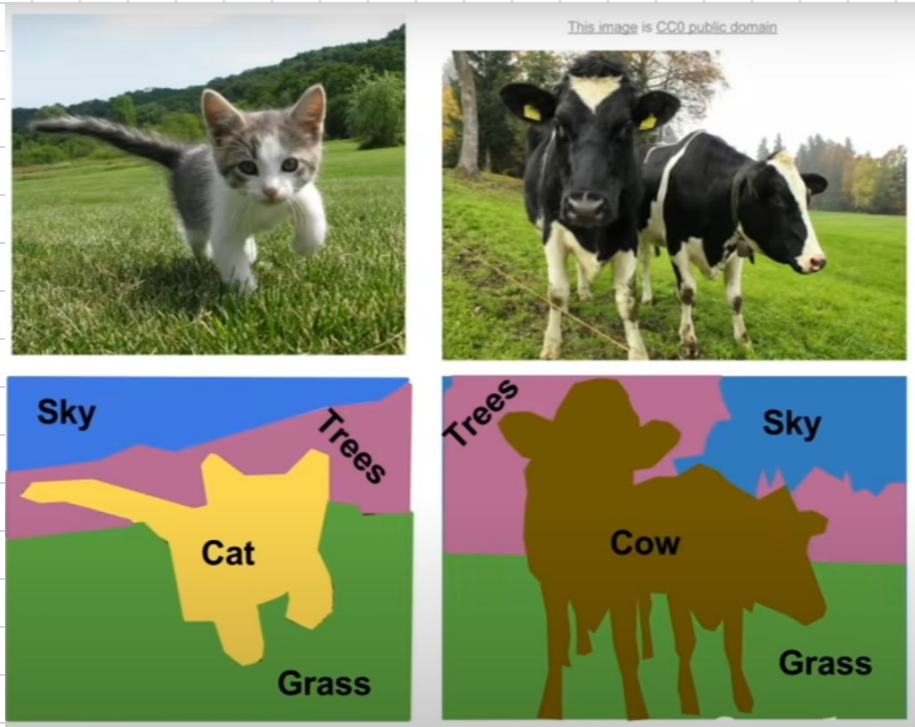


# Semantic segmentation:

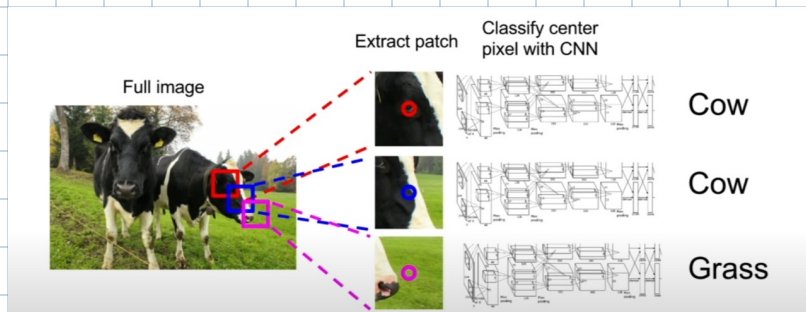
↳ Outputs a decision of a category for every pixel in that image.

↳ Does not differentiate between instances



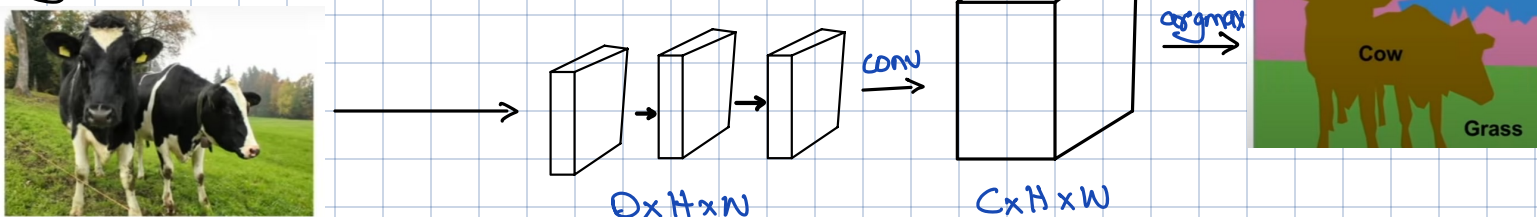
## Sliding window:

↳ Extract patches from the image  
↳ Apply CNN on the patch  
↳ Classify it

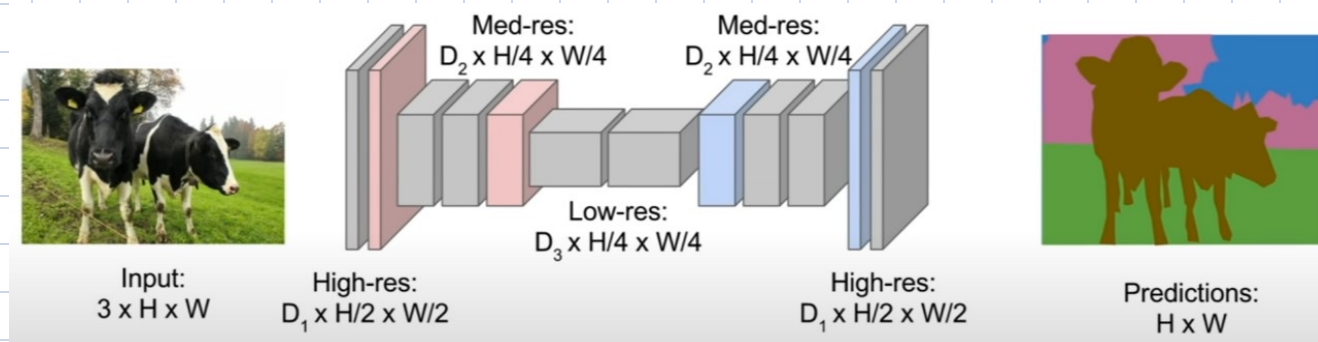


↳ Computationally expensive

## Fully convolutional:



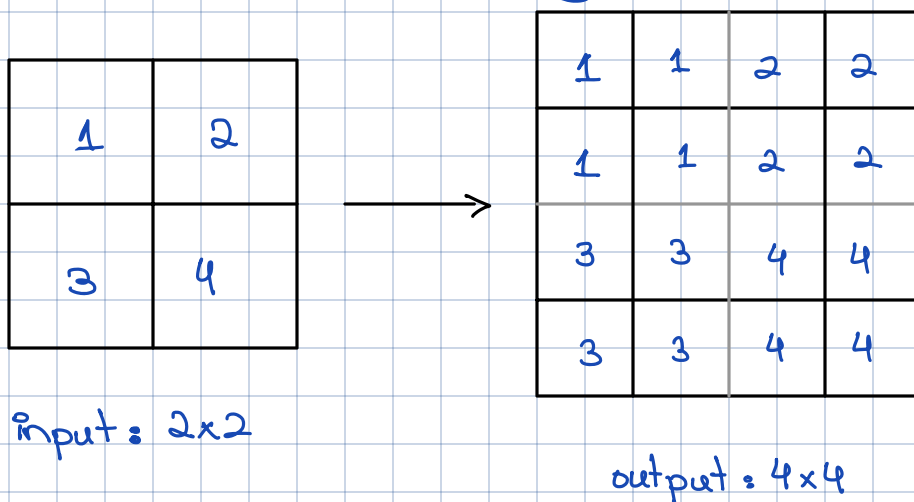
- ↳ A bunch of conv layers
- ↳ Final conv layer gives classification scores for every pixel
- ↳ Too computationally expensive
- ↳ You use upsampling & downsampling



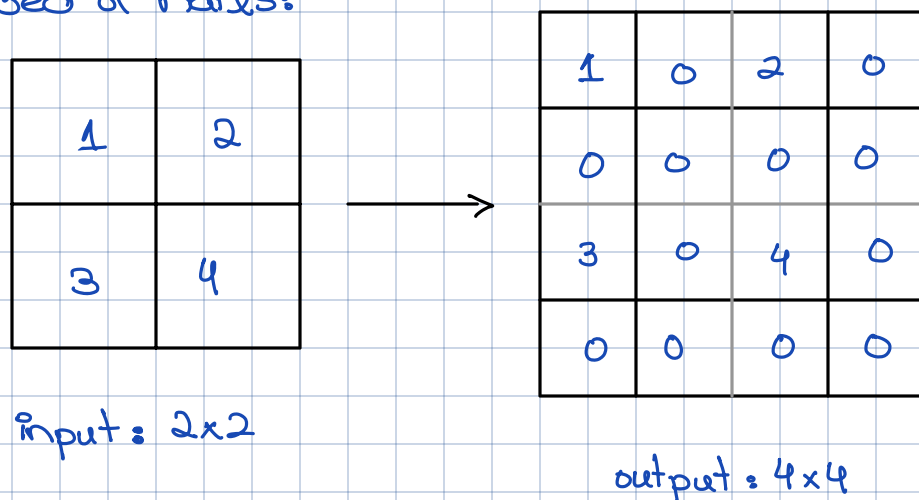
↳ Conv  $\rightarrow$  Down sampling  $\rightarrow$  Upsample  $\rightarrow$  Predictions

