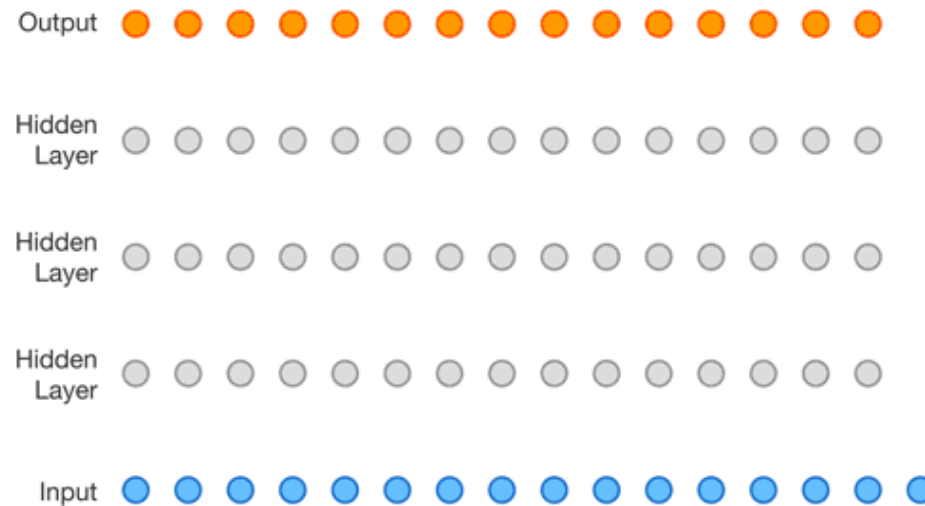


- One way to deal with large contexts in feedforward networks
  - dilated convolutions

# Dilated convolutions

In each layer, add more spacing between elements

- increase receptive field
- prevent explosion in number of parameters

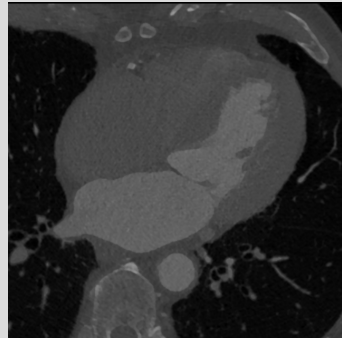


# Back to images

## Classification

## Regression

Image



CNN

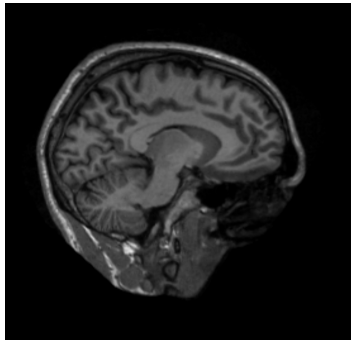
Heart



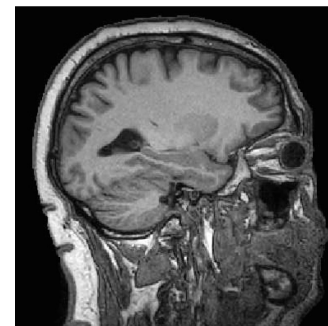
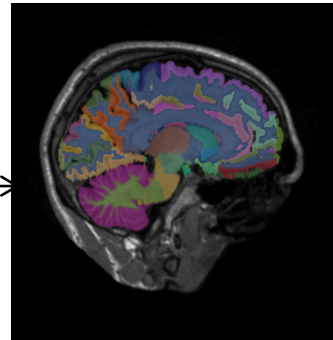
CNN

