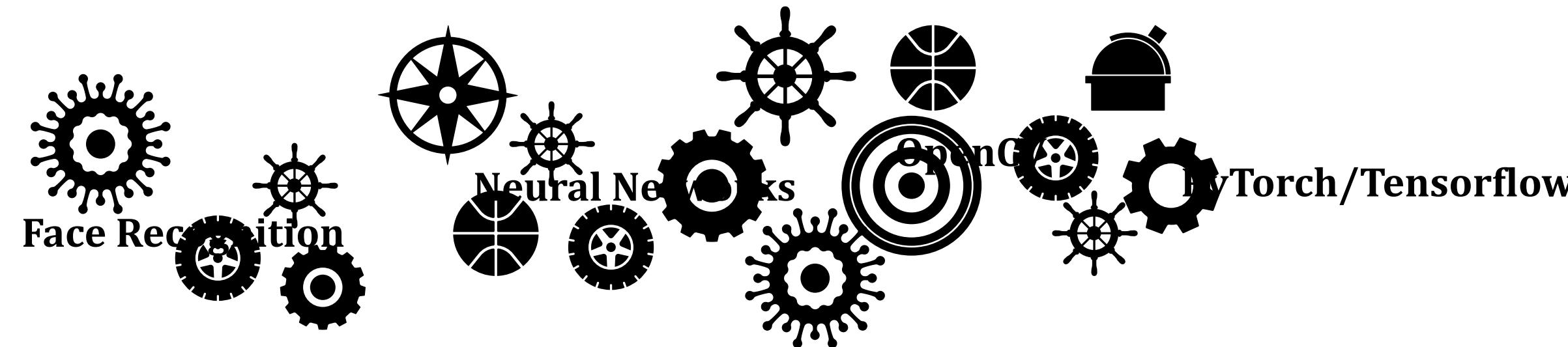
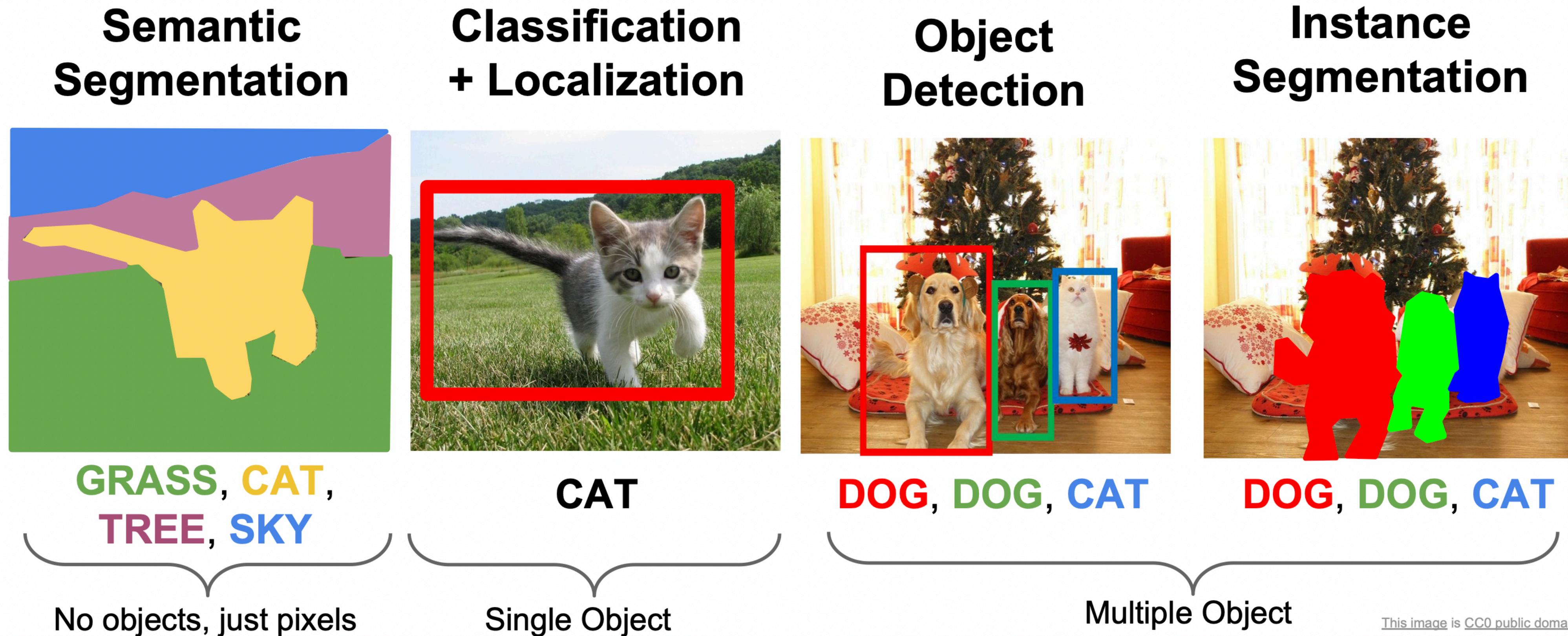


