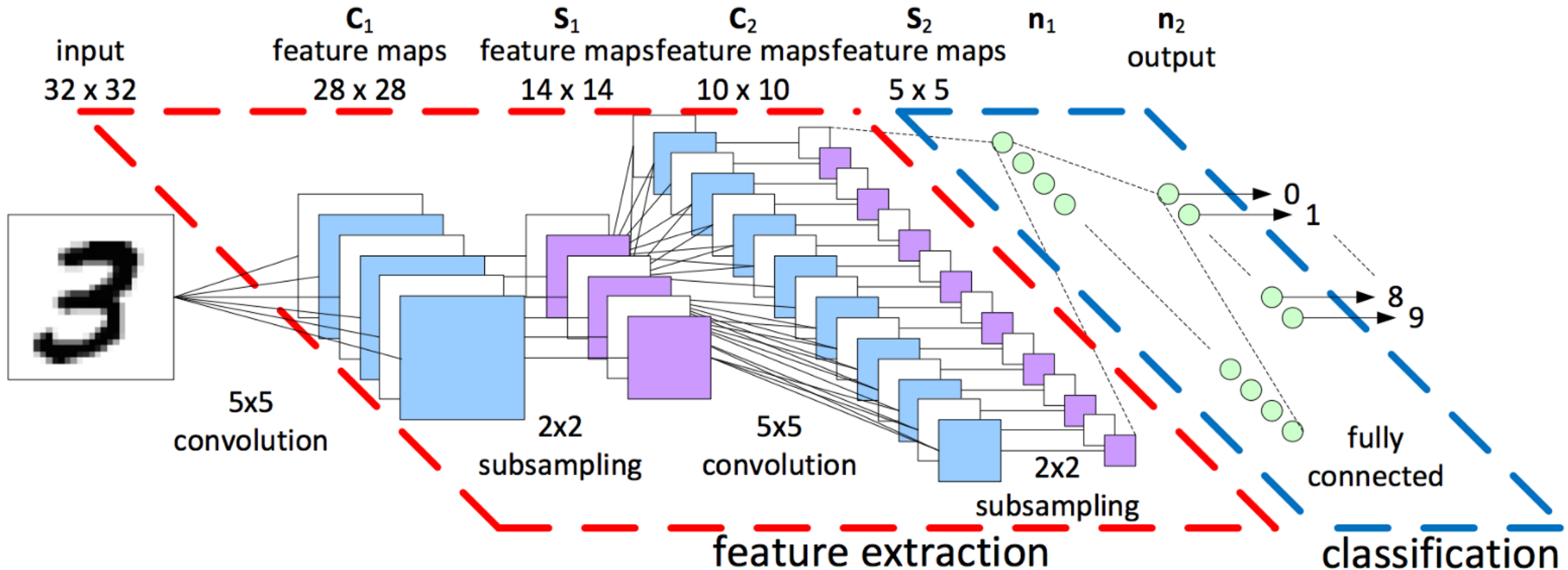


Introduction to Image Segmentation using Deep Learning

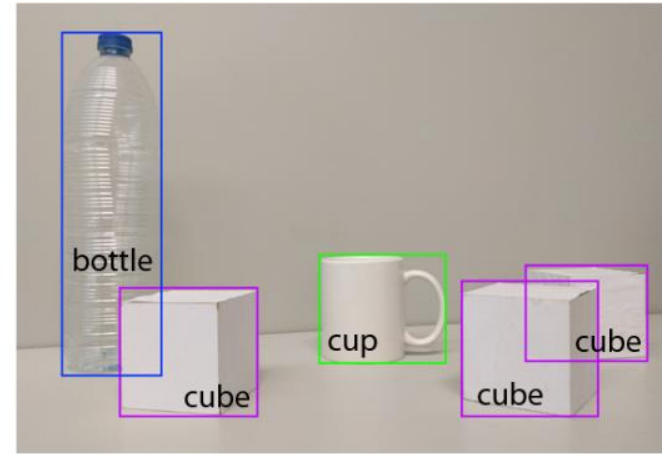
So far image classification



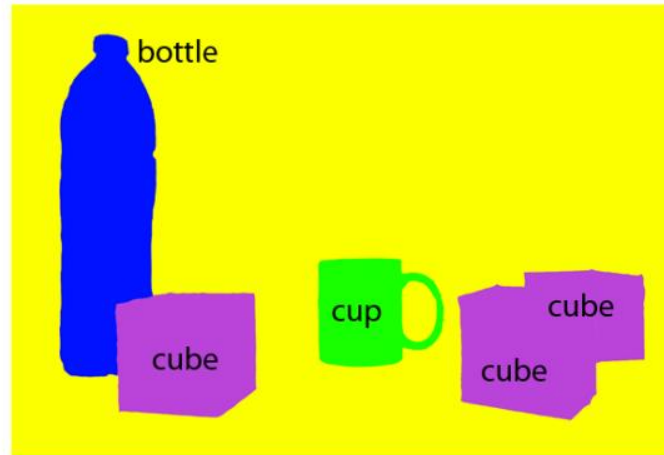
Other Computer Vision Tasks



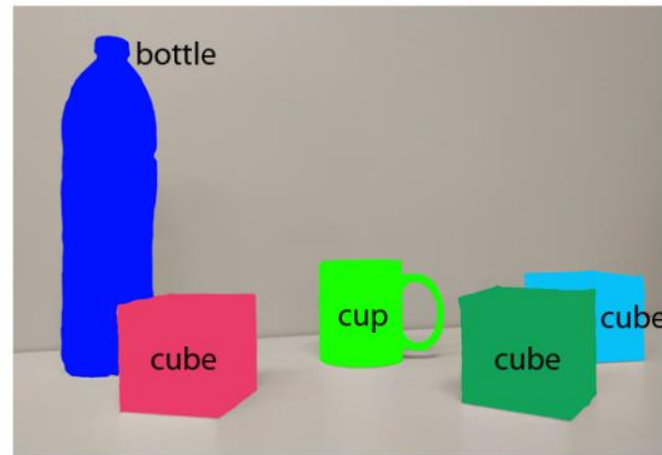
(a) Image classification



(b) Object localization



(c) Semantic segmentation



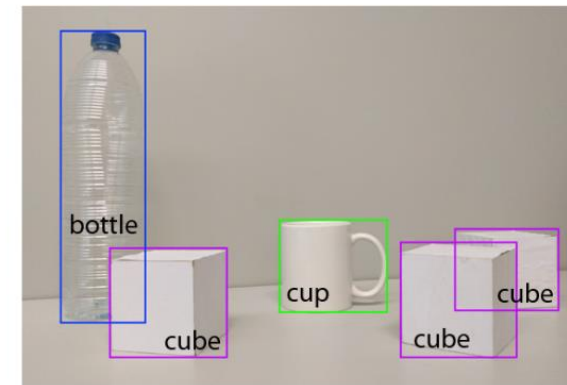
(d) Instance segmentation

Other Computer Vision Tasks

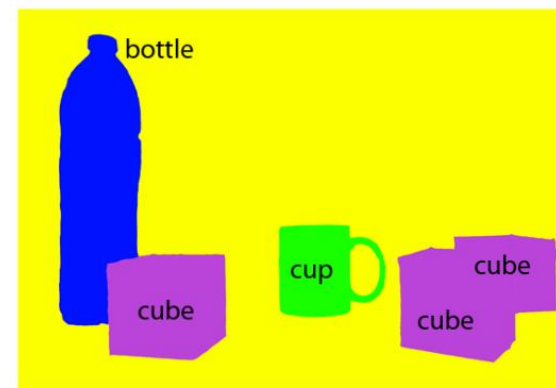
- **Object Localization or object detection** is to identify all the class relevant objects in an image
- **Segmentation** is a partition of an image into several parts by giving pixel label, but without any attempt at understanding what these parts represent.
- **Semantic segmentation** attempts to partition the image into semantically meaningful parts, and to classify each part into one of the pre-determined classes. But, still can't differ among objects/instances in the same class
- **Instance segmentation** can differ objects/instances in the same class while performing the semantic segmentation