300 ml

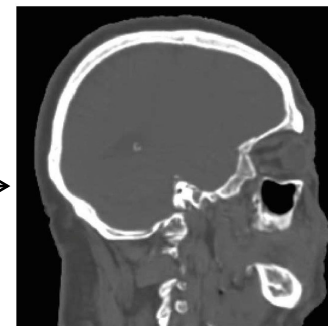
Voxel



CNN

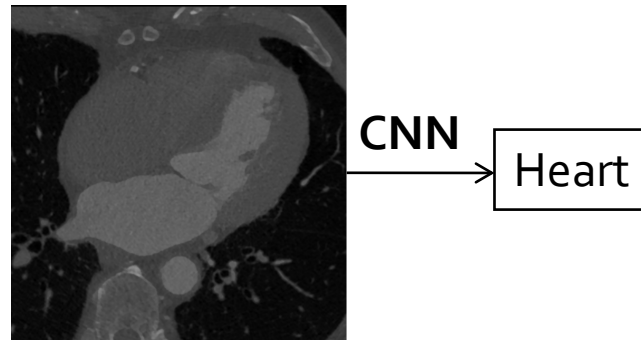


CNN

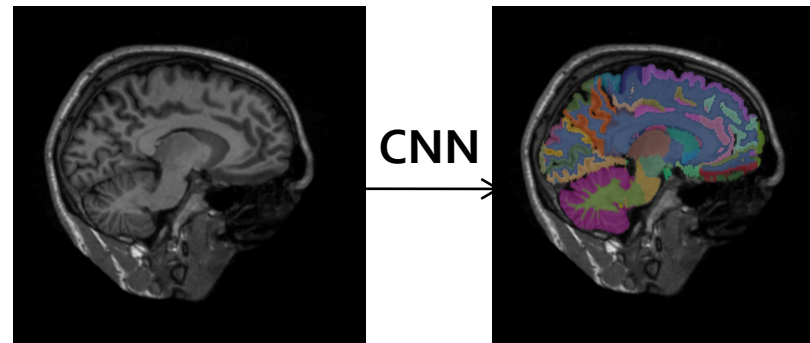


# CNNs for pixelwise prediction

LeNet, AlexNet, VGG-Net, GoogLeNet all predict one value per **image**

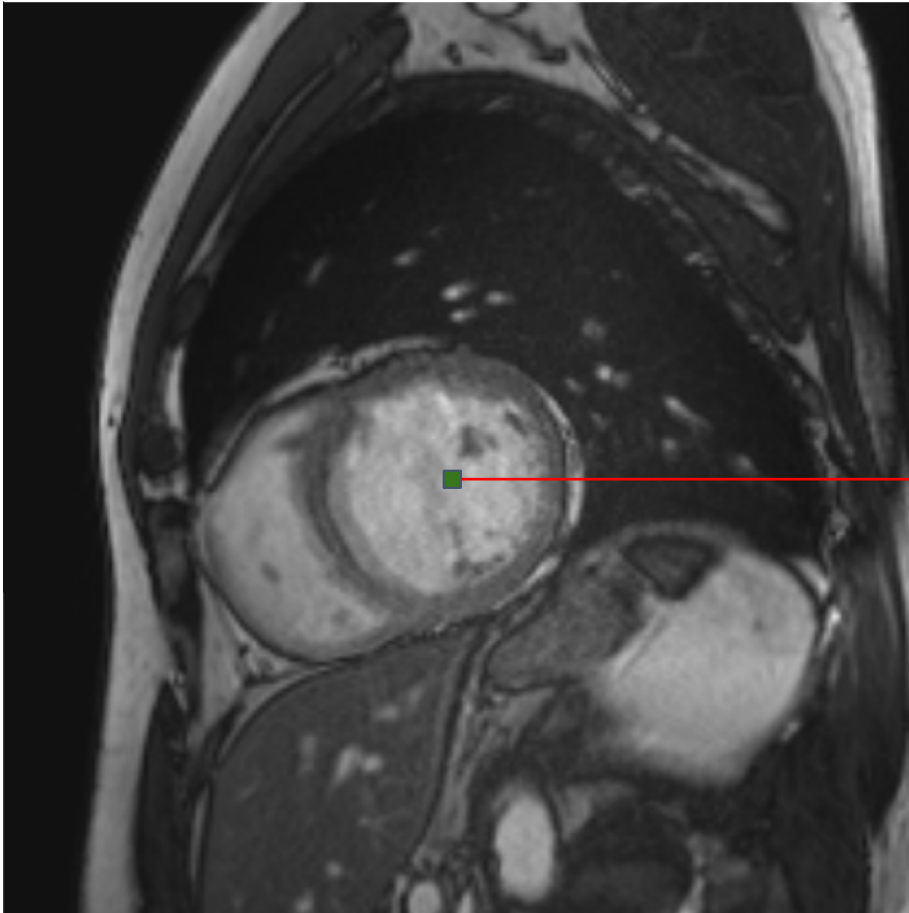


Often, we want to predict one value per **pixel/voxel**