# Computer Vision



Object Detection      Deep Learning      Segmentation  
Image Classification

# Semantic Segmentation

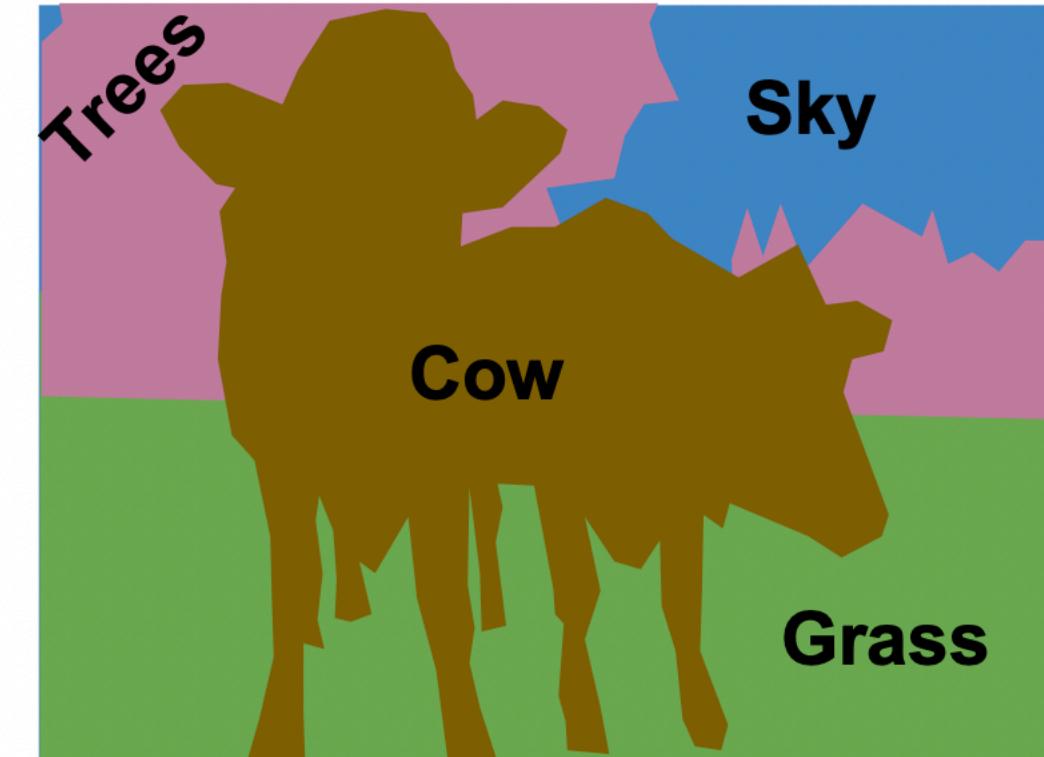
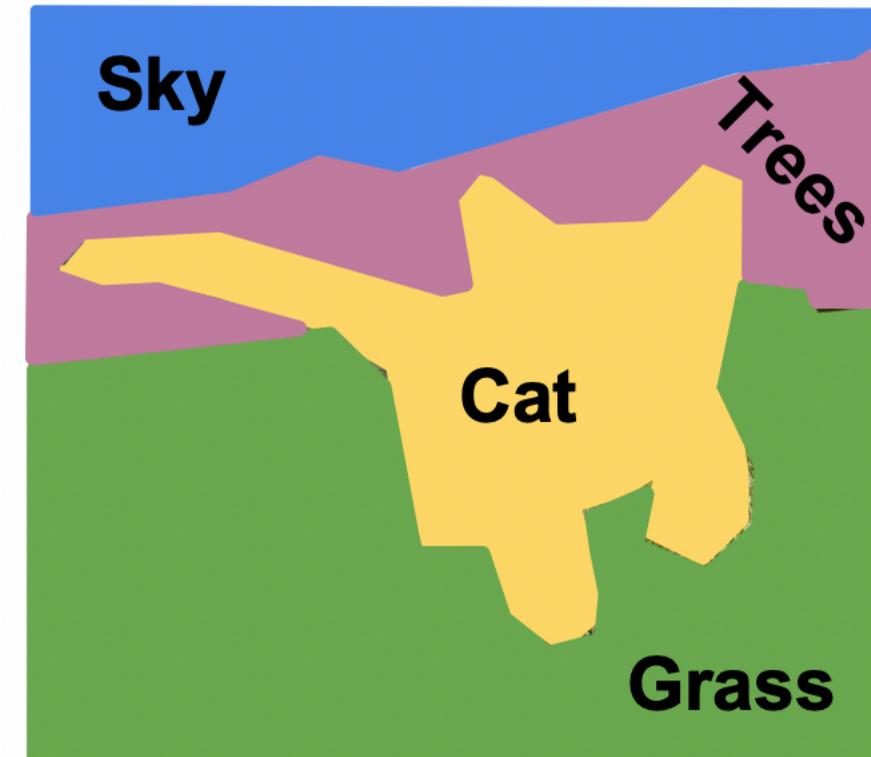


# Semantic Segmentation

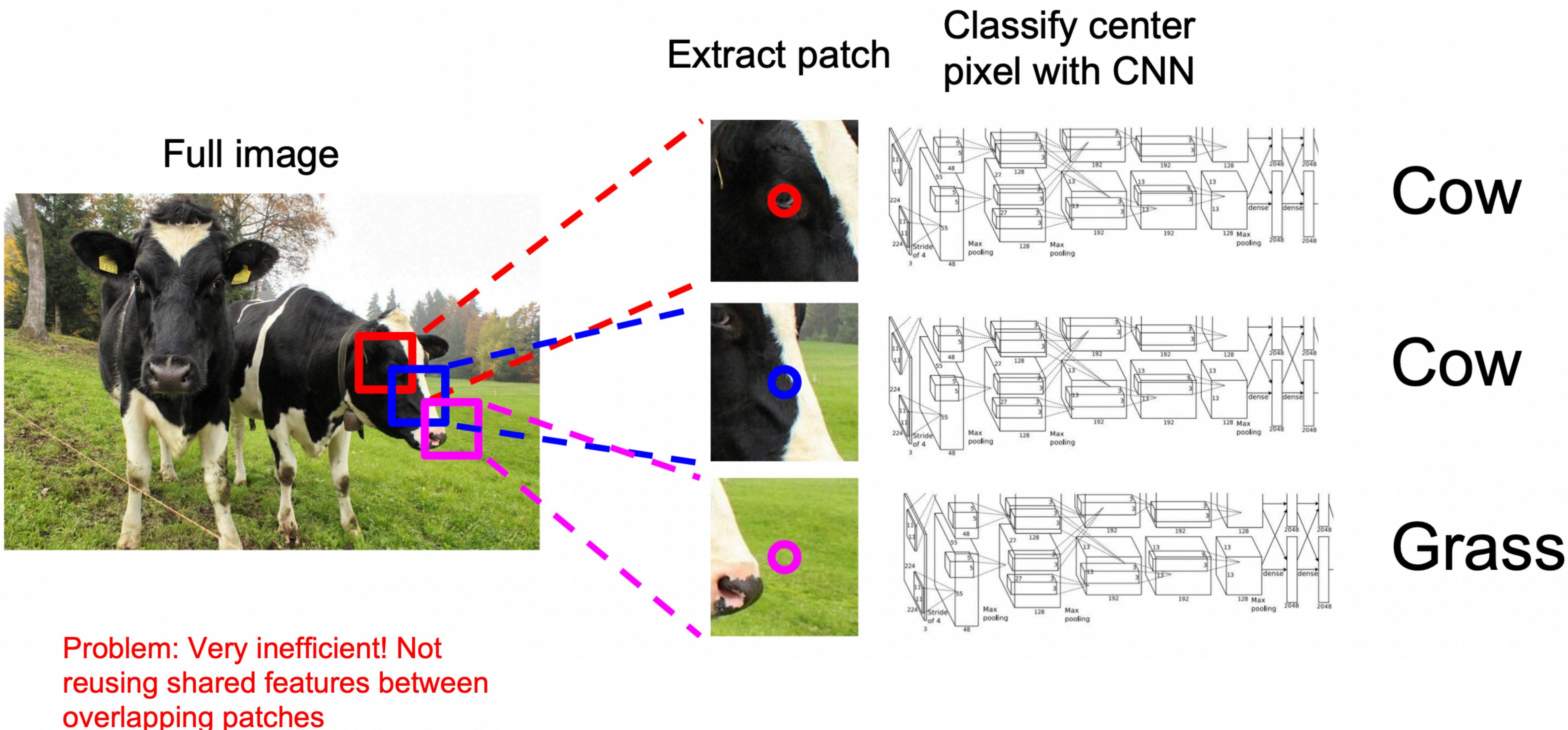
## Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