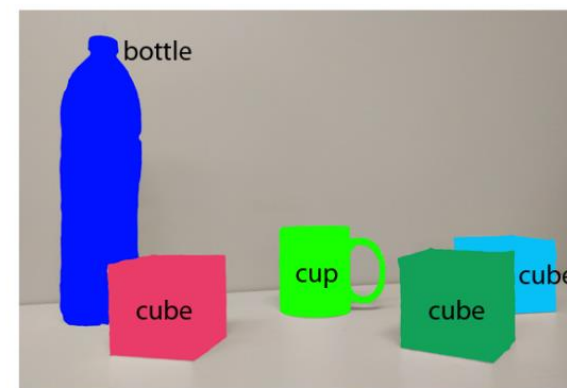
(a) Image classification



(b) Object localization



(c) Semantic segmentation

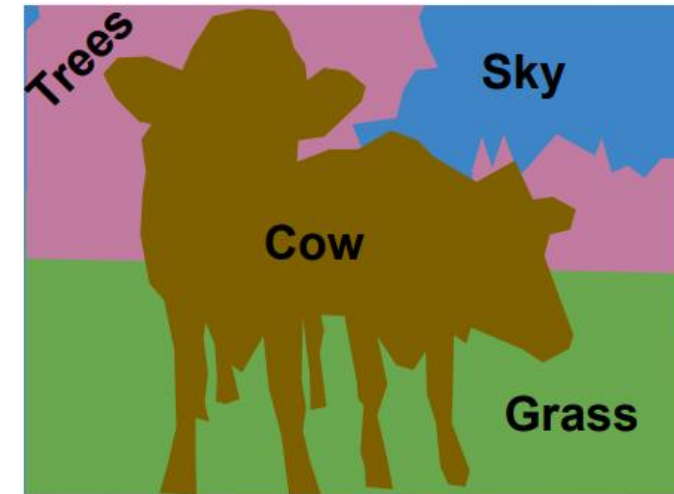
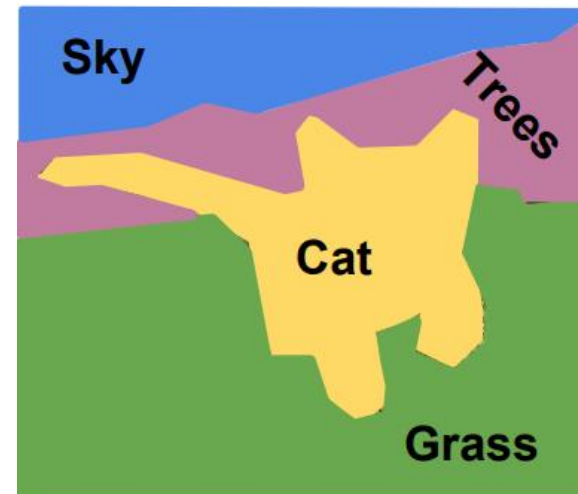


(d) Instance segmentation

A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, "A review on deep learning techniques applied to semantic segmentation,"

Semantic Segmentation

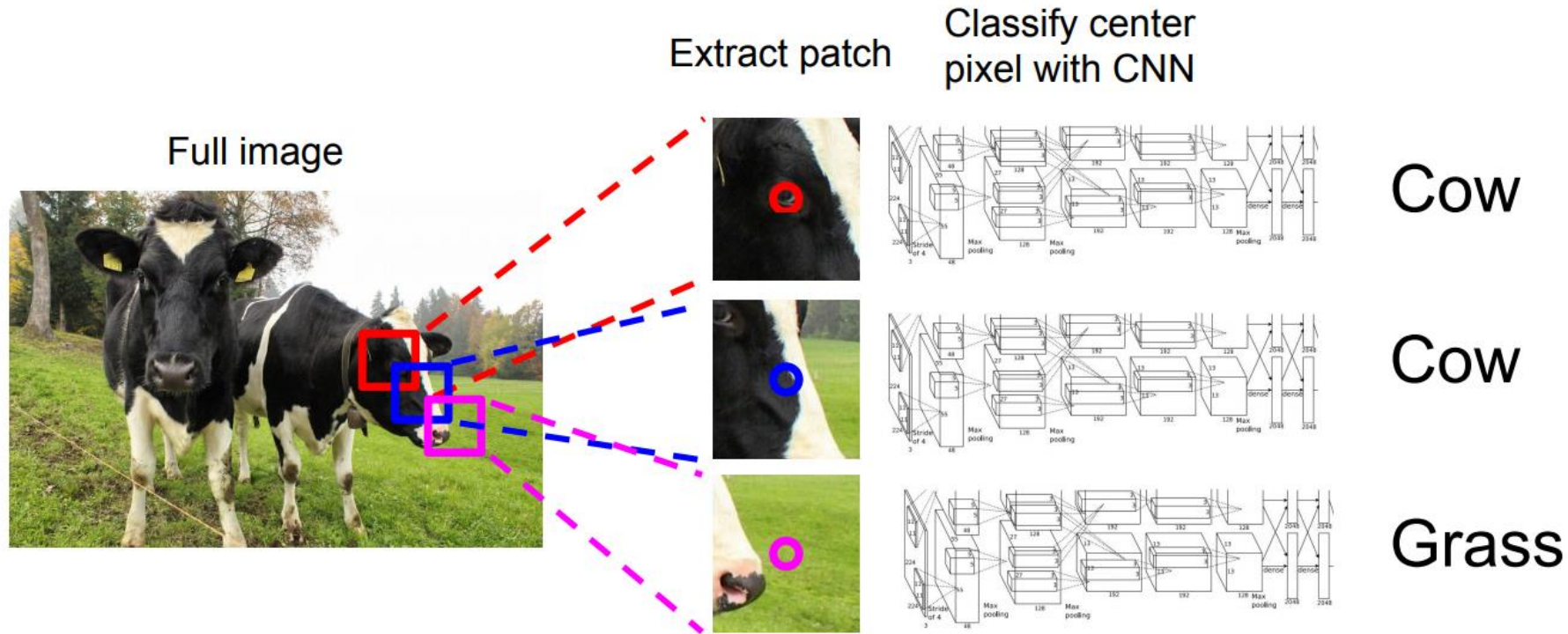
- Purely supervise learning
- Label each pixel in the image with a category label
- But do not differentiate instances, only care about pixels
- Dataset preparation is very expensive
- Online-tools are there for manual labeling
https://en.wikipedia.org/wiki/List_of_manual_image_annotation_tools
- Training could be slower than image classification



