

# Segmentation

---

Deep Learning

Aziz Temirkhanov  
Lambda, HSE

Fast recap

---

# Statistical learning theory

$\{(x_i, y_i)\}_{i=1}^{\ell}$  - training set       $x \in \mathbb{R}^d$  - data features  
 $y \in \mathbb{R}$  - target

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \mathcal{L}(\theta; x, y) \quad \mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(g_{\theta}(x_i), y_i)$$

# STL

1. Problem is formulated in business language (e.g. “Highlight a tumor within MRI scan”)
2. Reformulate it formally in terms of math and STL (“given an image  $x$ , predict segmentation mask  $y$ ”)
3. Recall a universal approximation theorem
4. Remember that you have a prior knowledge about your data distribution (besides the normality assumption)
5. Fit your function  $f$  (your NN) to approximate target function  $g$  (your PDF or other law)
6. During fitting, you have to minimize the loss function (ERM). Once it minimized, and other formalities are satisfied, you are done!

# Image segmentation

---

# Computer Vision Tasks

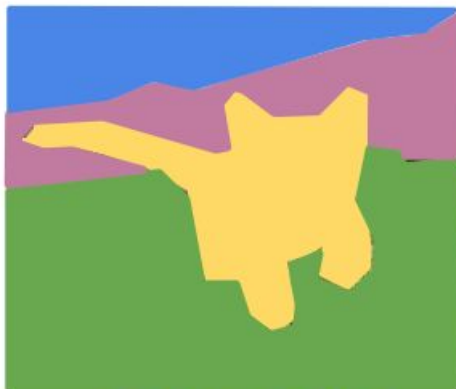
## Classification



**CAT**

No spatial extent

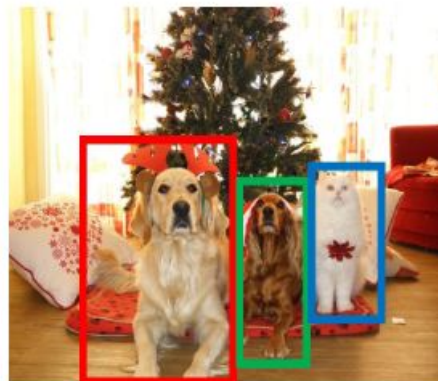
## Semantic Segmentation



**GRASS, CAT, TREE, SKY**

No objects, just pixels

## Object Detection



**DOG, DOG, CAT**

## Instance Segmentation



**DOG, DOG, CAT**

Multiple Object

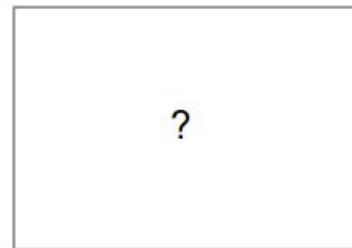
[This image is CC0 public domain](#)

# Formulation



GRASS, CAT,  
TREE, SKY, ...

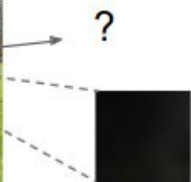
Paired training data: for each training image,  
each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

# Problem

Full image

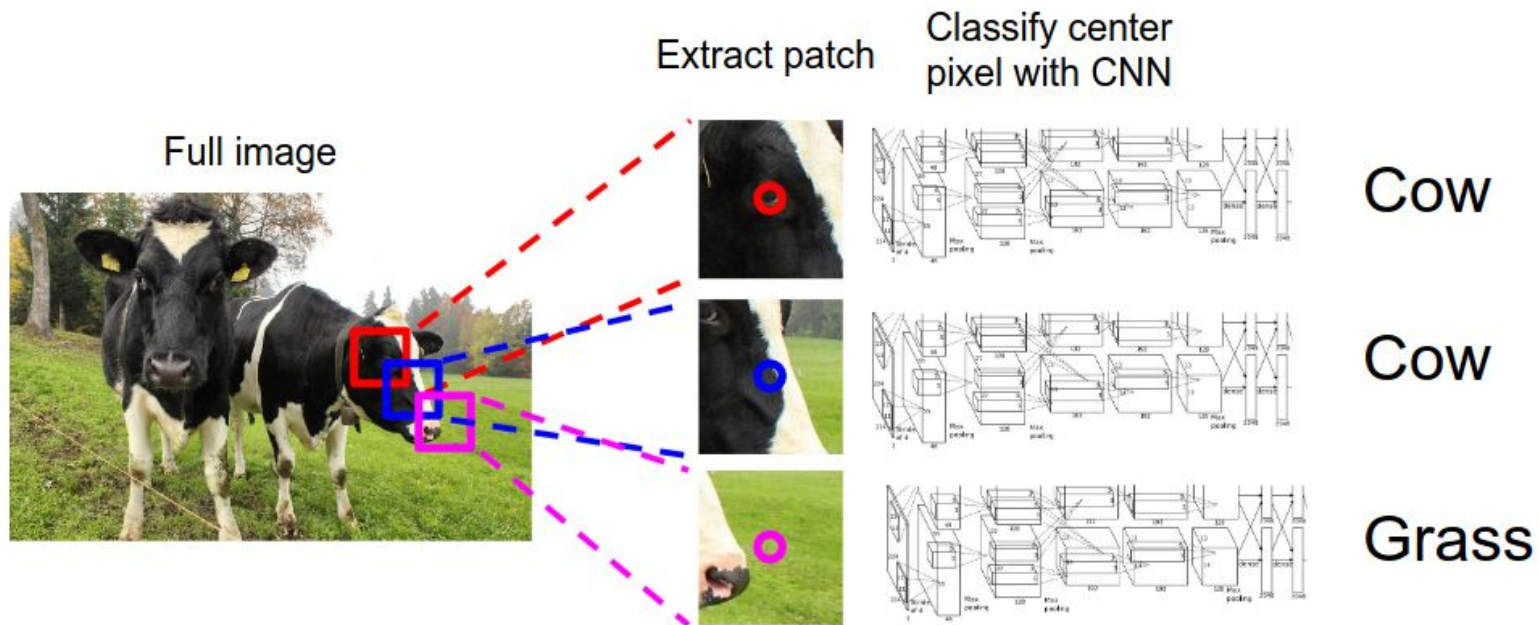


Impossible to classify without context

Q: how do we include context?



