

Semantic Segmentation

Xiaolong Wang

Last Class

- Finetuning with CNN
- The developments and insights of CNN architectures

This Class

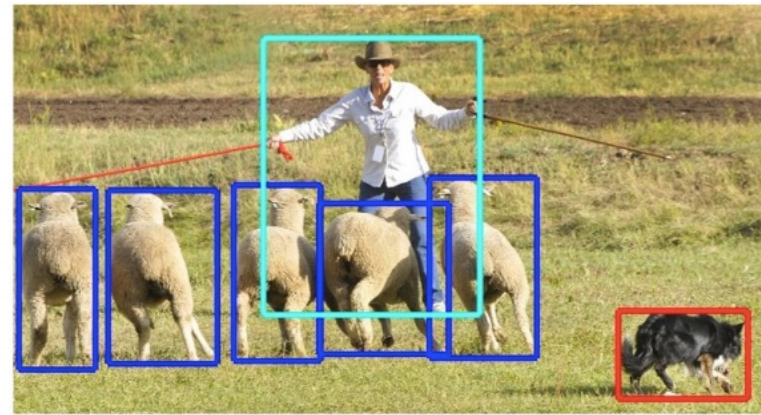
- Naïve FCN model for Image Segmentation
- Transpose Convolution
- Skip Connection, Hypercolumn

Segmentation Problem and FCN

The problem



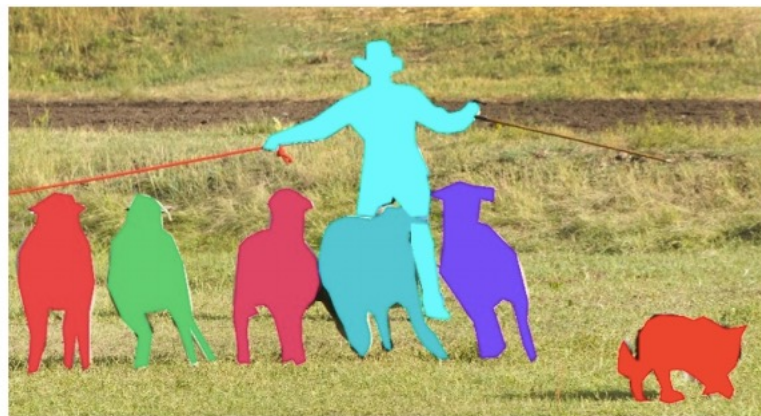
image classification



object detection

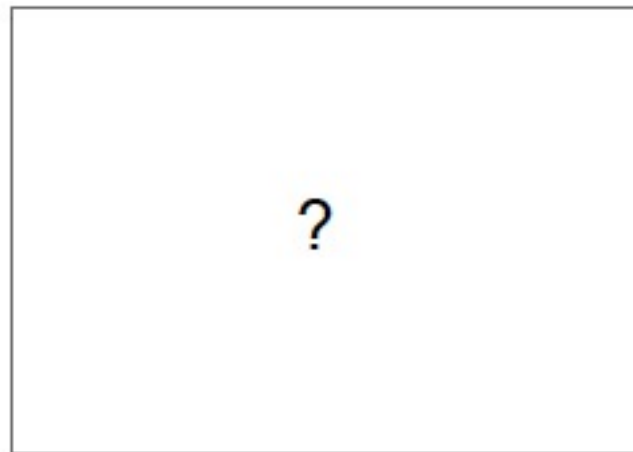
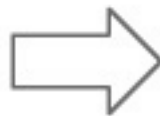


semantic segmentation



instance segmentation

The problem



Semantic Segmentation

Full image

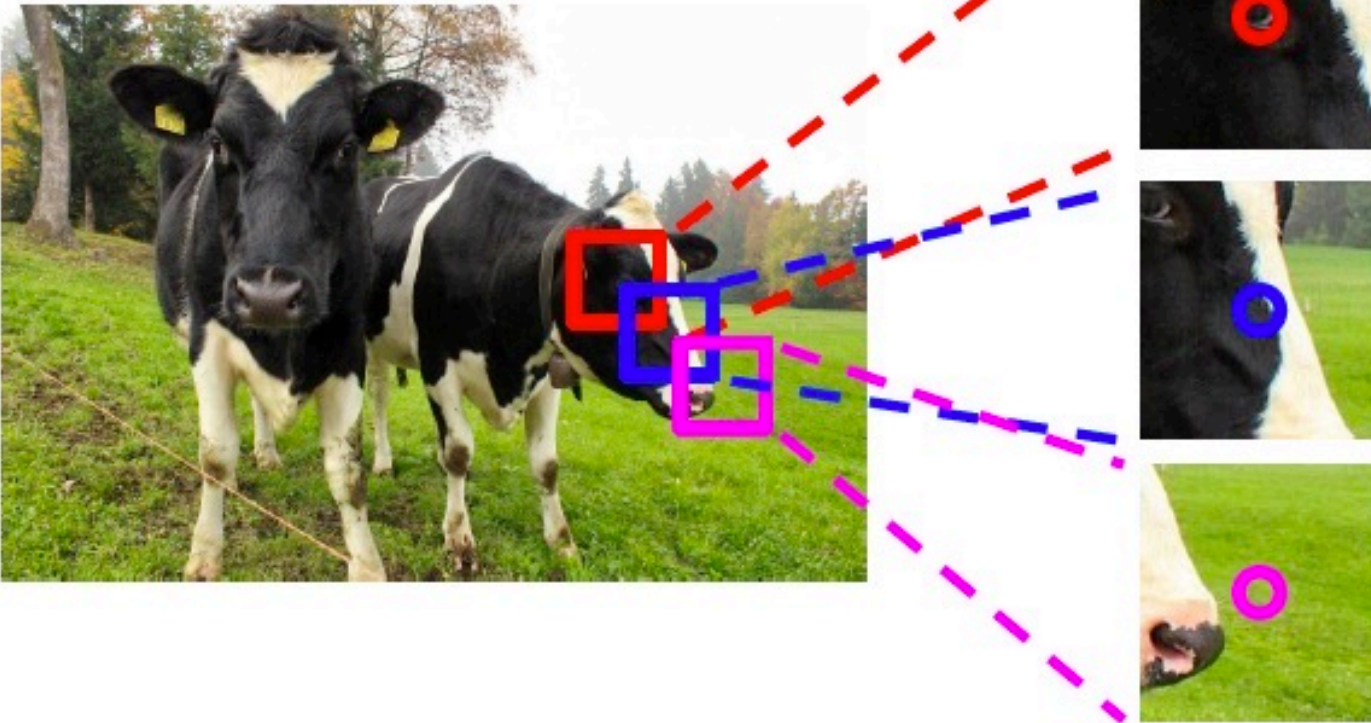


Simply staring one pixel is impossible to do the classification

Let's put in some context!