Semantic Segmentation methods

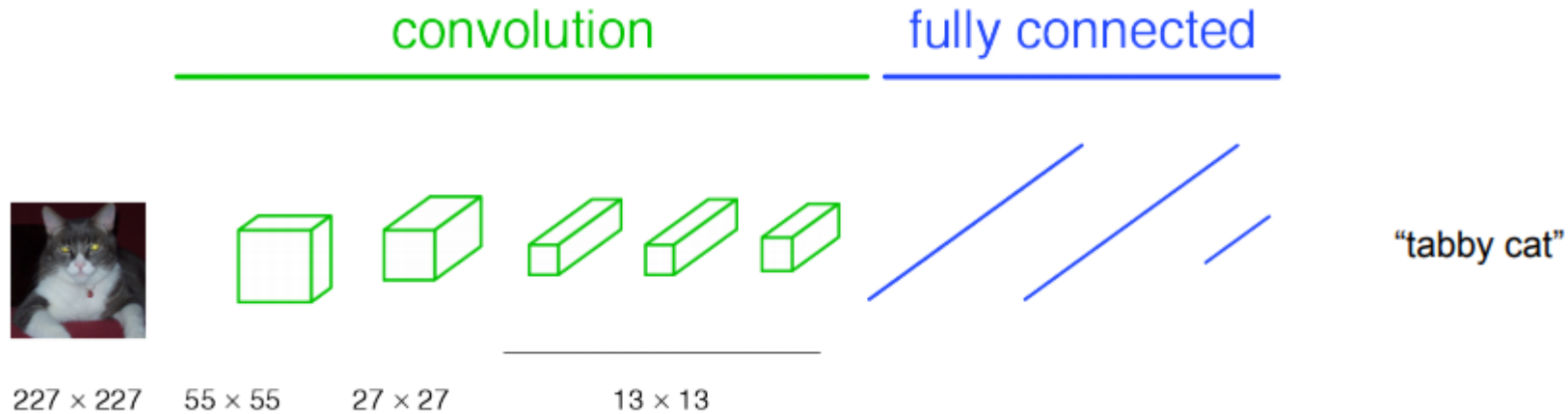
- Sliding Window approach
 - Farabet et al, “Learning Hierarchical Features for Scene Labeling,” TPAMI 2013
 - Pinheiro and Collobert, “Recurrent Convolutional Neural Networks for Scene Labeling”, ICML 2014
- Fully convolutional network approach
 - Long, Shelhamer, and Darrell, “Fully Convolutional Networks for Semantic Segmentation”, CVPR 2015
 - Noh et al, “Learning Deconvolution Network for Semantic Segmentation”, ICCV 2015

Semantic Segmentation: sliding window

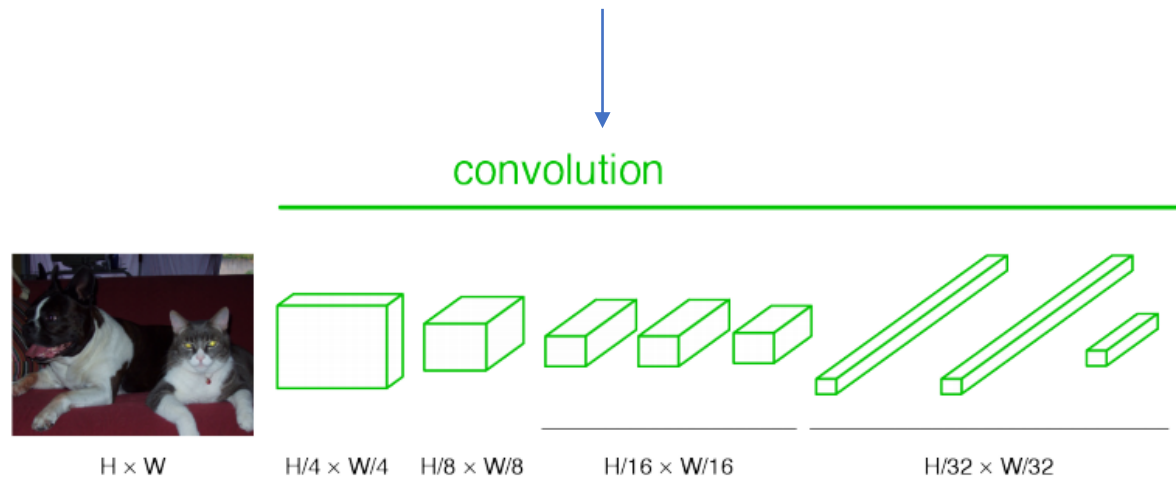
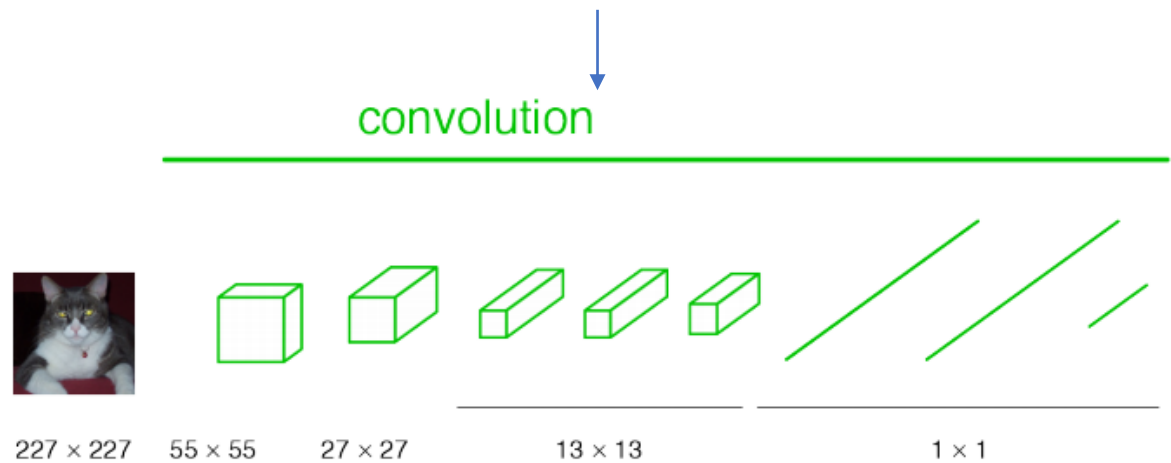
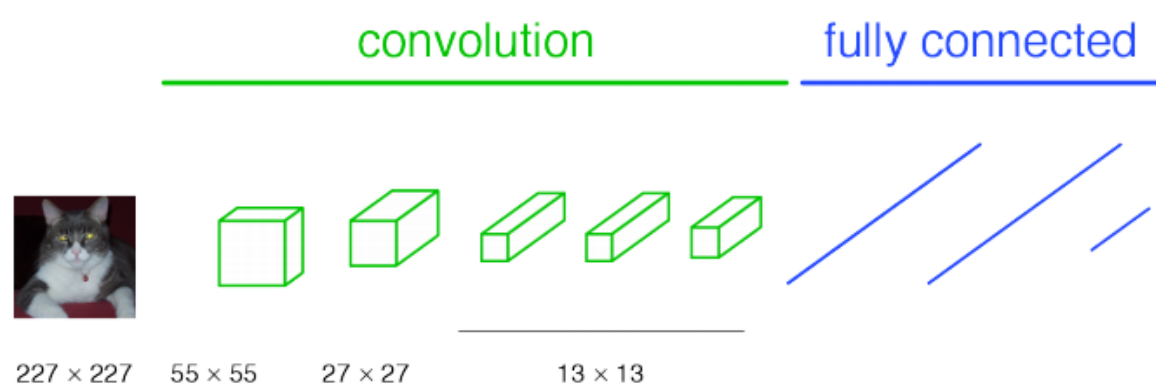


- Predict the middle pixel of the window while sliding it throughout the image
- But this method is very inefficient and less accurate
- Since method using a sliding window, for a single image there will be huge number of forward passes and backward passes.
- Another reason is that its not reusing shared features between overlapping patches