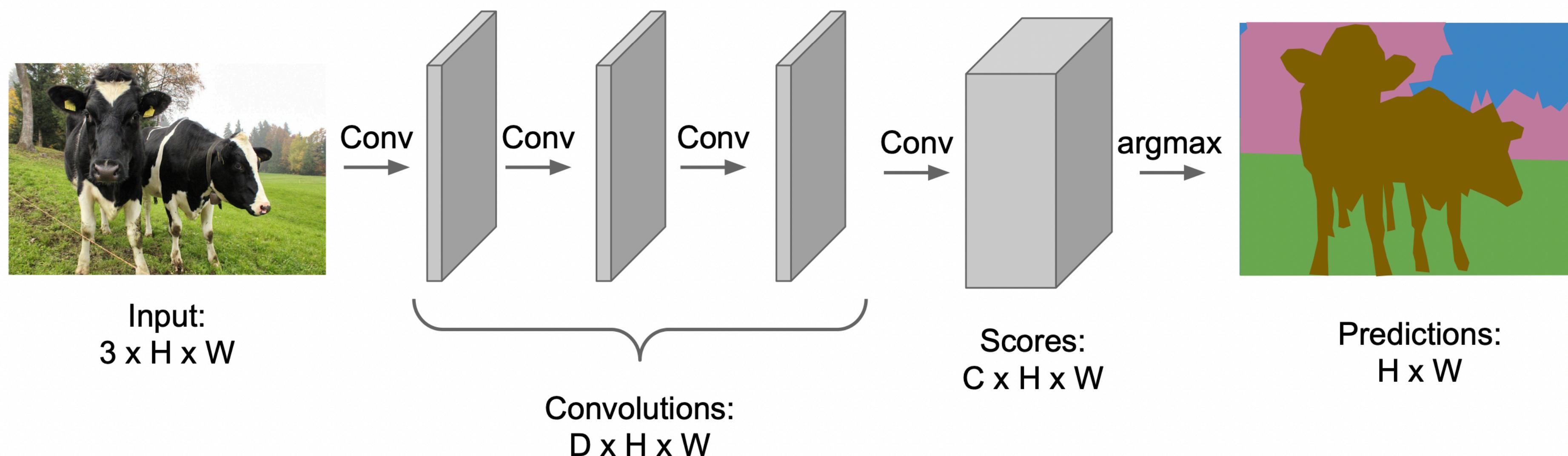
# Semantic Segmentation Idea: Sliding Window



# Semantic Segmentation

## Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers  
to make predictions for pixels all at once!



Problem: convolutions at  
original image resolution will  
be very expensive ...

# Semantic Segmentation

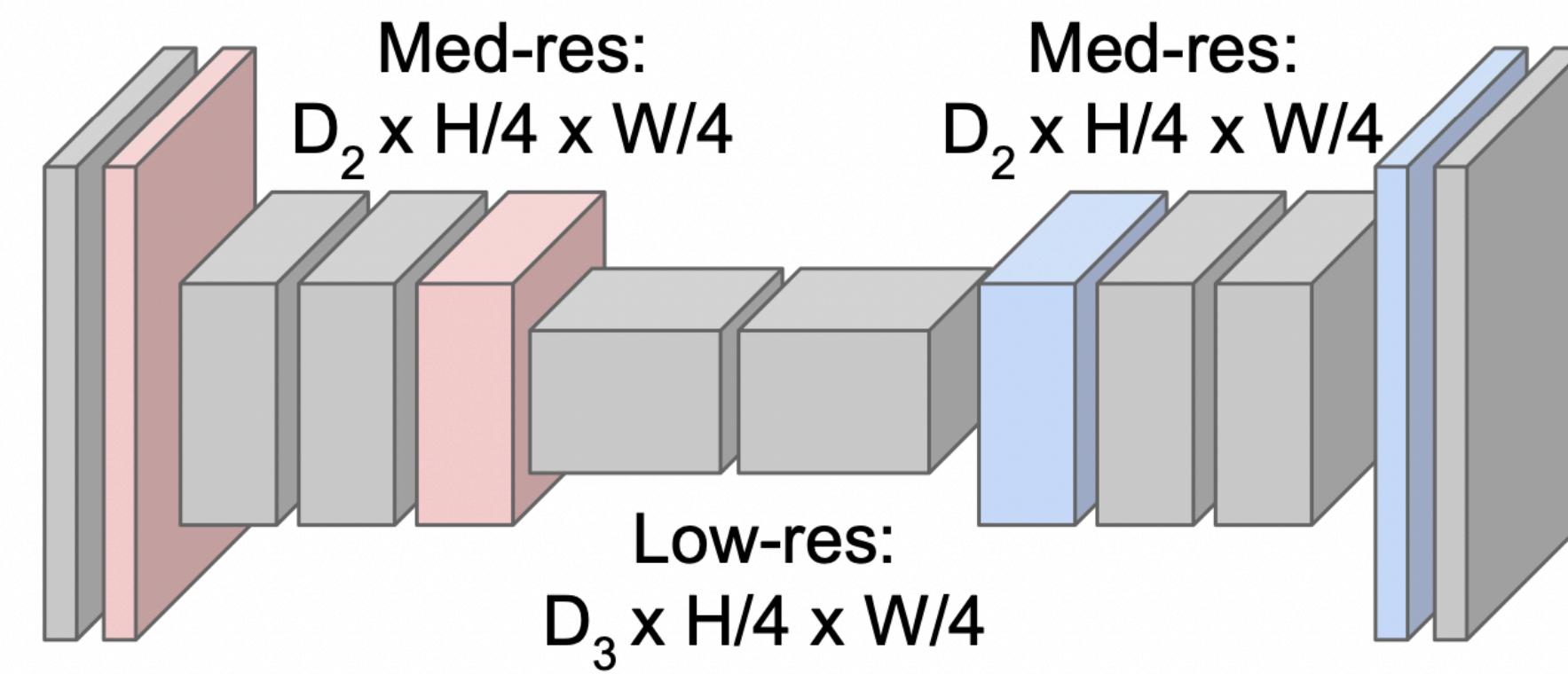
## Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution



Input:  
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with  
**downsampling** and **upsampling** inside the network!



**Upsampling:**  
???



Predictions:  
 $H \times W$

# Semantic Segmentation

## In-Network upsampling: “Unpooling”

**Nearest Neighbor**

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

**“Bed of Nails”**

1	2
3	4



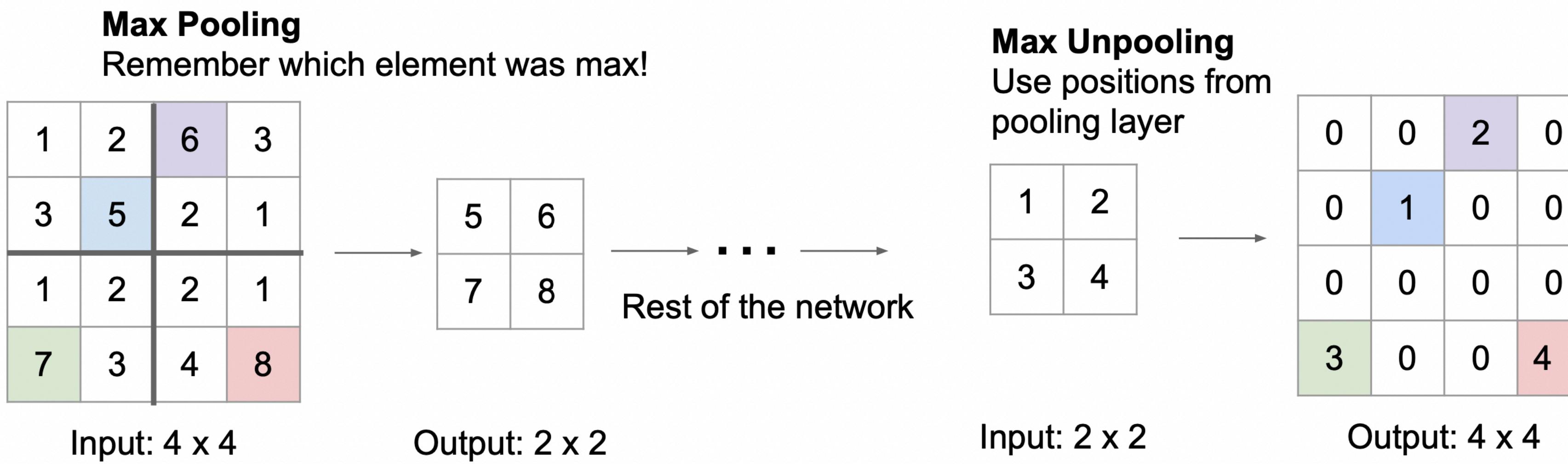
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

