

Deep Learning II: Advanced Neural Network Architectures

8DM40 Machine Learning in Medical Imaging and Biology

Jelmer Wolterink

02-10-2019



Deep Learning II

Me

- Postdoctoral researcher @ Amsterdam UMC – Location AMC
- Deep learning for 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

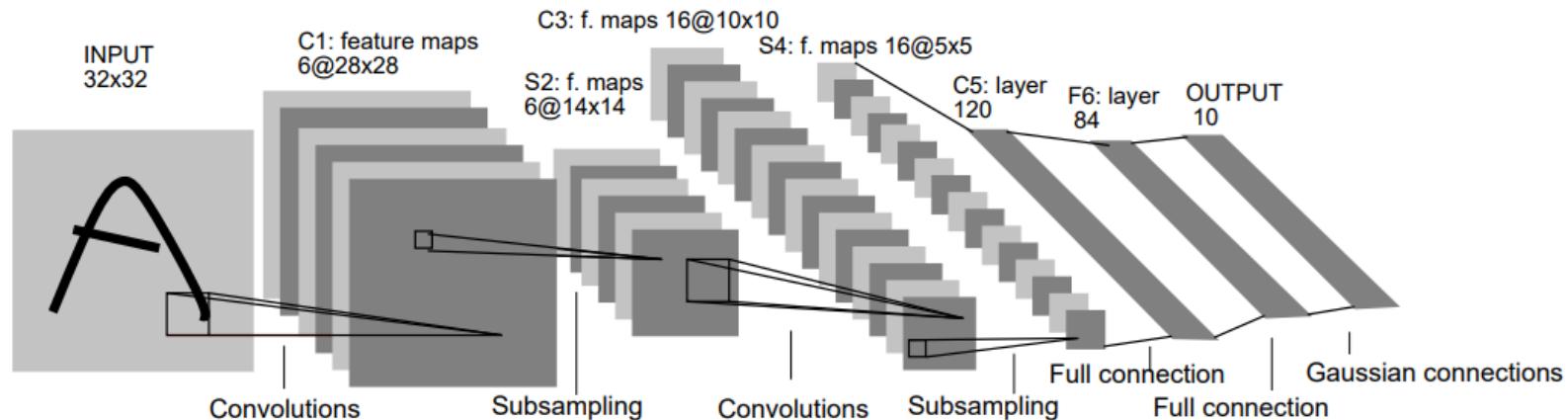
- Recurrent neural networks
- Long short term memory units

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)

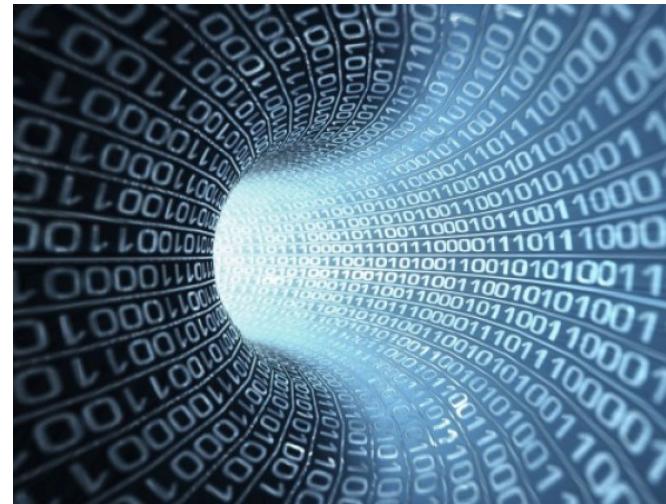




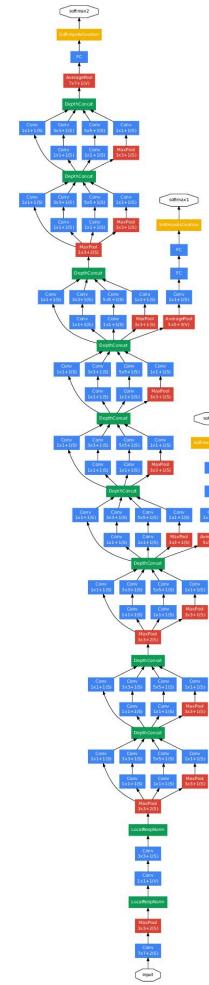
What happened between 1998 and now?



Compute



Data



Algorithms



Data: ImageNet challenge

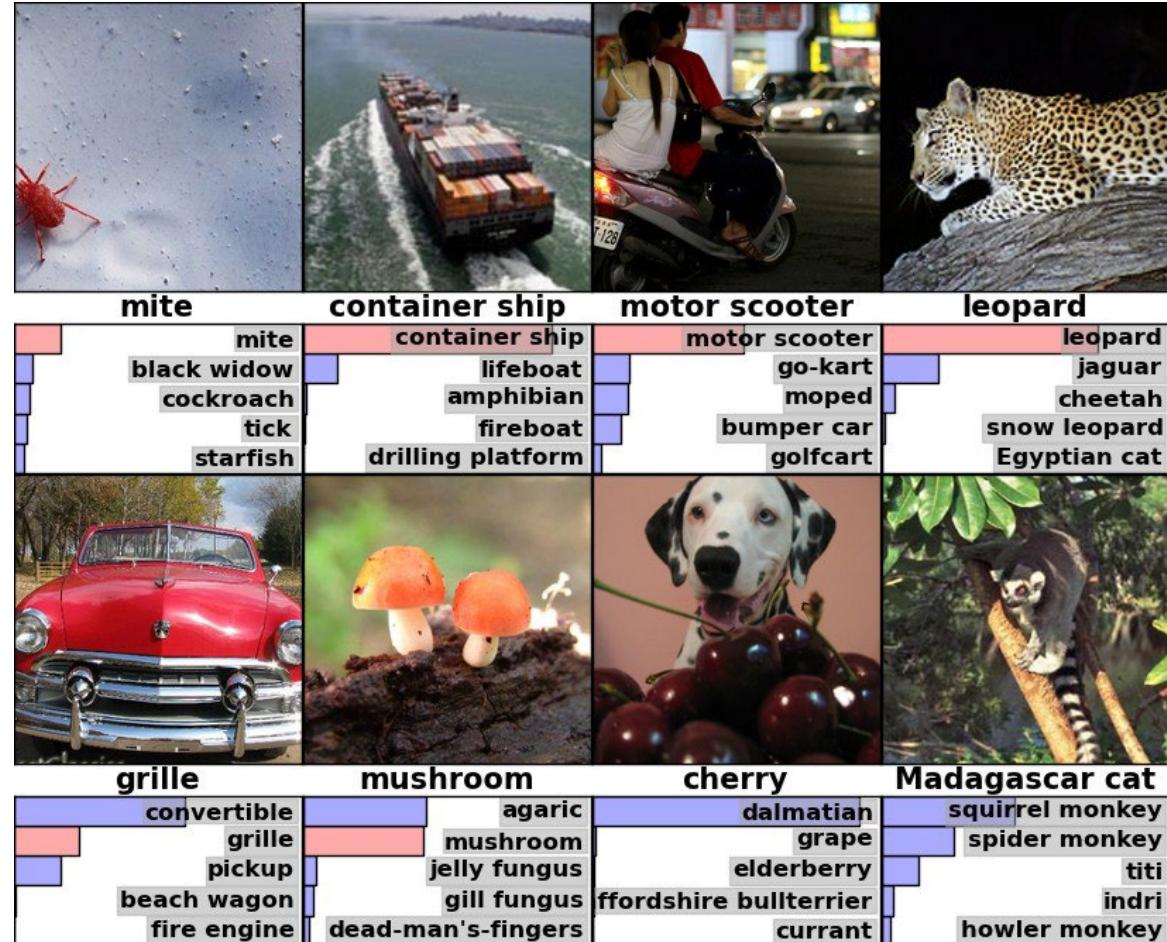
Benchmark for image classification/object detection

Data

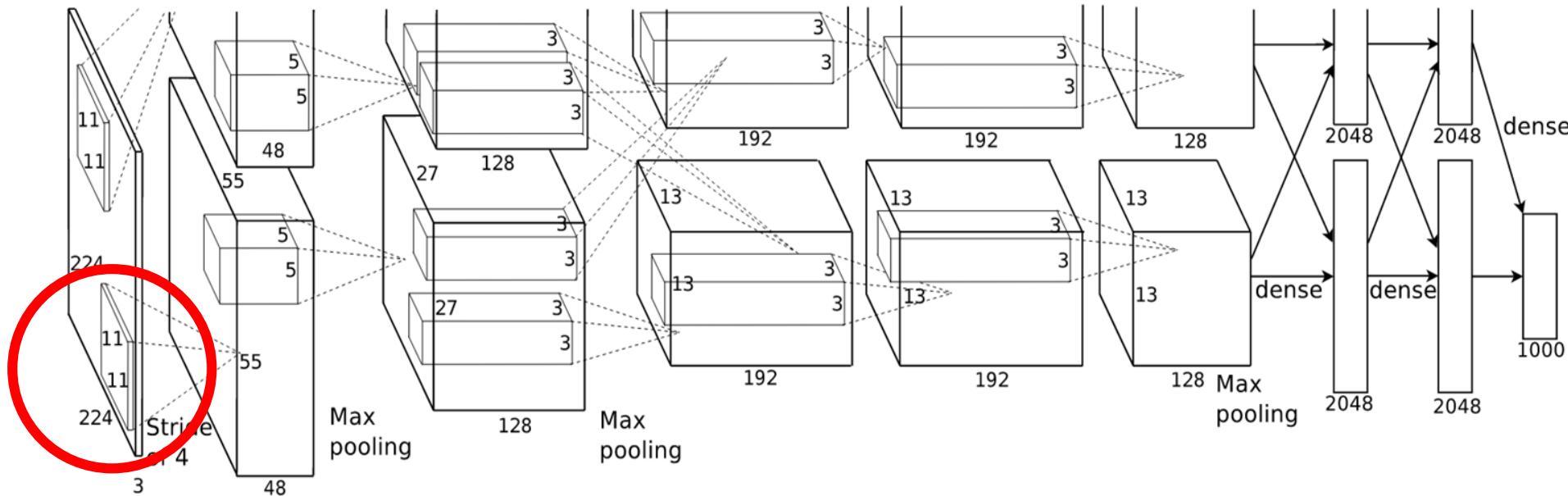
- > 1,200,000 RGB images
- Images show one of **1000** classes