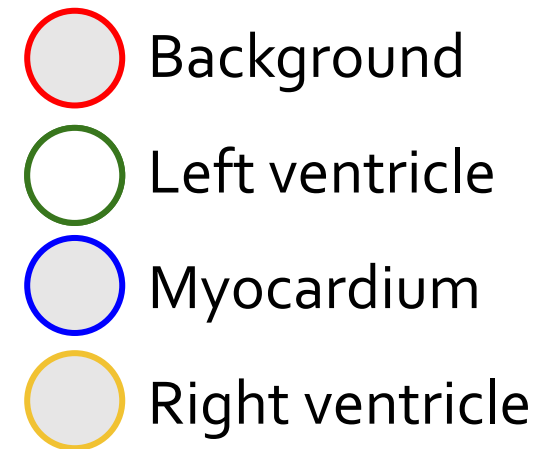




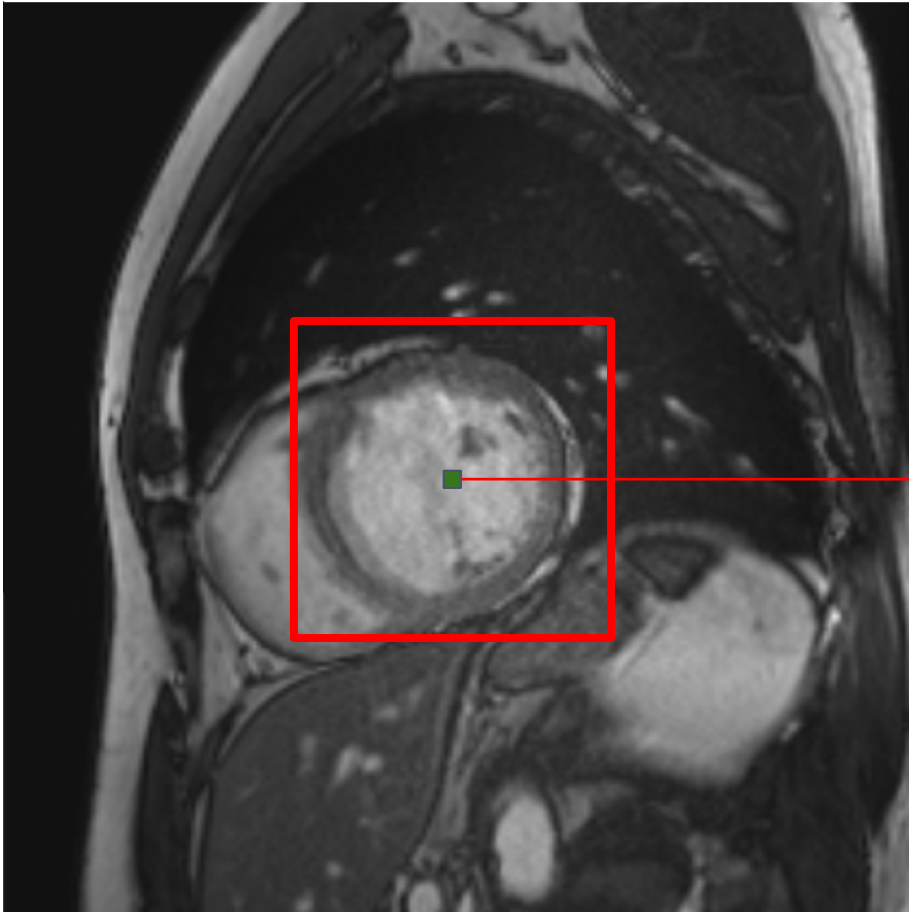
# Sliding window



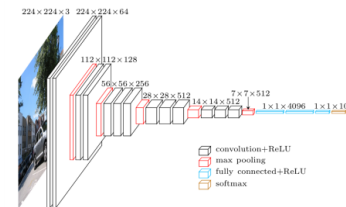
- **Image** = patch centered at voxel
- **Label** = class of center voxel



# Sliding window



- **Image** = patch centered at voxel
- **Label** = class of center voxel

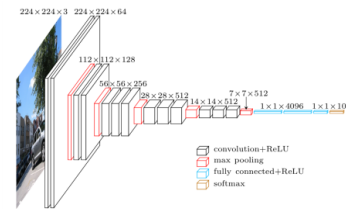
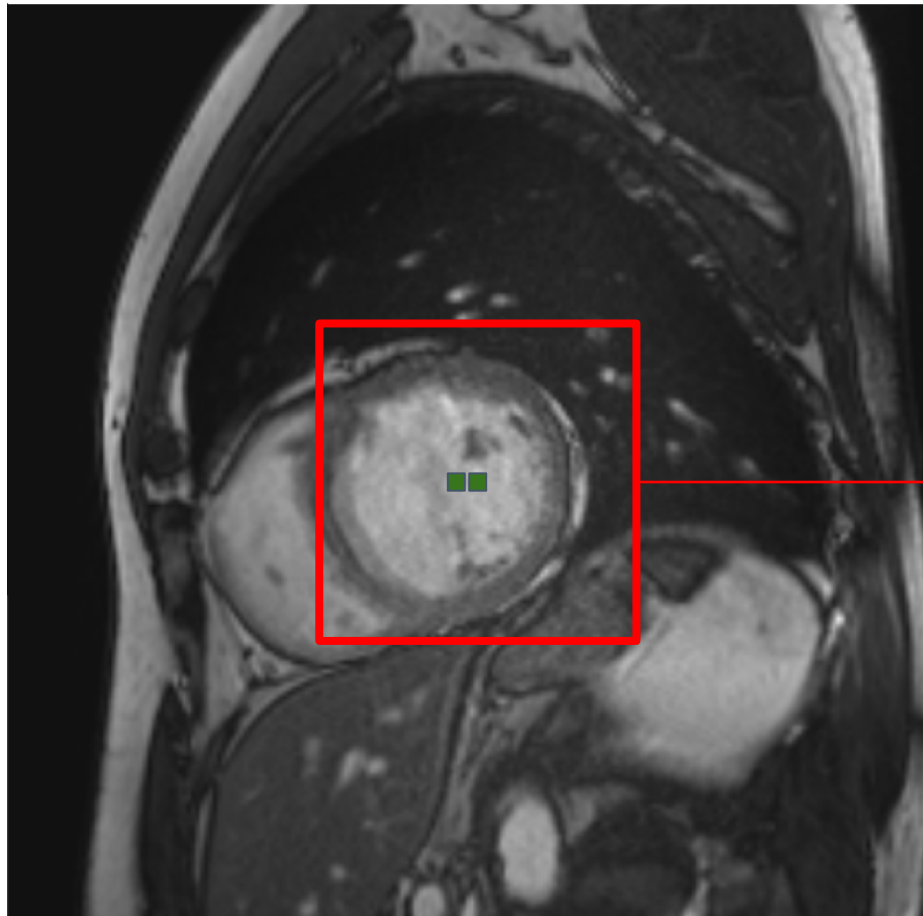