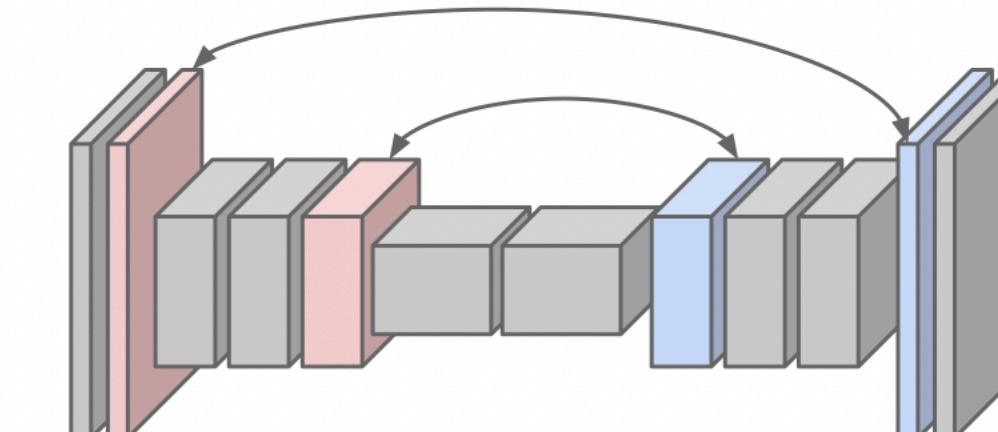
Output: 4 x 4

# Semantic Segmentation

## In-Network upsampling: “Max Unpooling”



Corresponding pairs of  
downsampling and  
upsampling layers

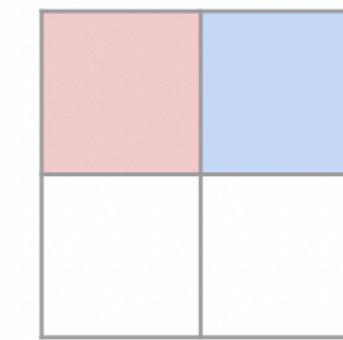


# Semantic Segmentation

## Learnable Upsampling: Transpose Convolution

**Other names:**

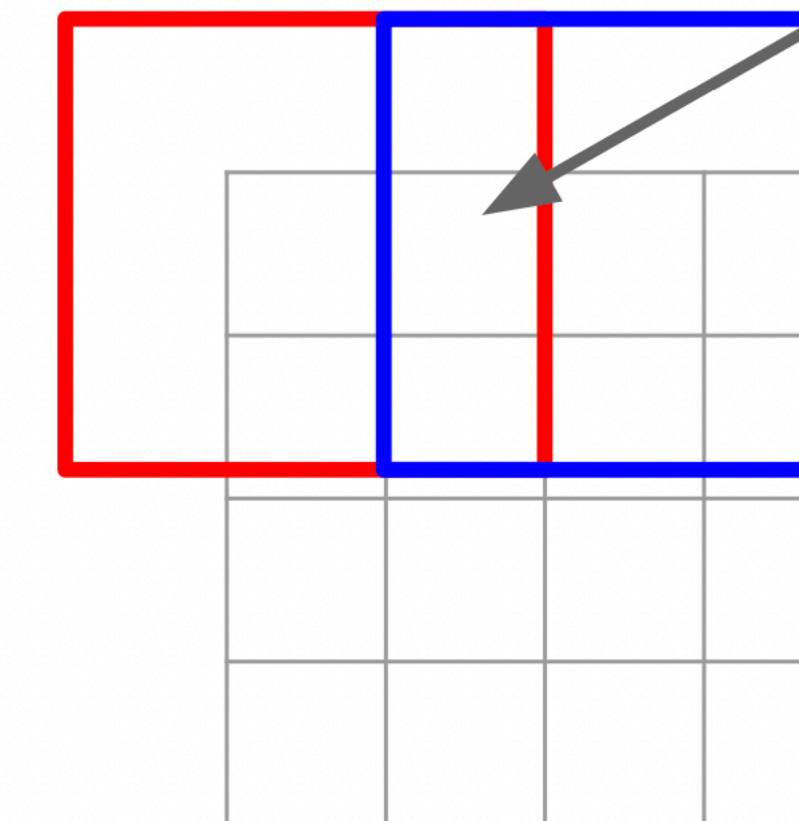
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution



Input: 2 x 2

**3 x 3 transpose convolution, stride 2 pad 1**

Input gives weight for filter



Output: 4 x 4

Sum where output overlaps

Filter moves 2 pixels in the output for every one pixel in the input

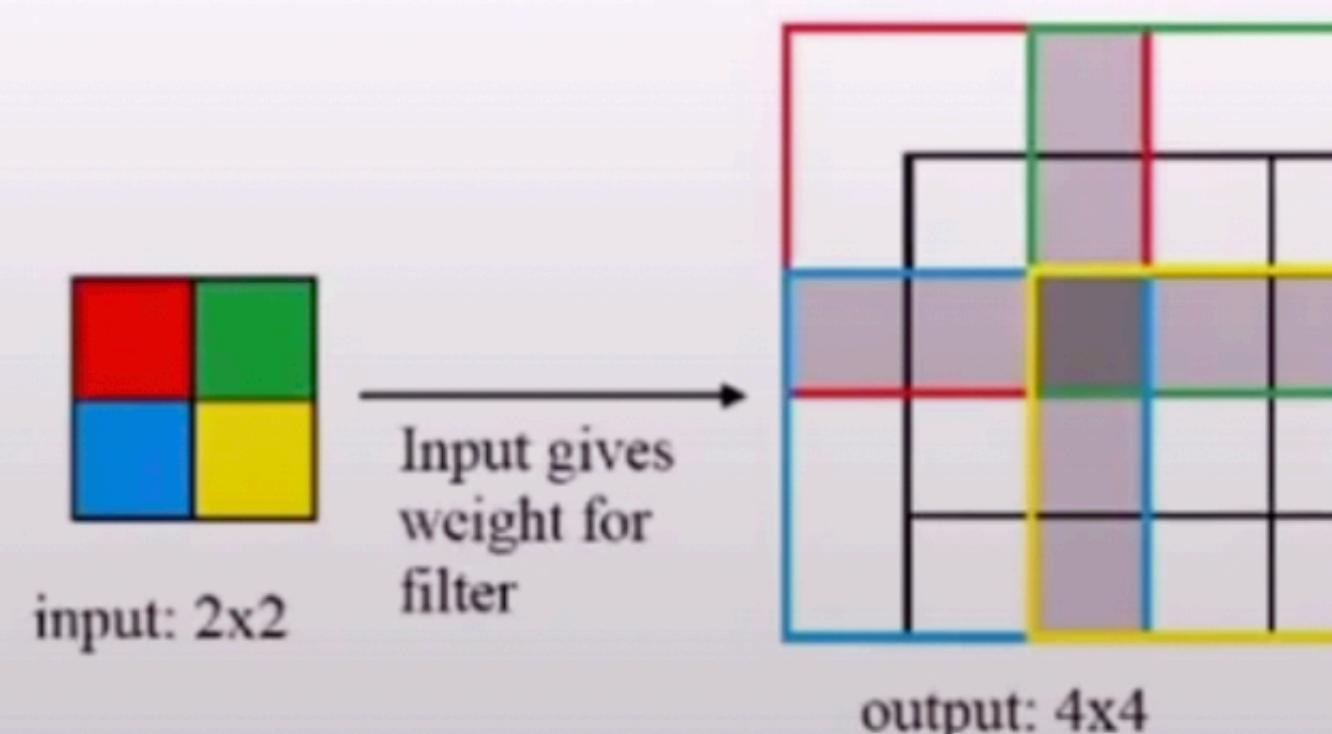
Stride gives ratio between movement in output and input