Task

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



AlexNet



- 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



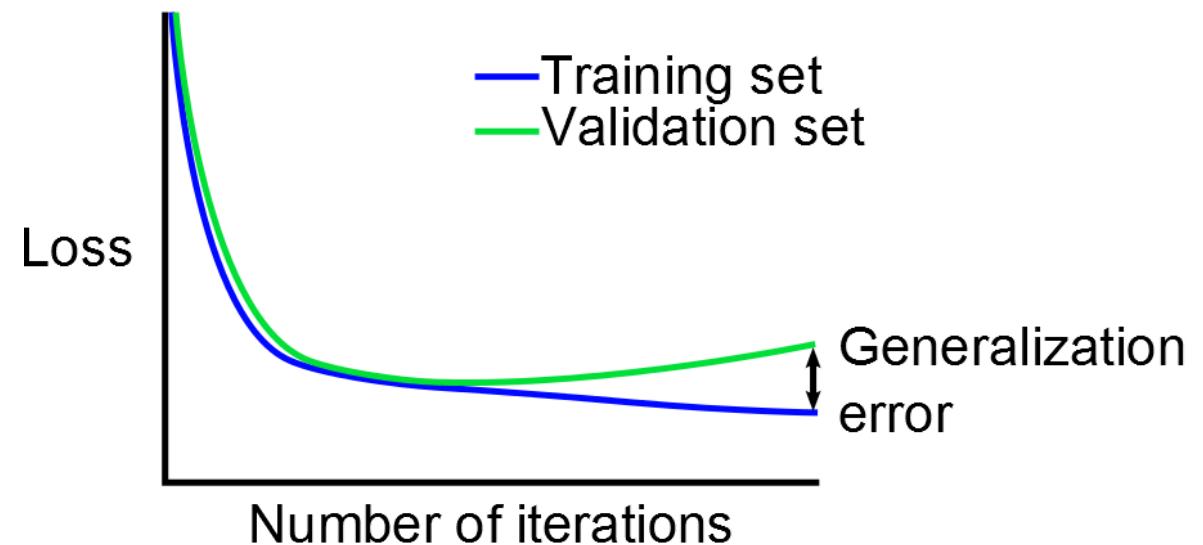
Overfitting

Reasons

- Too many parameters
- Not enough data

Solution

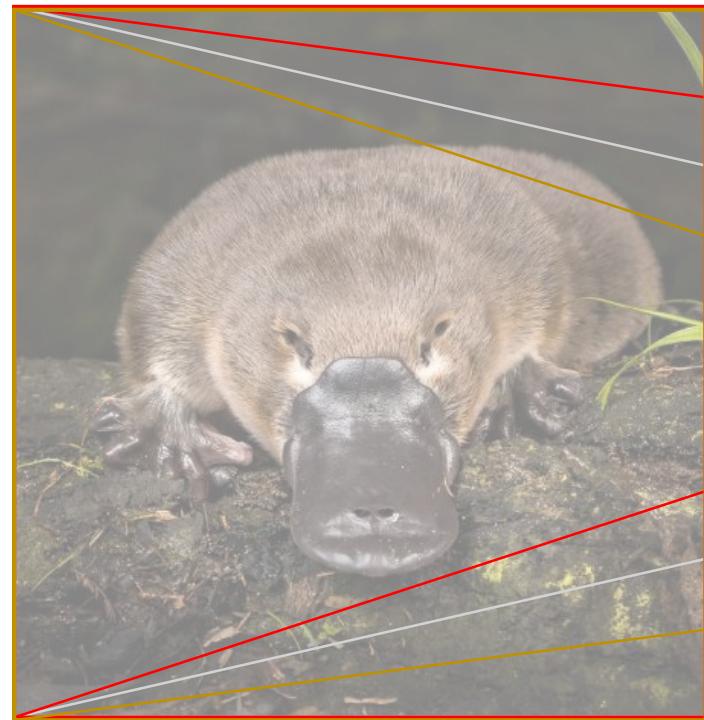
- Reduce number of parameters





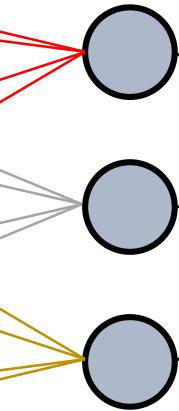
Kernel size

100



10,000 weights

100



Platypus

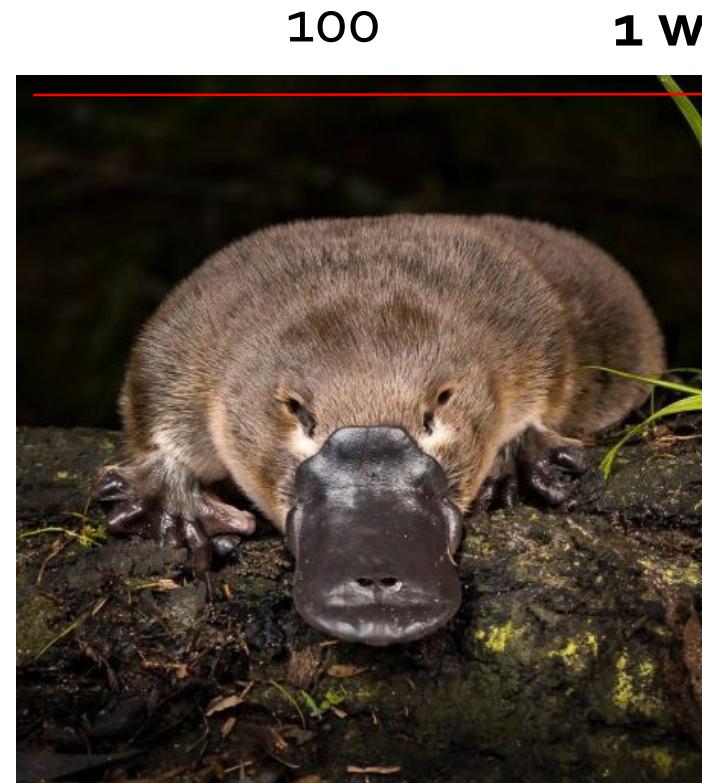
Input

Hidden

Output



Kernel size



100

1 weight



100

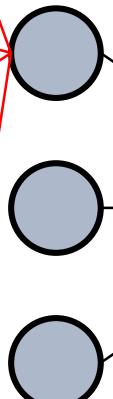
Input

Platypus

Hidden

Hidden

Output

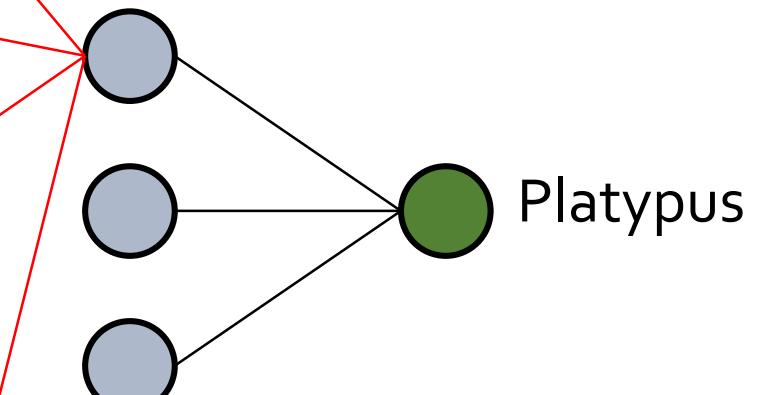
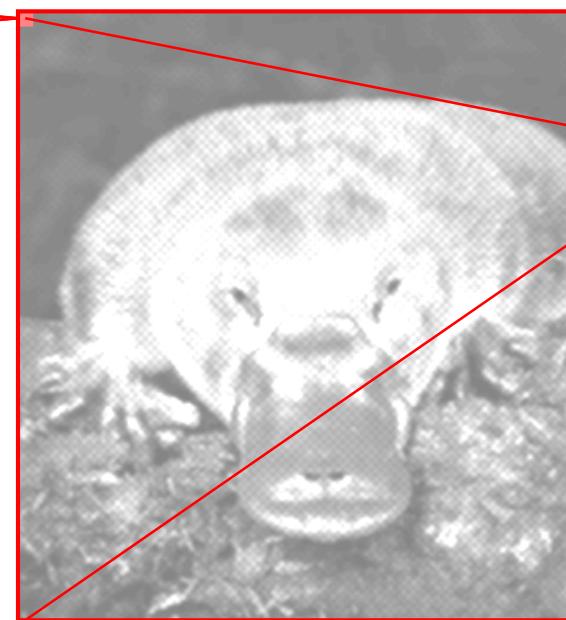
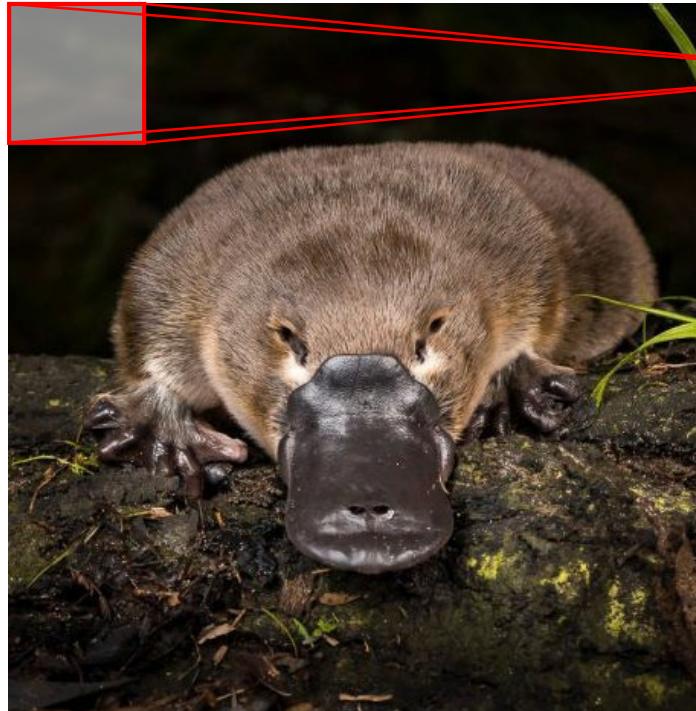




Kernel size

100

11 x 11 = 121 weights



Input

Hidden

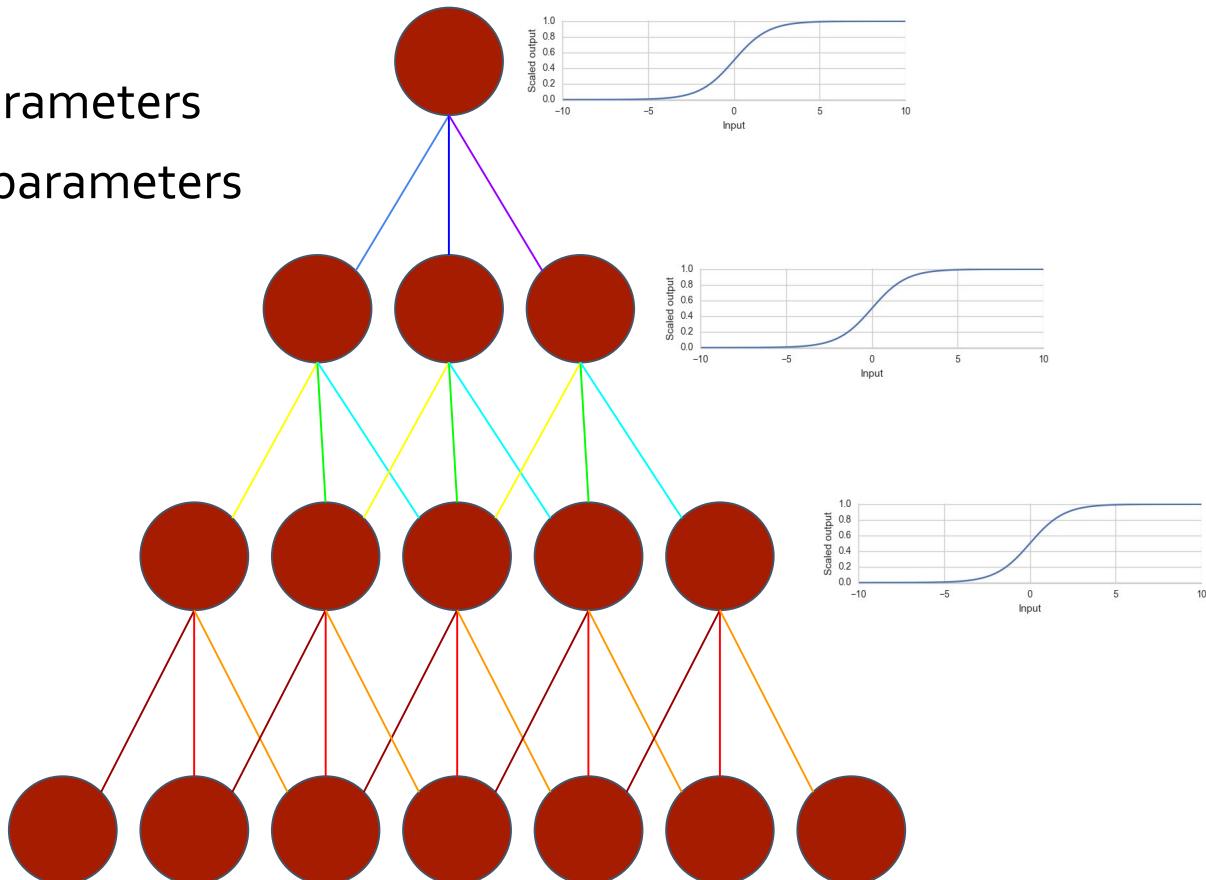
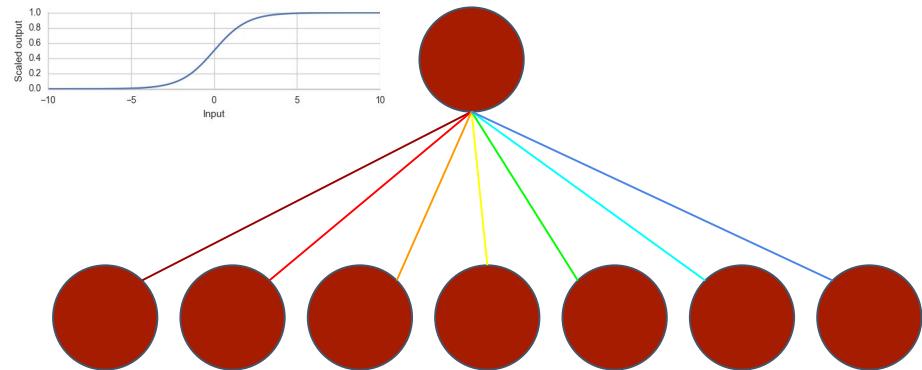
Hidden

