

CENG 506 Deep Learning

Lecture 8 - Segmentation

Slides were prepared using the course material of
Stanford's CNN Course (CS231n by Fei-Fei, Johnson, Yeung)

Previously we have seen

Classification



CAT

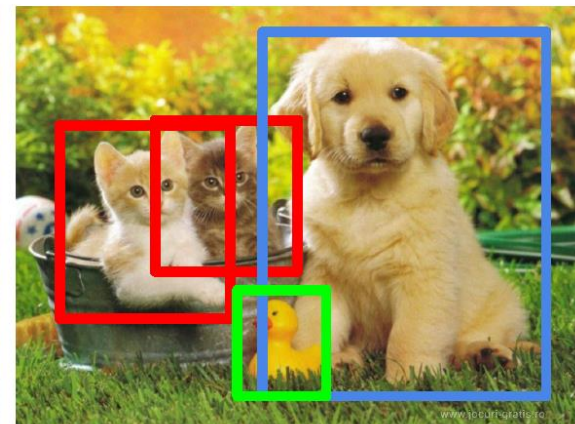
Classification + Localization



CAT

Single object

Object Detection



CAT, DOG, DUCK

Multiple objects

Today

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Instance Segmentation

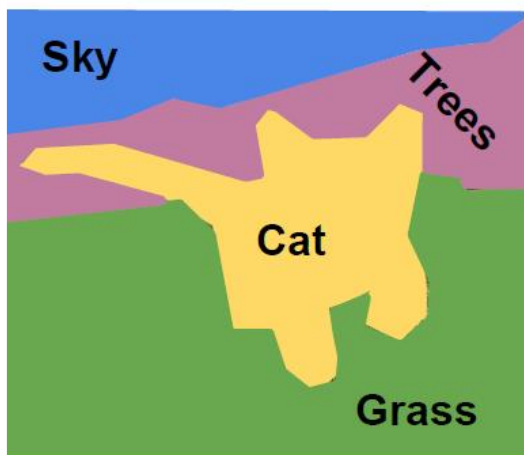


DOG, DOG, CAT

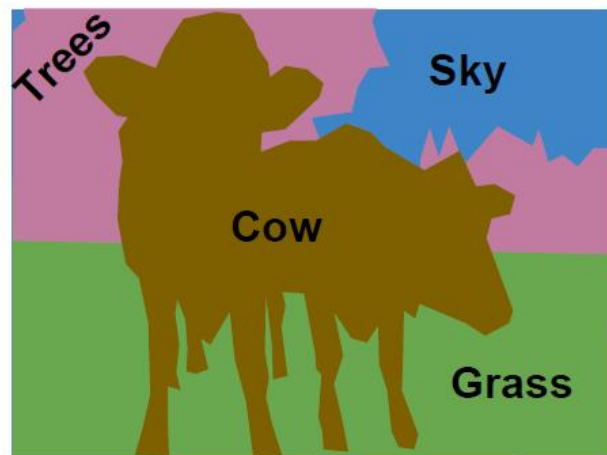
Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

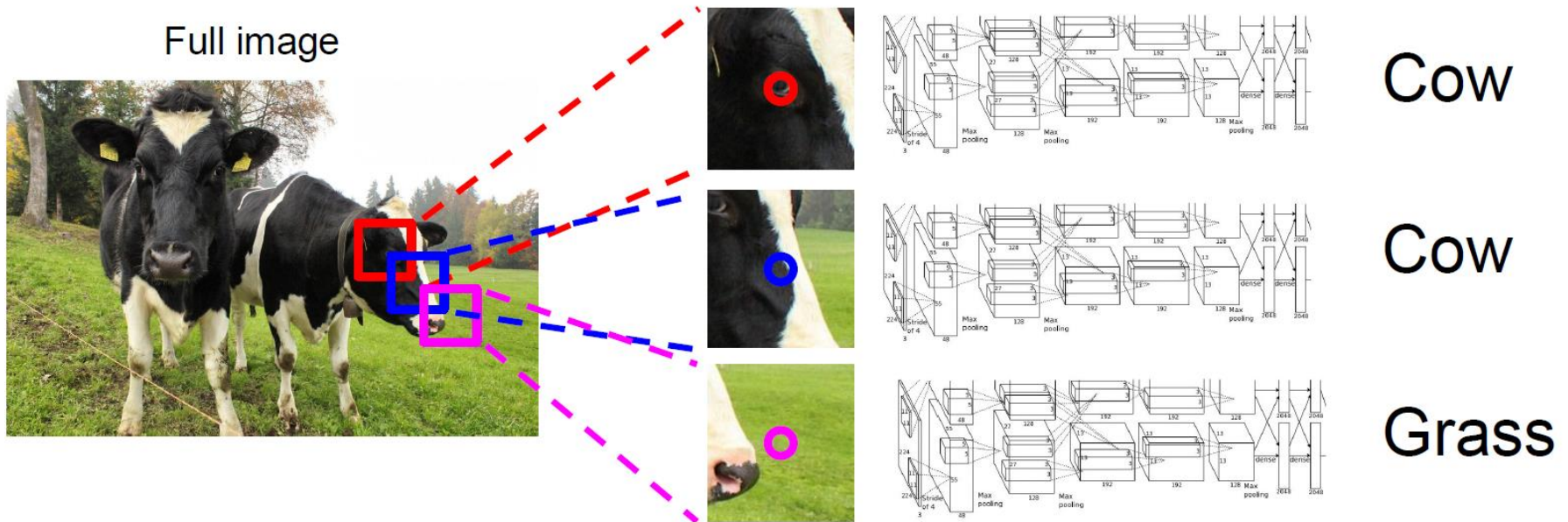


This image is CC0 public domain



Idea 1: Sliding window

Extract patch Classify center pixel with CNN



Problem: Very inefficient!