Semantic Segmentation: Fully Convolutional Networks

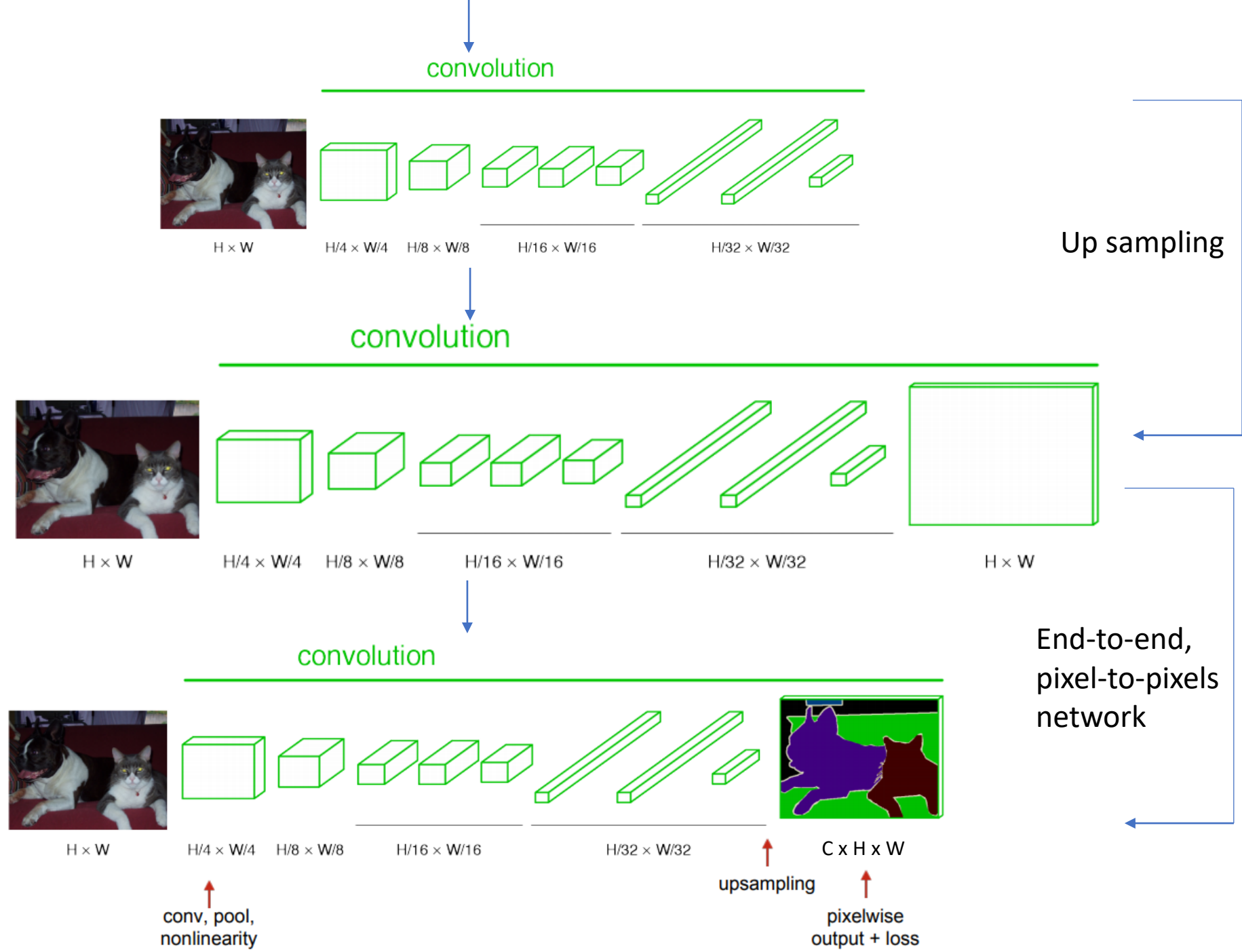


- Normal classification convolutional networks have a fully connected network for classification
- Since the fully connected networks have a fixed size, input image need to be resize into a fix size input
- But when the image is resize, we loss some pixels or generate new pixels which can cause accuracy problems in the segmentation
- There for the solution is to remove the fully connected part



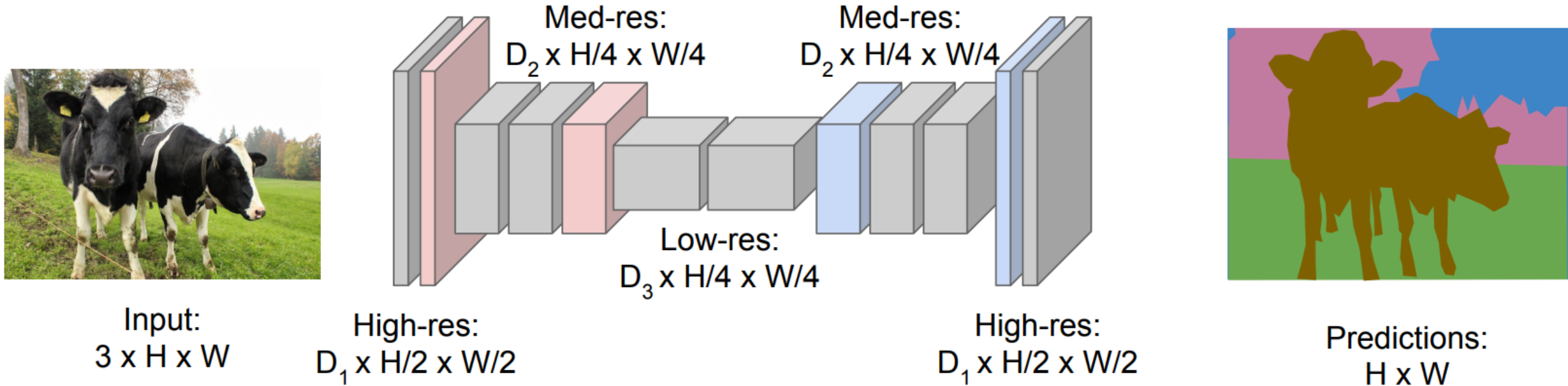
Making all convolutional filters

Making all convolutional filter sizes change with the input size



Designing a Fully Convolutional Networks

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



- Down sample the high resolution image to increase the computational efficiency
 - Strided convolution
 - Polling (max, average)
- Up sample the convolutional filters up to image resolution to compare it with the image size segmentation label.

Upsampling

- In-Network upsampling: “Unpooling”
 - Nearest Neighbor
 - Bed of Nails
 - Max Unpooling
- Learnable Upsampling: Transpose Convolution

In-Network upsampling: “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

In-Network upsampling: “Max Unpooling”

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4



5	6
7	8

Output: 2 x 2

...

Rest of the network

Max Unpooling

Use positions from pooling layer

1	2
3	4

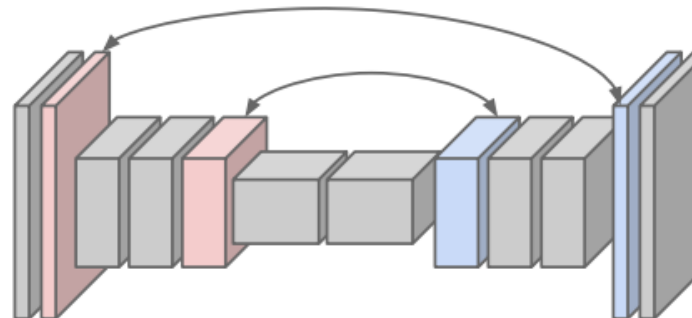
Input: 2 x 2



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

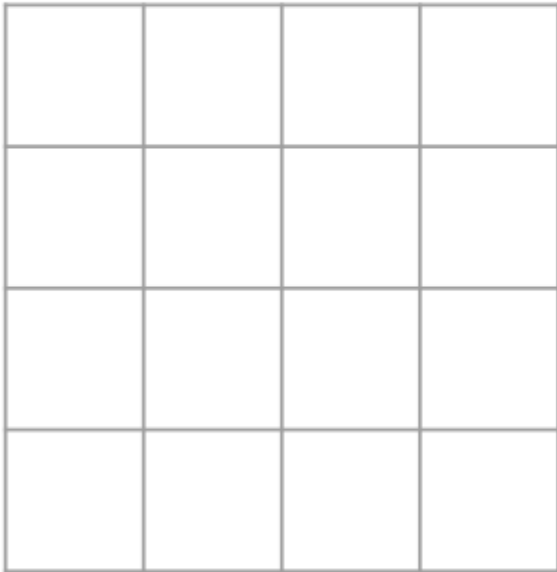
Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers

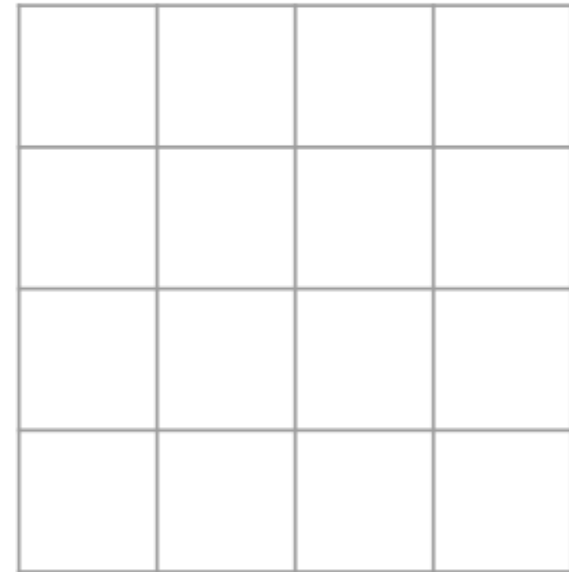


Learnable Upsampling: Transpose Convolution

Typical 3 x 3 convolution, stride 1 pad 1



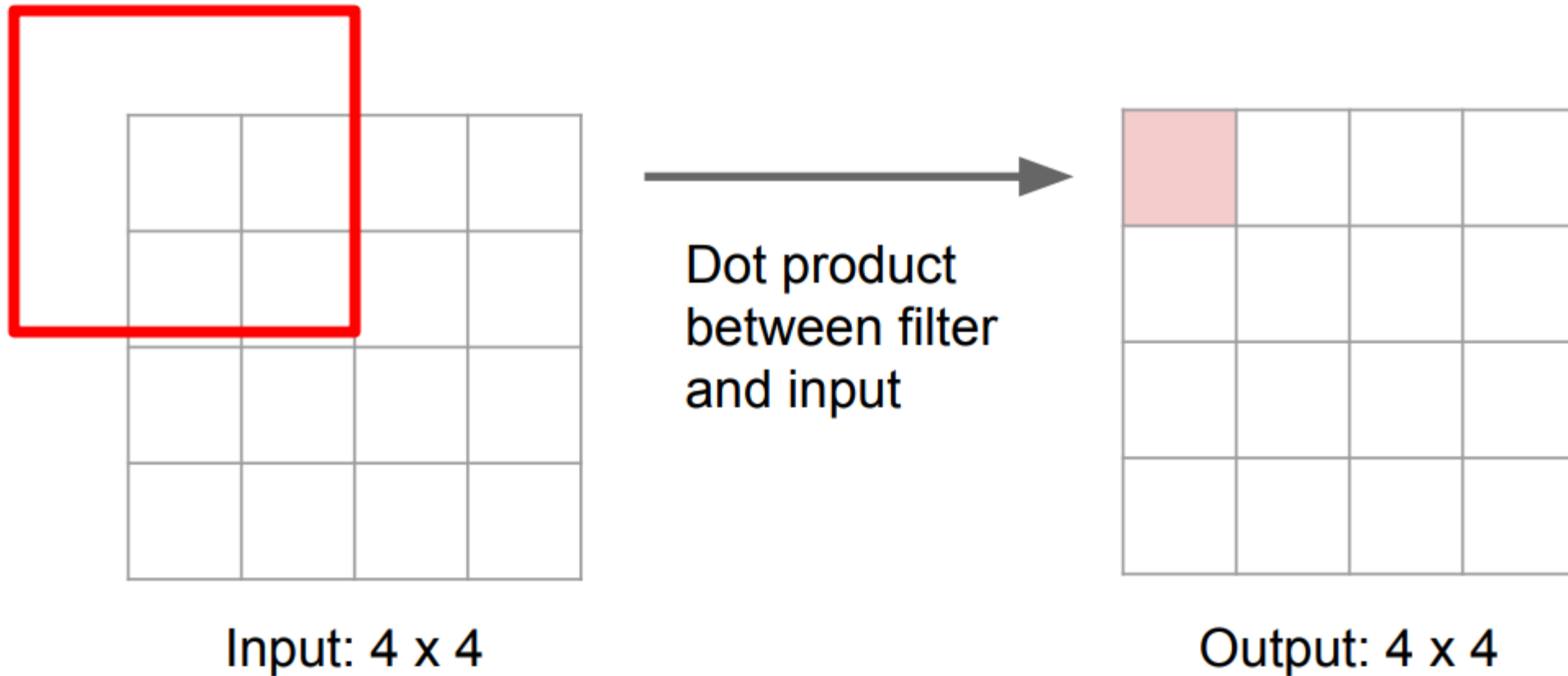
Input: 4 x 4



Output: 4 x 4

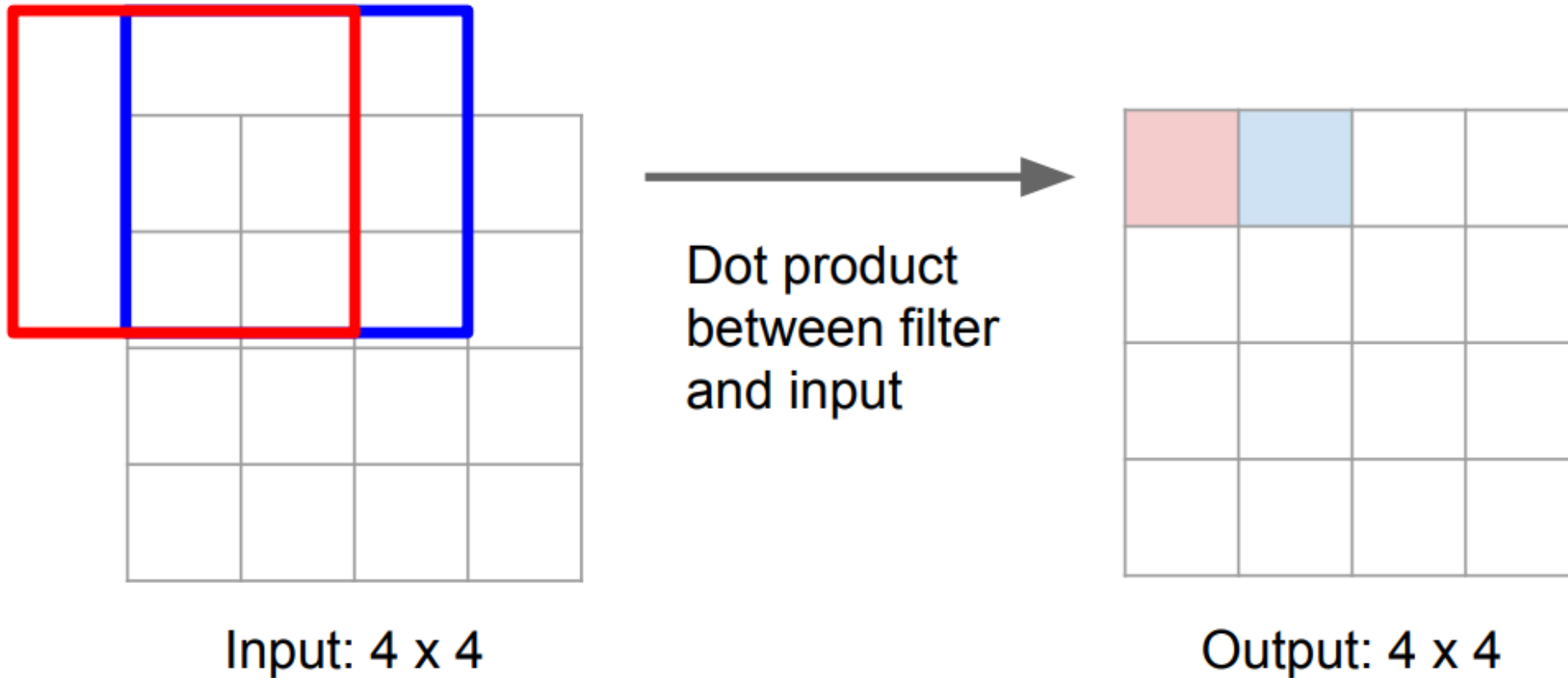
Learnable Upsampling: Transpose Convolution