Output



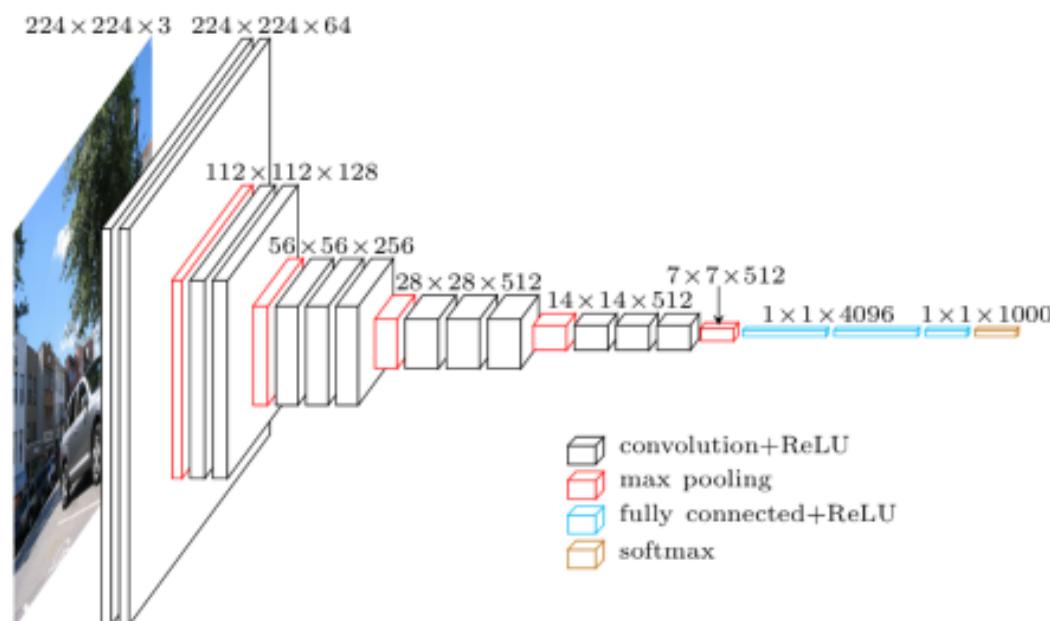
Using smaller kernels

- 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





VGG-Net

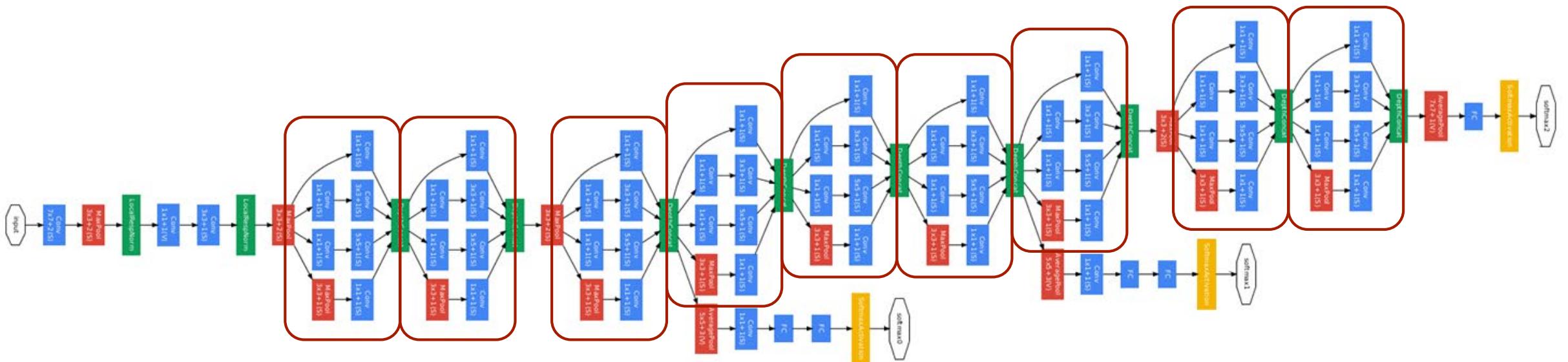


ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096	FC-4096	FC-4096	FC-4096	FC-4096	FC-4096
FC-1000					
soft-max					



GoogLeNet (Inception v1)

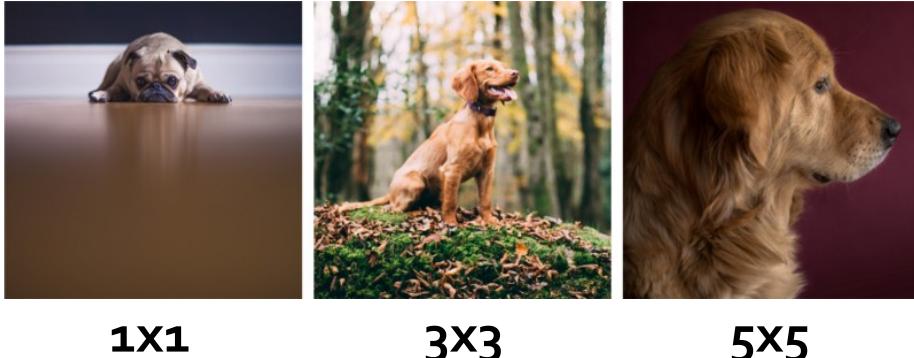
- 22 layer-network
 - Very deep compared to LeNet/AlexNet
 - SOTA on ImageNet (when published)



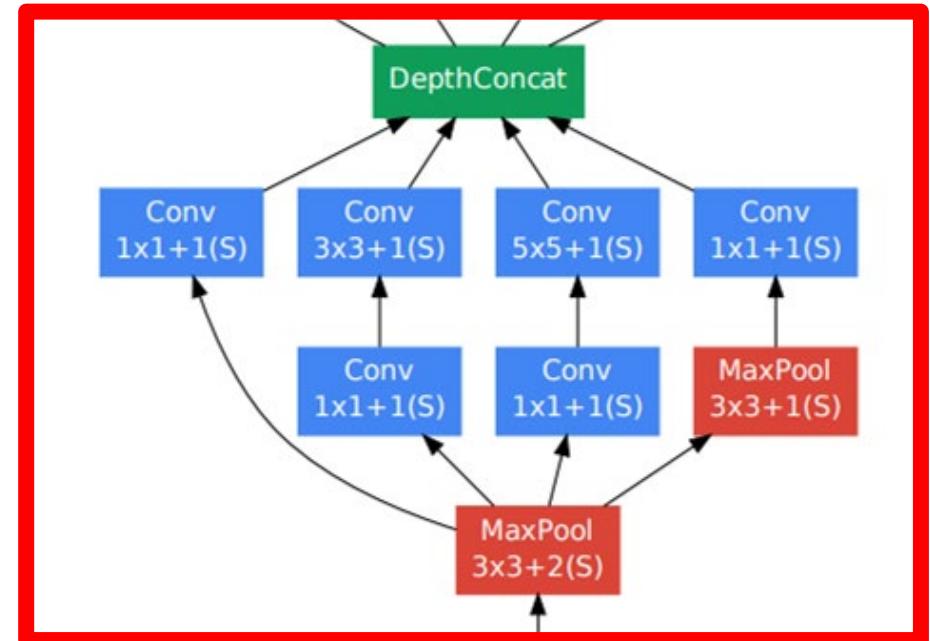


Inception module

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



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

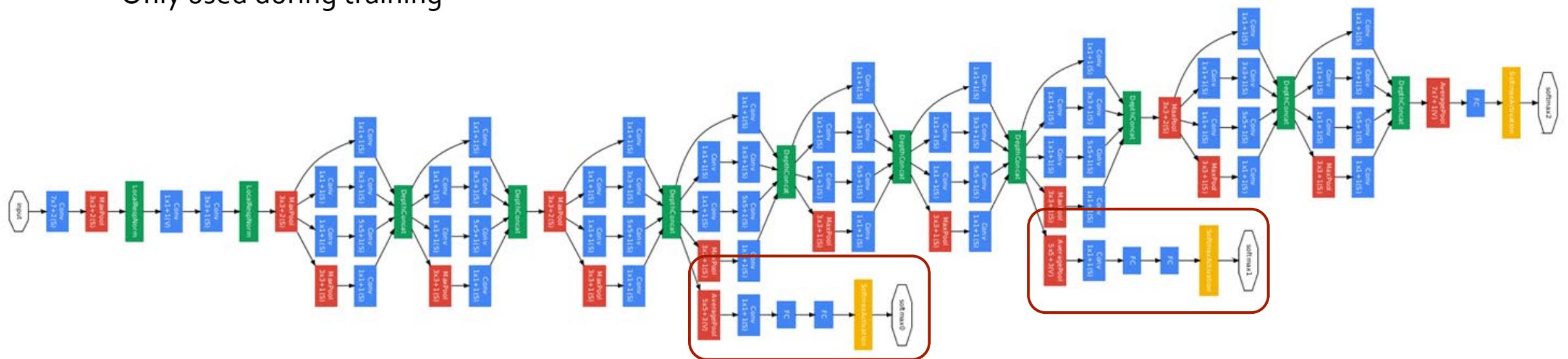




Auxiliary classifiers

Auxiliary classifiers provide extra supervision

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





Residual connections

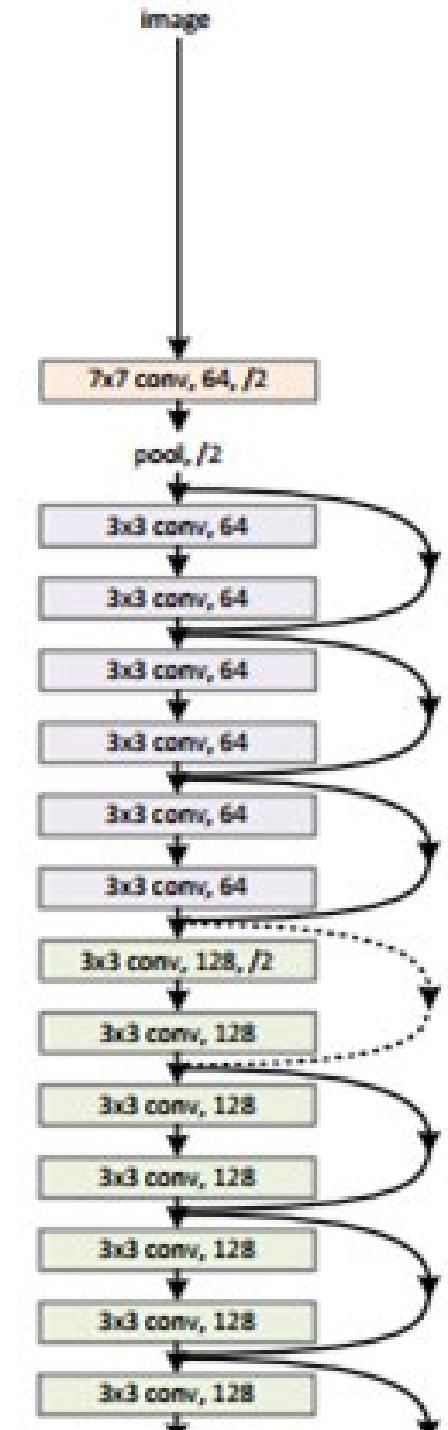
Very deep networks

- 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 its shallower variant

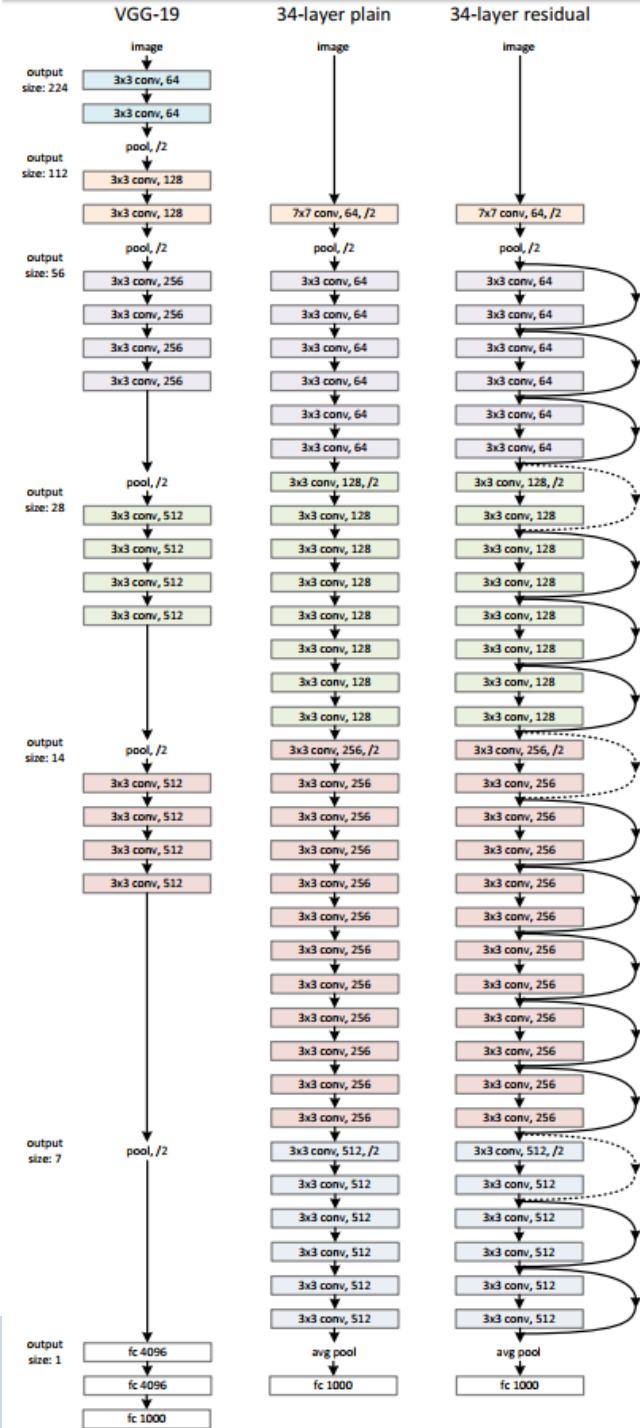
- If adding layers doesn't help, just use the skip connection