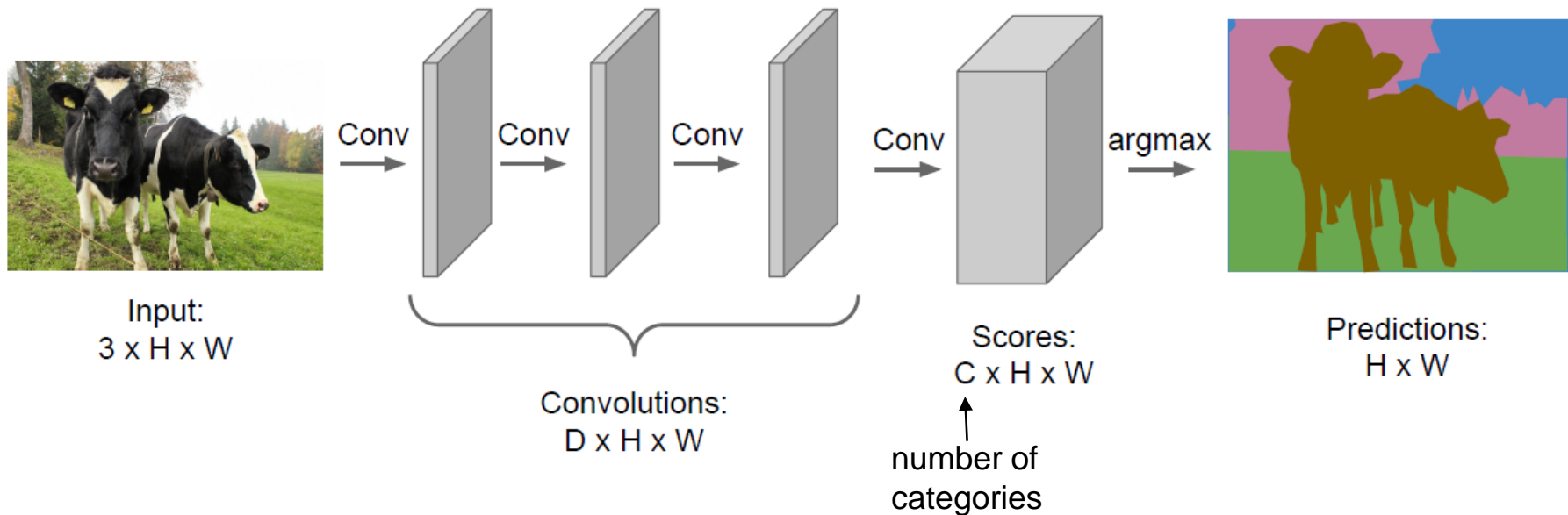
Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Idea 2: 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 ...

Idea 2: Fully convolutional

Design network as a bunch of convolutional layers, with

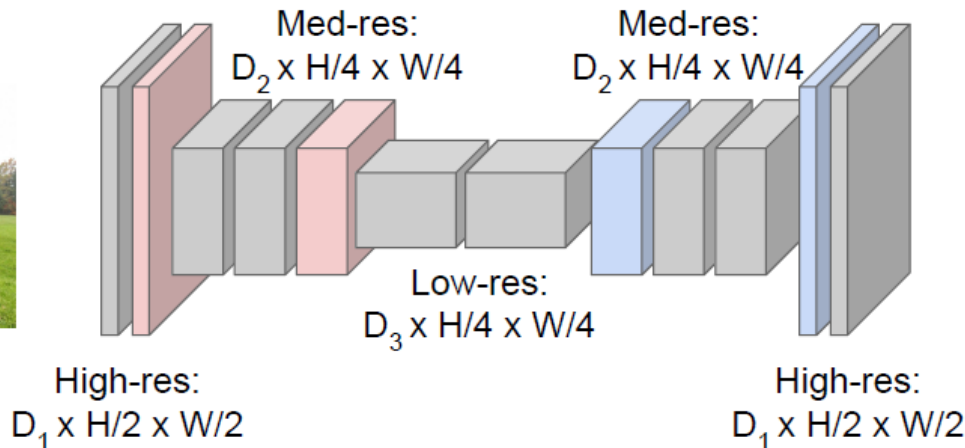
downsampling and **upsampling** inside the network!

Downsampling:
pooling,
convolution

Upsampling ?



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

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

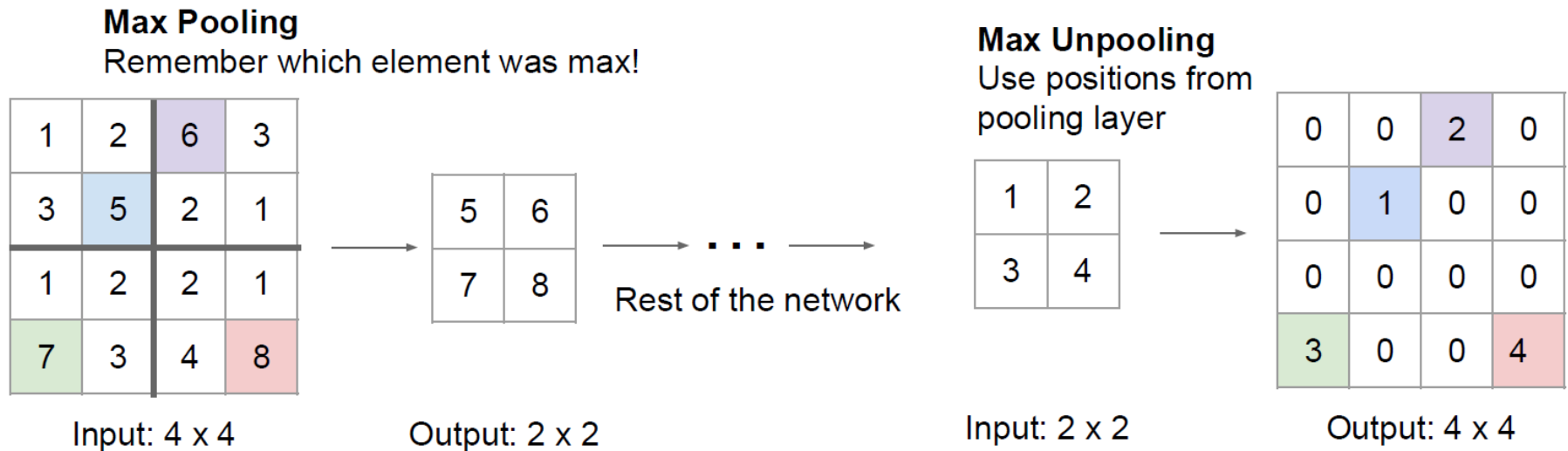


1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

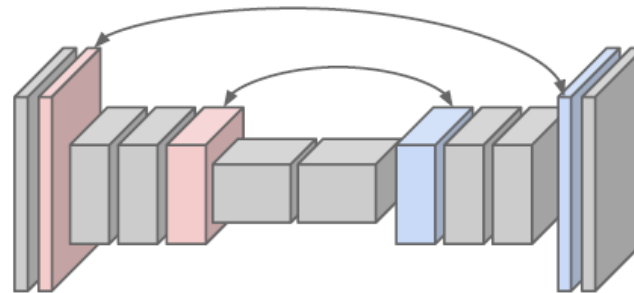
Input: 2 x 2

Output: 4 x 4

In-Network upsampling: “Max Unpooling”