Unpooling:

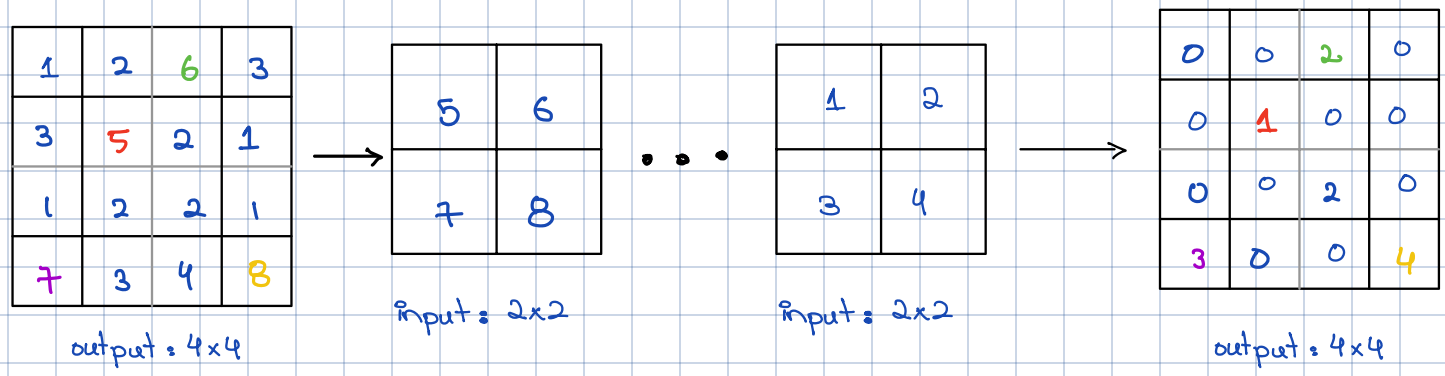
↳ Nearest neighbours unpooling



↳ Bed of nails:



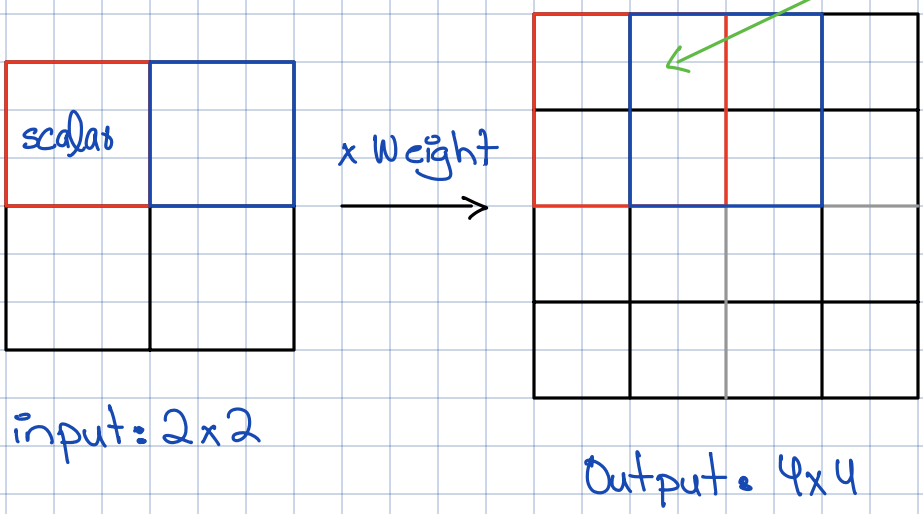
# Max unpooling:



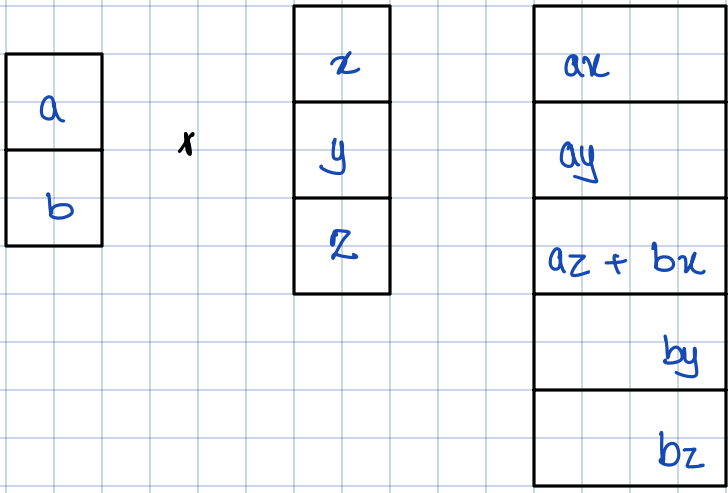
Transpool Convolution/De/up/fractionally/Backward strided conv

↳ Learnable Upsampling

sum where they overlap



Eg: 1D:

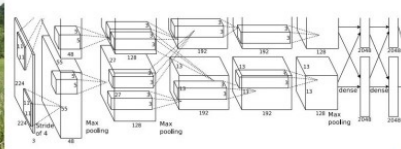


# Classification + localization:

↳ Classify image & draw boundary box around that object



This image is CC0 public domain



Fully Connected:  
4096 to 1000

Class Scores

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

correct label

softmax loss

Vector:  
4096

Fully Connected:  
4096 to 4

Box Coordinates  
(x, y, w, h)

L2 loss

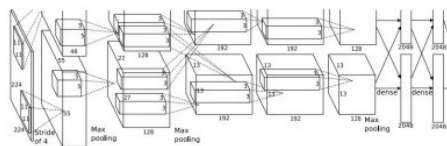
Correct box  
(x, y, w, h)

↳ Take weighted sum of the loss

↳ Hypersparameterise

## Human pose estimation:

↳ Represent by 14 joint locations



Vector:  
4096

Left foot: (x, y) → L2 loss

Right foot: (x, y) → L2 loss

Head top: (x, y) → L2 loss

Correct left foot: (x', y')

Correct head top: (x', y')

+

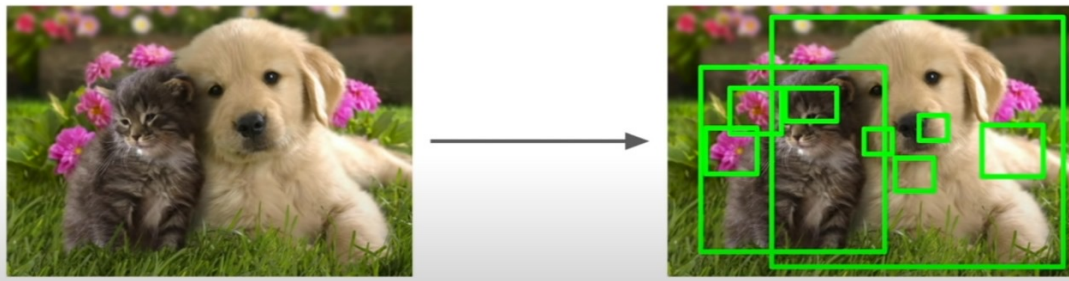
Loss

## Object detection:

↳ Can't use above due many objects thus many coordinates  
↳ Can't use sliding window because "How to choose the crop?".  
Objects can anywhere, at any size & too computationally expensive

## Region proposals:

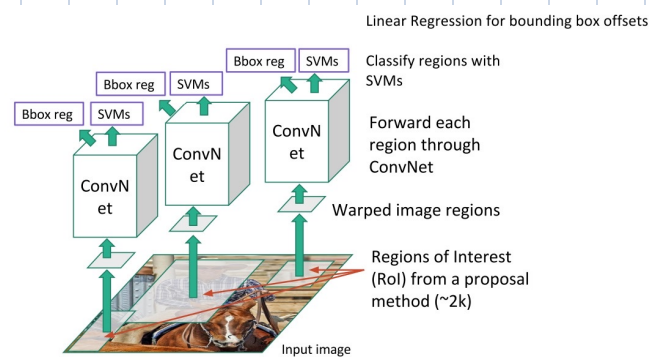
↳ gives regions where object might be present



↳ Selective search gives 2000 regions where objects might be found

## R-CNN:

- ↳ Take proposal regions
- ↳ Change them into same size
- ↳ Pass them through CNN
- ↳ Use SVM
- ↳ Predicts bounding box ( $x, y, w, h$ )