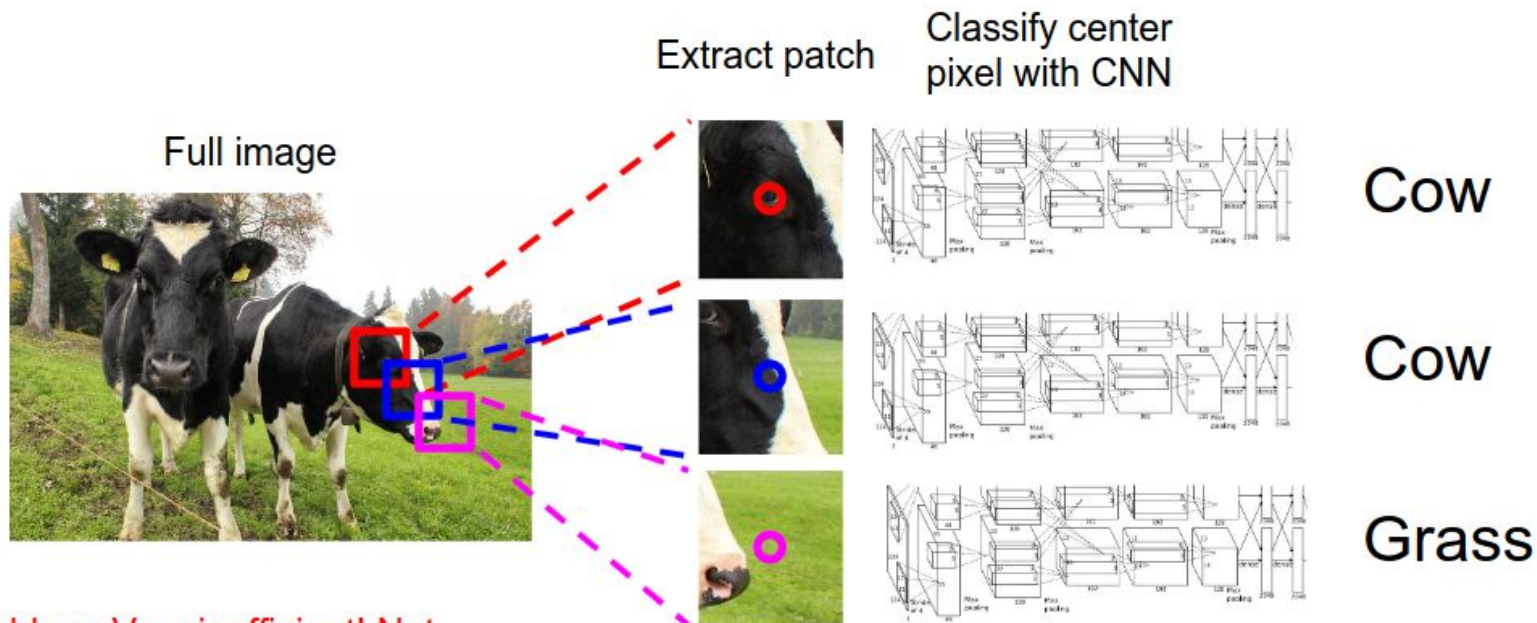
# Context



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

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

# Inefficiency

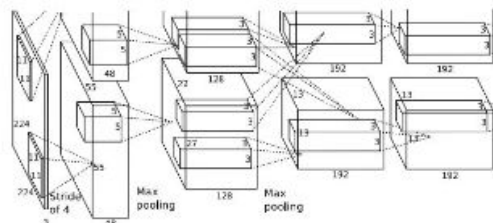


Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013  
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Shared Convolutional features