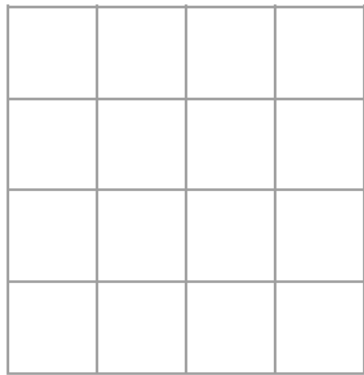


Corresponding pairs of
downsampling and
upsampling layers

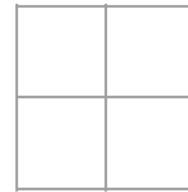


Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



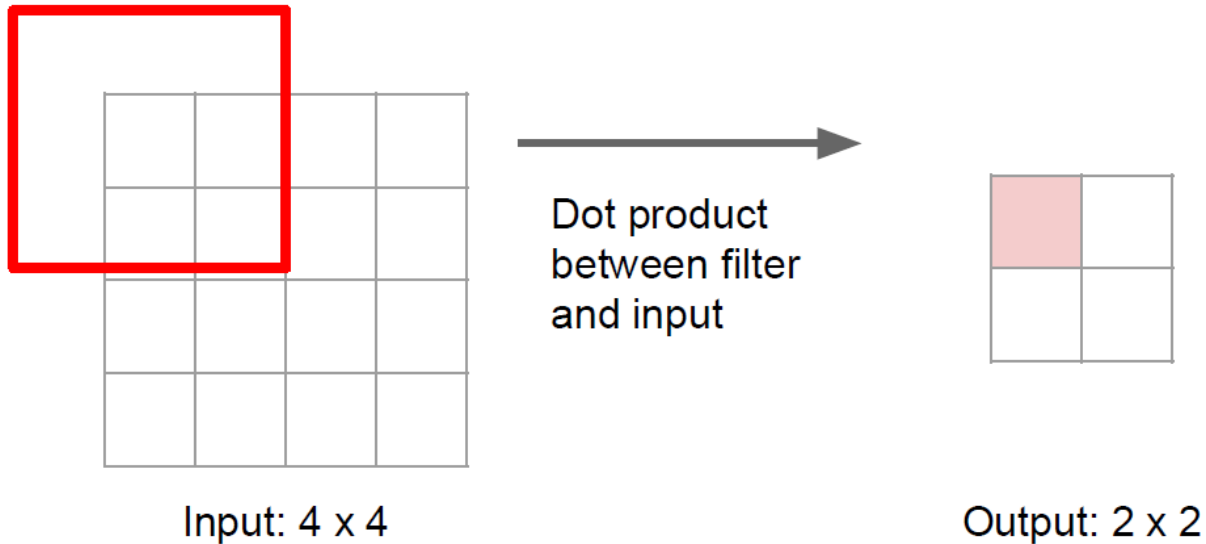
Input: 4 x 4



Output: 2 x 2

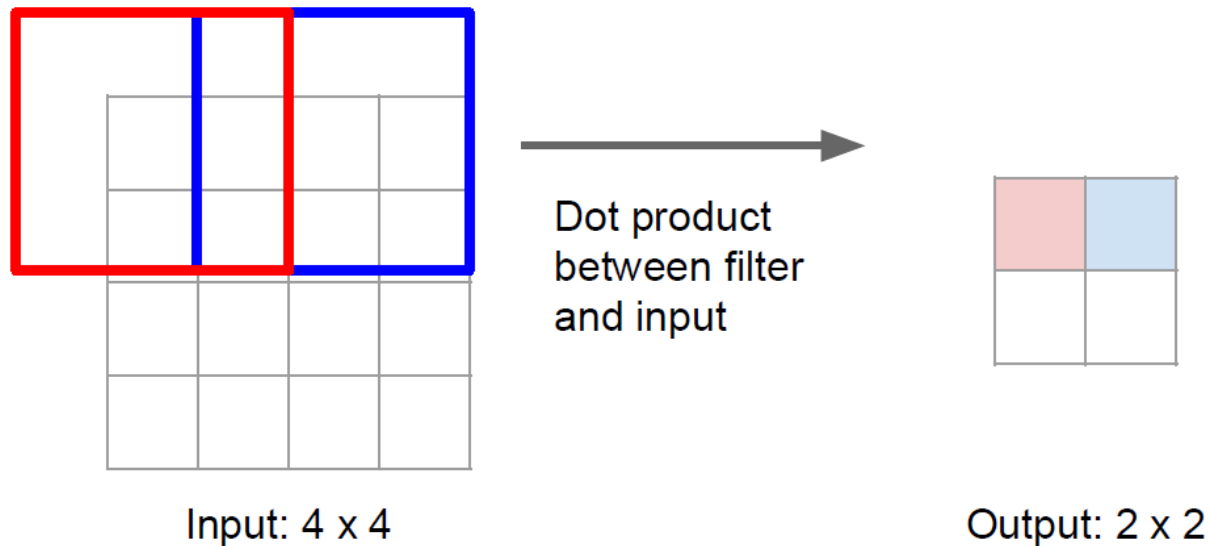
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Learnable Upsampling: Transpose Convolution

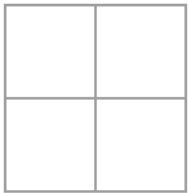
Recall: Normal 3 x 3 convolution, stride 2 pad 1



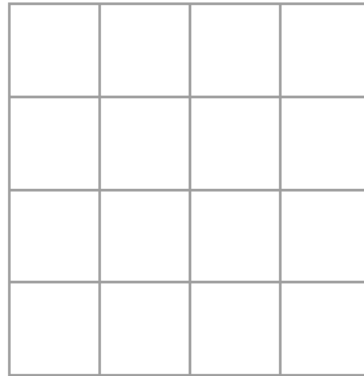
Filter moves 2 pixels in input for every one pixel in output
Stride gives ratio between movement in input and output