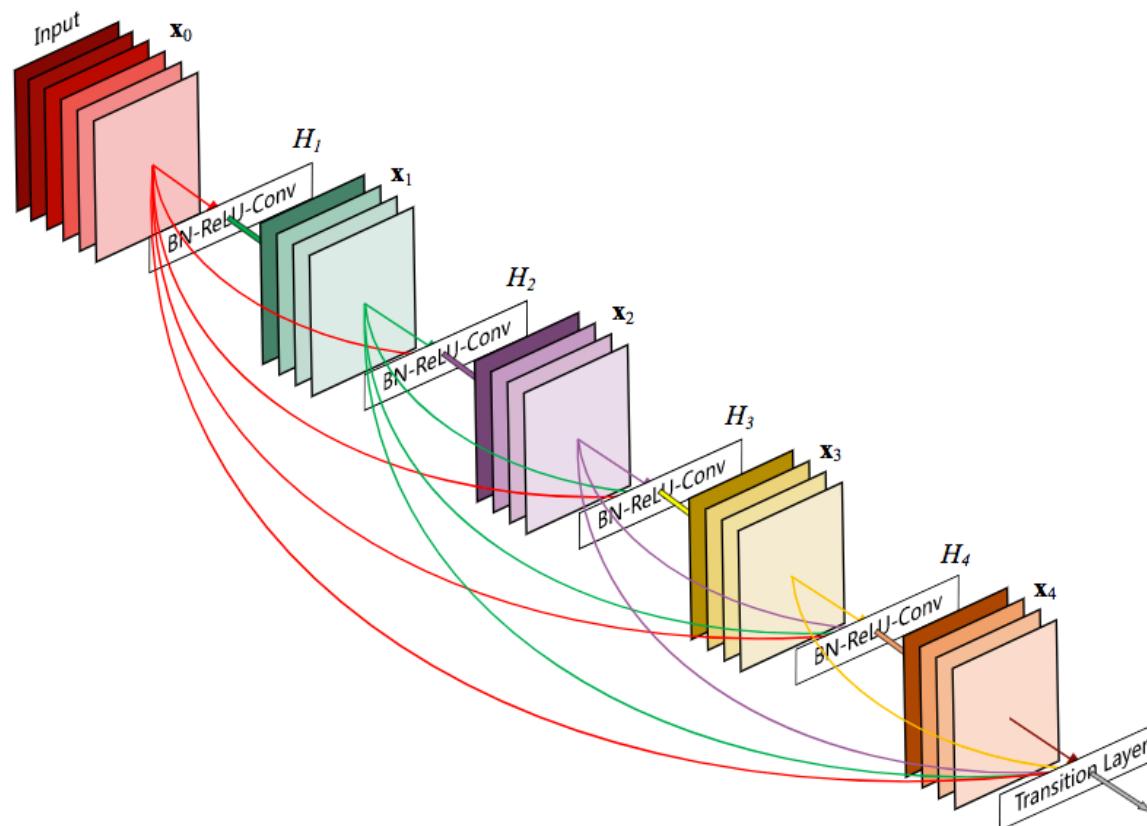
Residual network (ResNet)

- 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





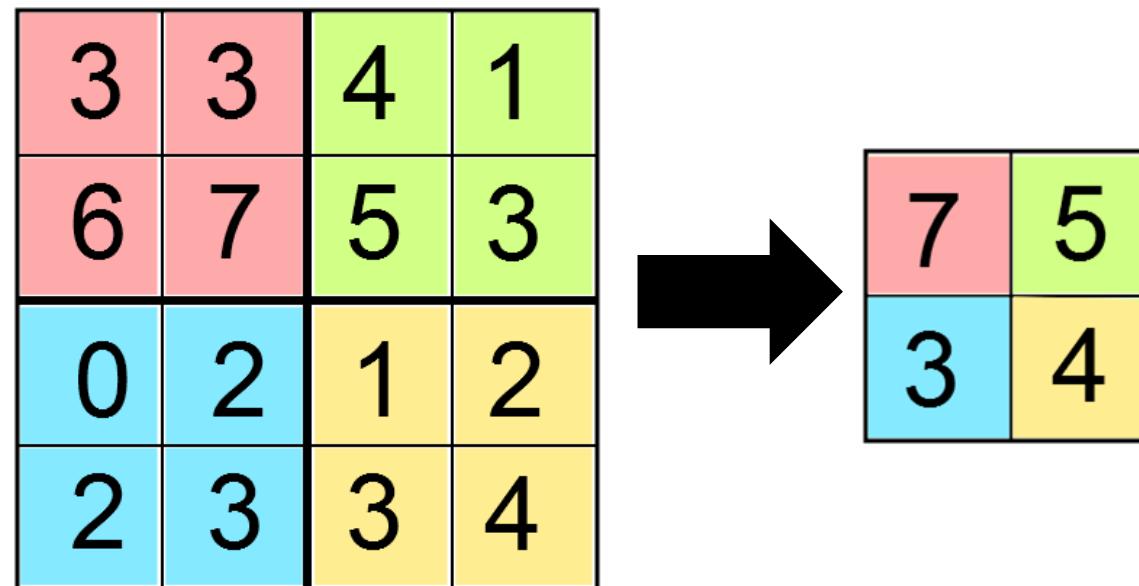
Dense networks





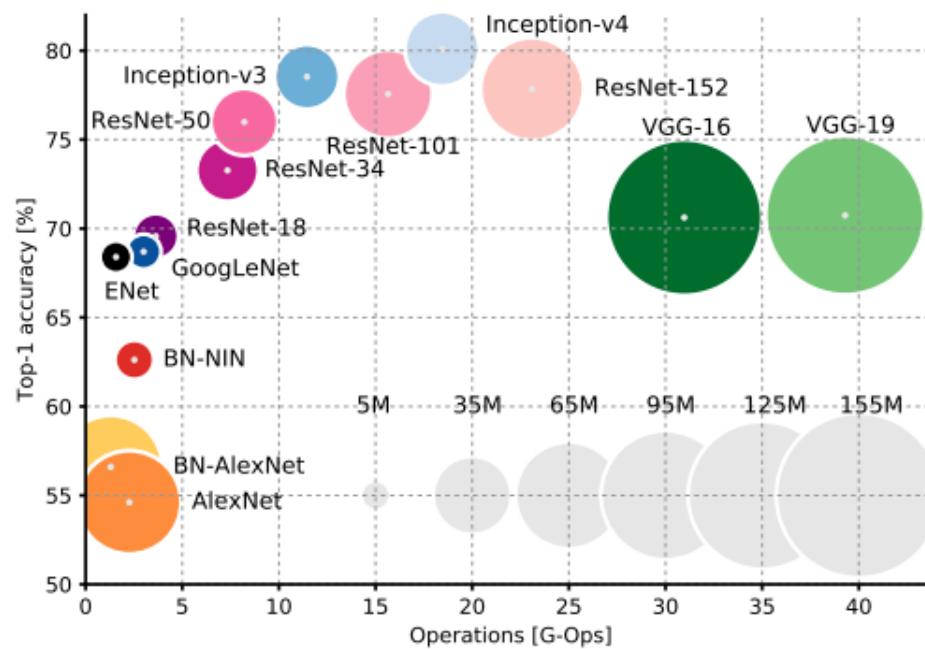
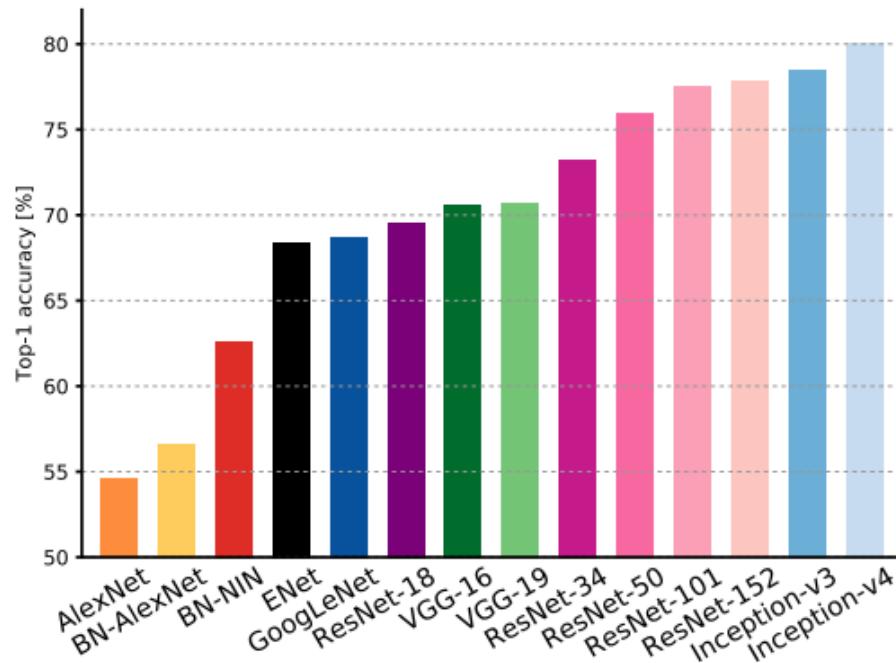
Downsampling

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





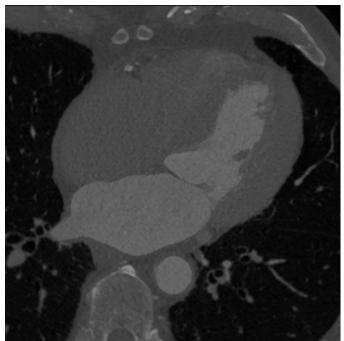
Complexity vs. accuracy



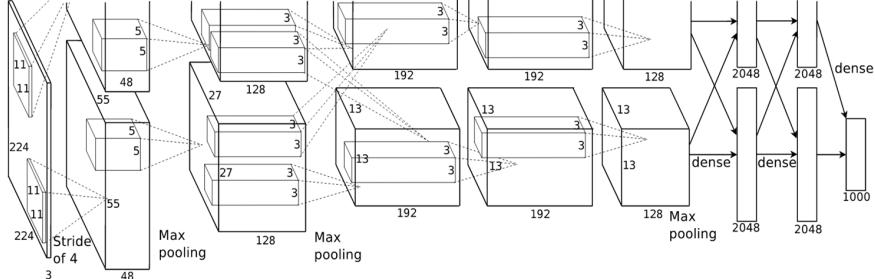


Example: Organ localization in CT

2D image



AlexNet

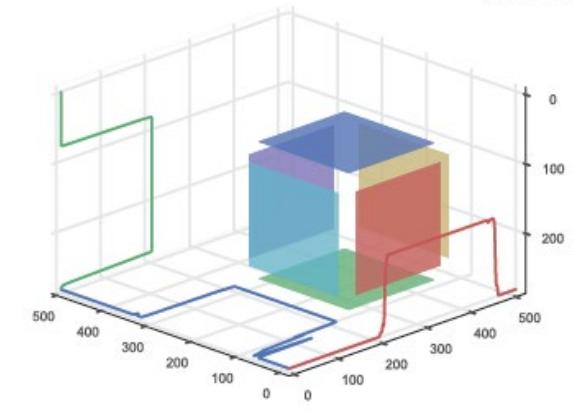
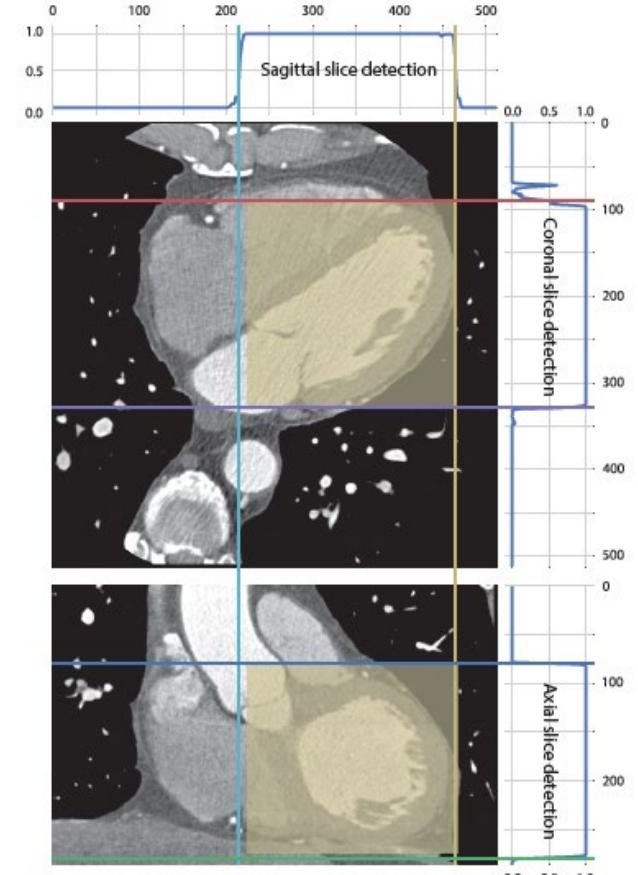


Ventricle?

- Yes
- No

→

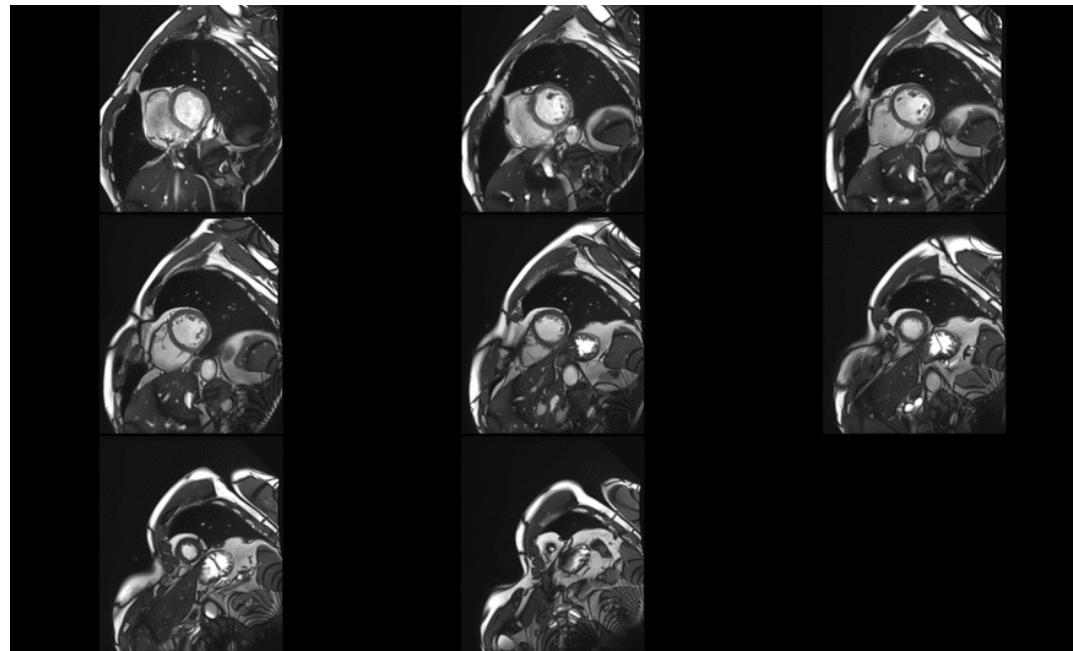
For each image slice





Sequences

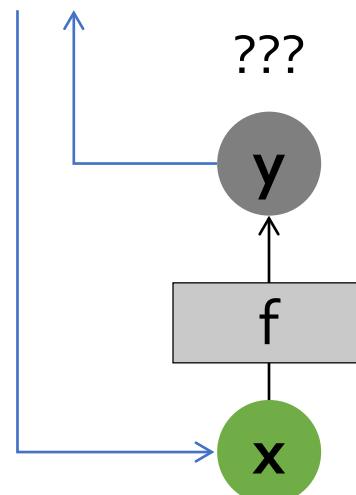
- A lot of data is sequential
- E.g. videos, audio, text, ECG, medical images, ...
- Can we use this in our neural network?