# Semantic Segmentation

## Upsampling Techniques: Checkerboard Artifacts

- Transpose convolution has uneven overlap when the kernel size is not divisible by the stride
- The uneven overlaps on the two axes multiply, creating a characteristic checkerboard artifact
- In principle, network could learn weights to avoid this but ...
- ... in practice, networks struggle to avoid it completely

3x3 transpose convolution, Stride 2



# Semantic Segmentation

## Checkerboard Artifacts

### How to avoid

- Choose appropriate kernel size: use a kernel size that is divisible by the stride, avoiding the overlap issue
- Separate upsampling from convolution to compute features
- For example:
  - Resize the image (e.g., using NN or bilinear interpolation)
  - Add a convolutional layer

# Semantic Segmentation

## Semantic Segmentation Idea: Fully Convolutional

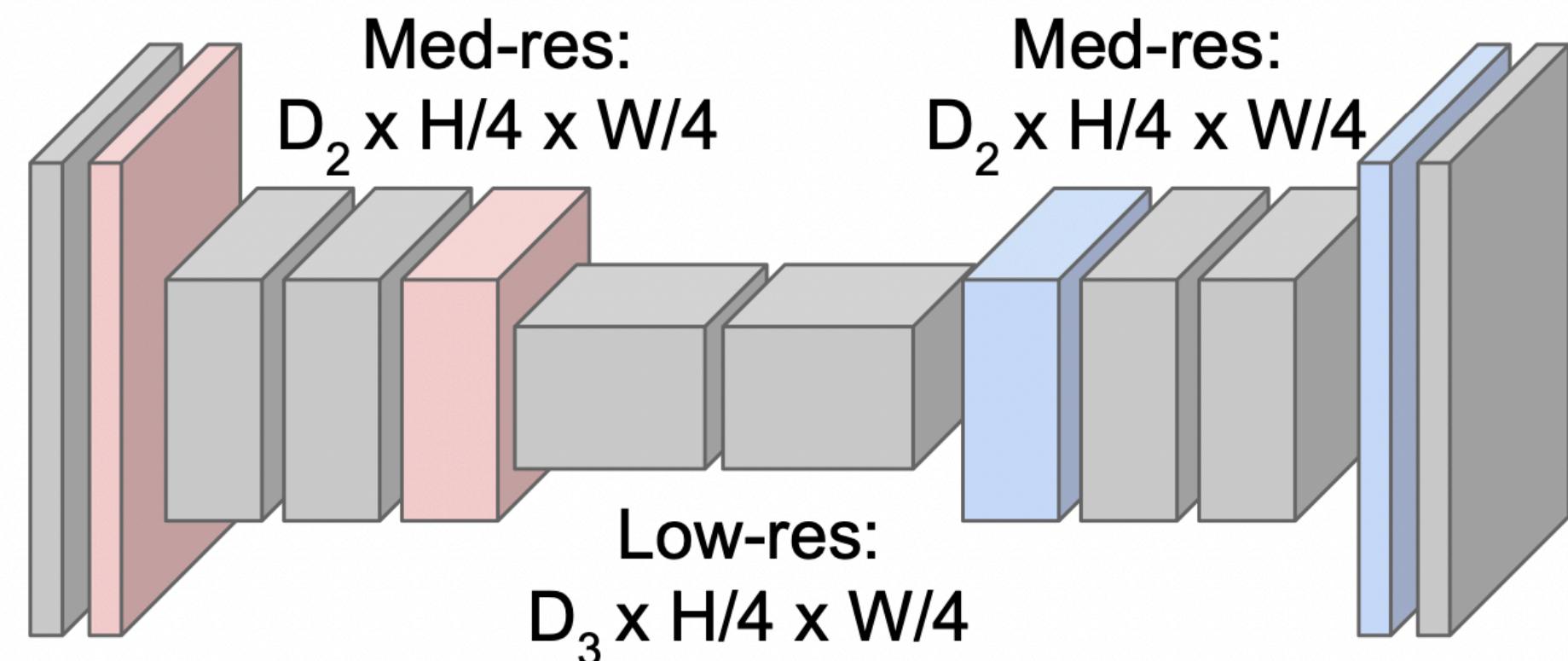
**Downsampling:**  
Pooling, strided convolution



Input:  
 $3 \times H \times W$

High-res:  
 $D_1 \times H/2 \times W/2$

Design network as a bunch of convolutional layers, with  
**downsampling** and **upsampling** inside the network!



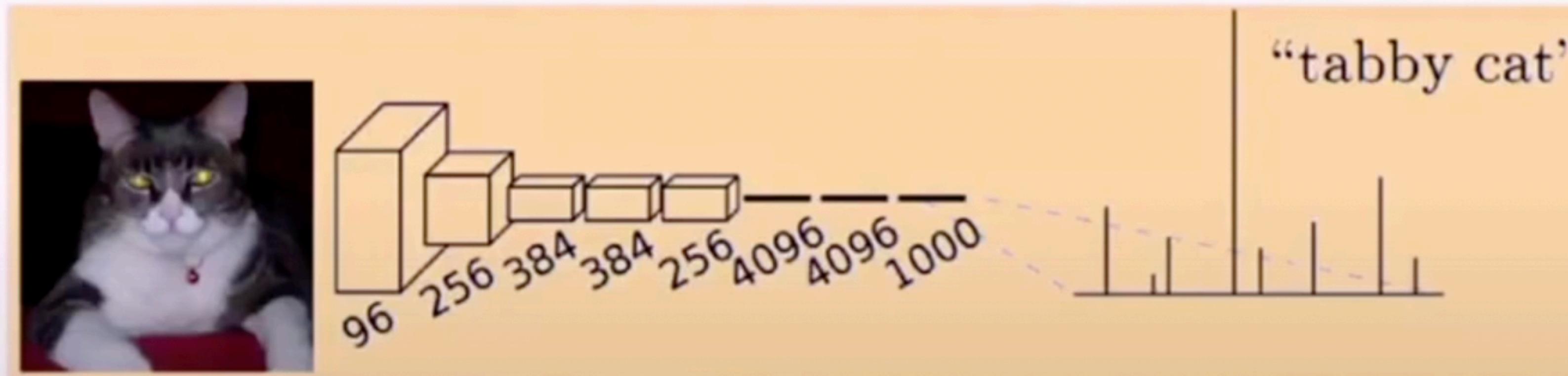
**Upsampling:**  
Unpooling or strided transpose convolution