Full image

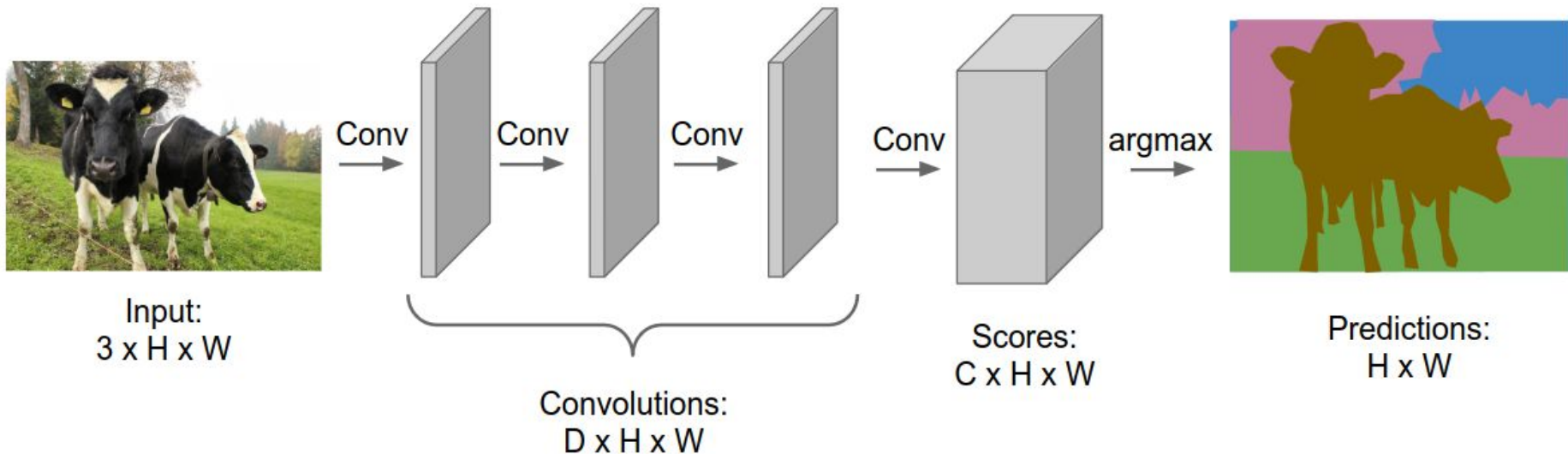


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

**Problem:** classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

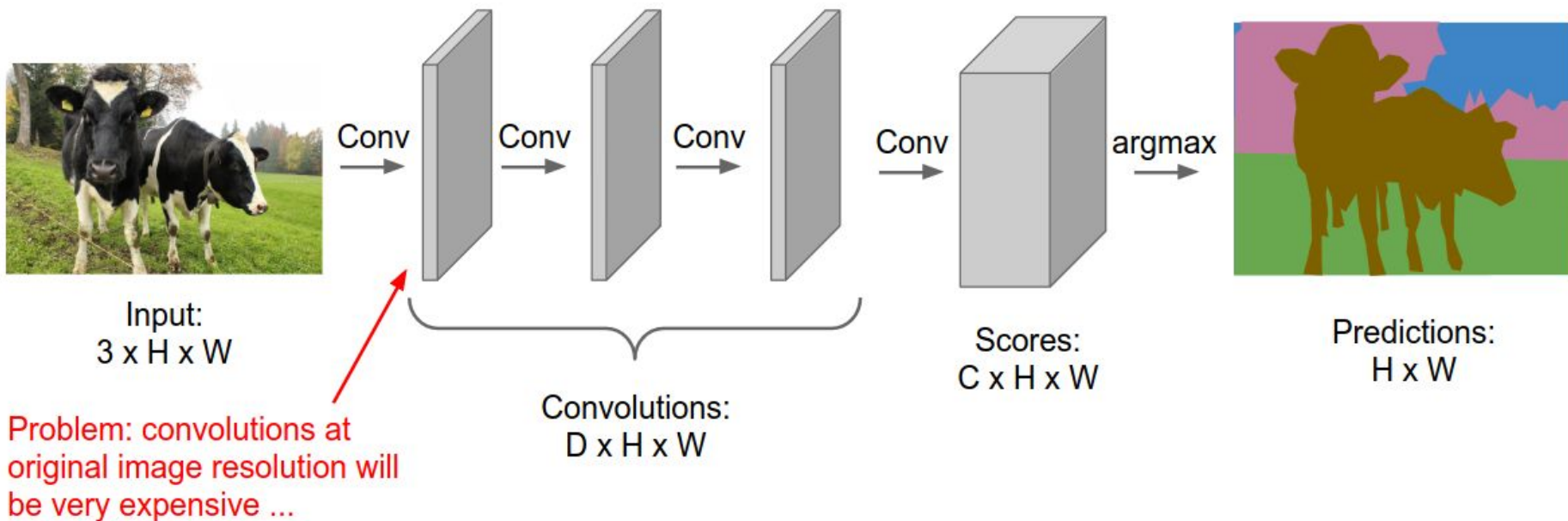
# CNN Segmentation

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



# CNN Segmentation

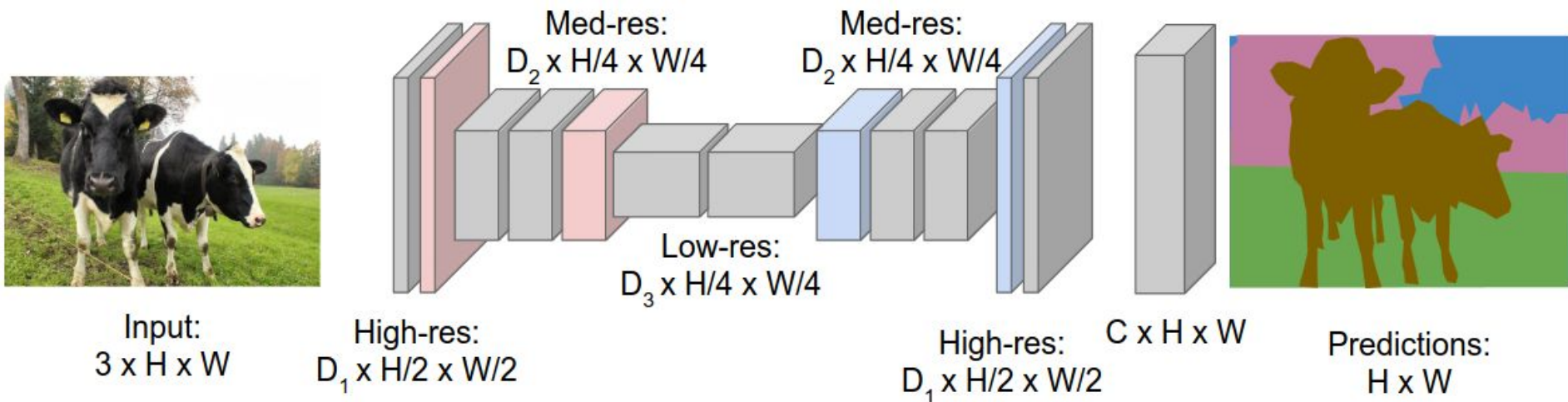
Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!





# CNN Segmentation

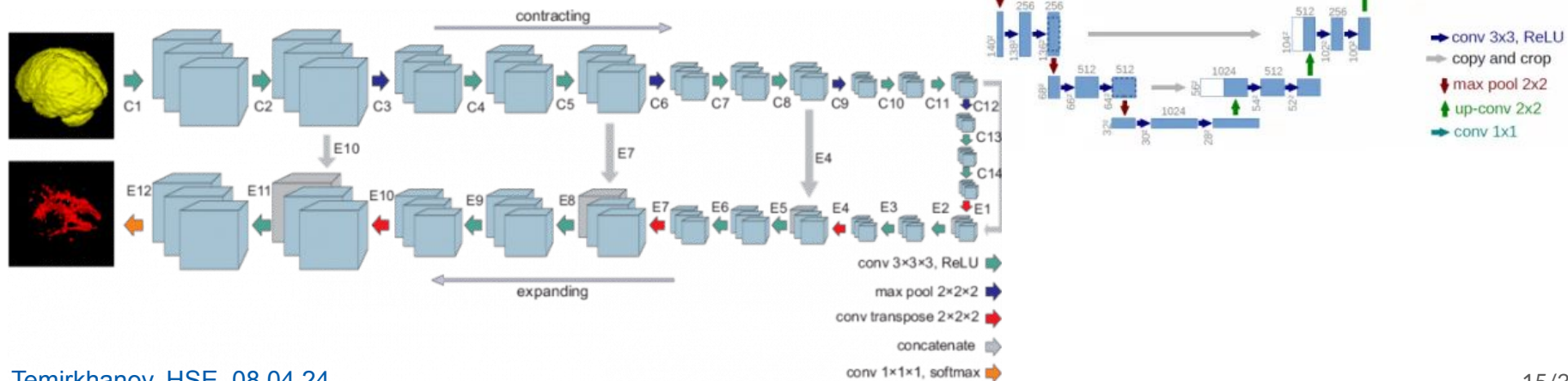
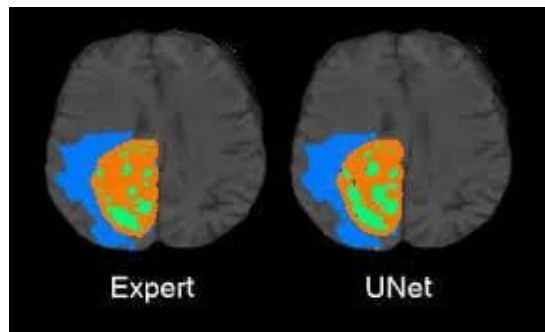
Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# UNet