- Background
- Left ventricle
- Myocardium
- Right ventricle

# Sliding window

Combination of thousands of image classification tasks

- **Image** = patch centered at voxel
- **Label** = class of center voxel

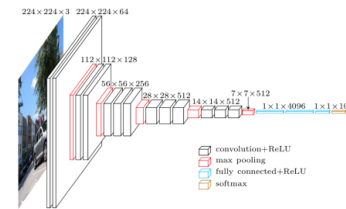
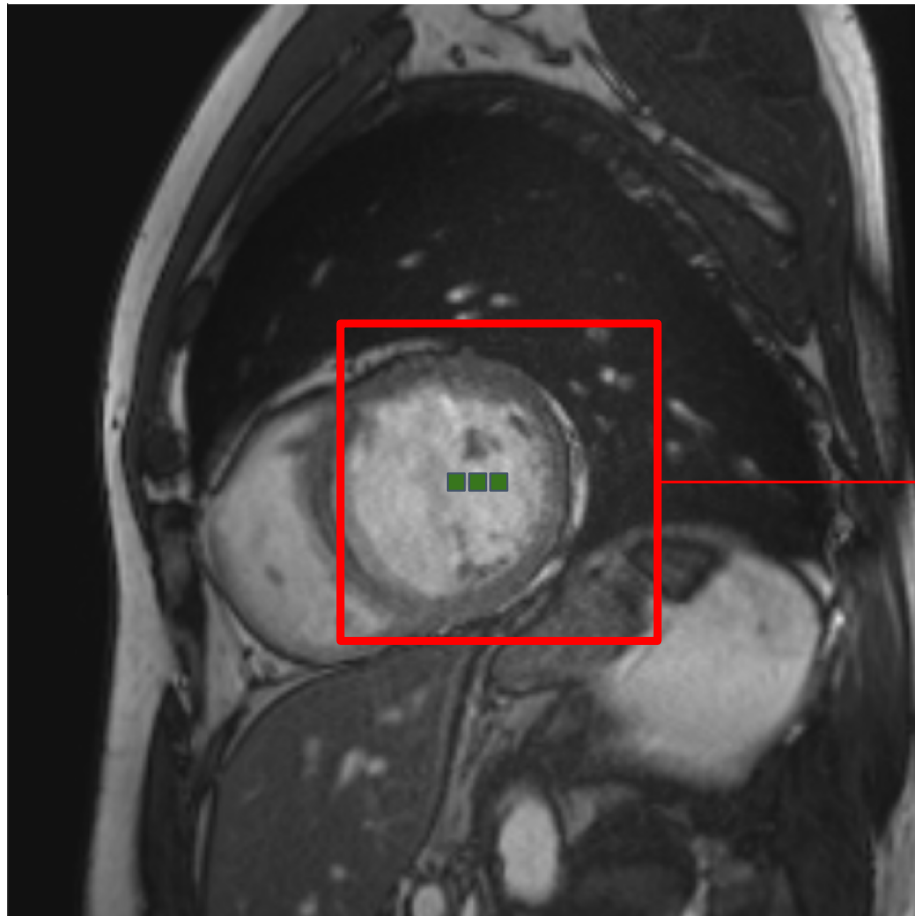


