# Next Topic: Fully-Convolutional Networks
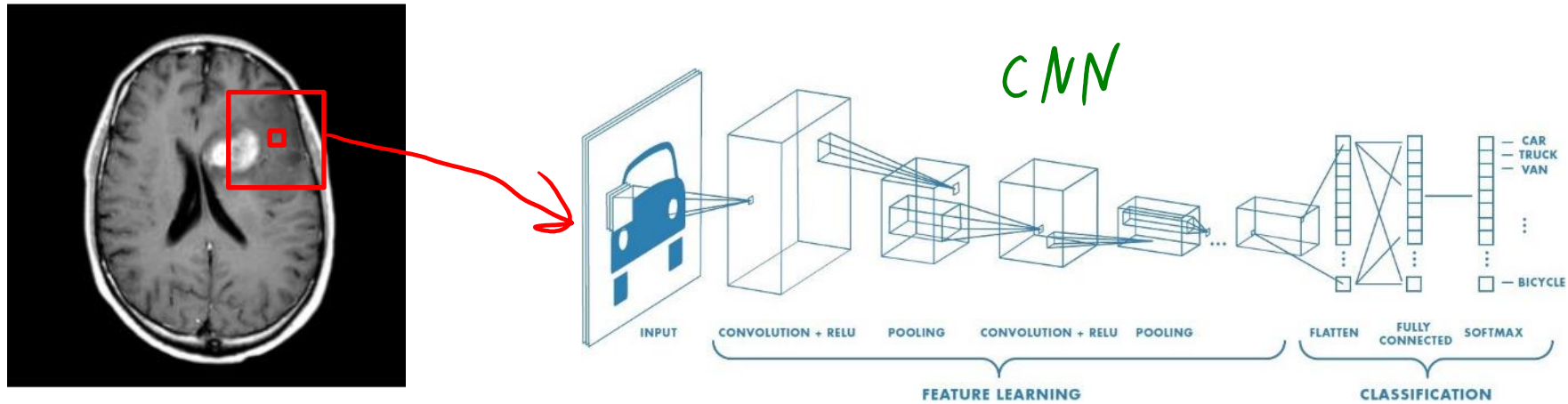
# Motivation: Pixel Classification

- Suppose we want to make a prediction at each pixel in an image:
  - Segmenting tumours, segmenting street scenes, depth estimation.



- How can we use CNNs for this problem?
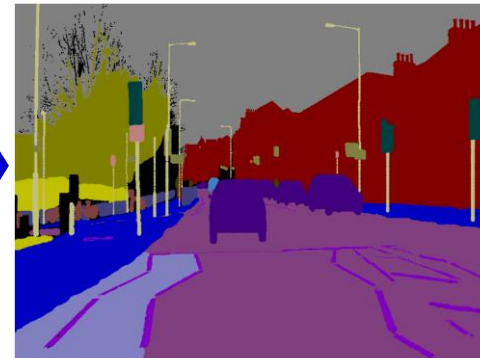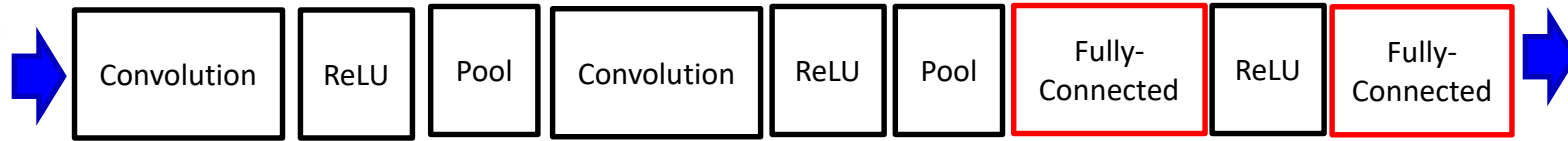
# Naïve Approach: Sliding Window Classifier

- Train a CNN that predicts pixel label given its neighbourhood.