Recurrent neural networks (RNNs)

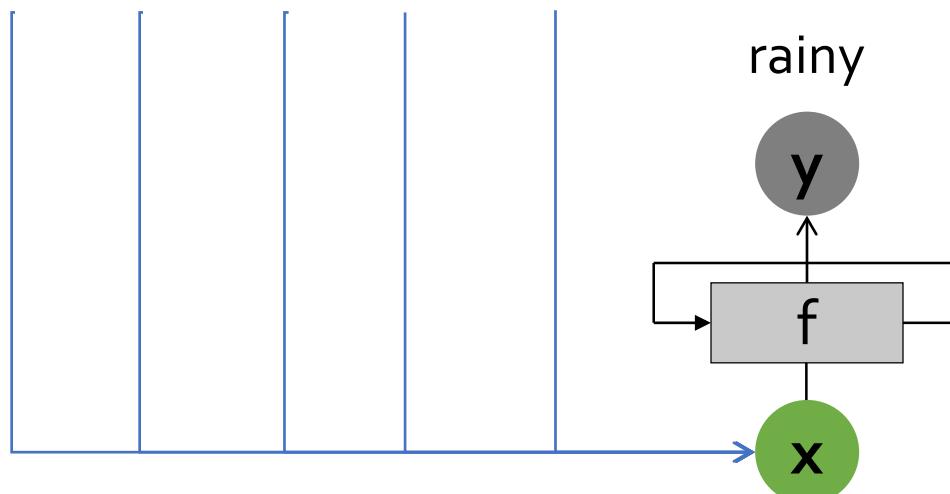
- It would be good to use information that came before
- A feedforward neural network has no ‘memory’
- Consider training a neural network to predict the next word
 - “It’s October, the weather is ...”





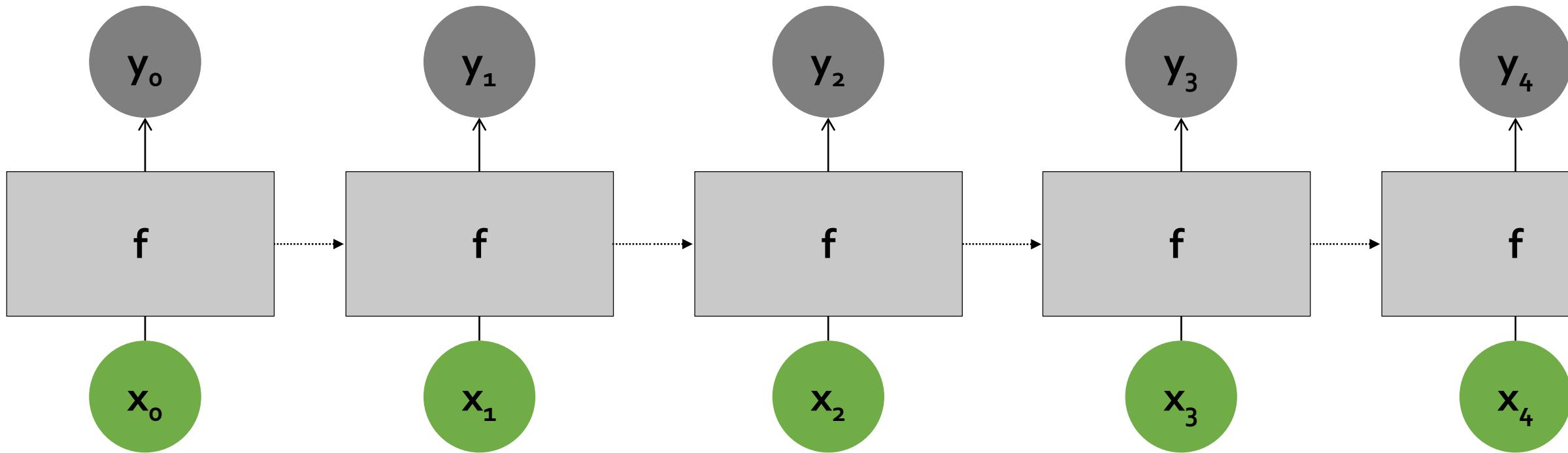
Recurrent neural networks (RNNs)

- 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 October, the weather is ...”





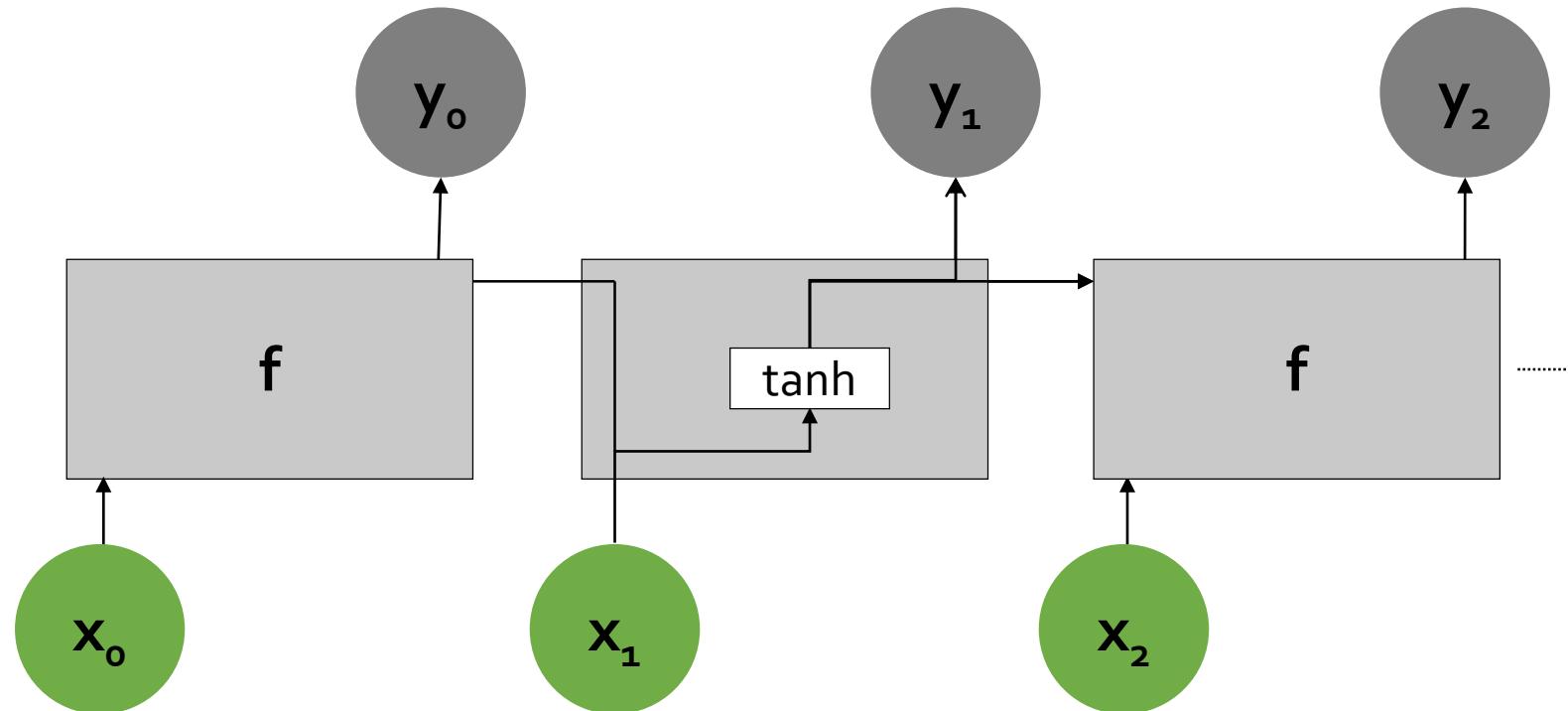
Unrolling a recurrent neural network





Unrolling a recurrent neural network

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



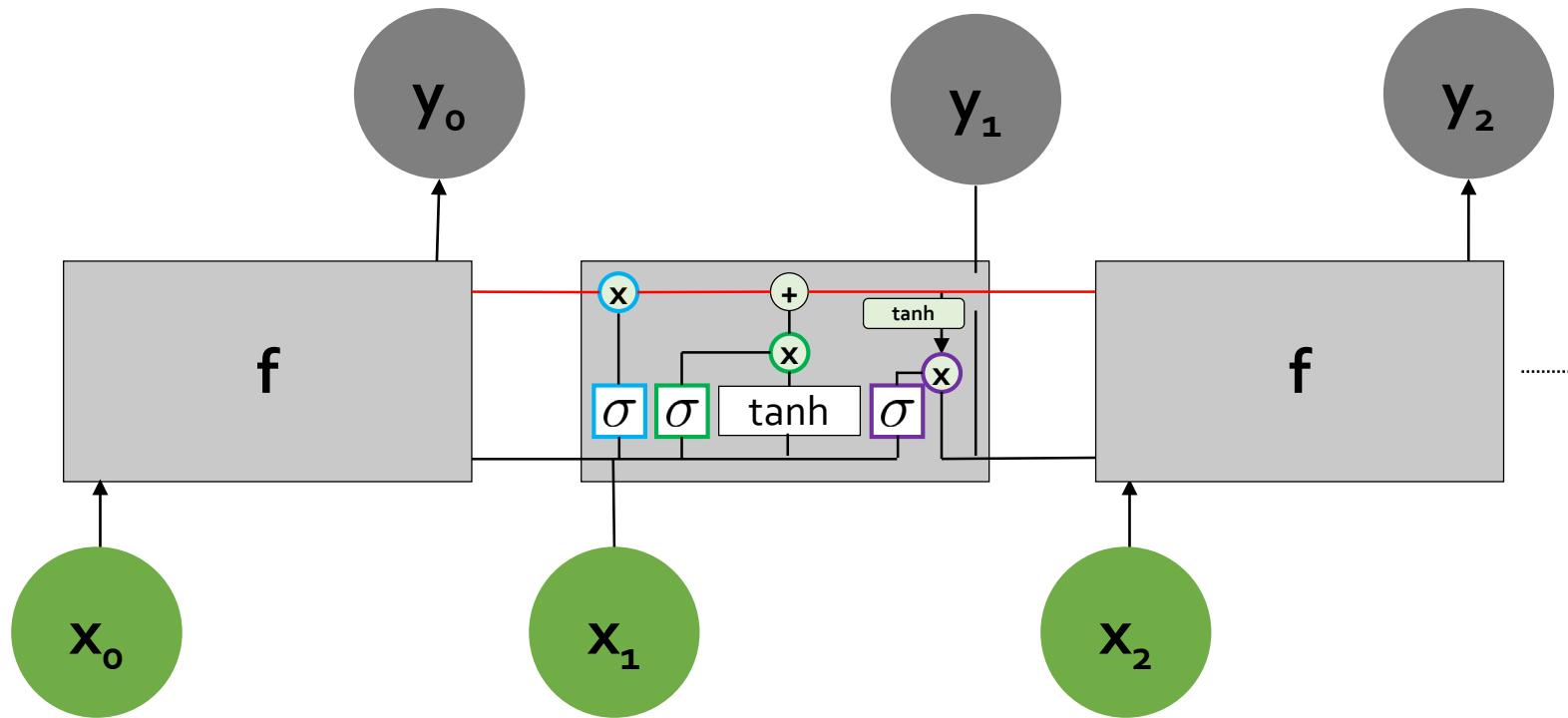


Long short-term memory (LSTM)

1. Cell state

2. Gates

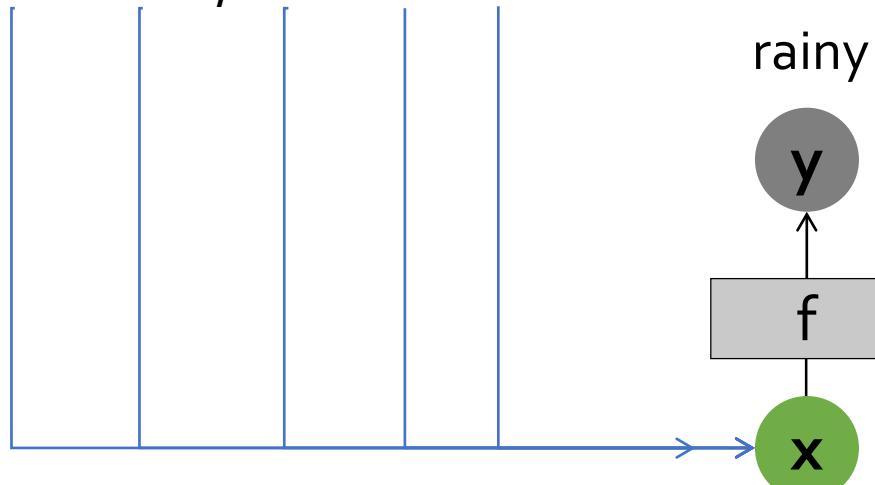
- Forget gate
- Input gate
- Output gate