Normal 3 x 3 convolution, stride 1 pad 1



Learnable Upsampling: Transpose Convolution

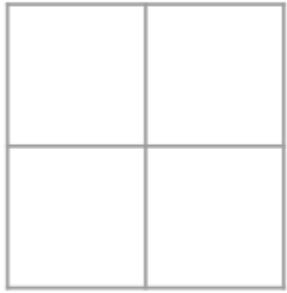
Normal 3 x 3 convolution, stride 1 pad 1



Learnable Upsampling: Transpose Convolution

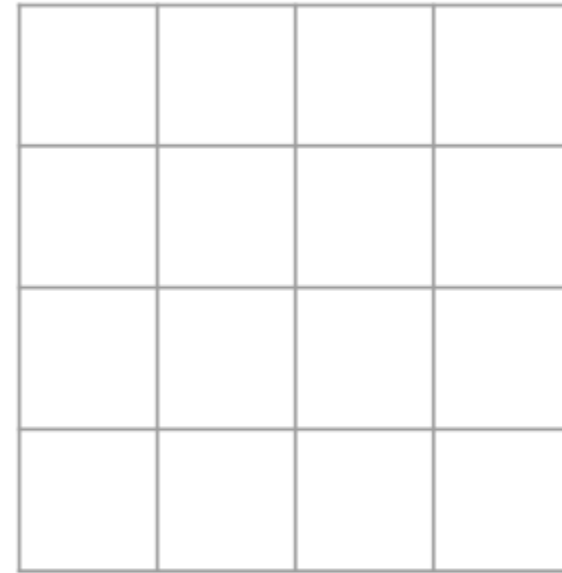
Now the objective is increase the filter size for the decoder to match with the label

3 x 3 transpose convolution, stride 2 pad 1



Input: 2 x 2

Transpose convolution
→

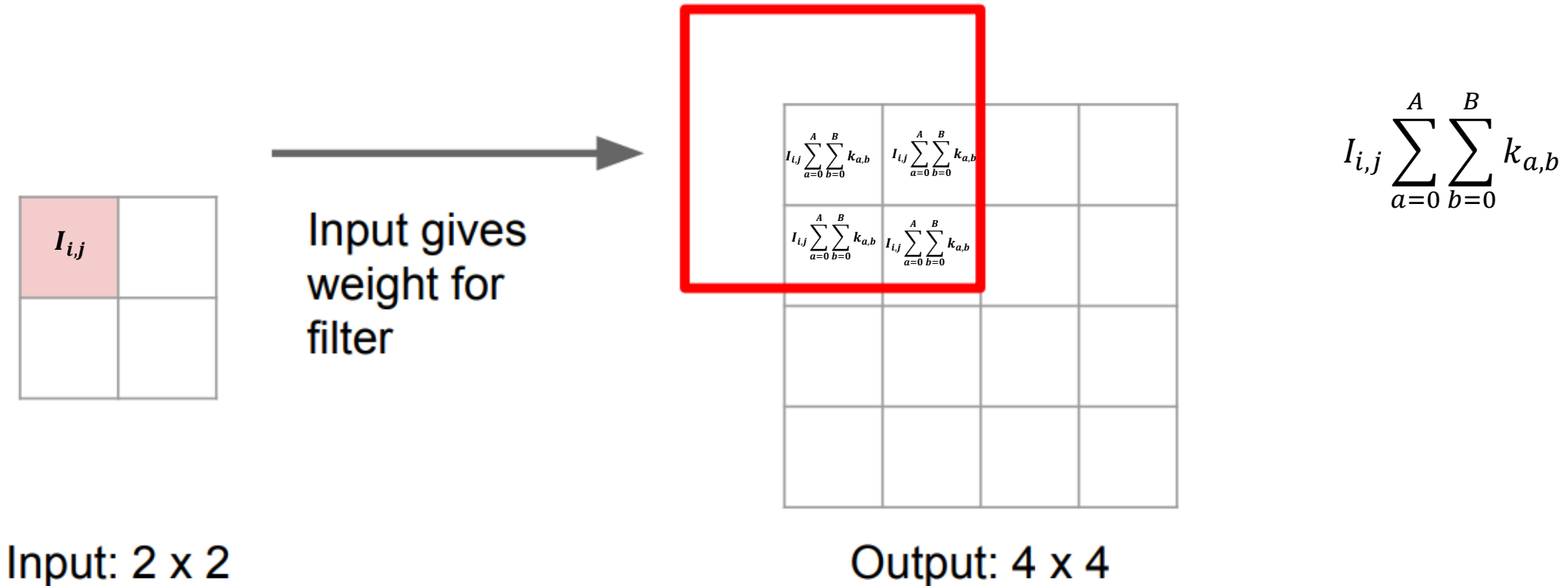


Output: 4 x 4

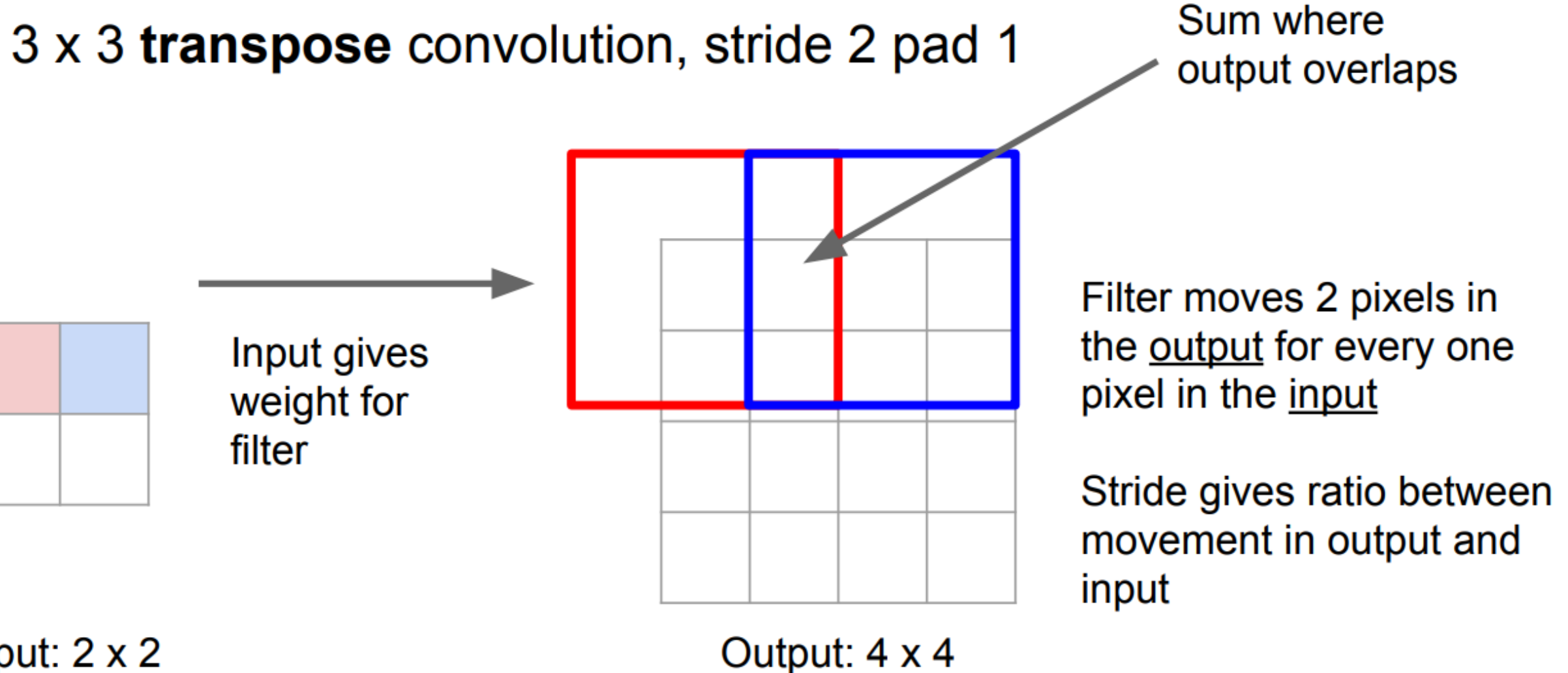
Learnable Upsampling: Transpose Convolution