- To label all pixels, apply the CNN to each pixel.
  - Advantages:
    - Turns pixel labeling into image classification.
    - Can be applied to images of different sizes.
  - Disadvantage: this is slow.
    - (Cost of applying CNN) * (number of pixels in the image).
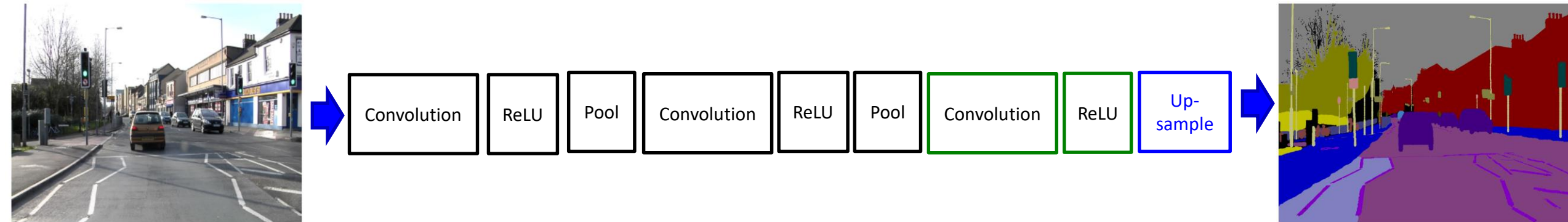
# Encoding-Decoding for Pixel Classification

- Similar to multi-label, could use CNN to generate an image encoding.
  - With output layer making a prediction at each pixel.



- Much-faster classification.
  - Small number of "global" convolutions, instead of repeated "local" convolutions.
- But, the encoding has mixed up all the space information.
  - Due to fully-connected layers.
  - Fully-connected layer needs to learn "how to put the image together".
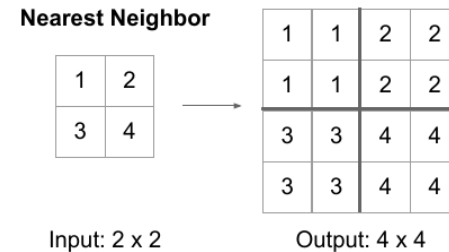    - And all input images must be the same size.

# Fully-Convolutional Networks

- Fully-convolutional networks (FCNs):
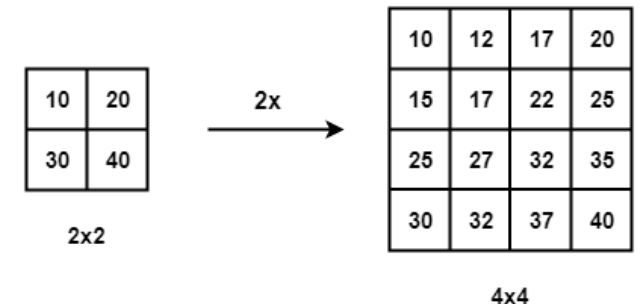  - CNNs with no fully-connected layers (only convolutional and pooling).



- Maintains fast classification of the encoding-decoding approach.
- Same parameters used across space at all layers.
  - This allows using the network on inputs of different sizes.
- Needs upsampling layer(s) to get back to image size.
- FCNs quickly achieved state of the art results on many tasks.

# Traditional Upsampling Methods

- In upsampling, we want to go from a small image to a bigger image.
  - This requires interpolation: guessing "what is inbetween the pixels".
- Classic upsampling operator is nearest neighbours interpolation:
  - But this creates blocky/pixelated images.



- Another classic method is bilinear interpolation:
  - Weighted combination of corners:
  - More smooth methods include "bicubic" and "splines".



- In FCNs, we learn the upsampling/interpolation operator.