- Background
- Left ventricle
- Myocardium
- Right ventricle

# Sliding window

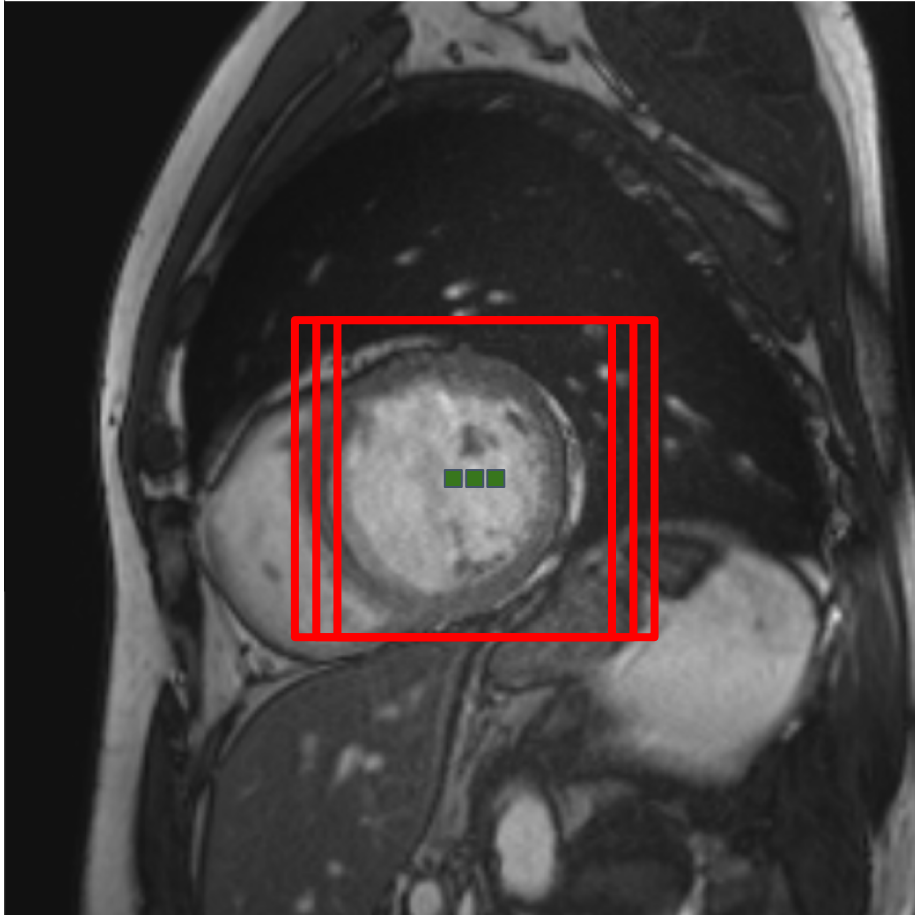
Combination of thousands of image classification tasks

- **Image** = patch centered at voxel
- **Label** = class of center voxel



- Background
- Left ventricle
- Myocardium
- Right ventricle

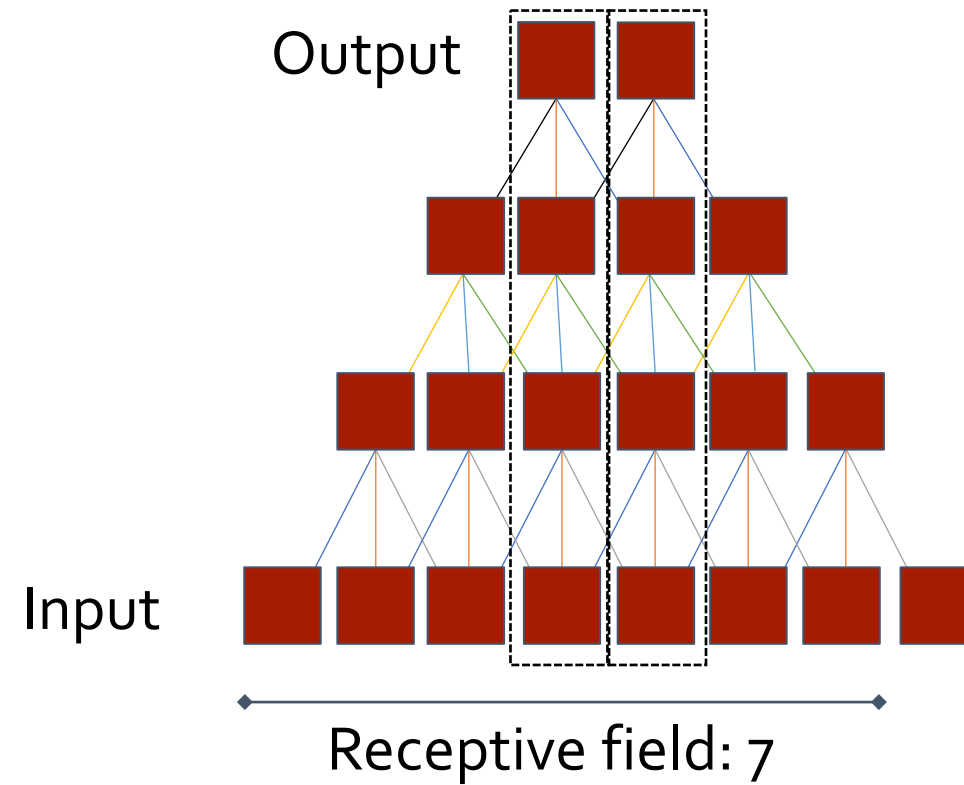
# Sliding window



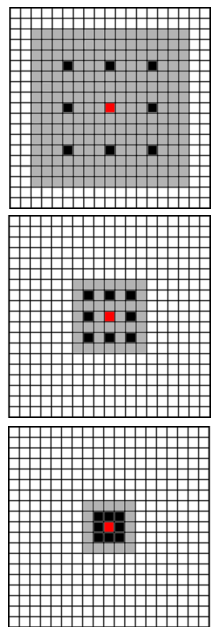
## Sliding window approaches are inefficient