# Unpooling

**Nearest Neighbor**

1	2
3	4

Input: 2 x 2



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

Output: 4 x 4

**“Bed of Nails”**

1	2
3	4

Input: 2 x 2

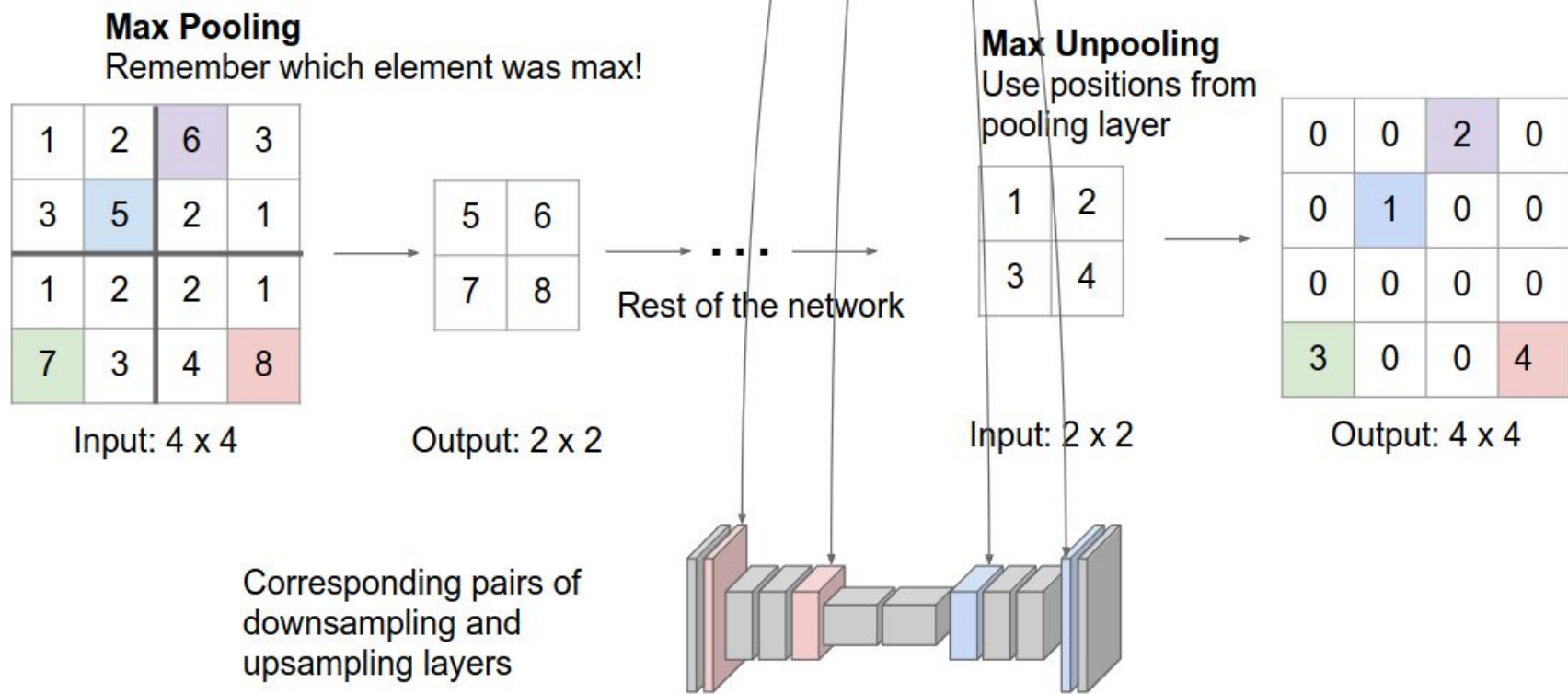


1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output: 4 x 4

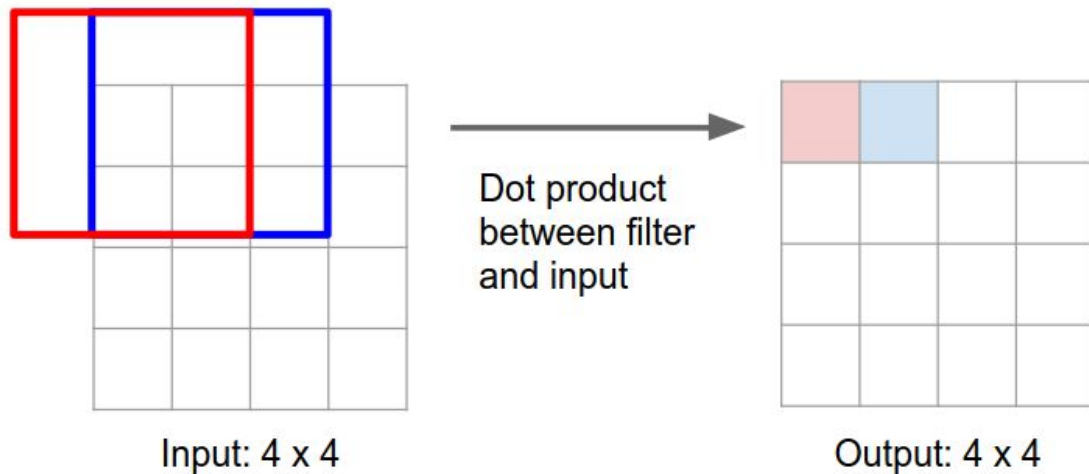


# Max unpooling



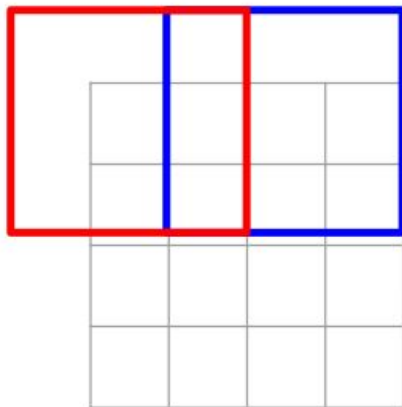
# Transpose Convolution

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



# Transpose Convolution

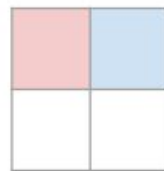
**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product  
between filter  
and input