Learnable Upsampling: Transpose Convolution

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



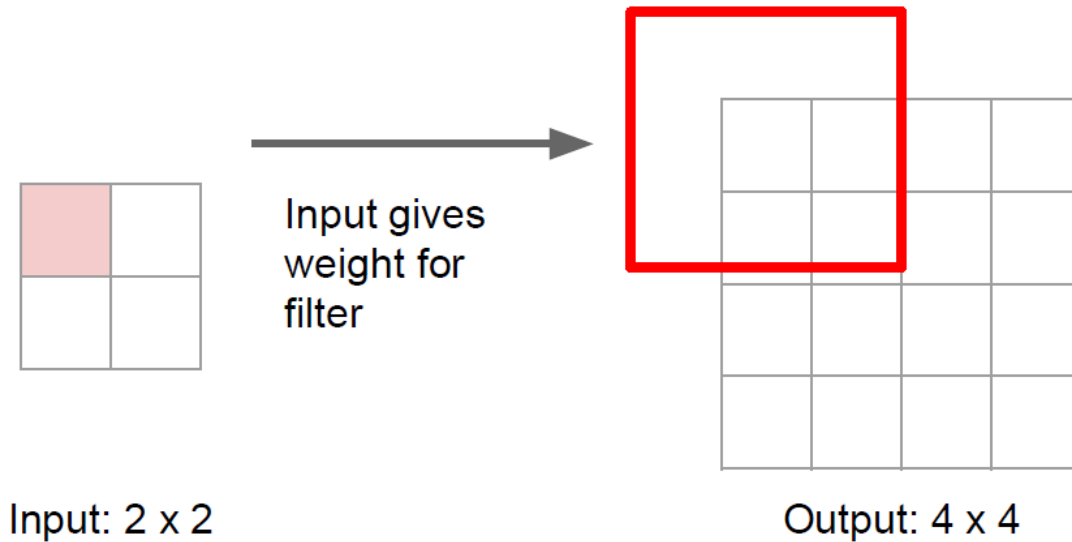
Input: 2 x 2



Output: 4 x 4

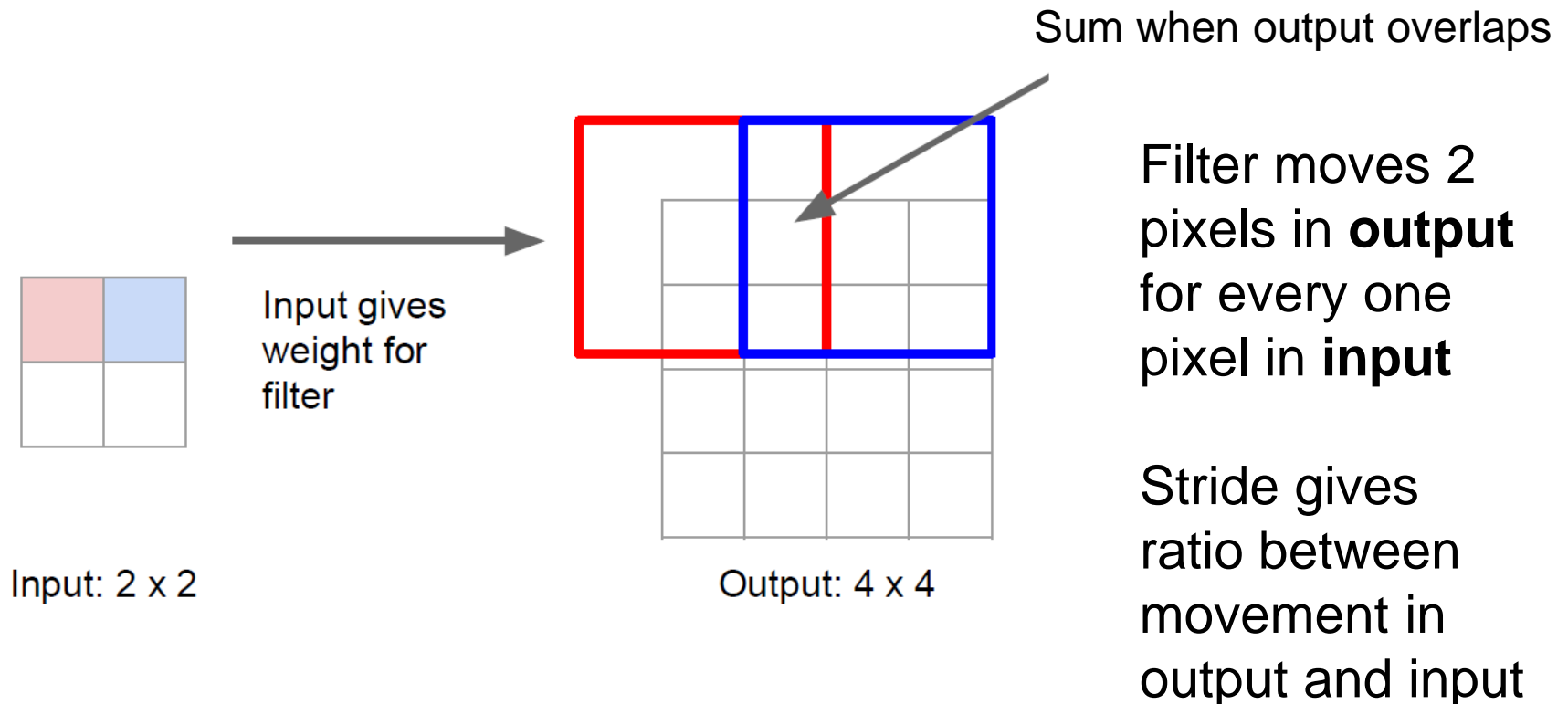
Learnable Upsampling: Transpose Convolution