# Upsampling with Transposed Convolution

- FCN Upsampling layer is implemented with a transposed convolution.
  - Sometimes called "deconvolution" in ML or "fractionally-strided convolution".
    - But not related to deconvolution in signal processing.
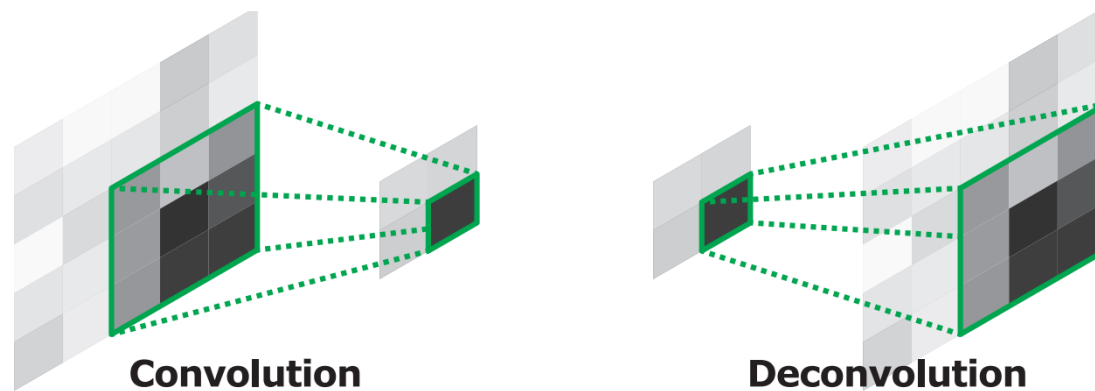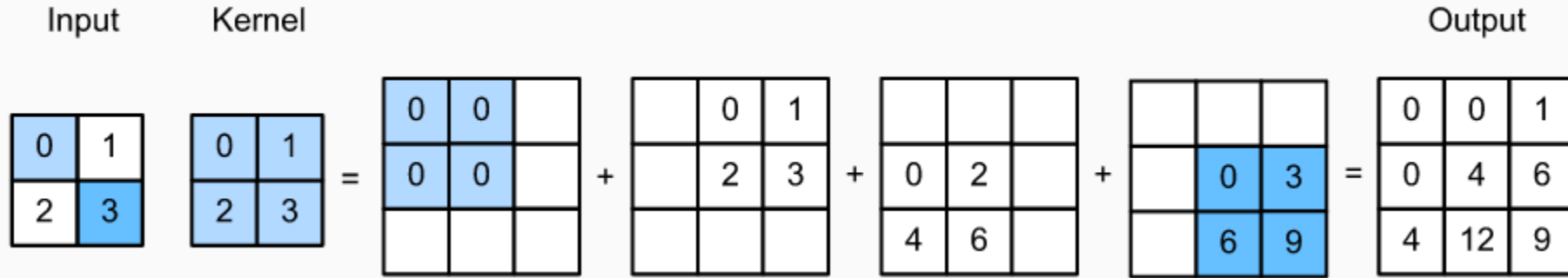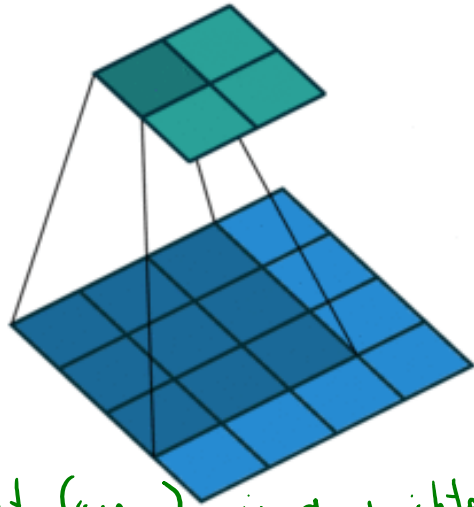


**Convolution**          **Deconvolution**

Figure 3. Illustration of deconvolution and unpooling operations.

- Convolution generates 1 pixel by taking weighted combination of several pixels.
  - And we learn the weights.
- Transposed convolution generates several pixels by weighting 1 pixel.
  - And we learn the weights.
  - This generates overlapping regions, which get added together to make final image.
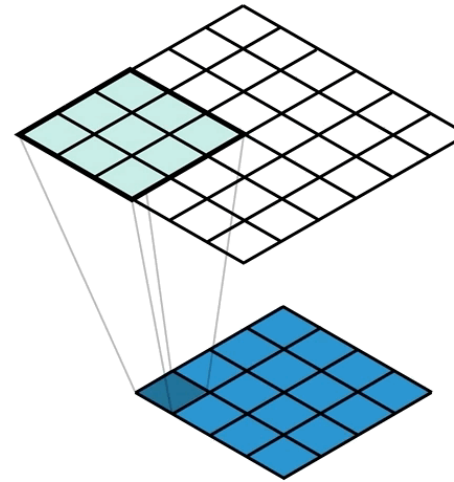
# Upsampling with Transposed Convolution



Convolution

Each output (green) is a weighted combination of one 3x3 region

Transposed convolution

Each input (blue) gives a 3x3 region. These regions overlap, so we take a weighted combination at each pixel to give final result.

- Animations [here](#) and [here](#).

# Why is it called "transposed" convolution?

- We can write the convolution operator as a matrix multiplication, z=Wx.



- In transposed convolution, non-zero pattern of 'W' is transposed from convolution.
  - You can implement transposed convolution as a convolution.



- In this example the filter is the same, but does not need to be:
  - Transposed convolution is not the "reverse" of convolution (it only "reverses" the size).