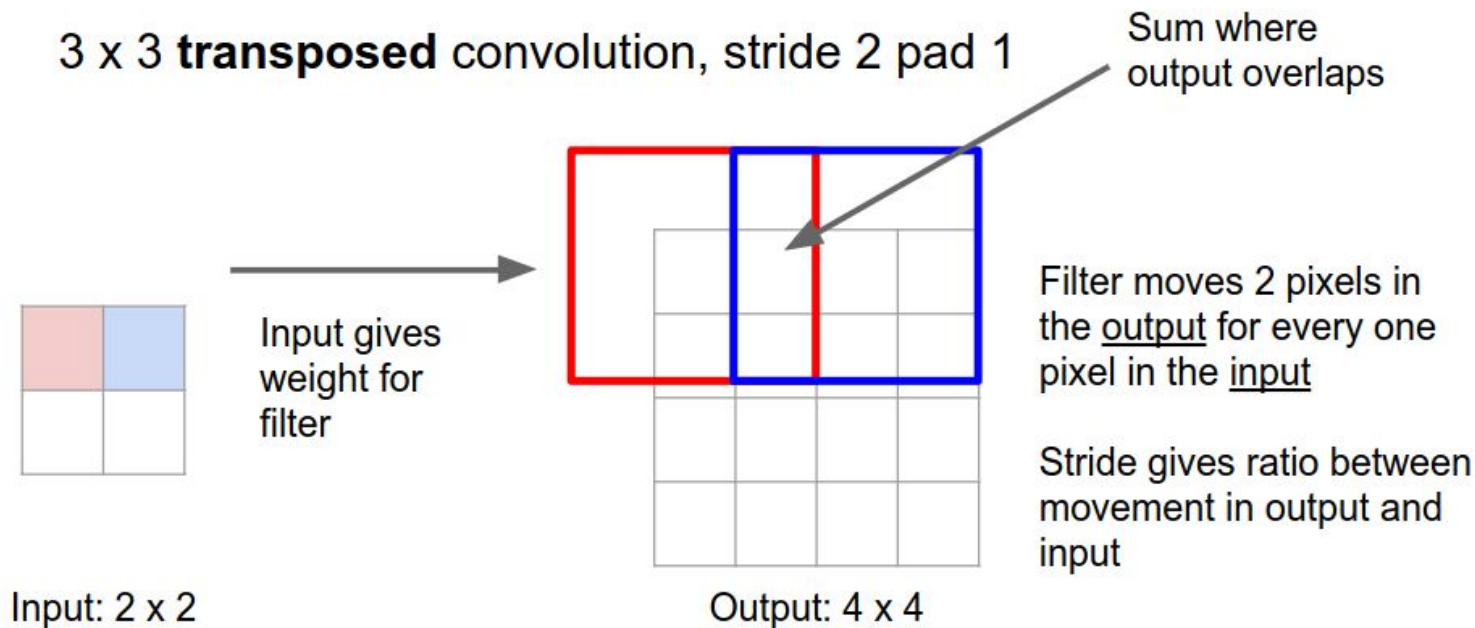
Output: 2 x 2

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

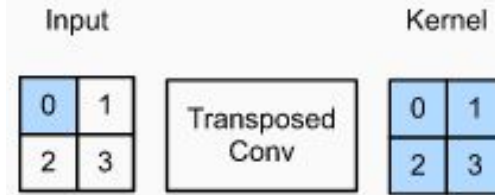
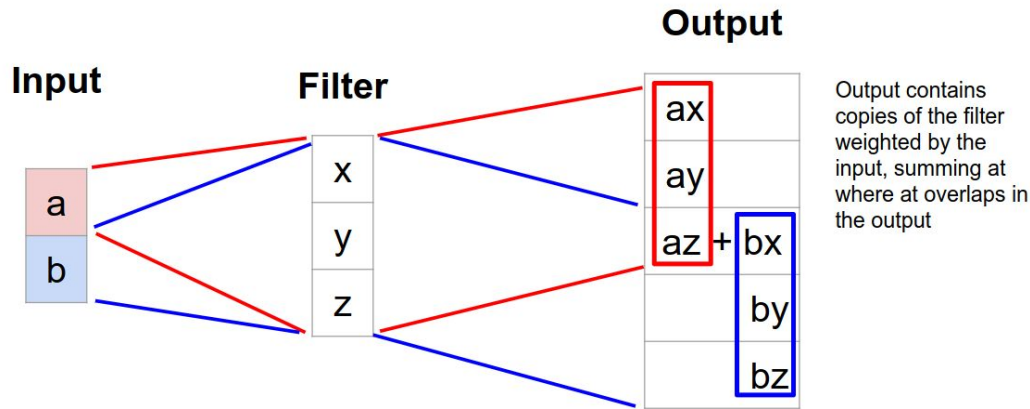
Stride gives ratio between movement in input and output

We can interpret strided convolution as “learnable downsampling”.

# Transposed Convolution



# Transposed Convolution



**Output**

$$\begin{bmatrix} 0 & 0 & \\ 0 & 0 & \\ & & \end{bmatrix} + \begin{bmatrix} & 0 & 1 \\ & 2 & 3 \\ & & \end{bmatrix} + \begin{bmatrix} & & \\ 0 & 2 & \\ 4 & 6 & \end{bmatrix} + \begin{bmatrix} & & \\ & 0 & 3 \\ & 6 & 9 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 4 & 6 \\ 4 & 12 & 9 \end{bmatrix}$$