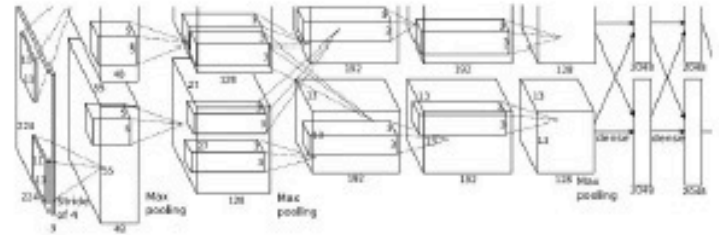
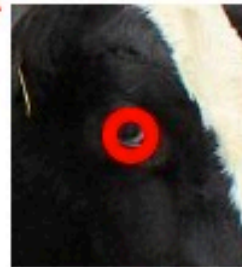
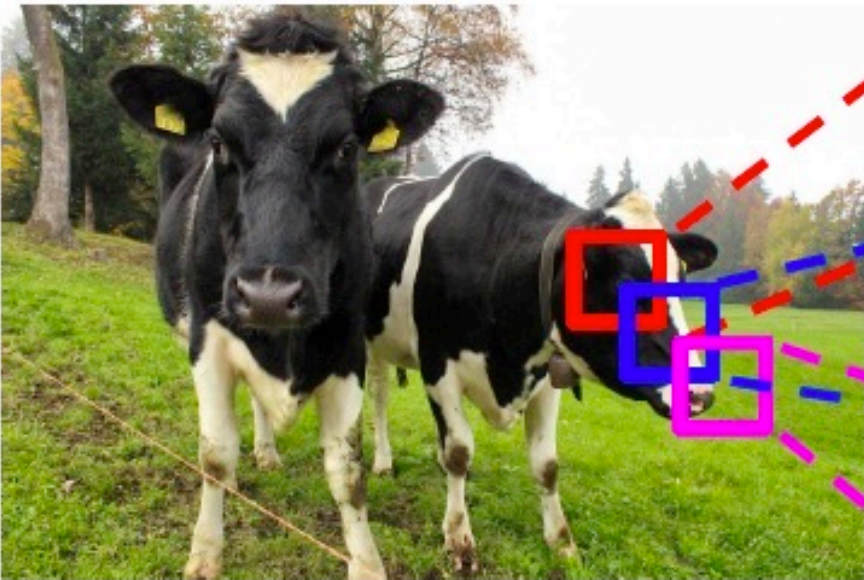
Semantic Segmentation

Full image

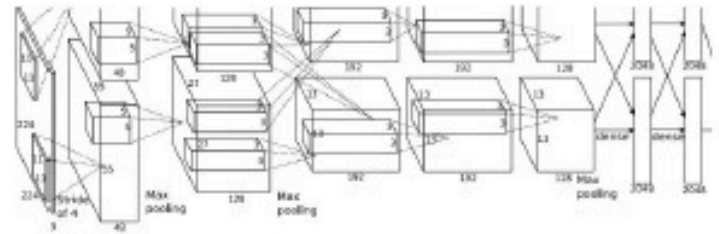
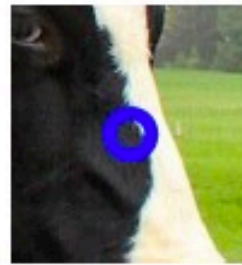


Semantic Segmentation

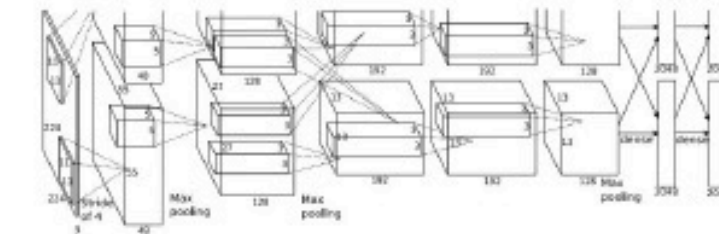
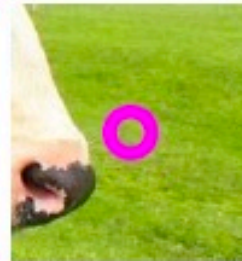
Full image



Cow



Cow

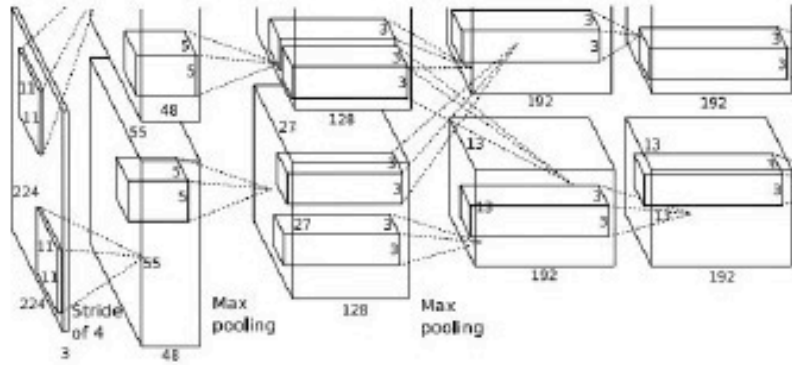


Grass

Time Consuming!

Can we process the whole image at one time?

Full image



AlexNet input:
 $227 \times 277 \times 3$

AlexNet Conv5:
 $13 \times 13 \times 128$

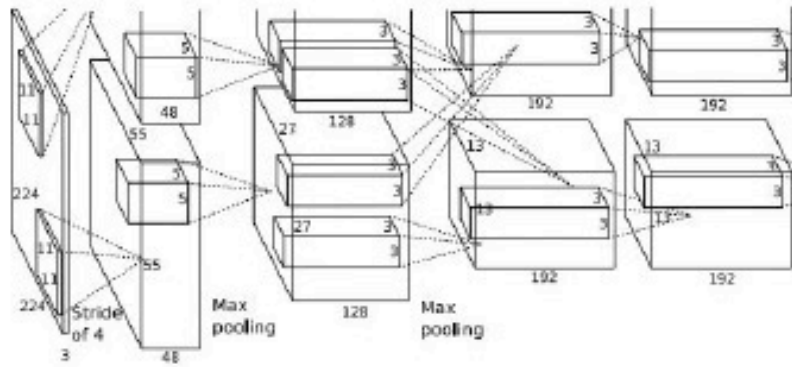


Output ?
 $13 \times 13 \times 21$

Output is too small!