- Each patch is processed separately
- Lots of redundant operations
- We would like to re-use/share operations

# All-convolutional network



# Dilated convolutions



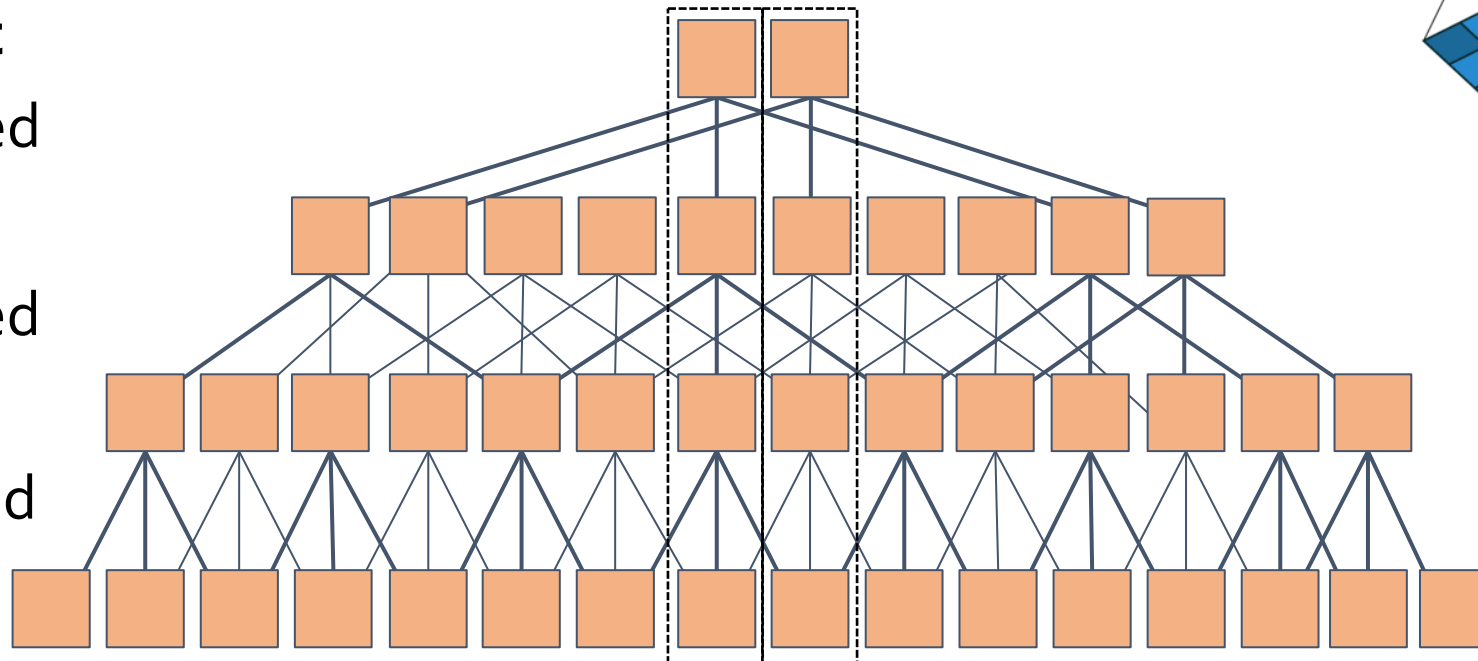
**Output**

4-dilated

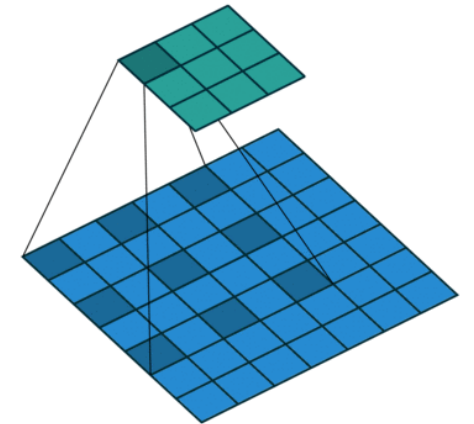
2-dilated

1-dilated

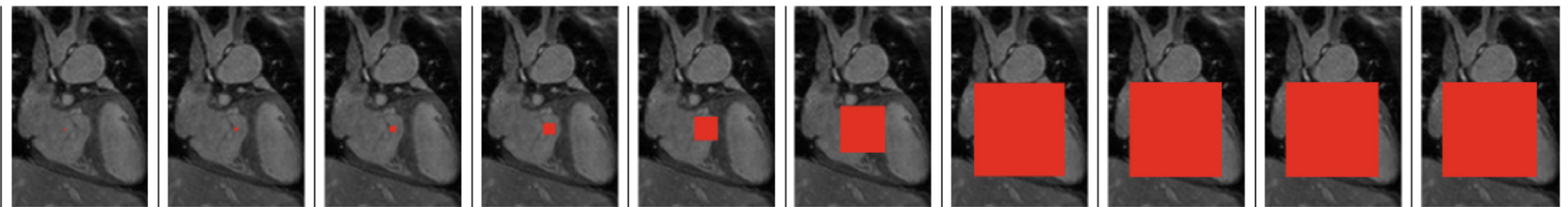
**Input**



Receptive field: 15



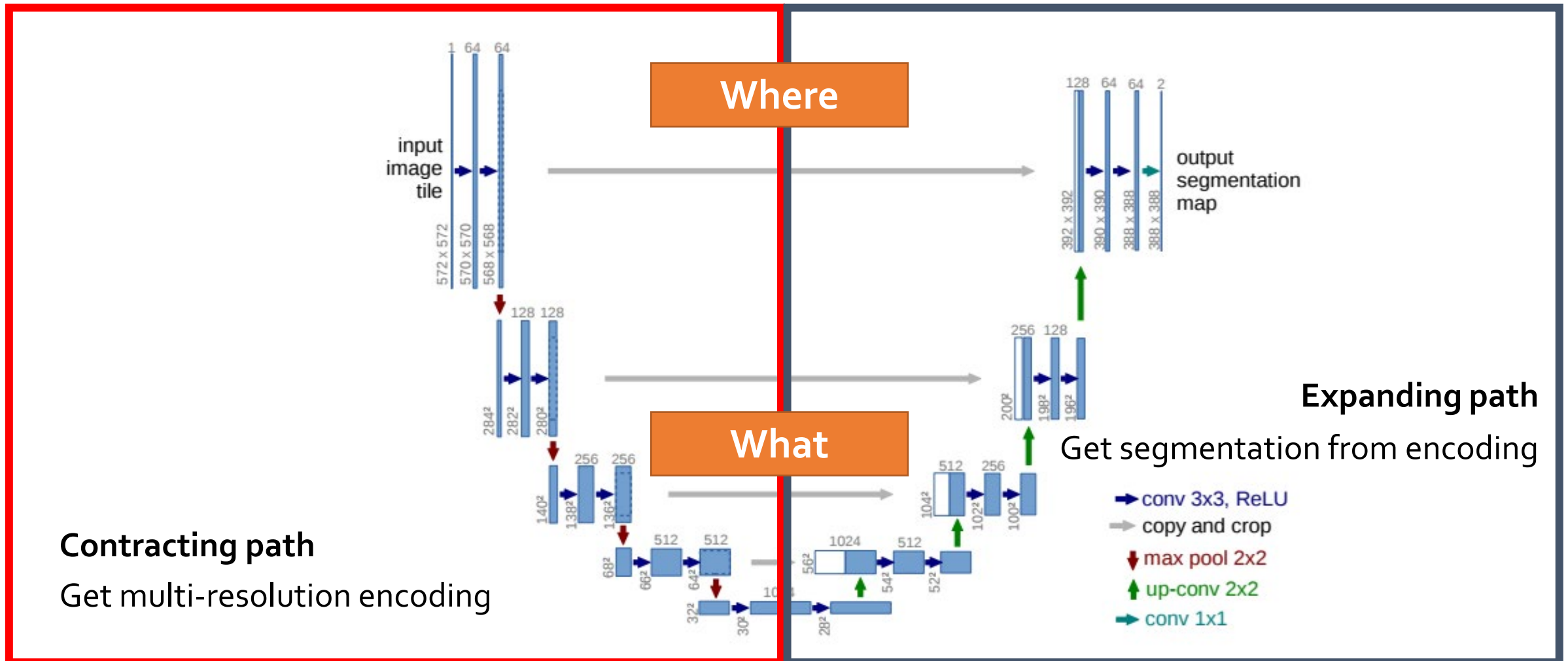