Recurrent vs. feedforward

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



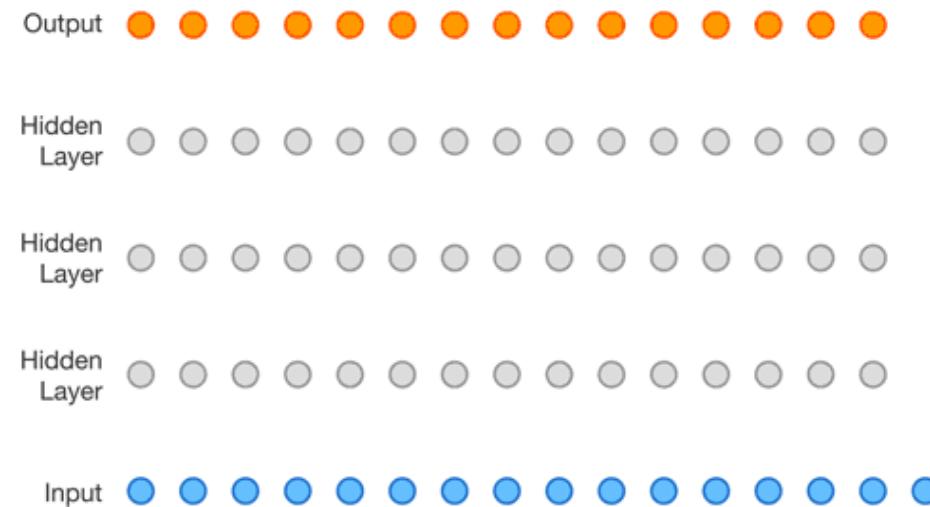
- One way to deal with large contexts in feedforward networks
 - dilated convolutions



Dilated convolutions

In each layer, add more spacing between elements

- increase receptive field
- prevent explosion in number of parameters



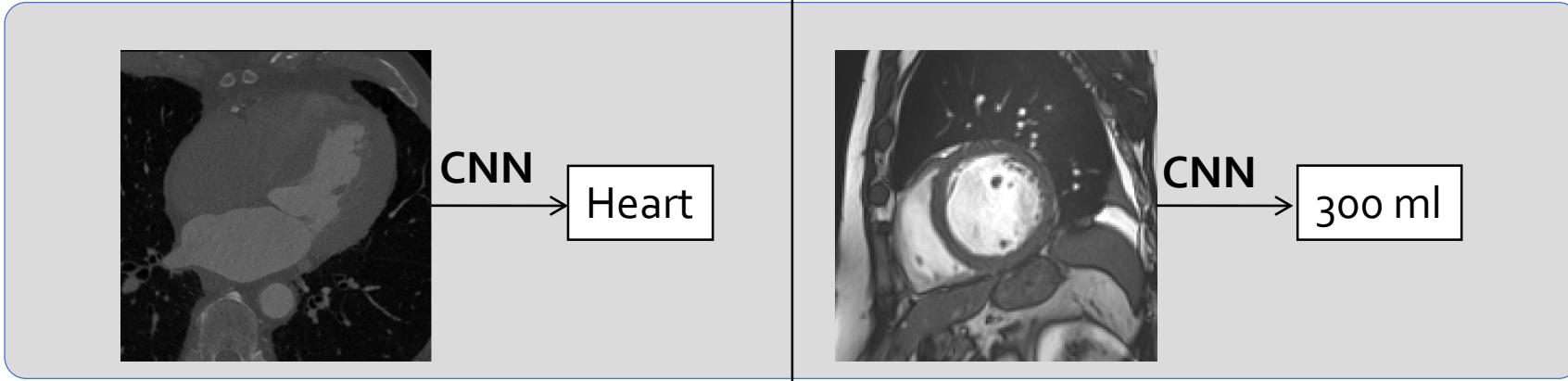


Back to images

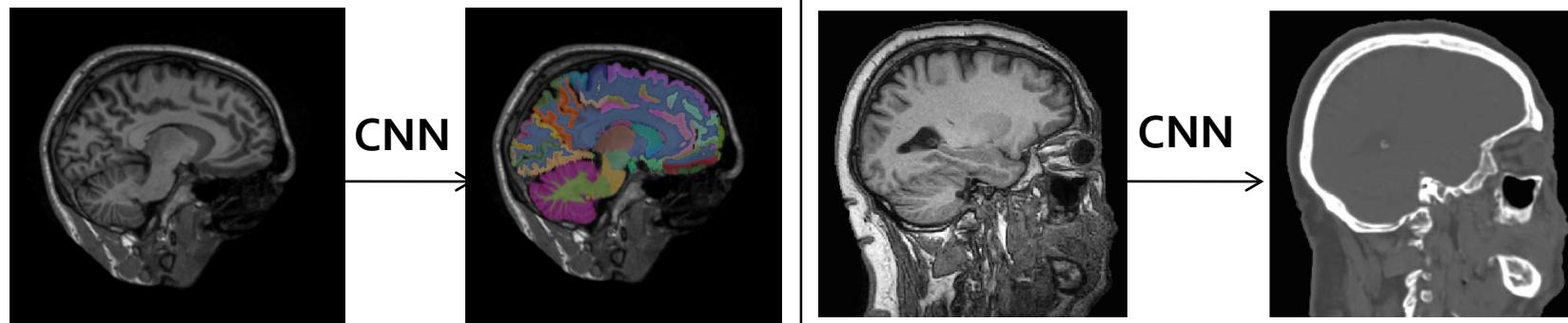
Classification

Regression

Image



Voxel



de Vos, Bob D., et al. "ConvNet-based localization of anatomical structures in 3-D medical images." *IEEE Trans Med Imaging* 36.7 (2017): 1470-1481.

Luo, Gongning, et al. "Multi-views fusion CNN for left ventricular volumes estimation on cardiac MR images." *IEEE Transactions on Biomedical Engineering* 65.9 (2018): 1924-1934.

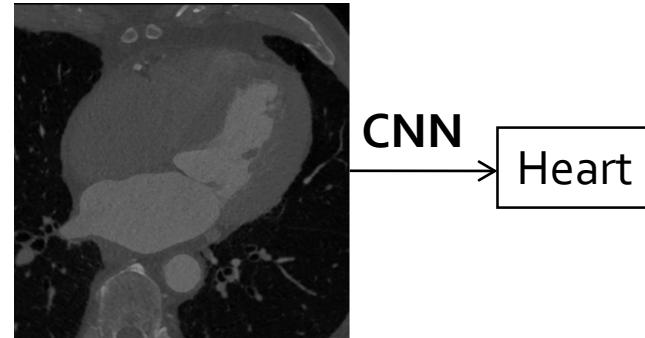
Moeskops, Pim, et al. "Automatic segmentation of MR brain images with a convolutional neural network." *IEEE transactions on medical imaging* 35.5 (2016): 1252-1261.

Wolterink, Jelmer M., et al. "Deep MR to CT synthesis using unpaired data." *International Workshop on Simulation and Synthesis in Medical Imaging*. Springer, Cham, 2017.

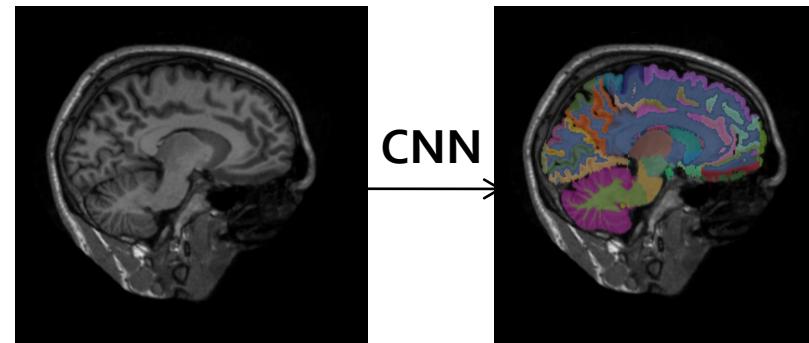


CNNs for pixelwise prediction

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

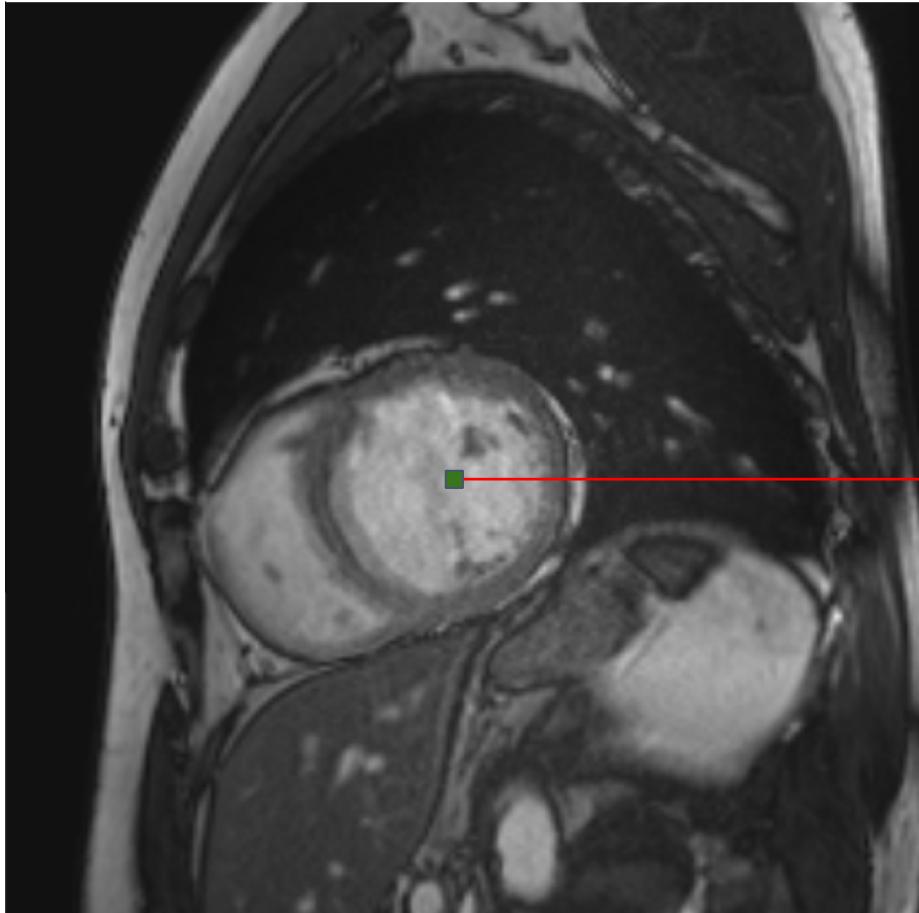


Often, we want to predict one value per **pixel/voxel**





Sliding window

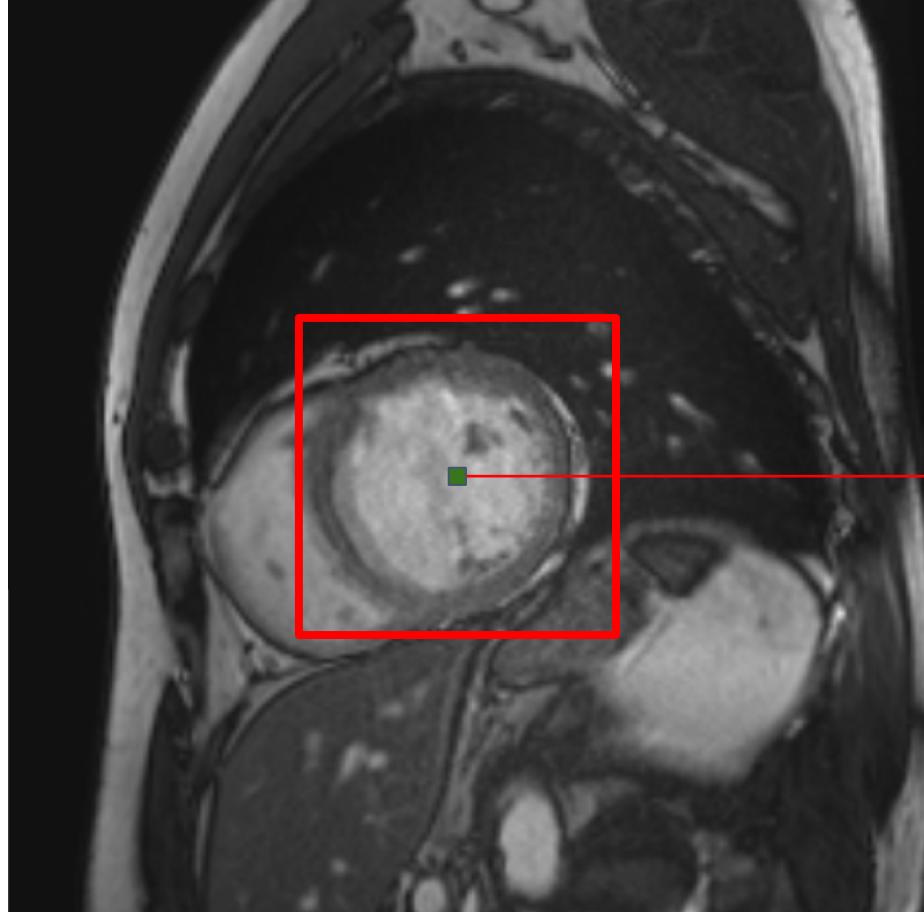


- **Image** = patch centered at voxel
- **Label** = class of center voxel

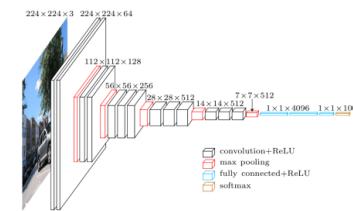
- Background
- Left ventricle
- Myocardium
- Right ventricle