Can we process the whole image at one time?

Full image



AlexNet input:
 $227 \times 277 \times 3$

AlexNet Conv5:
 $13 \times 13 \times 128$

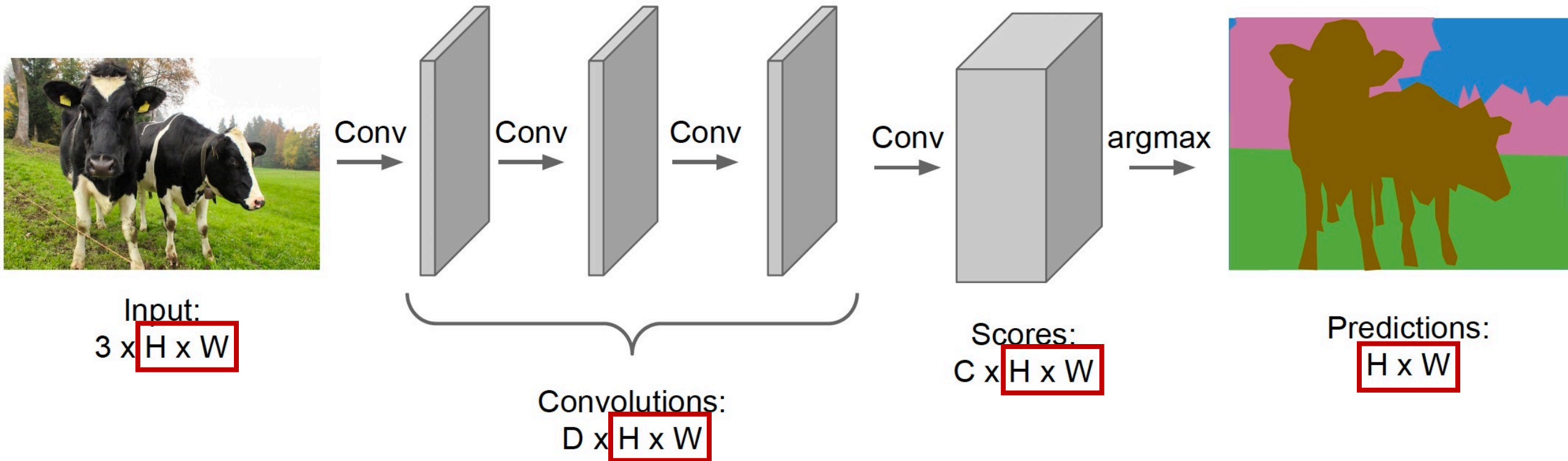
FC +
reshape



Output ?
 $227 \times 227 \times 21$

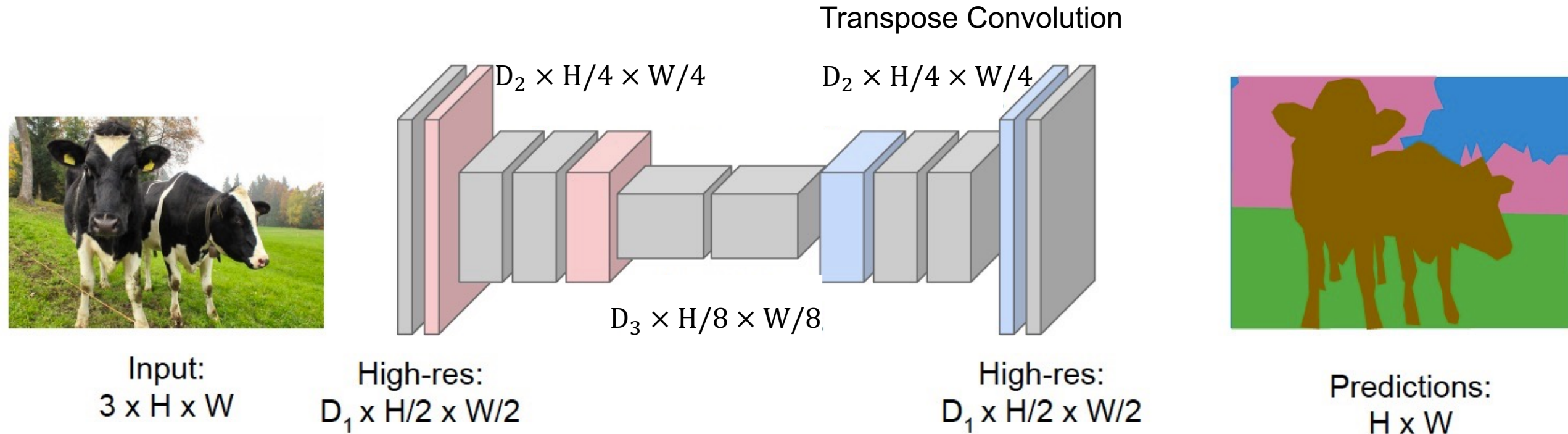
Huge number of Parameters for FC: $13 \times 13 \times 128 \times 227 \times 227 \times 21$

Fully Convolutional Network



Convolution at original image resolution has high computation cost.