# Increasing Resolution: FCN Skip Connections

- Convolutions and pooling lose a lot of information.
- Original FCN paper considered adding skip connections to help upsampling:
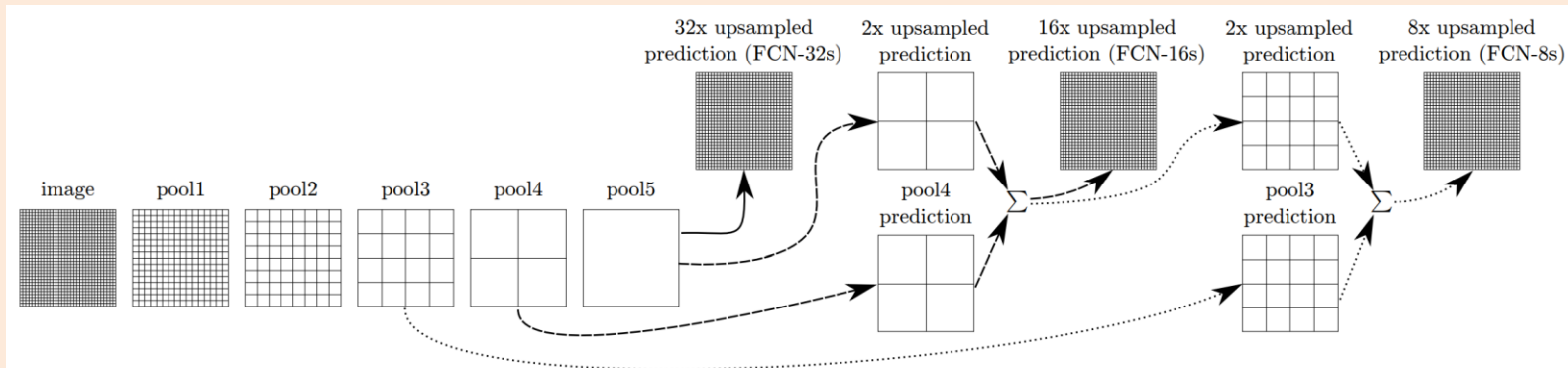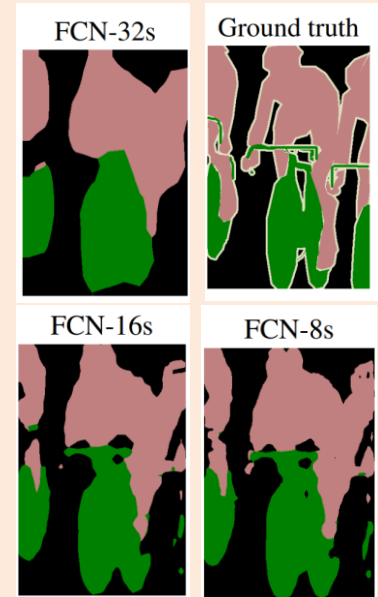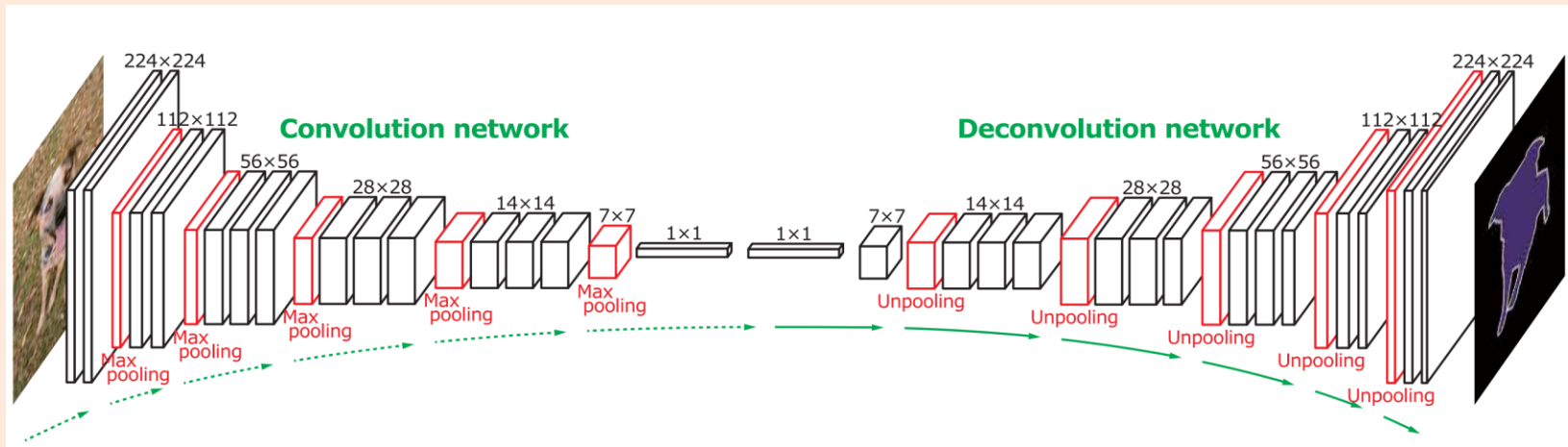