Predictions:  
 $H \times W$

# Semantic Segmentation

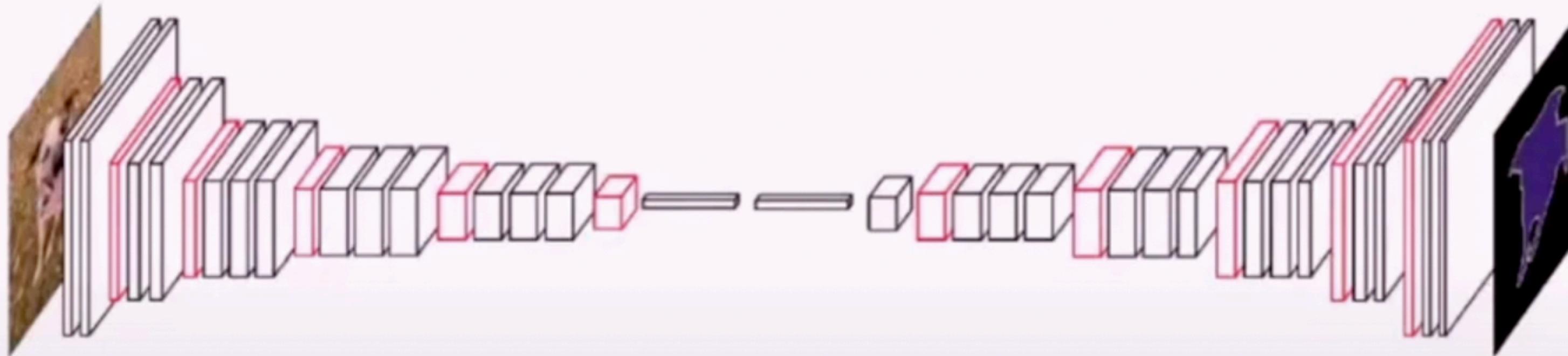
- Convolutional neural networks with **fully connected layers** used for classification:
  - CNN layers for feature extraction
  - Output dimensions reduced because of subsampling
  - **Fully connected layers** have fixed inputs and remove spatial information
  - Output: vector encoding class probabilities



# Semantic Segmentation

## Encoder and Decoder

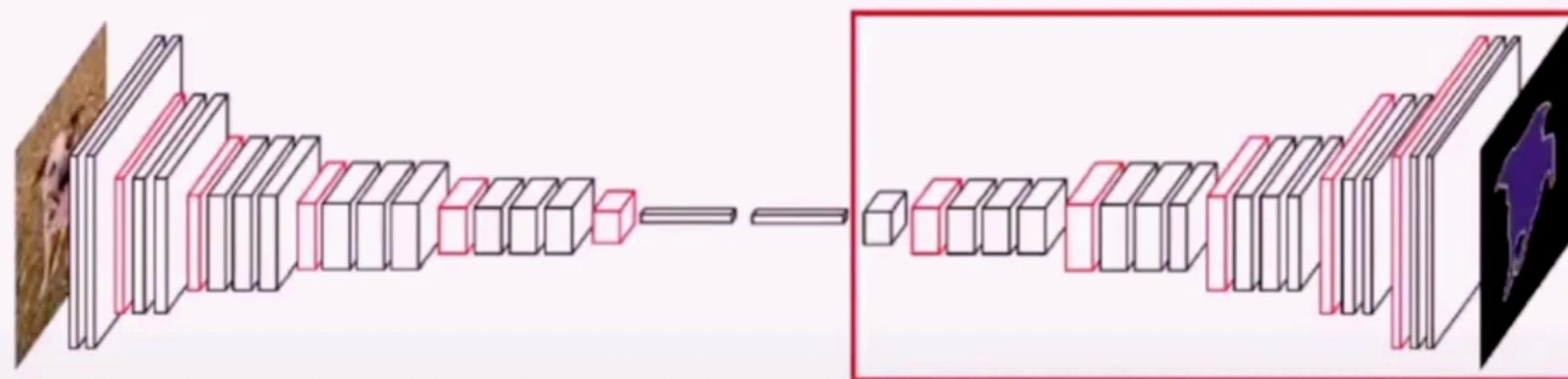
- Most methods use a classification network as basis
- If it can pickup the class, it should be able to learn the correct features



# Semantic Segmentation

## Encoder and Decoder (cont.)

- Network has to learn to **decode** (map) the output of the encoder to pixel-wise predictions
- **Decoder** part of the network



Example networks with different decoders:

- Long et al.'s Fully Convolutional Networks [13]
- SegNet [1]
- U-net [21]

# Semantic Segmentation

## Integrating Context Knowledge: Combining Where and What

- Semantic segmentation requires integration of information from various spatial scales
  - Need to balance local and global information
- Local information crucial to achieve good pixel-level accuracy
- Global context of the image enables to resolve local ambiguities
- CNNs struggle with this balance
- Different approaches to integrate local & global information