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

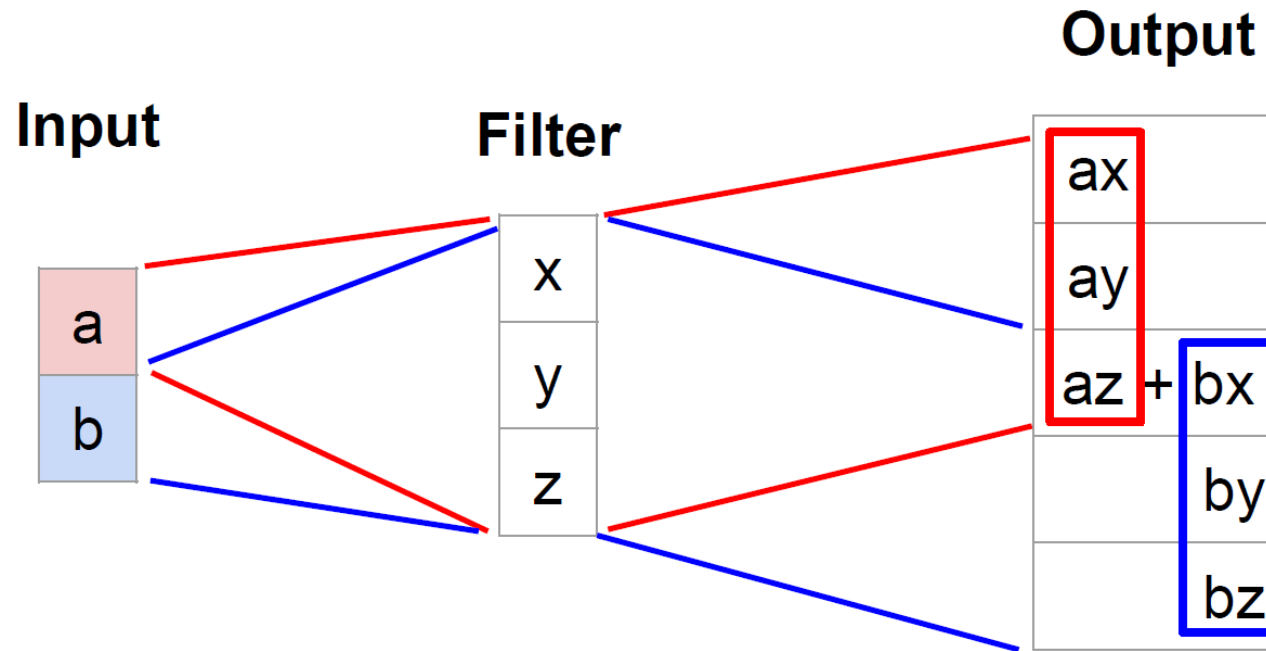


Learnable Upsampling: Transpose Convolution

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



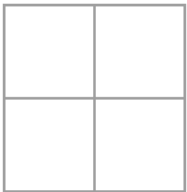
Transpose Convolution: 1D Example



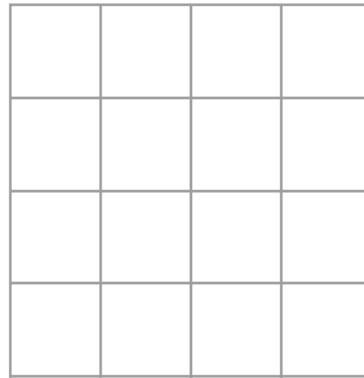
Output contains copies of the filter weighted by the input, summing where they overlap in the output

Need to crop one pixel from output to make output exactly 2x input

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



Input: 2 x 2



Output: 4 x 4

Learnable Upsampling

Other Names for Transpose Convolution:

- Deconvolution
- Upconvolution

Idea 2: Fully convolutional

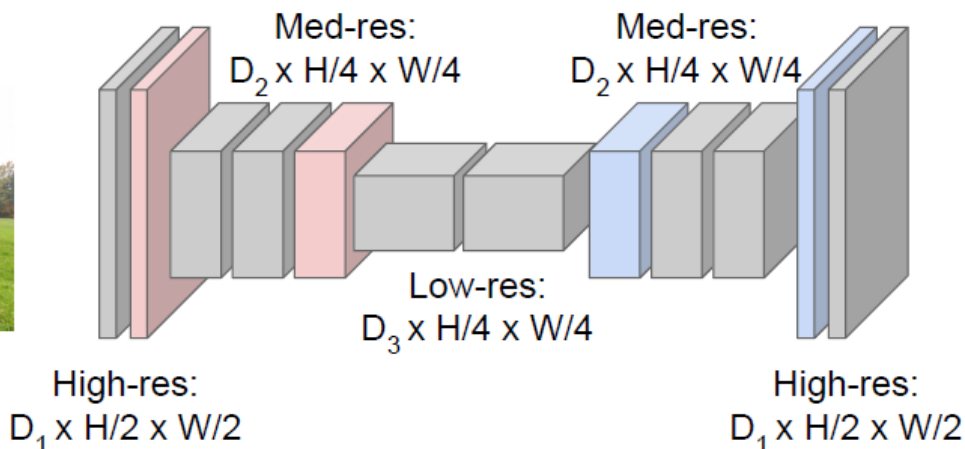
Downsampling:
pooling,
convolution

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

Upsampling:
Unpooling and
transpose
convolution
(deconvolution)

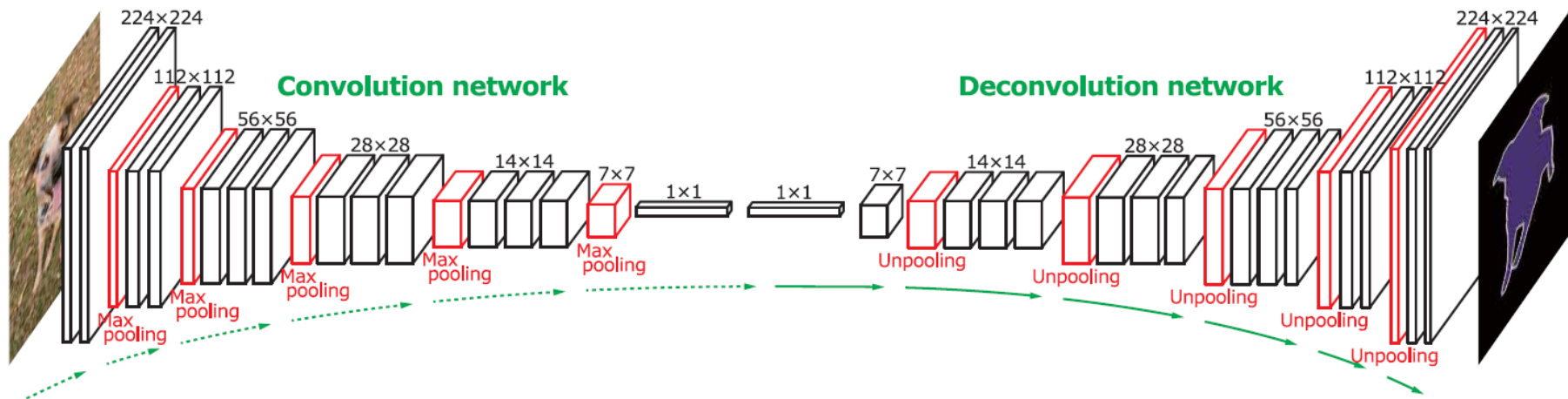


Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Idea 2 Example: Noh et al, 2015



- Convolution network: VGG 16-layer net.
- Given a feature vector obtained from the convolution network, a deconvolution network produces pixel-wise segmentation map.
- Deconvolution layers densify the sparse activations by unpooling and multiple learned deconvolution filters.