# Example



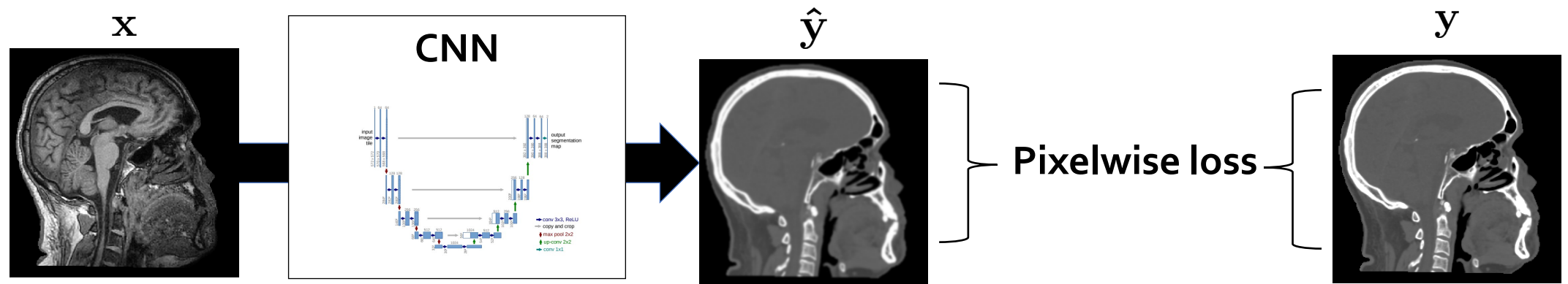
Layer	1	2	3	4	5	6	7	8	9	10
Convolution	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$1 \times 1$	$1 \times 1$
Dilation	1	1	2	4	8	16	32	1	1	1
Field	$3 \times 3$	$5 \times 5$	$9 \times 9$	$17 \times 17$	$33 \times 33$	$65 \times 65$	$129 \times 129$	$131 \times 131$	$131 \times 131$	$131 \times 131$
Channels	32	32	32	32	32	32	32	32	192	3
Parameters	320	9248	9248	9248	9248	9248	9248	9344	6912	579