Sliding window



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

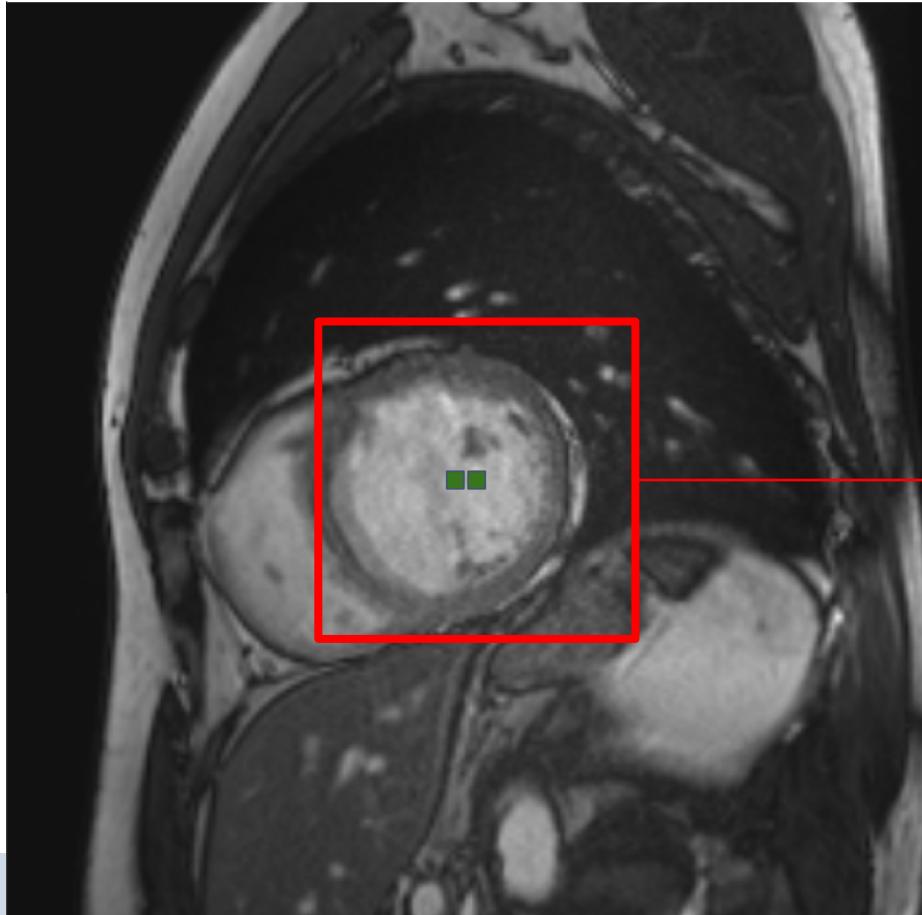


- Background
- Left ventricle
- Myocardium
- Right ventricle

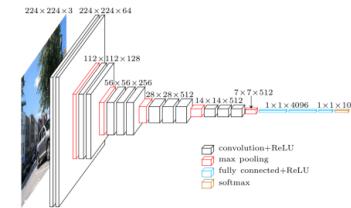


Sliding window

Combination of thousands of image classification tasks



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

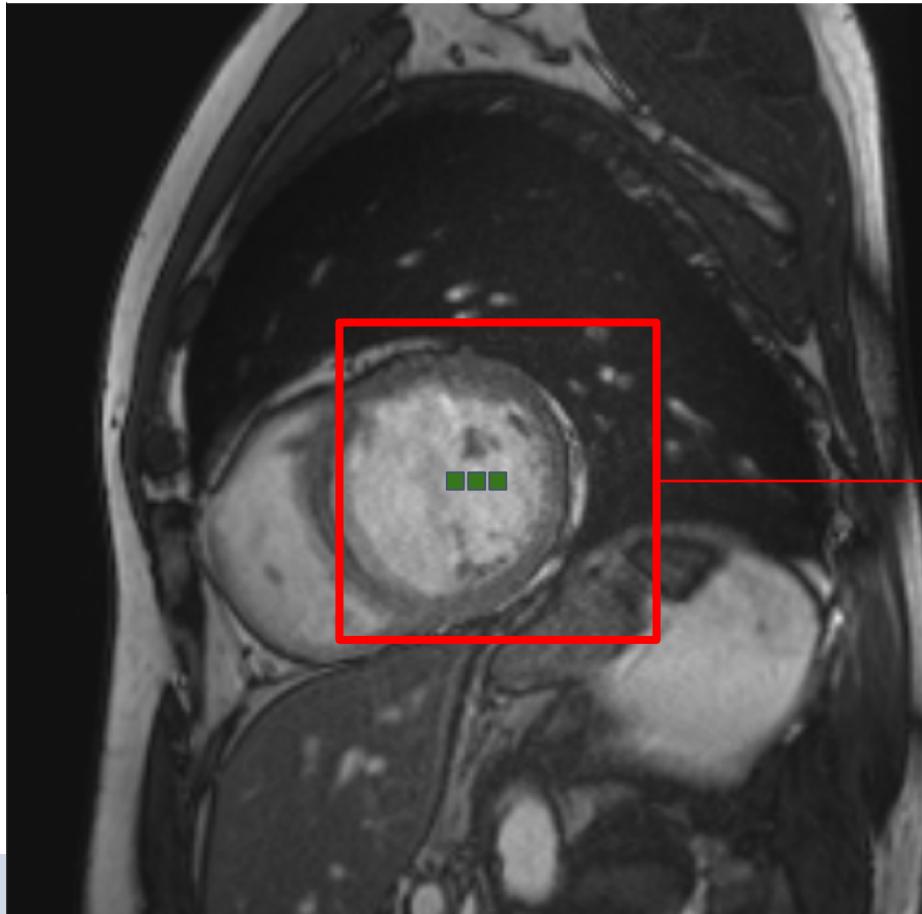


- Background
- Left ventricle
- Myocardium
- Right ventricle

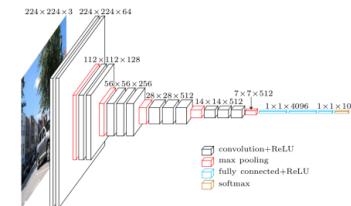


Sliding window

Combination of thousands of image classification tasks



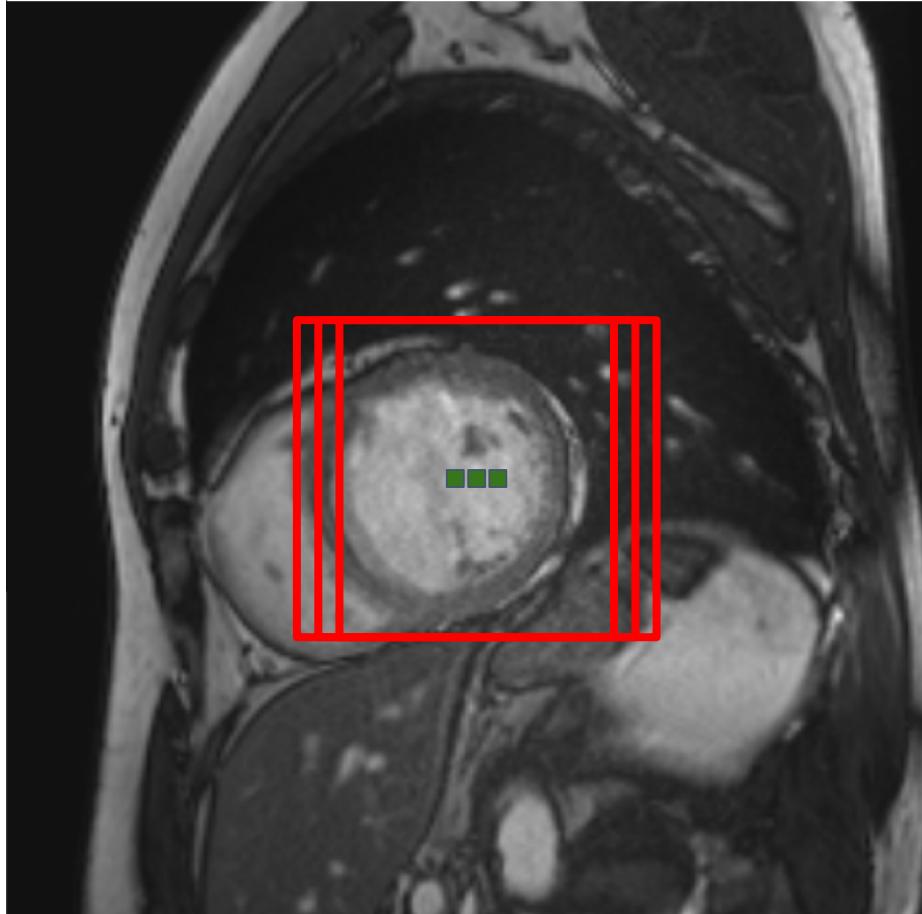
- **Image** = patch centered at voxel
- **Label** = class of center voxel