Idea 2 Example: Noh et al, 2015

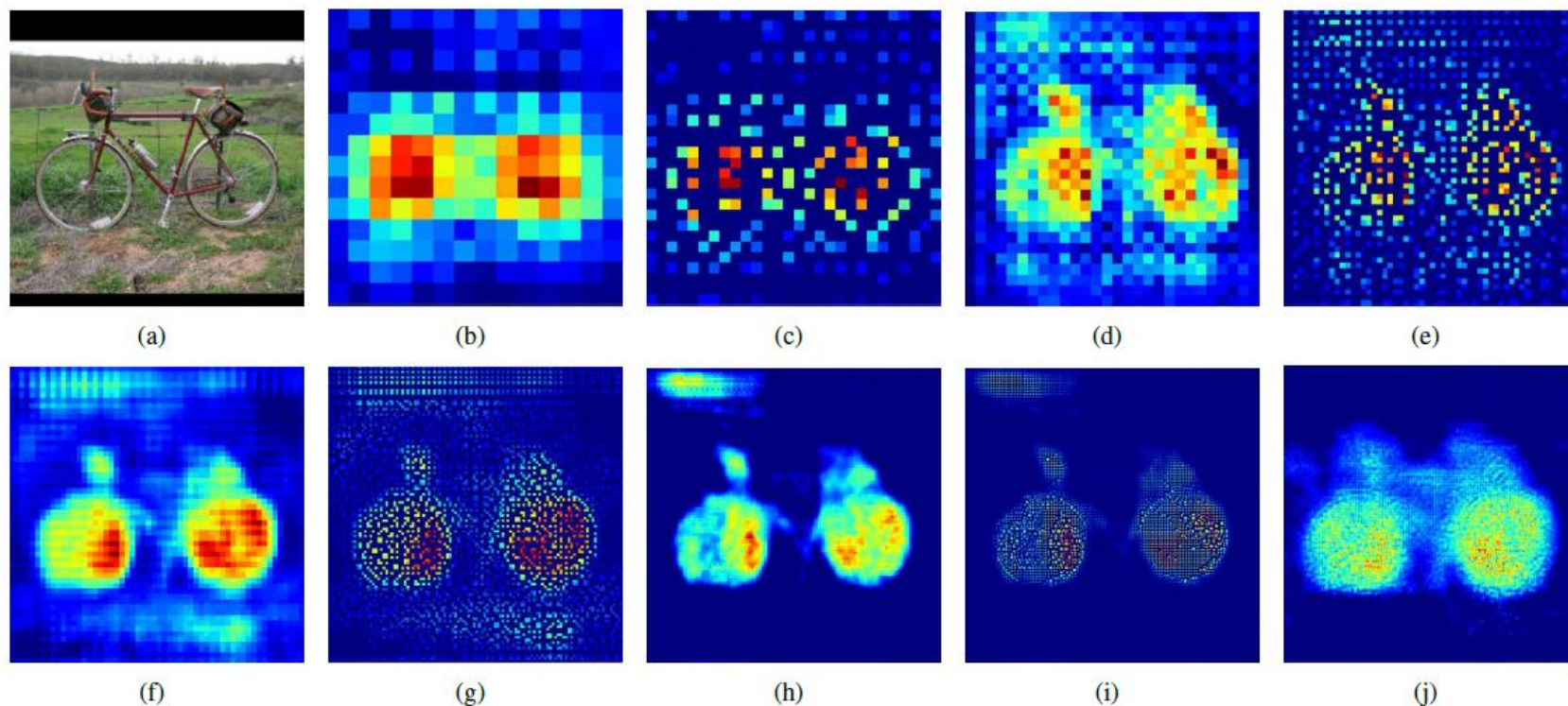
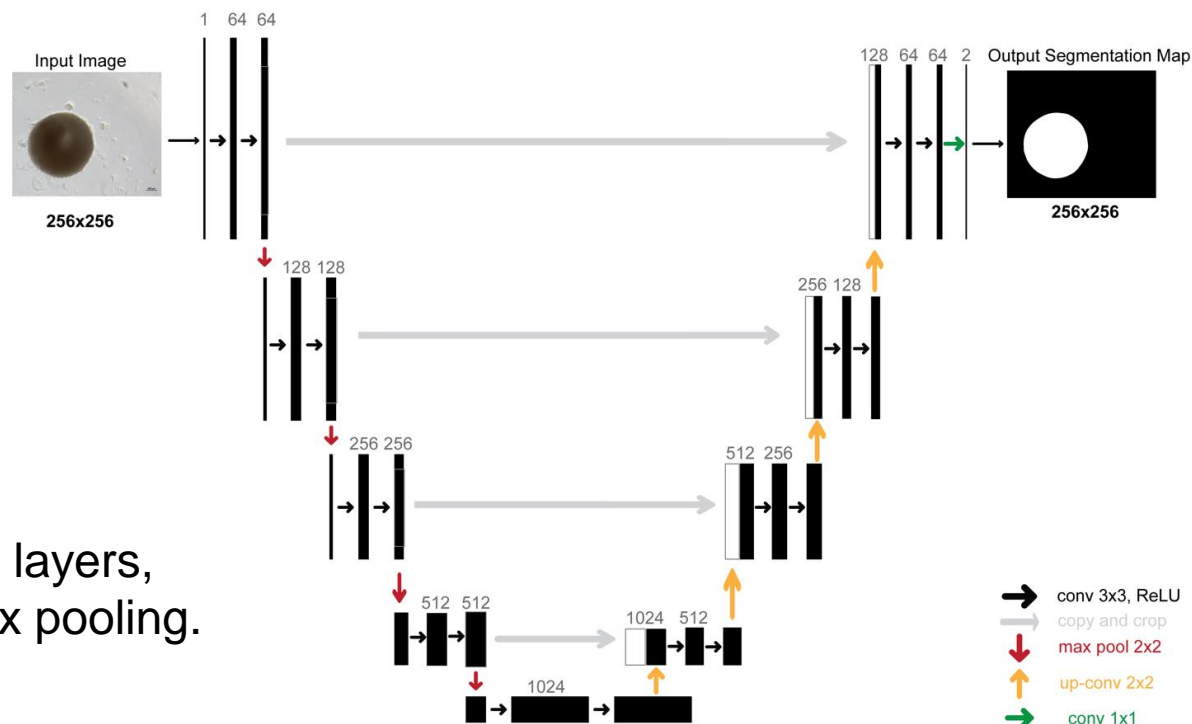


Figure 4. Visualization of activations in our deconvolution network. The activation maps from (b) to (j) correspond to the output maps from lower to higher layers in the deconvolution network. We select the most representative activation in each layer for effective visualization. The image in (a) is an input, and the rest are the outputs from (b) the last 14×14 deconvolutional layer, (c) the 28×28 unpooling layer, (d) the last 28×28 deconvolutional layer, (e) the 56×56 unpooling layer, (f) the last 56×56 deconvolutional layer, (g) the 112×112 unpooling layer, (h) the last 112×112 deconvolutional layer, (i) the 224×224 unpooling layer and (j) the last 224×224 deconvolutional layer. The finer details of the object are revealed, as the features are forward-propagated through the layers in the deconvolution network. Note that noisy activations from background are suppressed through propagation while the activations closely related to the target classes are amplified. It shows that the learned filters in higher deconvolutional layers tend to capture class-specific shape information.