Now the objective is increase the filter size for the decoder to match with the label

3 x 3 transpose convolution, stride 2 pad 1



Learnable Upsampling: Transpose Convolution

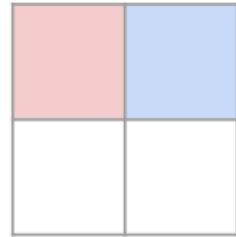


Learnable Upsampling: Transpose Convolution

Other names:

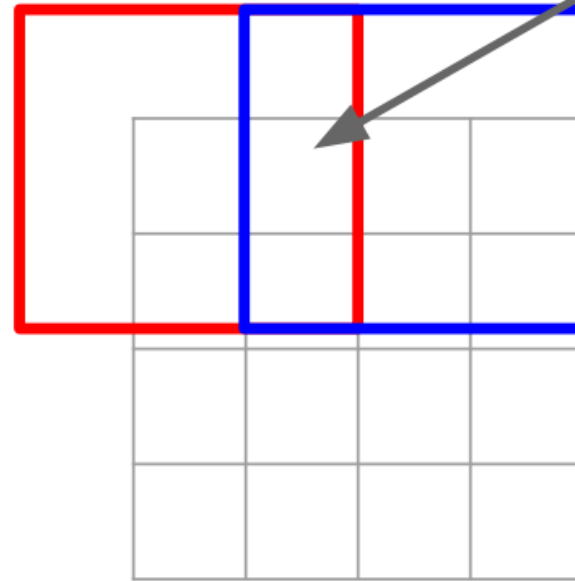
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

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



Input: 2 x 2

Input gives weight for filter



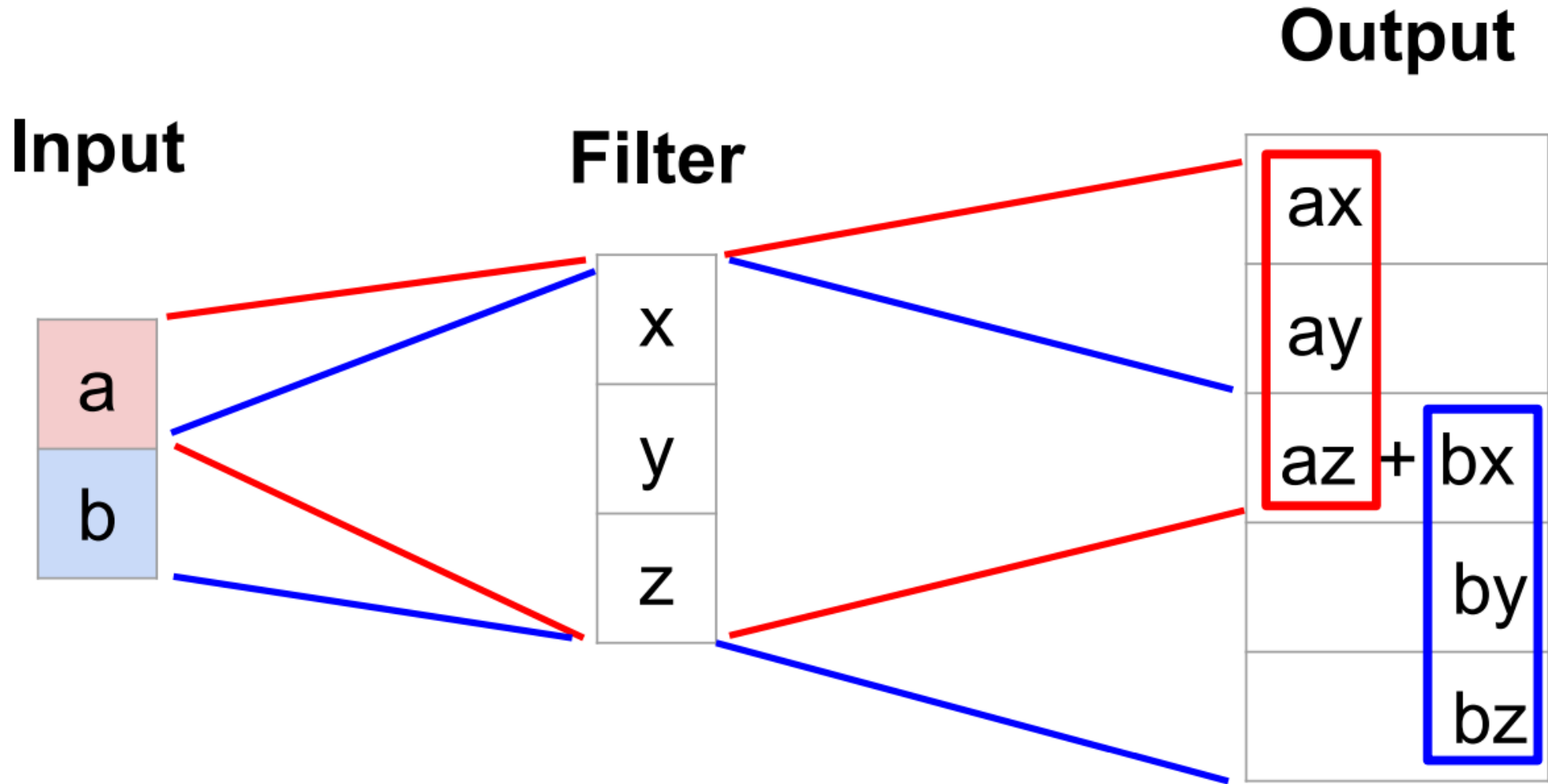
Output: 4 x 4

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

Transpose Convolution: 1D Example



Output contains copies of the filter weighted by the input, summing at where it overlaps in the output

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

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X \vec{a}$$

kernel size = 3

stride=1

Padding=1

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel
size=3, stride=1, padding=1

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X \vec{a}$$

Transpose Convolution



$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

When stride=1, convolution transpose is just a regular convolution (with different padding rules)

Transpose Convolution 1D Example in Segmentation

$$\vec{x} * \vec{a} = X \vec{a}$$

The diagram illustrates a 1D transpose convolution operation. On the left, a 2x6 grid represents the kernel matrix X . The first row contains $x, y, z, 0, 0, 0$ and the second row contains $0, 0, x, y, z, 0$. A blue bracket above the first three columns is labeled "kernel size = 3". A blue arrow below the first two rows is labeled "stride=2". To the right of the grid is a vertical column vector $\vec{a} = \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix}$. A blue arrow points from the label "Padding=1" to the top and bottom zeros of this vector. An equals sign follows the vector, leading to a 2x2 grid representing the output: $\begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$.