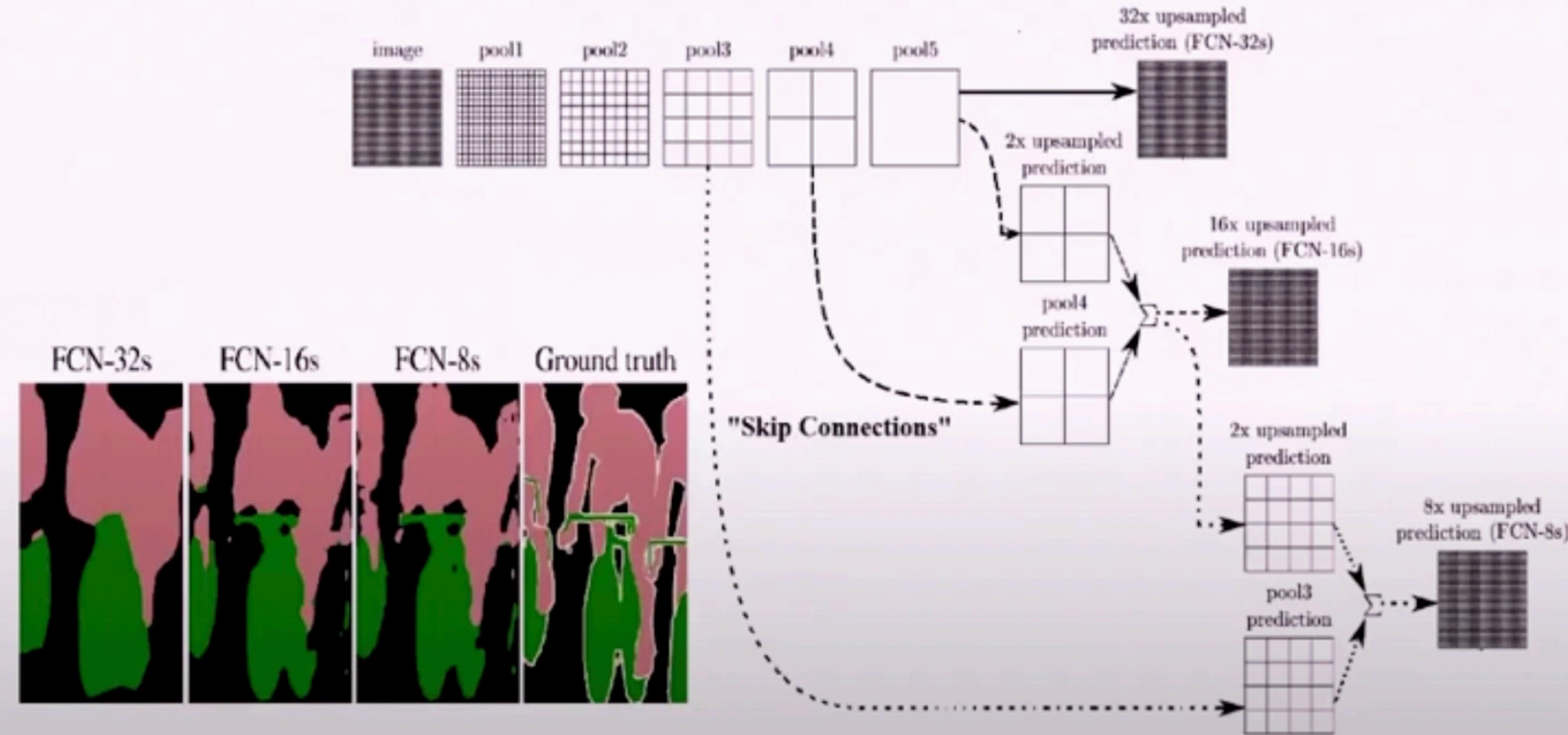
# Semantic Segmentation

## Long et al.'s Fully Convolutional Networks [13]

- Upsampling using learnable transposed convolutions
- Idea: Add links combining the final prediction and previous (lower) layers with finer strides
- Additional  $1 \times 1$  convolution after pooling layer, predictions are added up  
→ Make local predictions with global structure
- Network topology is a directed acyclic graph (DAG), with “**skip connections**” from lower to higher layers
- Refinement of coarse segmentation

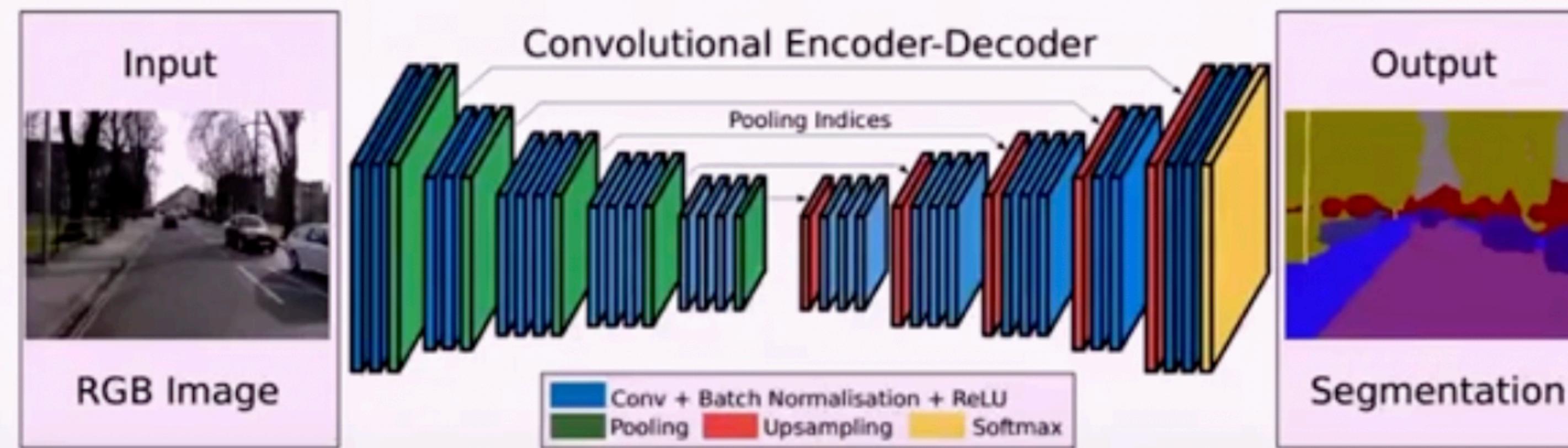
# Semantic Segmentation

## Long et al.'s Fully Convolutional Networks [13]



# Semantic Segmentation

## SegNet [1]:



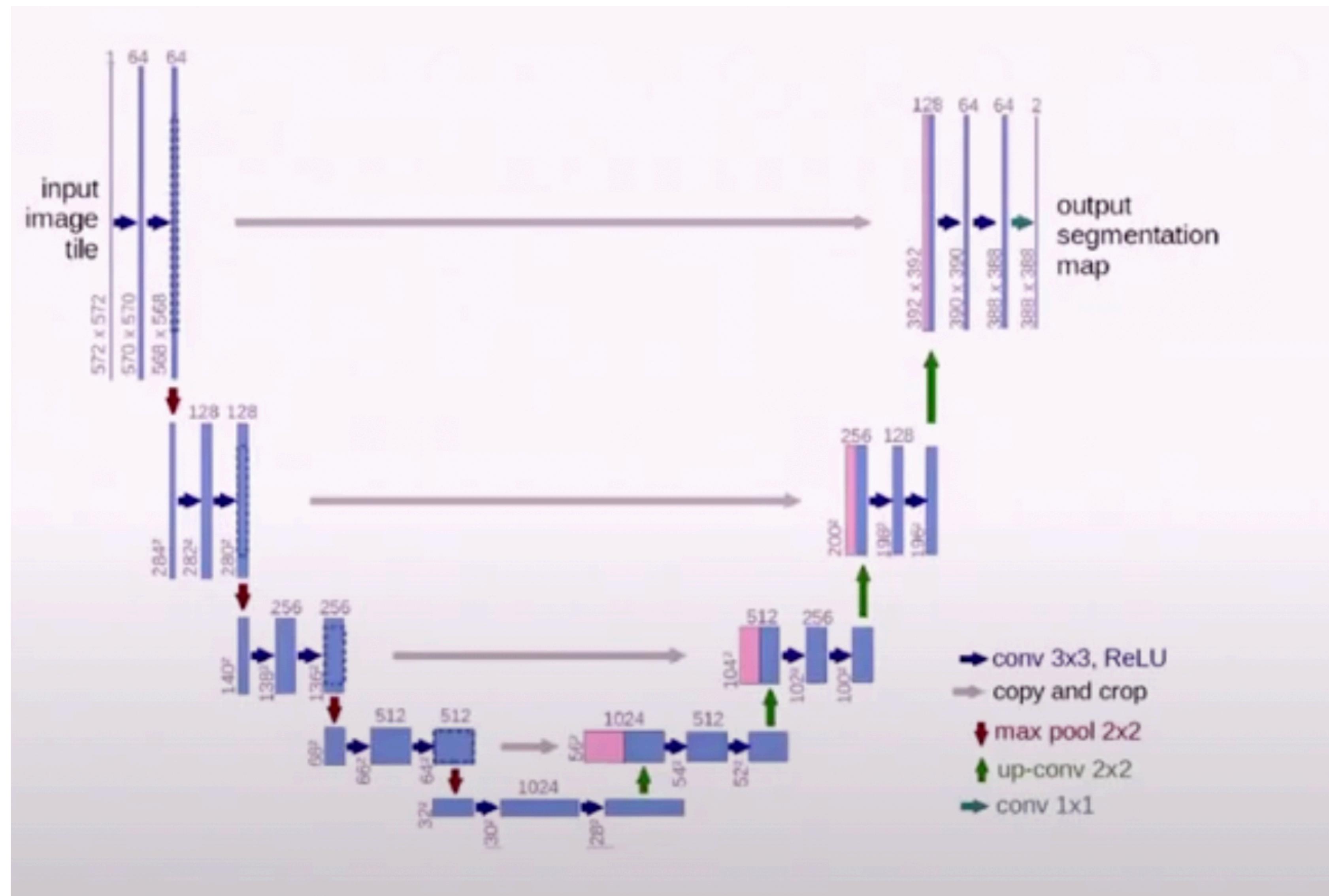
- Decoder: Upsampling & convolution layers followed by softmax
- Each upsampling layer corresponds to a **max-pooling layer** in encoder
- Upsampling reuses max-pooling indices from feature maps in encoder
- provides context information