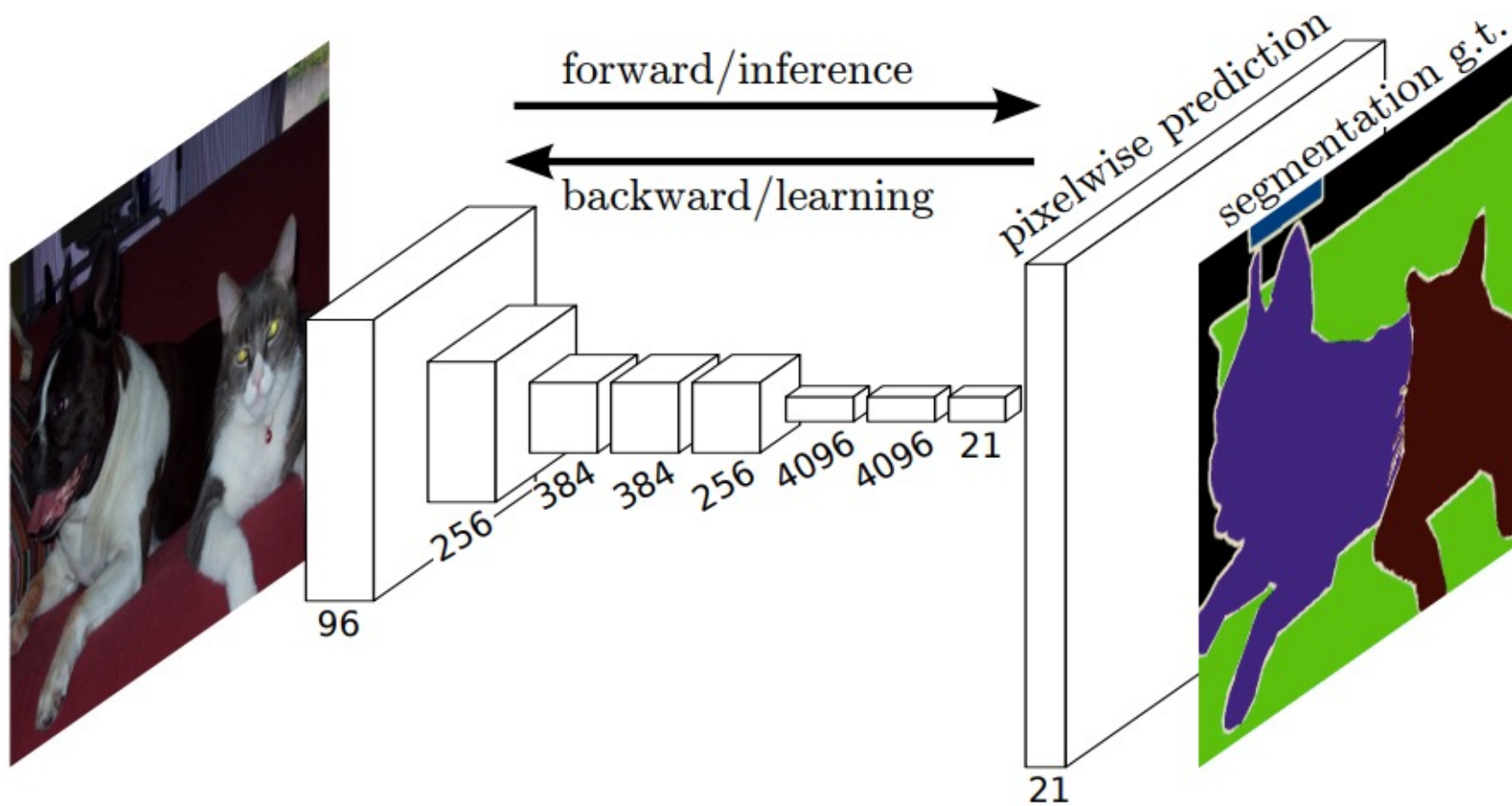
Fully Convolutional Network



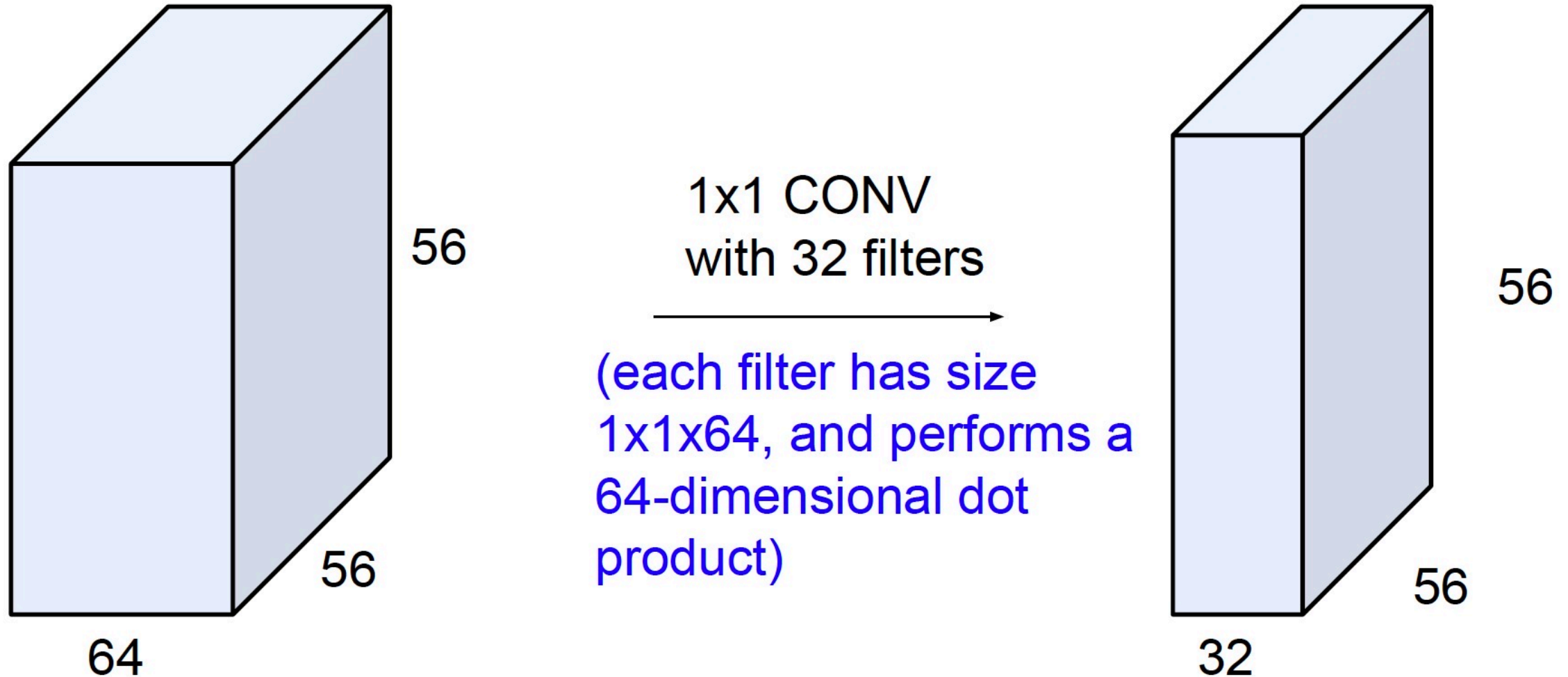
Make the feature map small increases the receptive field