Figure 3. Our DAG nets learn to combine coarse, high layer information with fine, low layer information. Layers are shown as grids that reveal relative spatial coarseness. Only pooling and prediction layers are shown; intermediate convolution layers (including our converted fully connected layers) are omitted. Solid line (FCN-32s): Our single-stream net, described in Section 4.1, upsamples stride 32 predictions back to pixels in a single step. Dashed line (FCN-16s): Combining predictions from both the final layer and the pool4 layer, at stride 16, lets our net predict finer details, while retaining high-level semantic information. Dotted line (FCN-8s): Additional predictions from pool3, at stride 8, provide further precision.

- Allows using high-resolution information from earlier layers.
- They first trained the low-resolution FCN-32, then FCN-16, then FCN-8.
  - "First learn to encode at a low resolution", then slowly increase resolution.
  - Parameters of transposed convolutions initialized to simulate "bilinear interpolation".

# Increasing Resolution: Deconvolution Networks

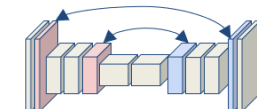- Alternate resolution-increasing method is deconvolution networks:



- Includes transposed convolution layers and unpooling layers.
  - Store the max pooling argmax values.
  - Restores "where" activation happened.
    - Still loses the "non-argmax" information.

# Increasing Resolution: U-Nets

- Another popular variant is u-nets:



- Decoding has connections to same-resolution encoding step.

# High Resolution Predictions with FCNs

- "AutoPortrait": automatic high-resolution photo re-touching.



Figure 1: Overview of AutoRetouch: (1) input portraits cropped to the head, (2) the images pass through the fully convolutional architecture of AutoRetouch, (3) the output portraits are retouched while preserving fine textures. (4) Sample before-after patches scaled up. We empirically show that our final models have desirable generalization properties.

# High Resolution Predictions with FCNs

- "AutoPortrait": automatic high-resolution photo re-touching.

# Source of Labels

- Labeling every pixel takes a <span style="color:red">long time</span>.

- Where we get the labels from?
  - Domain expert (medical images).
  - Grad students or paid labelers (ImageNet).
  - Simulated environments.
    - <span style="color:green">High number</span> of <span style="color:red">lower-quality</span> examples.
    - Often a net gain with fine-tuning on real images.
    - Can get data at night, in fog, or dangerous situations.



Video game



Google street view

# Source of Labels

- Recent works recognize you <span style="color:red">do not need to label every pixel.</span>
  - You can evaluate loss/gradient on a subset of labeled pixels.
  - Could have labeler click on a few pixels inside objects, and a few outside.
    - Many variations are possible, that let you label a lot of images in a short time.
- Penguin counting based on "single pixel" labels in training data:
  - And some tricks to separate objects and remove false positives:

# Summary

- Autoencoders:
  - Neural network where the output is the input.
    - Non-linear generalization of PCA.
  - Encode data into a bottleneck layer, then decode predict original input.
  - Can be used for visualization, compression, outlier detection, pre-training.
- Denoising autoencoders train to uncorrupt/enhance images.
  - Useful for removing noise, adding colour, super-resolution, and so on.
- Encoding-Decoding approach to multi-label classification:
  - Have all classes shared the same hidden layer(s).
  - Reduces number of parameters.
  - Models dependencies between classes, while keeping inference easy.
- Pre-training:
  - Use parameters from model trained a on large diverse dataset, to initialize SGD for new dataset.
- Fully-convolutional networks (FCNs):
  - CNNs where every layer maintains spatial information.
  - Useful for handling images of different sizes.
  - Requires upsampling to be used for pixel classification.
  - Transposed convolutions learn upsampling operators.

- Next time: are we learning harmful biases?