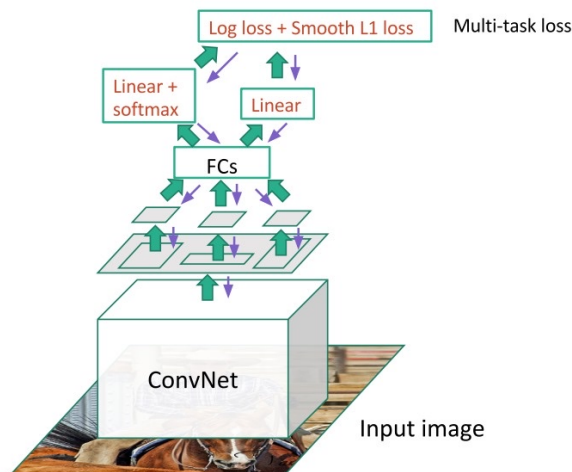


## Problems:

- ↳ Computationally expensive
- ↳ Region proposals NOT learned
- ↳ Training is slow

## Fast R-CNN

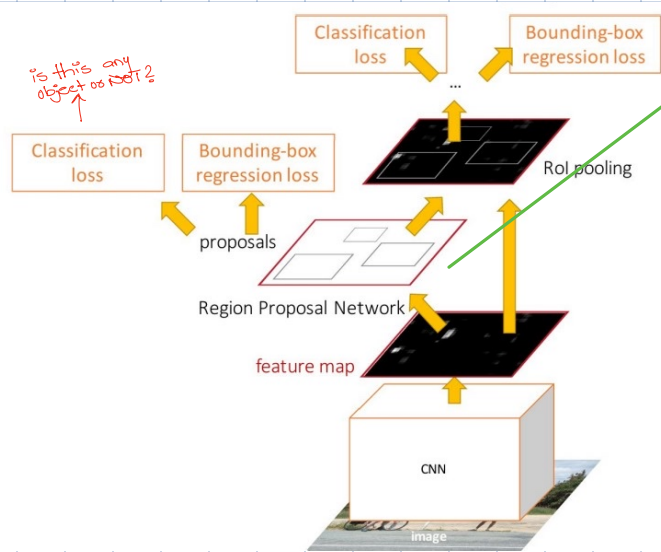
- ↳ Run image through CNN
- ↳ Crop patches corresponding to the feature map on the input



↳ This ends up getting slowed by the process of computing region proposals



# Faster R-CNN.



Any time region proposal has some overlap with any ground truth object then that's positive region proposal

You only look once  
↑

Detections without proposals: YOLO/SSD  $\rightarrow$  single shot detection

$\hookrightarrow$  Divide image into grids  $7 \times 7$

$\hookrightarrow$  Image a set of base bounding boxes,  $B=3$

$\hookrightarrow$  Within each grid cell:

$\hookrightarrow$  Predict offset from base bounding box

$\hookrightarrow$  Predict scores

$\hookrightarrow$  Output:  $7 \times 7 \times (5 \times C + C)$

Instance Segmentation:

$\hookrightarrow$  Which pixels corresponds to object

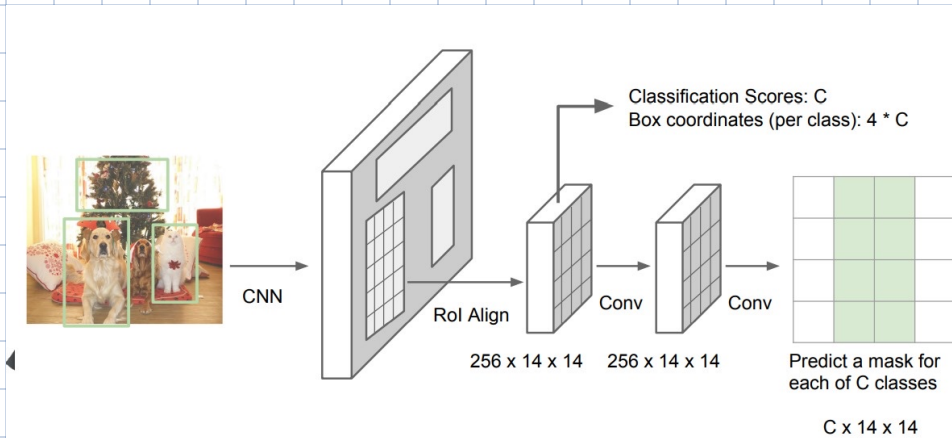
$\hookrightarrow$  Differentiates between multiply same instances



# Mask R-CNN:

↳ Same as faster Faster R-CNN

↳ Now, it labels the exact pixels in the region proposals



↳ Can also do pose detection