# Semantic Segmentation

## U-Net [21]

- The network consists of:
  - Encoder: A contracting path to capture context
  - Decoder: A **symmetric** expansion path for localization (decoder)
- Encoder follows typical architecture of a CNN
- Decoder (expansion path) consists of
  - **Upsampling of feature maps** followed by a  $2 \times 2$  convolution that halves the number of feature channels.
  - **Concatenation** with the corresponding cropped feature map from the contracting path
- The training strategy relies on use of data augmentation
  - Non-rigid deformation
  - Rotation & translation

# Semantic Segmentation



# Semantic Segmentation

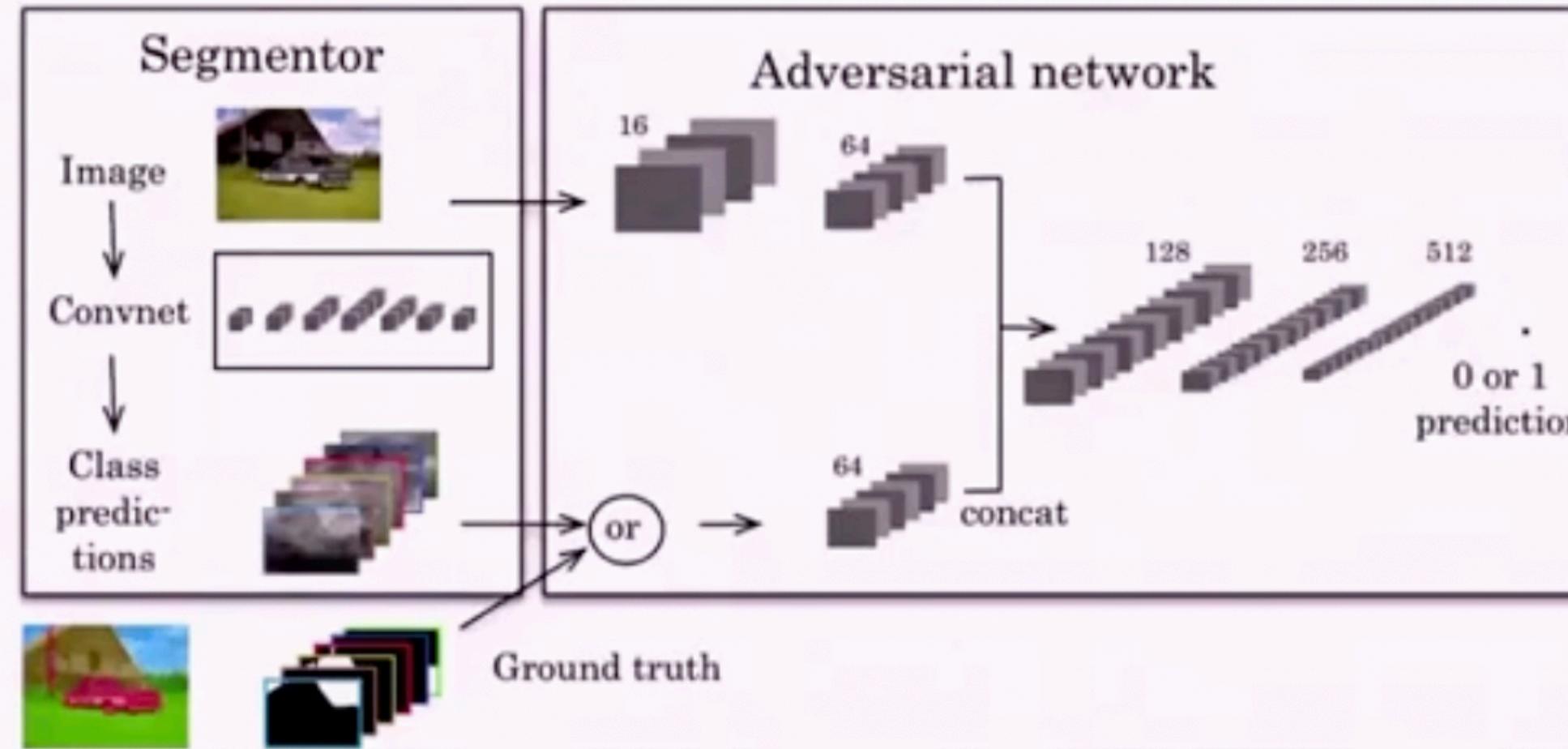
So far: Combination of feature maps from different levels + upsampling

Alternatives & extensions:

- Dilated convolutions
- Network stacks
- Multi-scale networks
- Defer context modeling to another network, e.g. an RNN
- Refinement as a post-processing step with Conditional Random Fields (CRFs)

# Semantic Segmentation

## Adversarial Networks for Segmentation [14]



Source: Luc et al., 2016

**Optimize Segmentor with hybrid loss function with two terms:**

- 1.) Usual pixel-wise multi-class cross-entropy for semantic segmentation
- 2.) Loss based on min-max game with Discriminator

# Semantic Segmentation

- Given a data set of  $N$  training images  $\mathbf{x}_n$  and a corresponding label maps  $\mathbf{y}_n$  the loss function defined as:

$$l(\theta_s, \theta_a) = \sum_{n=1}^N \underbrace{l_{mce}(s(\mathbf{x}_n), \mathbf{y}_n)}_{\text{multi-class cross-entropy}} - \lambda \underbrace{[l_{bce}(a(\mathbf{x}_n, \mathbf{y}_n), 1) + l_{bce}(a(\mathbf{x}_n, s(\mathbf{x}_n)), 0)]}_{\text{adversarial loss } \rightarrow \text{binary classification}}$$

- $(\mathbf{x}_n, \mathbf{y}_n), n \in \{1, \dots, N\}$ : Data set with ground truth labels
  - $\theta_s$ : Parameters of the Segmentor  $s(\mathbf{x}, \mathbf{y})$
  - $\theta_a$ : Parameters of the Discriminator  $a(\mathbf{x}, \mathbf{y})$
- Segmentor is trained both on the segmentation and on fooling the discriminator
- **Multi-task learning** with adversarial task