Another Idea 2 Example: U-Net



- Repeated double 3x3 conv layers, each followed by a 2x2 max pooling.
- After each downsampling, the no of feature channels are doubled.
- In the second half, each up-sampling step is followed by a 2x2 up-convolution and a concatenation with the corresponding level of the first half.
- Final layer uses a convolution to reduce the feature map depth to the desired number of classes.

Instance Segmentation

Instance Segmentation

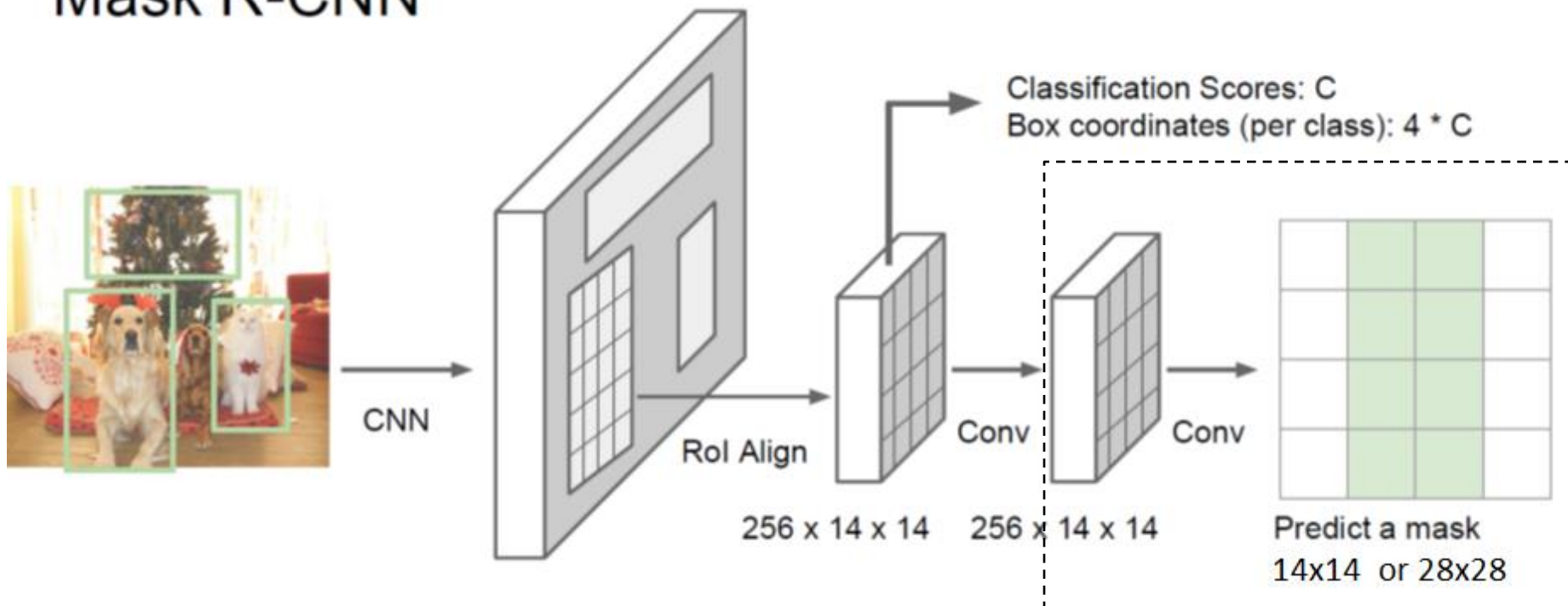


DOG, DOG, CAT

Instance Segmentation

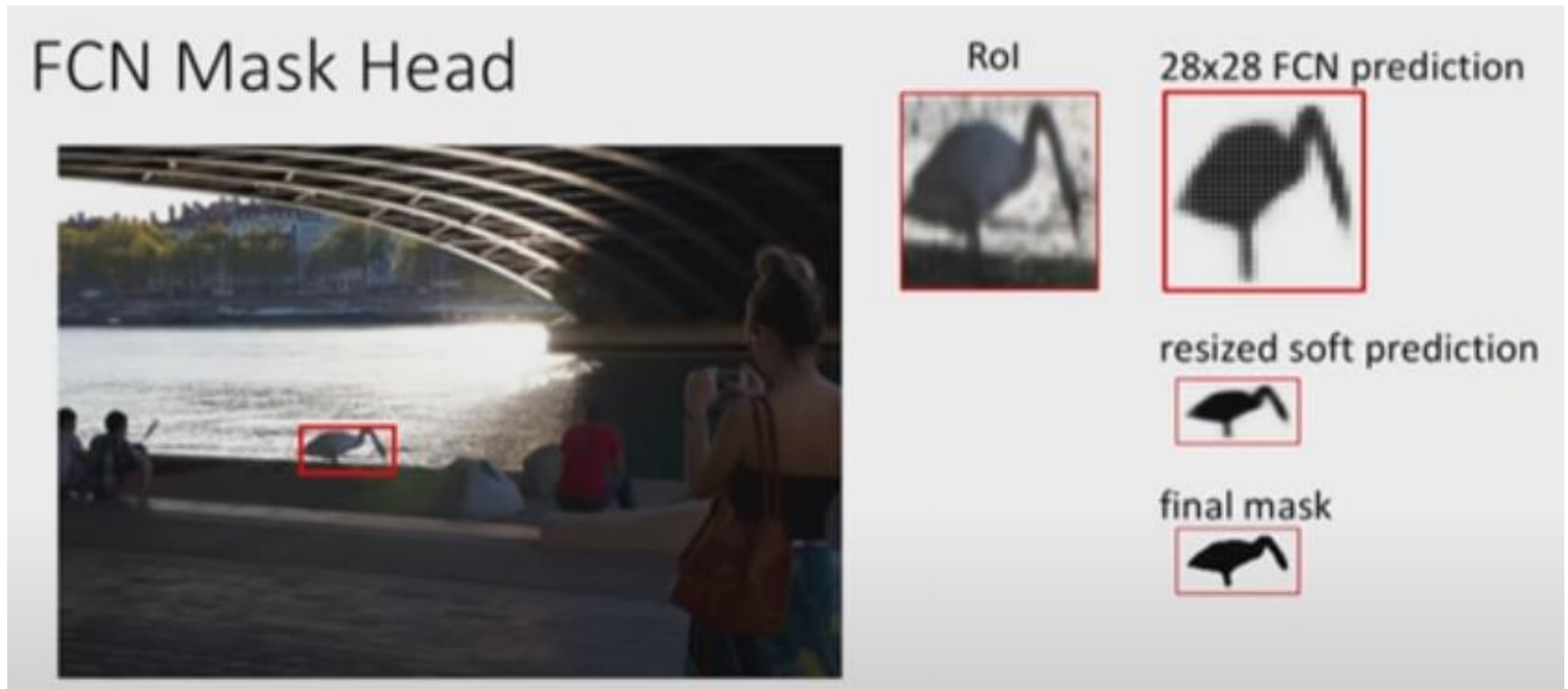
(Object detection + Semantic segmentation)