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, **stride=2**, padding=1

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X \vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix} = \begin{bmatrix} A \\ B \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Encoder

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} A \\ B \end{bmatrix} = \begin{bmatrix} Ax \\ Ay \\ Az + Ax \\ By \\ Bz \\ 0 \end{bmatrix}$$

Decoder

Semantic Segmentation Idea: Fully Convolutional

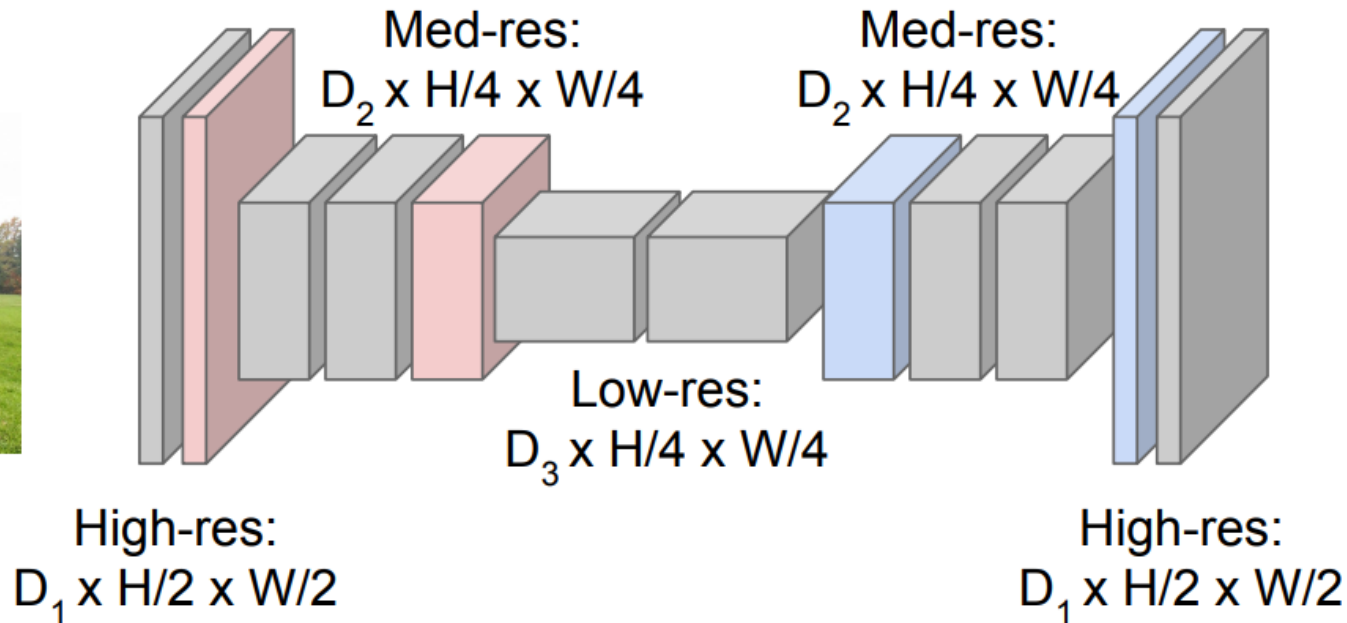
Downsampling:
Pooling, strided
convolution

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

Upsampling:
Unpooling or strided
transpose convolution



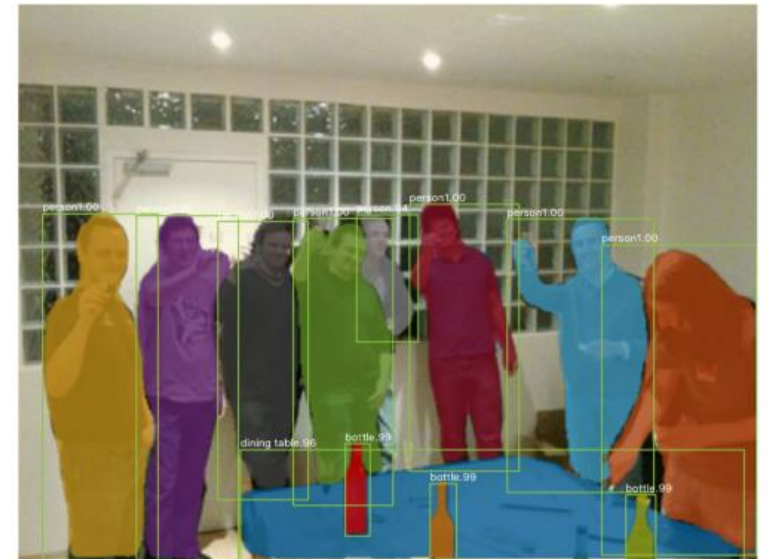
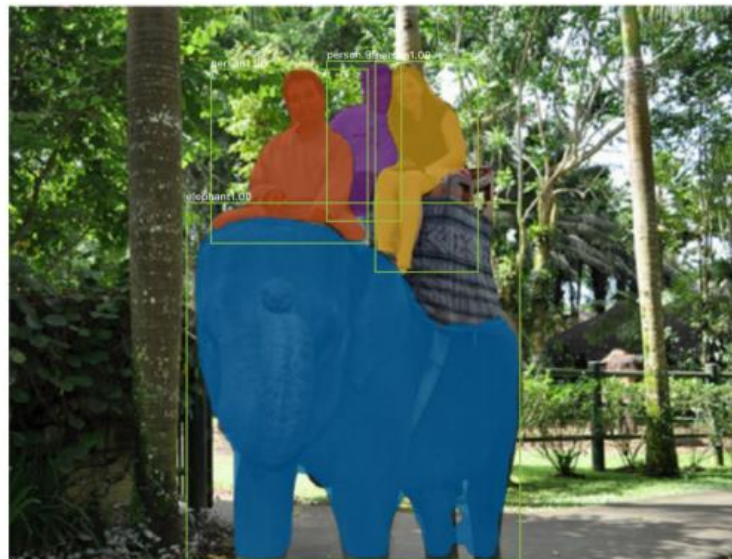
Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Instance Segmentation

- Instance segmentation techniques needs semantic segmentation theories as well as object detection theories
- (object localization) R-CNN → Fast-RCNN → Faster-RCNN → Mask R-CNN (instance segmentation)
- This can be a topic to a next meeting



References

- Stanford Vision : <http://cs231n.stanford.edu/>
- A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, *"A review on deep learning techniques applied to semantic segmentation,"*
- Farabet et al, *"Learning Hierarchical Features for Scene Labeling,"* TPAMI 2013
- Pinheiro and Collobert, *"Recurrent Convolutional Neural Networks for Scene Labeling,"* ICML 2014
- Long, Shelhamer, and Darrell, *"Fully Convolutional Networks for Semantic Segmentation,"* CVPR 2015
- Noh et al, *"Learning Deconvolution Network for Semantic Segmentation,"* ICCV 2015
- K. He, G. Gkioxari, P. Dollar, and R. Girshick. Mask R-² CNN. arXiv:1703.06870, 2017

Thank You