- Background
- Left ventricle
- Myocardium
- Right ventricle



Sliding window

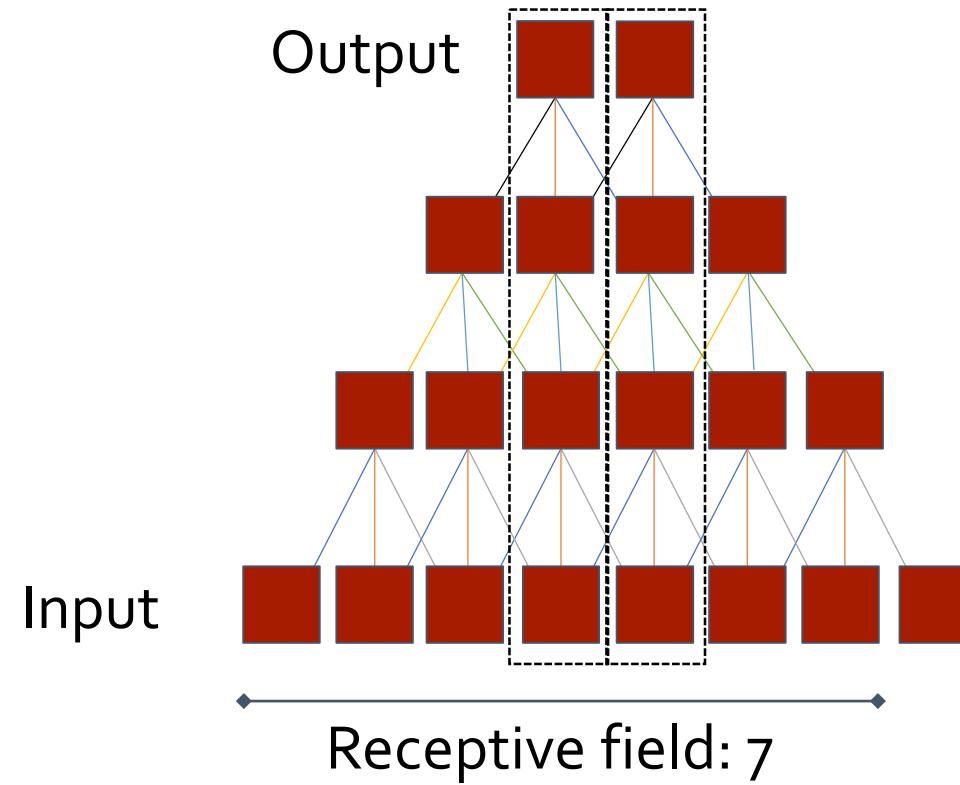


Sliding window approaches are inefficient

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



All-convolutional network



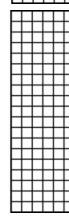


Dilated convolutions



Output

4-dilated

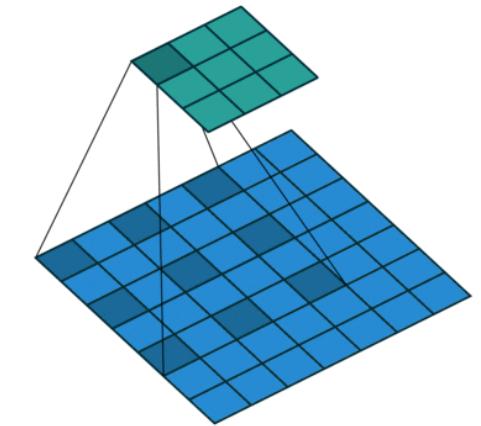
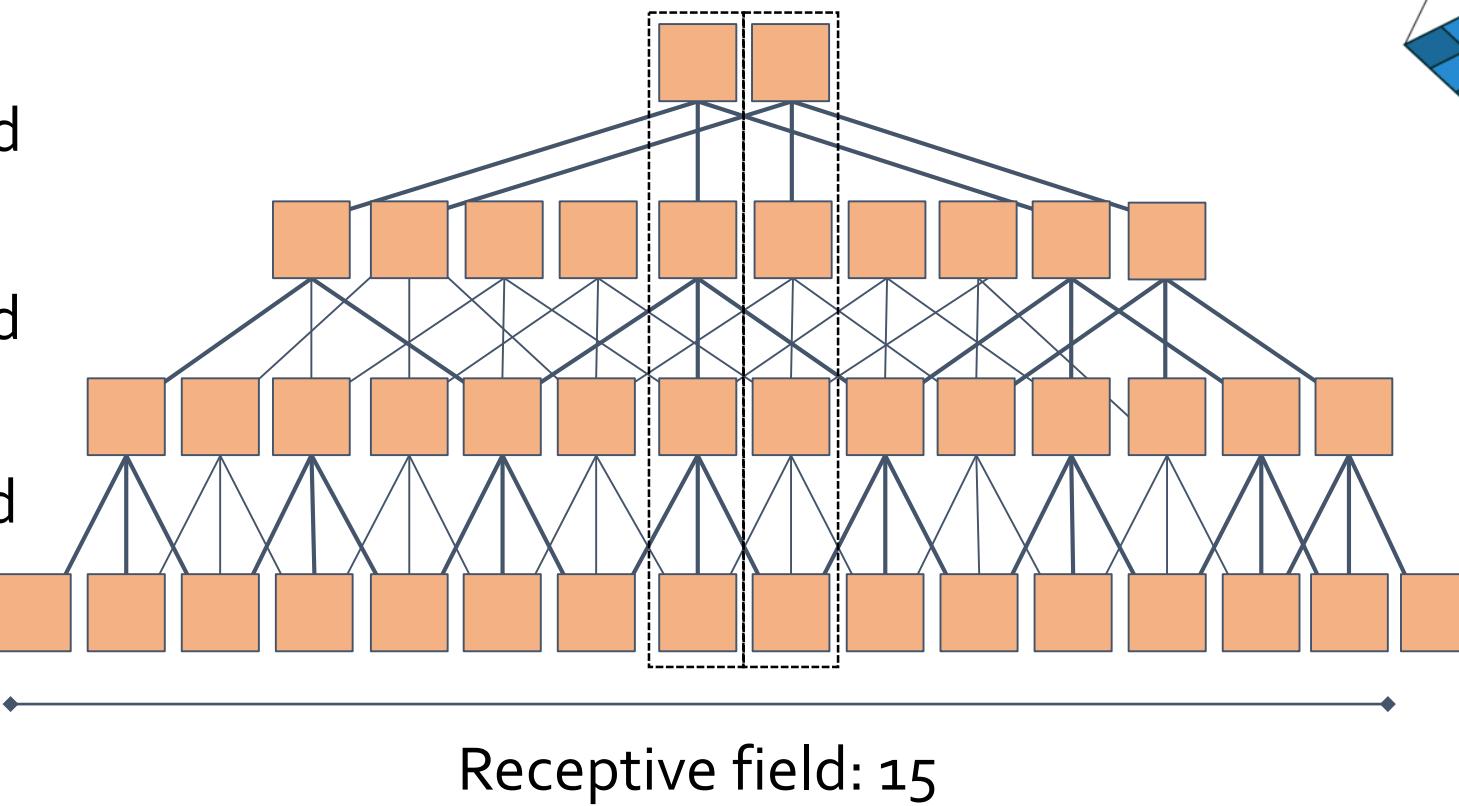


2-dilated



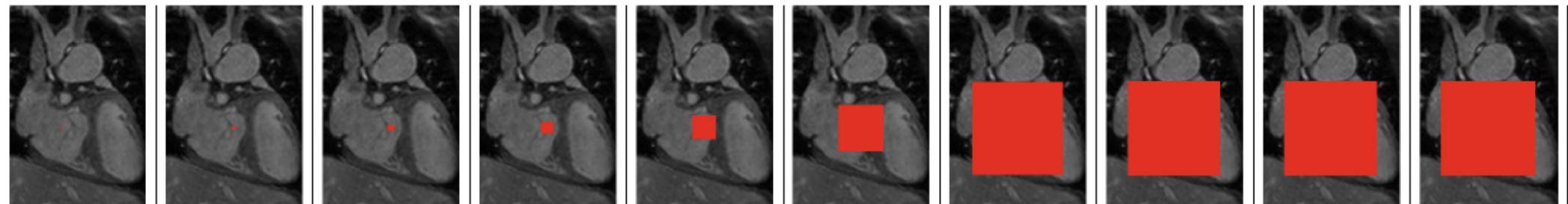
1-dilated

Input





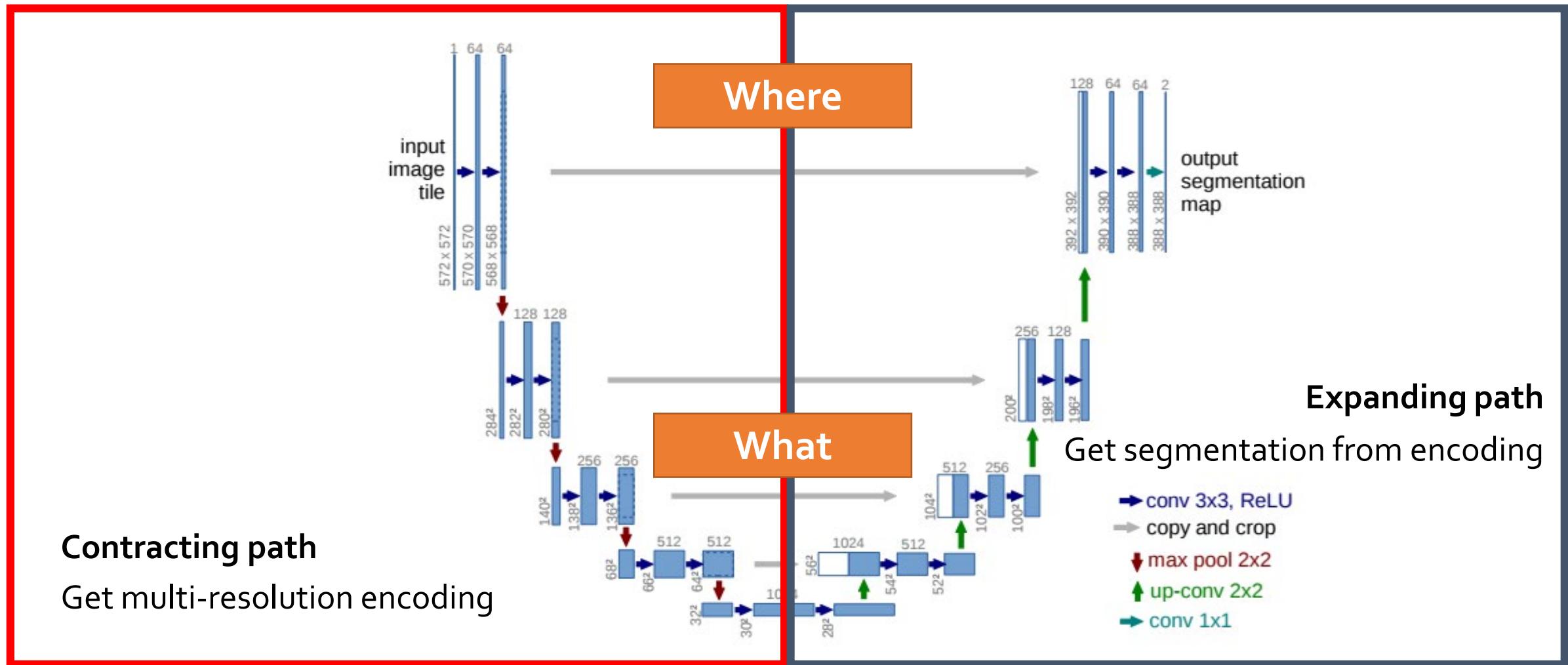
Example



Layer	1	2	3	4	5	6	7	8	9	10
Convolution	3×3	3×3	3×3	3×3	3×3	3×3	3×3	3×3	1×1	1×1
Dilation	1	1	2	4	8	16	32	1	1	1
Field	3×3	5×5	9×9	17×17	33×33	65×65	129×129	131×131	131×131	131×131
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



Contracting path

Get multi-resolution encoding

Where

What

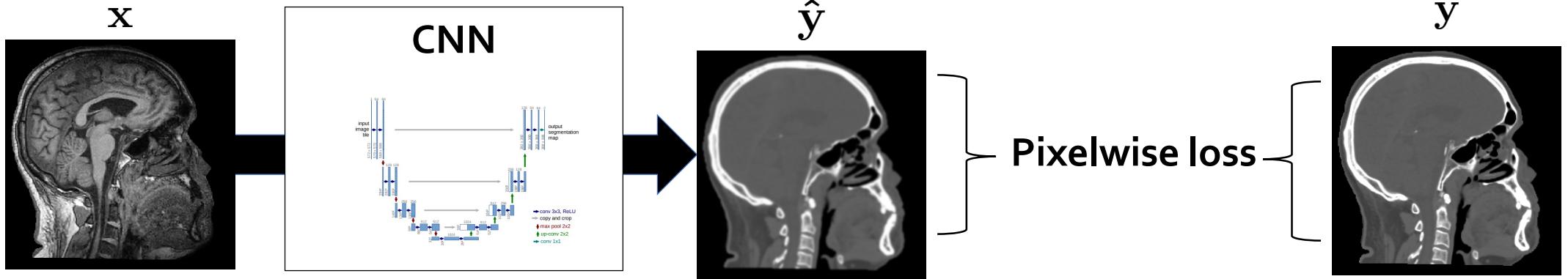
Expanding path

Get segmentation from encoding

- Blue arrow: conv 3x3, ReLU
- Grey arrow: copy and crop
- Red arrow: max pool 2x2
- Green arrow: up-conv 2x2
- Cyan arrow: conv 1x1



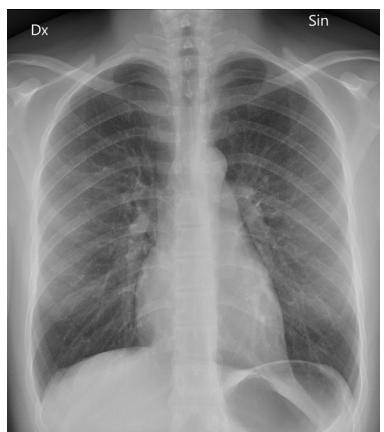
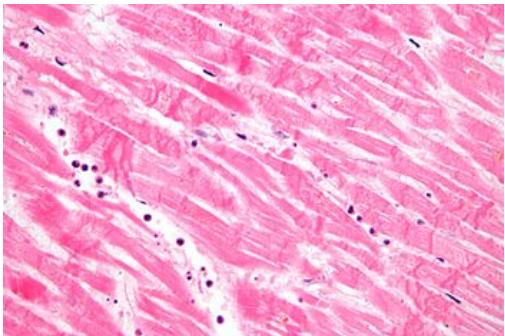
Training



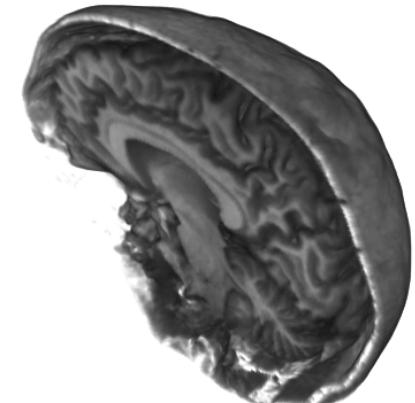
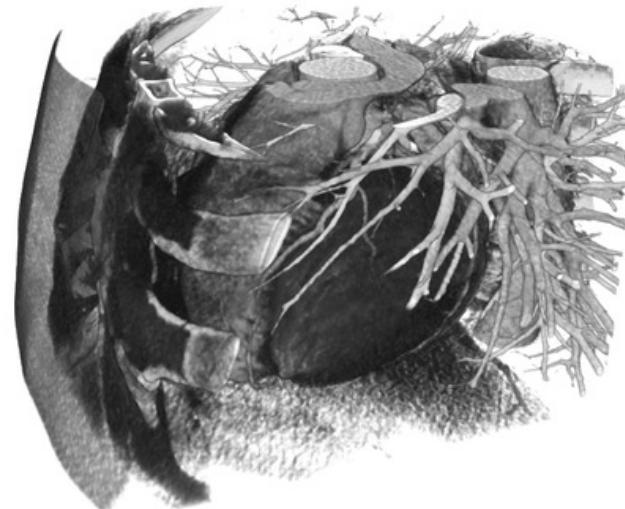


2D or 3D images

2D data



3D data





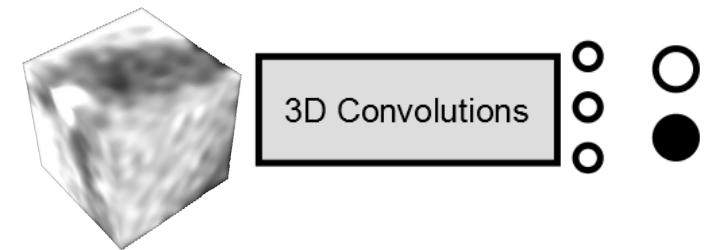
3D networks

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 FACE NSFW COLOR

MORE MODELS ▾

VIEW DOCS



General

LANGUAGE

English (en)

PREDICTED CONCEPT

coffee

0.994

cup

0.993

dawn

0.988

hot

0.987

no person

0.987

breakfast

0.987

caffeine

0.980

still life

0.978



TRY YOUR OWN IMAGE OR VIDEO