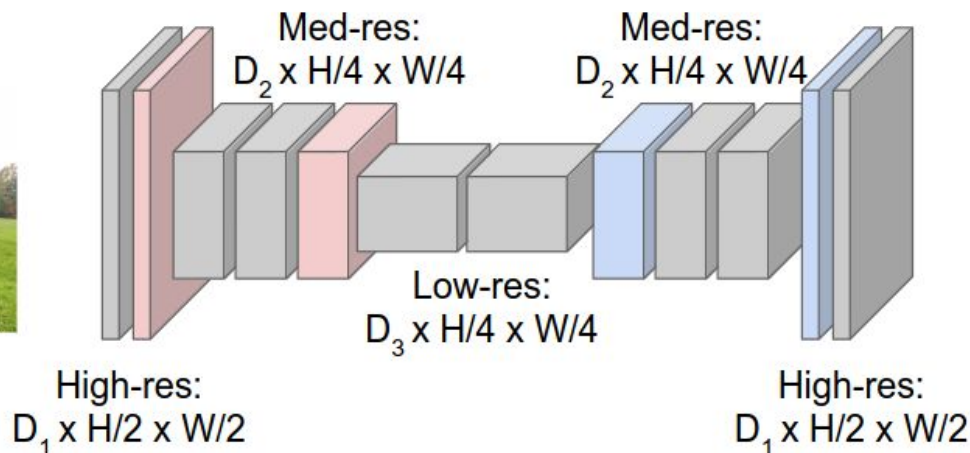
# Semantic Segmentation

**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:**  
Unpooling or strided  
transposed convolution

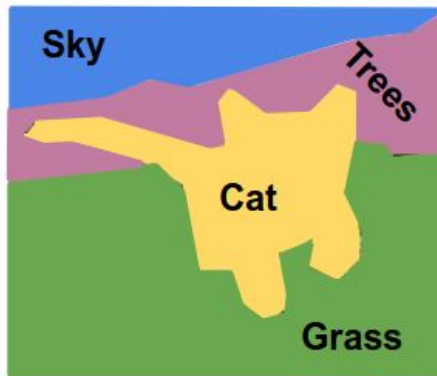


Predictions:  
 $H \times W$

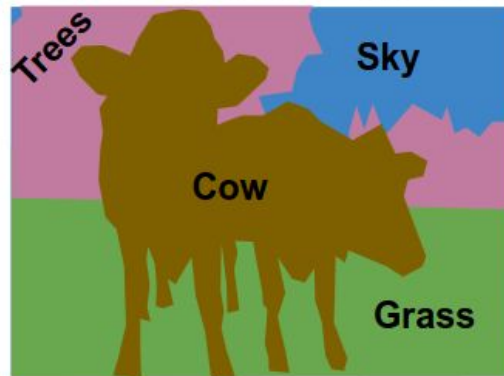
# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



[This image is CC0 public domain](#)