Make the feature map larger again increases the resolution

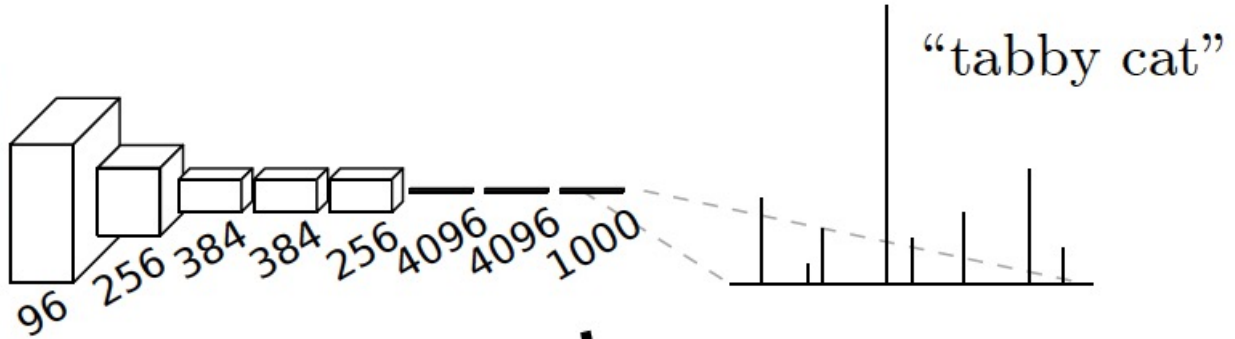
The FCN paper



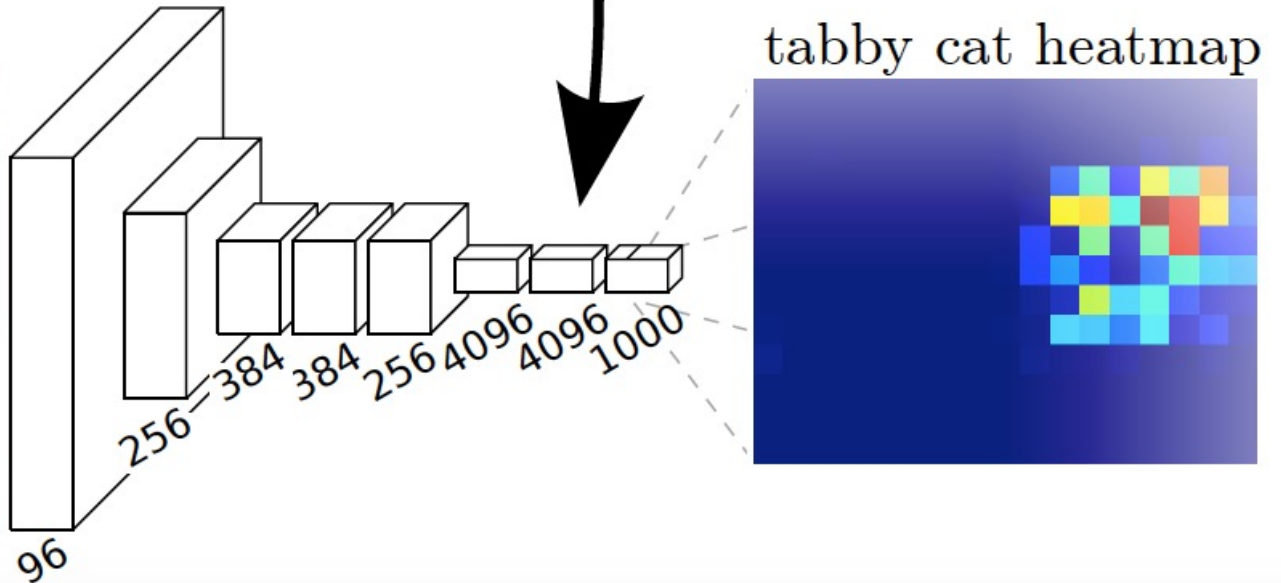
1 x 1 convolutions



The FCN paper



convolutionalization



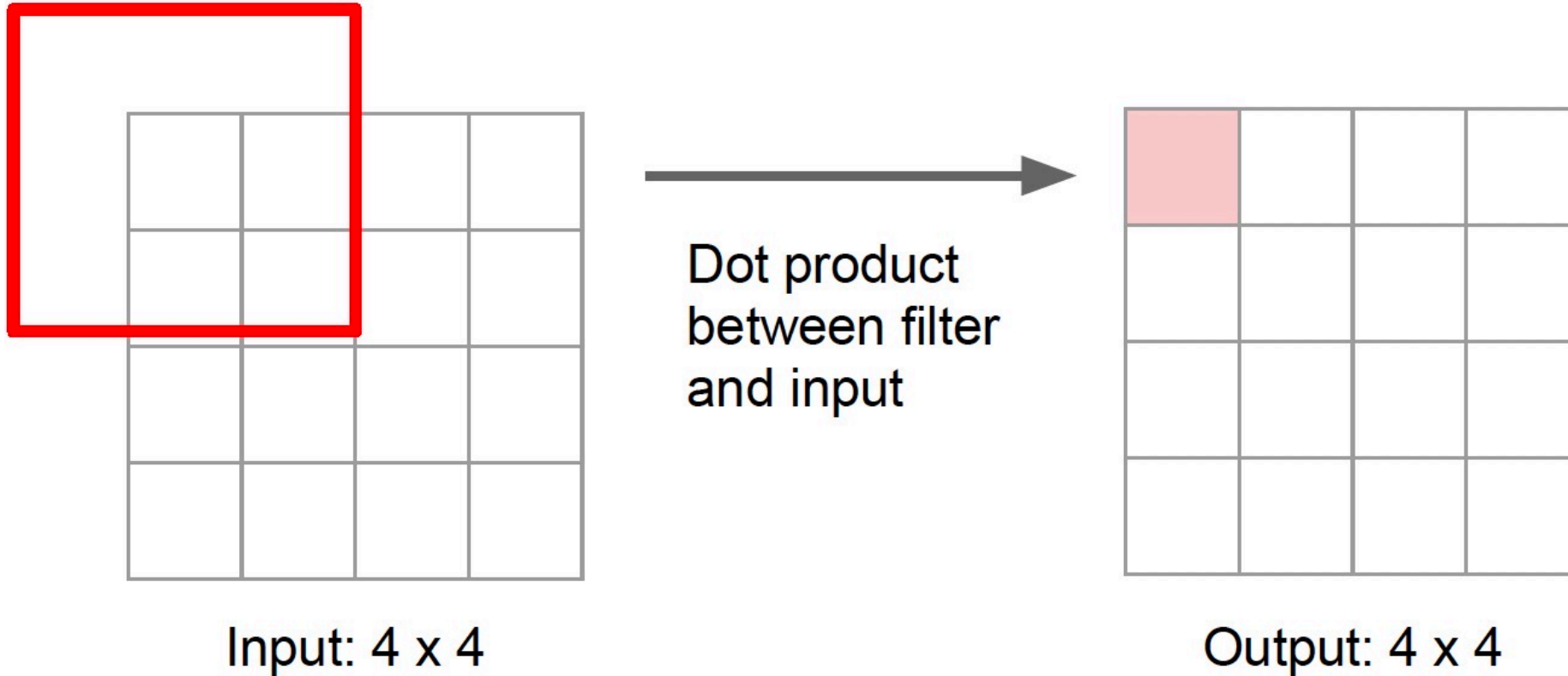
The upsampling

- Upsampling Layer