# Encoder-decoder architecture: U-Net

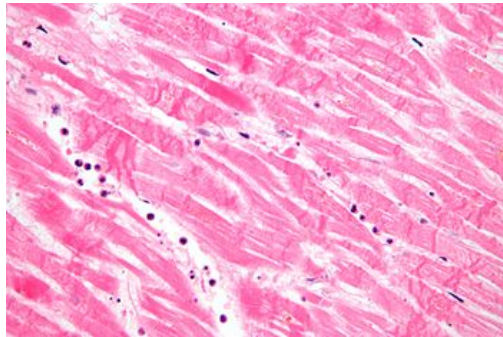


# Training

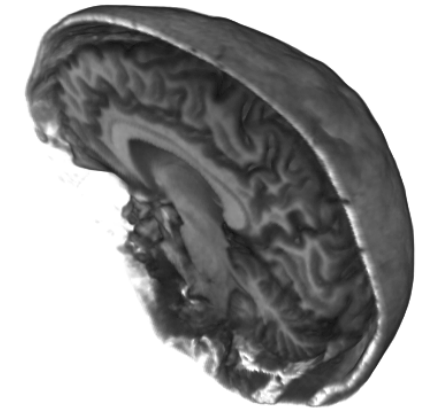
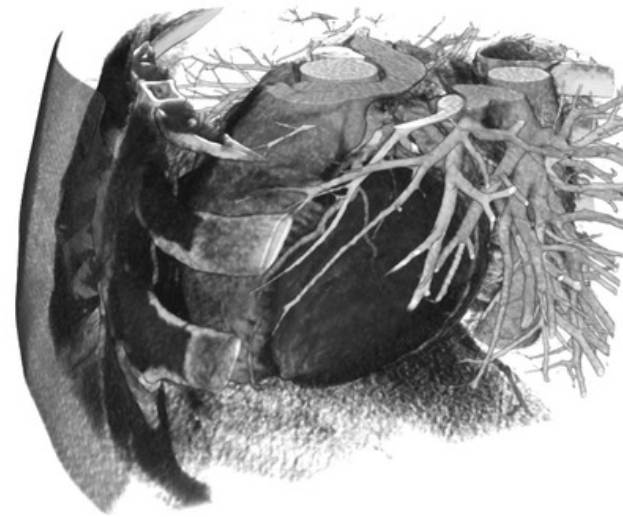


# 2D or 3D images

2D data



3D data



# 3D networks

Many medical images are 3D instead of 2D

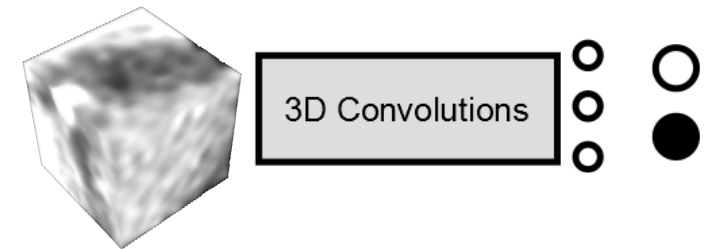
- MR images
- CT images

Can we just use 3D layers instead of 2D layers? Sure!

- 3D convolution layers in Keras, TensorFlow, PyTorch, etc.
- 3D network architectures (e.g. U-Net, V-Net)

But

- Is your data really 3D (think about acquisition)? Isotropy?
- Increase in memory consumption + operations + parameters



# Summary

## Advanced architectures

- AlexNet, GoogleNet, VGG-Net, ResNet
- Deeper, larger, better + some tricks

## Recurrent neural networks

- RNNs + LSTMs

## Per image prediction != per voxel prediction

- All-convolutional networks
- Encoder-decoder architectures

## 2D/3D data

- Most 2D neural networks extend to 3D



General

[VIEW DOCS](#)

LANGUAGE

English (en)

PREDICTED CONCEPT

PROBABILITY

coffee

0.994

cup

0.993

dawn

0.988

hot

0.987

no person

0.987

breakfast

0.987

caffeine

0.980

still life

0.980



TRY YOUR OWN IMAGE OR VIDEO