# Instance Segmentation

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Instance Segmentation

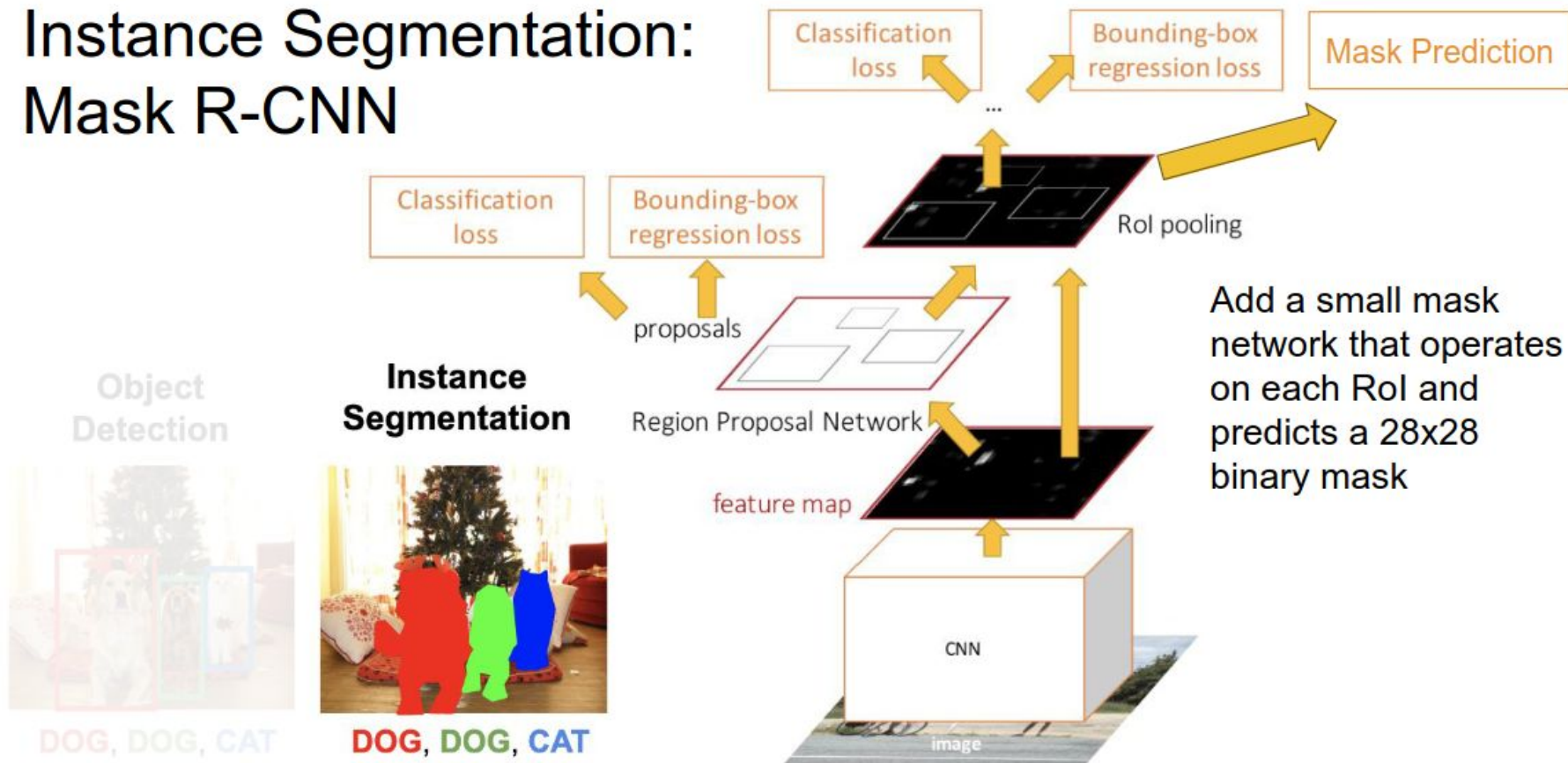


DOG, DOG, CAT

Multiple Object

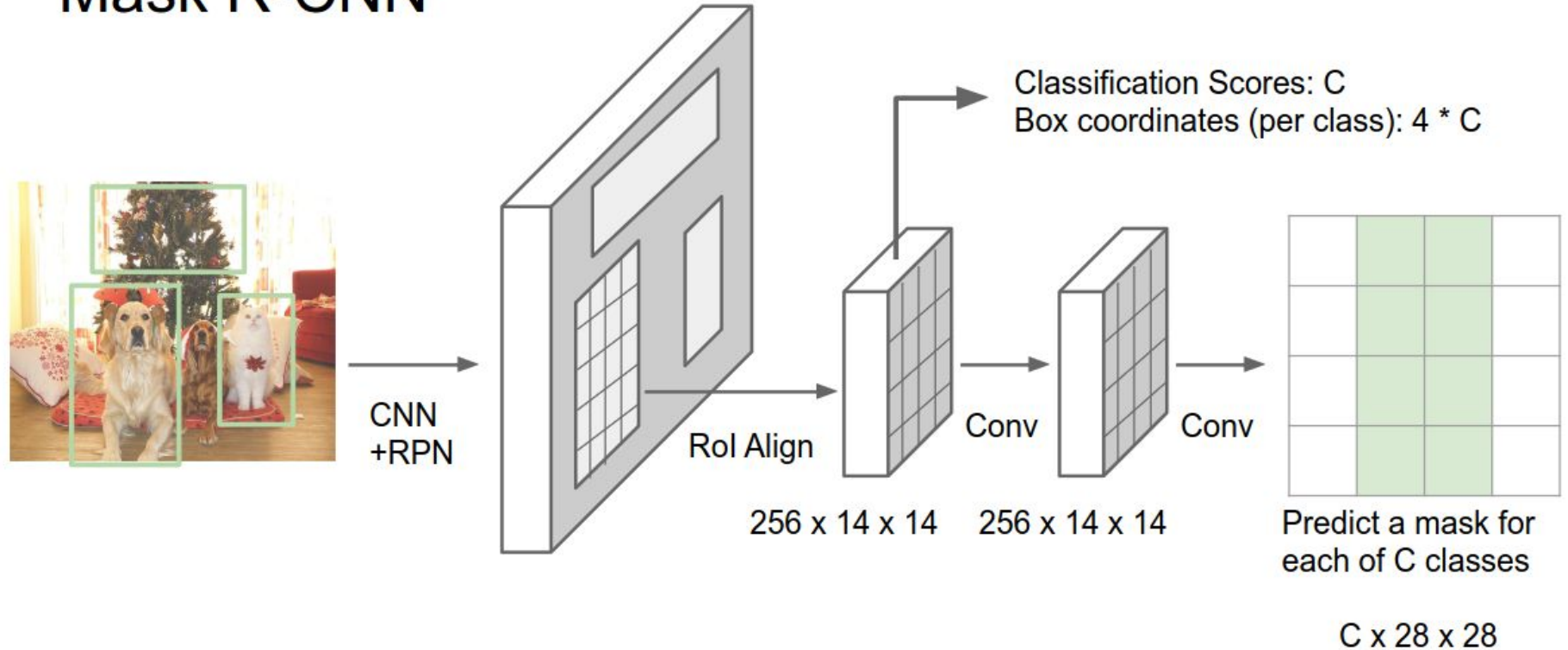


# Instance Segmentation: Mask R-CNN



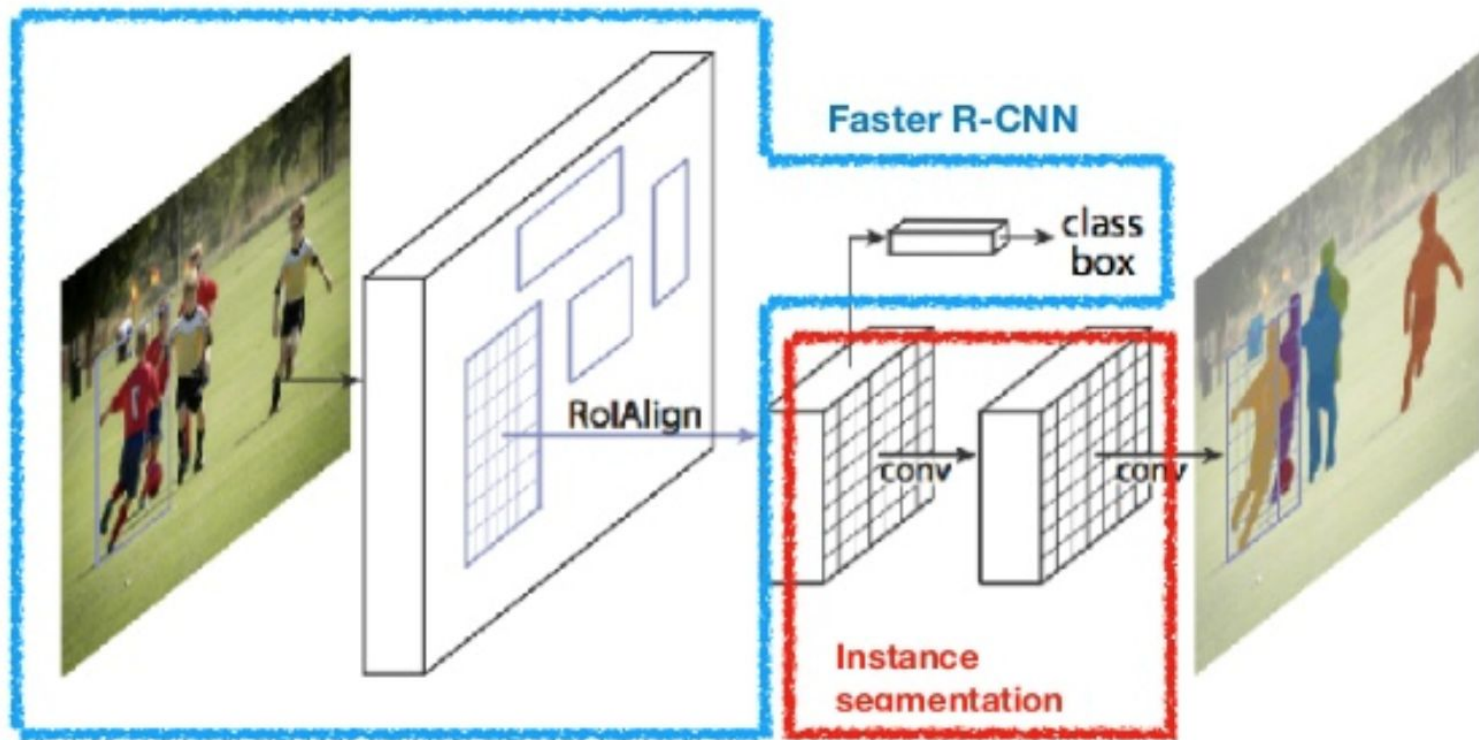
He et al, "Mask R-CNN", ICCV 2017

# Mask R-CNN

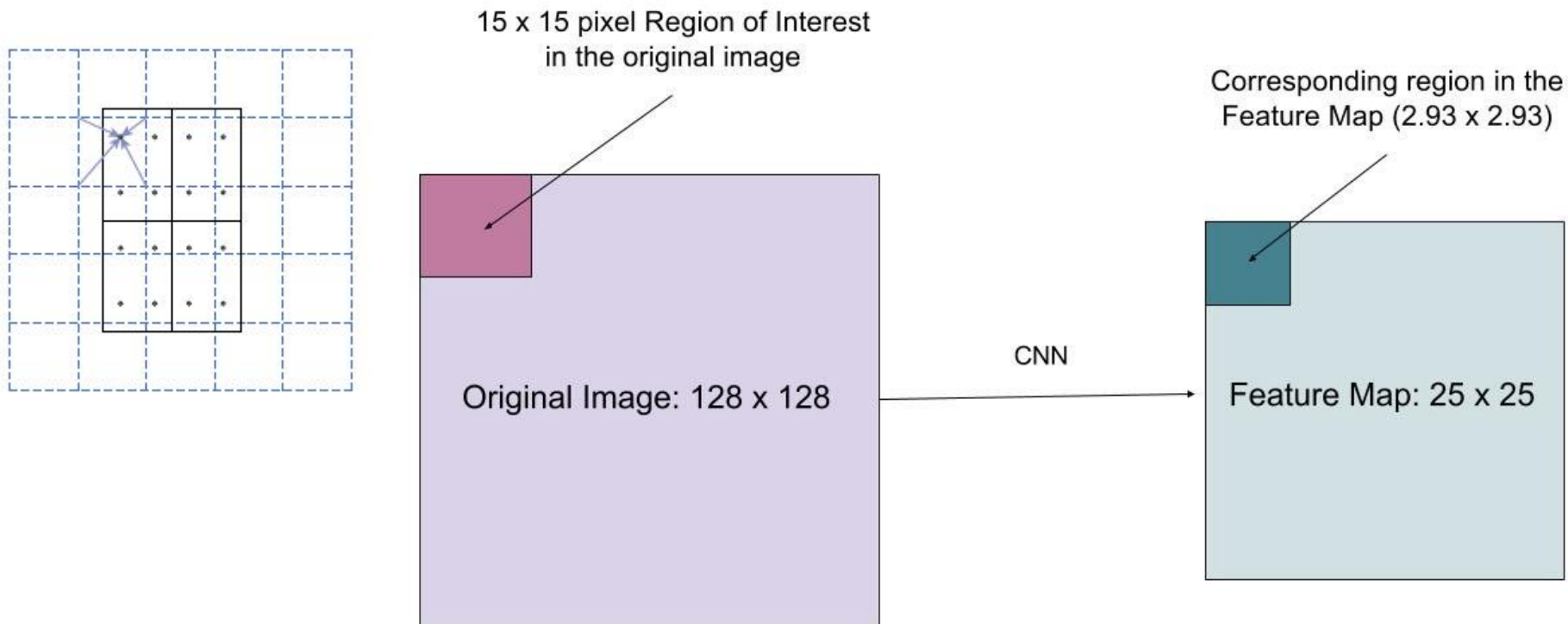


He et al, "Mask R-CNN", arXiv 2017

# Mask R-CNN



# ROI Alignment

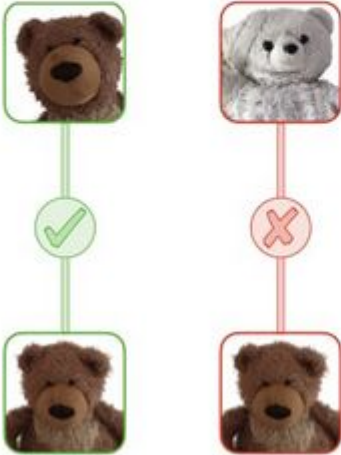
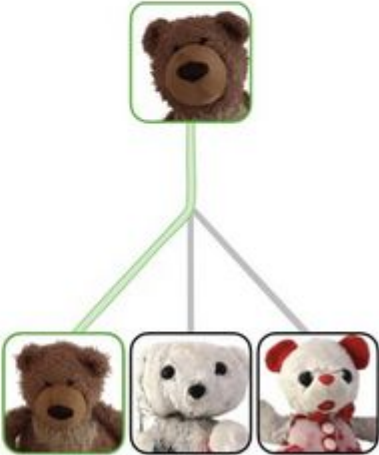


# Bonus part

---

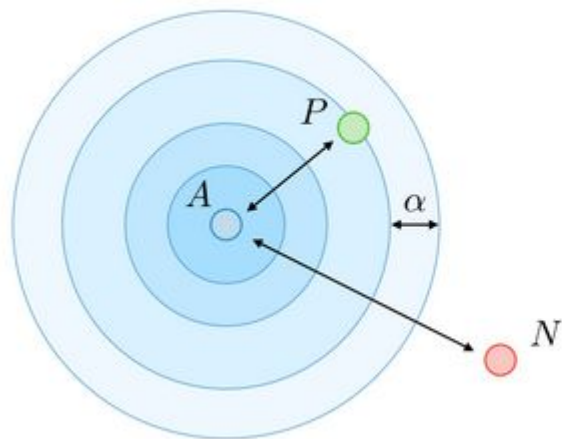
Addition vision tasks

# Face verification and recognition

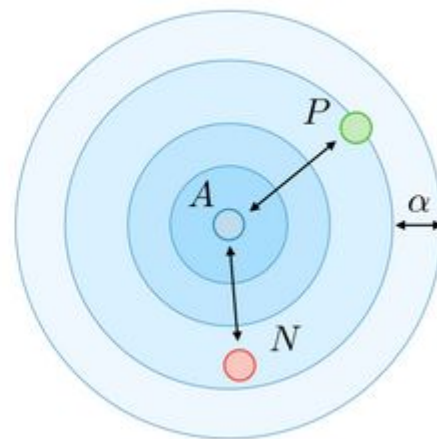
Face verification	Face recognition
<ul style="list-style-type: none"><li>• Is this the correct person?</li><li>• One-to-one lookup</li></ul>	<ul style="list-style-type: none"><li>• Is this one of the <math>K</math> persons in the database?</li><li>• One-to-many lookup</li></ul>
<p>Query</p>  <p>Reference</p> <p>The diagram illustrates the face verification process. It shows two vertical pairs of teddy bear images. The left pair consists of two identical brown teddy bears, connected by a green line with a green checkmark icon in the middle, indicating a successful match. The right pair consists of a grey teddy bear at the top and a brown teddy bear at the bottom, connected by a red line with a red 'X' icon in the middle, indicating a failed match.</p>	<p>Query</p>  <p>Database</p> <p>The diagram illustrates the face recognition process. At the top, a brown teddy bear is shown in a green-bordered box, labeled 'Query'. Below it, three lines branch out to three different teddy bears in a row, each in a black-bordered box, labeled 'Database'. The first bear is brown, the second is white, and the third is white with red ears. A green line connects the query bear to the first database bear, while grey lines connect it to the other two, representing a one-to-many lookup.</p>

# Triplet Loss

$$\ell(A, P, N) = \max(d(A, P) - d(A, N) + \alpha, 0)$$



$$\ell(A, P, N) = 0$$



$$\ell(A, P, N) > 0$$



# Style Transfer and Domain Adaptation



Content  $C$

+



Style  $S$

=



Generated image  $G$