Mask R-CNN



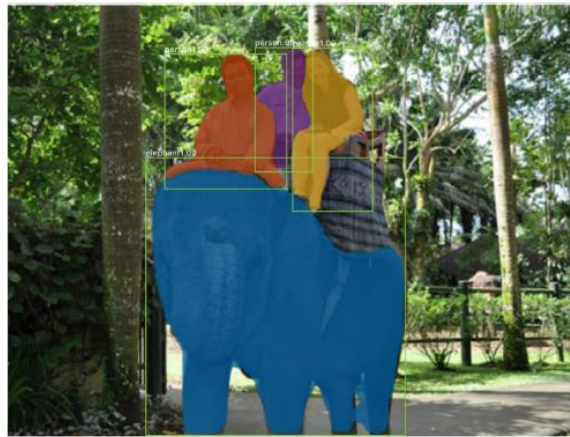
We add a third 'Mask' head to Faster R-CNN.

Mask R-CNN



‘Mask head’ is actually a fully convolutional network with some downsampling and upsampling.

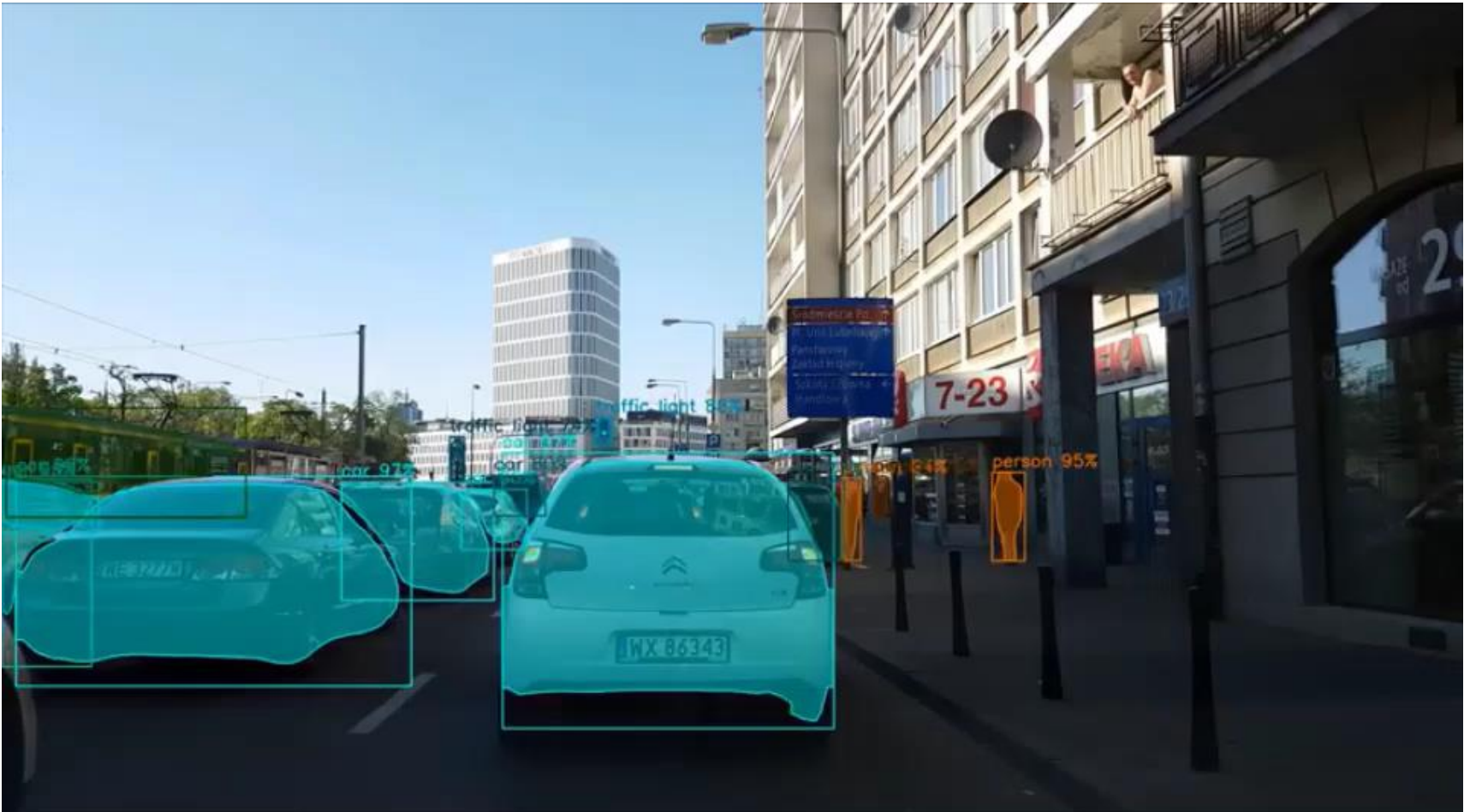
Mask R-CNN: Very Good Results!



He et al, "Mask R-CNN", ICCV 2017.

Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.

Mask R-CNN: Very Good Results!



See the full video at: <https://www.youtube.com/watch?v=OOT3UIXZztE>

Panoptic Segmentation

Semantic segmentation covers all pixels of the image. Instance segmentation concentrates on instances.

We can combine the results of separate semantic segmentation and instance segmentation tasks via a network.



Panoptic Segmentation: EfficientPS

EfficientPS consists of a shared backbone, a two-way Feature Pyramid Network (FPN), instance and semantic heads, and a panoptic fusion module.