- Deconvolution Layer
- Transpose Convolution Layer

Transpose Convolution

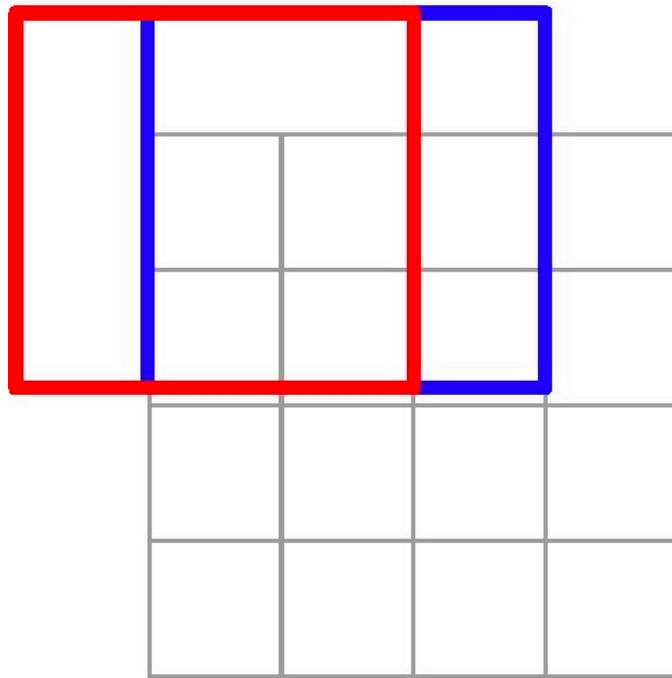
Recall Convolution

3 X 3 convolution with stride 1 and padding 1



Recall Convolution

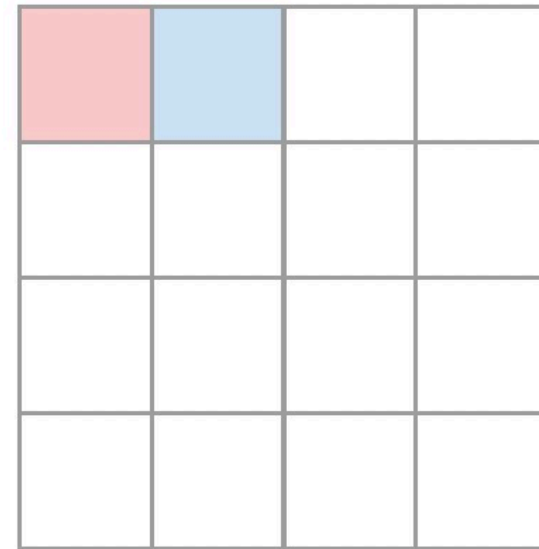
3 X 3 convolution with stride 1 and padding 1



Input: 4 x 4



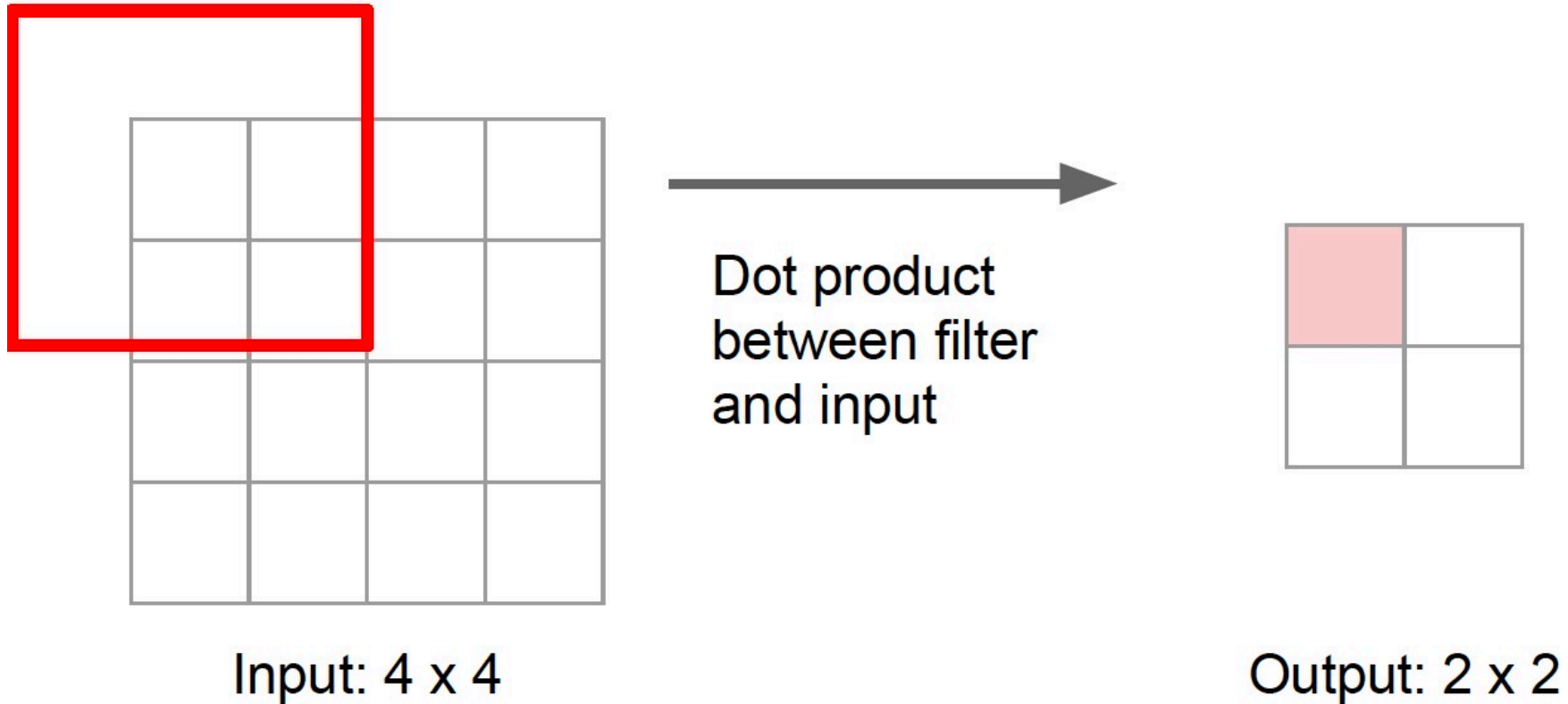
Dot product
between filter
and input



Output: 4 x 4

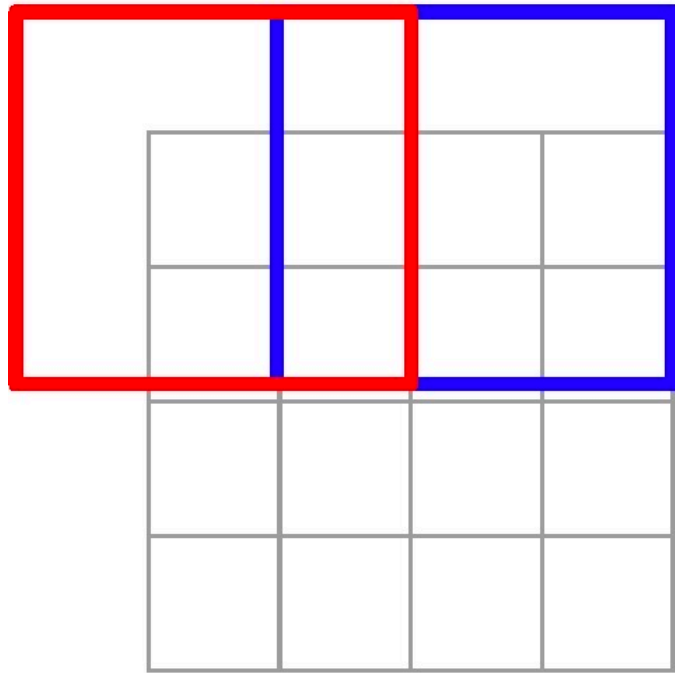
Recall Convolution

3 X 3 convolution with stride 2 and padding 1



Recall Convolution

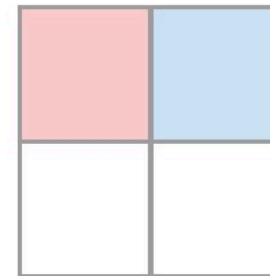
3 X 3 convolution with stride 2 and padding 1



Input: 4 x 4



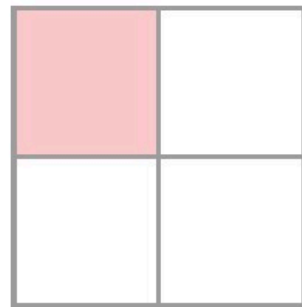
Dot product
between filter
and input



Output: 2 x 2

Transpose Convolution

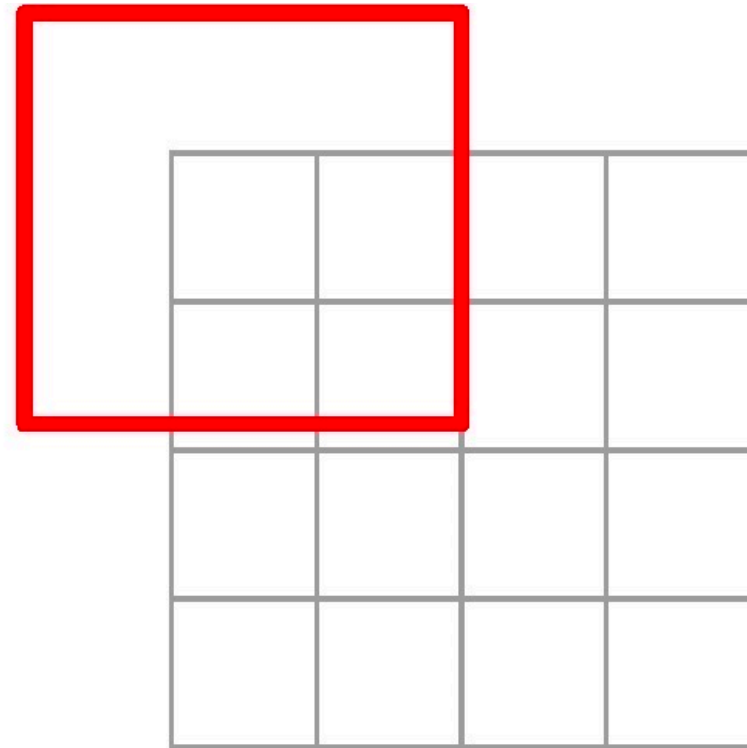
3 X 3 transpose convolution, stride 2 and padding 1



Input: 2 x 2



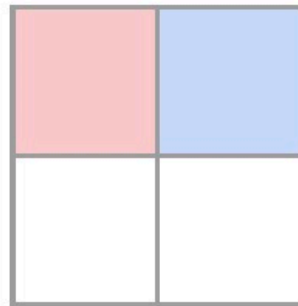
Input gives
weight for
filter



Output: 4 x 4

Transpose Convolution

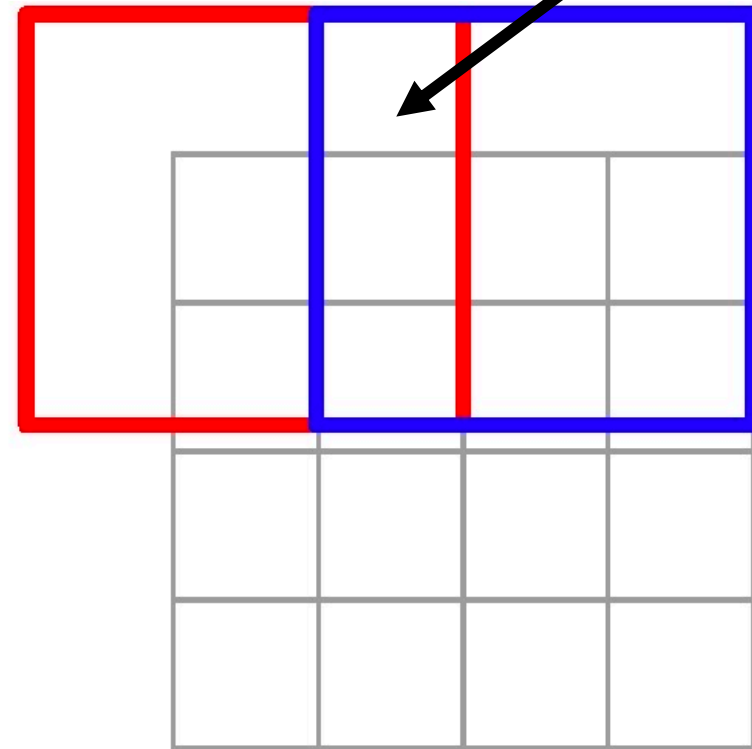
3 X 3 transpose convolution, stride 2 and padding 1



Input: 2 x 2



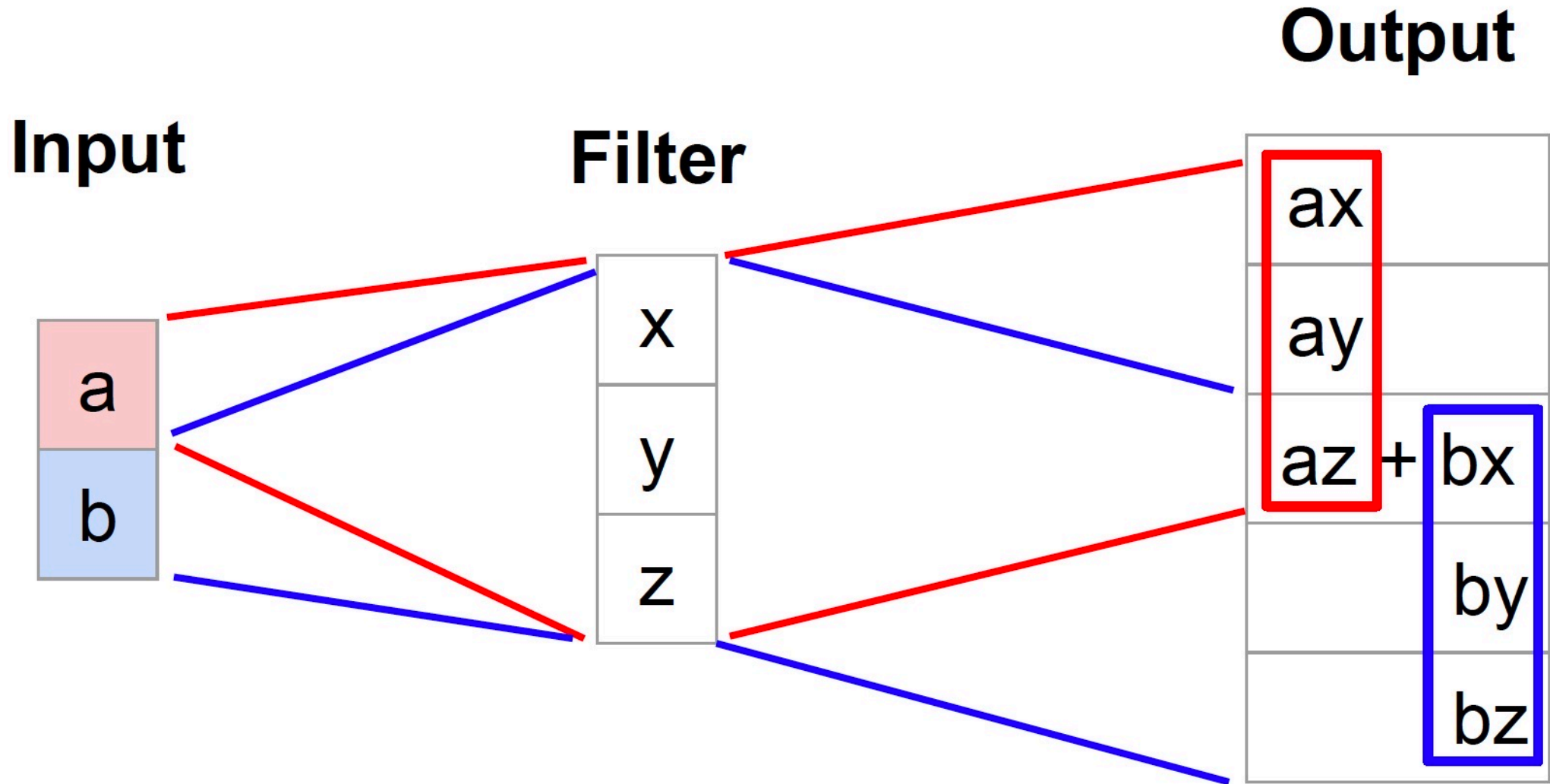
Input gives weight for filter



Output: 4 x 4

Sum over the overlapping region

1D Transpose Convolution Example



2D Convolution

Regular convolution (stride 1, pad 0)

$$\begin{array}{|c|c|c|c|} \hline x_{11} & x_{12} & x_{13} & x_{14} \\ \hline x_{21} & x_{22} & x_{23} & x_{24} \\ \hline x_{31} & x_{32} & x_{33} & x_{34} \\ \hline x_{41} & x_{42} & x_{43} & x_{44} \\ \hline \end{array}
 \quad * \quad
 \begin{array}{|c|c|c|} \hline w_{11} & w_{12} & w_{13} \\ \hline w_{21} & w_{22} & w_{23} \\ \hline w_{31} & w_{32} & w_{33} \\ \hline \end{array}
 =
 \begin{array}{|c|c|} \hline z_{11} & z_{12} \\ \hline z_{21} & z_{22} \\ \hline \end{array}$$

Matrix-vector form:

$$\begin{pmatrix} w_{11} & w_{12} & w_{13} & 0 & w_{21} & w_{22} & w_{23} & 0 & w_{31} & w_{32} & w_{33} & 0 & 0 & 0 & 0 & 0 \\ 0 & w_{11} & w_{12} & w_{13} & 0 & w_{21} & w_{22} & w_{23} & 0 & w_{31} & w_{32} & w_{33} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & w_{11} & w_{12} & w_{13} & 0 & w_{21} & w_{22} & w_{23} & 0 & w_{31} & w_{32} & w_{33} & 0 \\ 0 & 0 & 0 & 0 & 0 & w_{11} & w_{12} & w_{13} & 0 & w_{21} & w_{22} & w_{23} & 0 & w_{31} & w_{32} & w_{33} \end{pmatrix}
 \begin{pmatrix} x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \\ \vdots \\ x_{44} \end{pmatrix}
 =
 \begin{pmatrix} z_{11} \\ z_{12} \\ z_{21} \\ z_{22} \end{pmatrix}$$

4x4 input, 2x2 output

Transpose Convolution

z_{11}	z_{12}
z_{21}	z_{22}

$*T$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

=

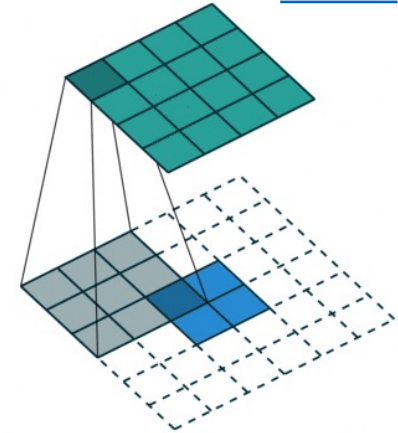
x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}

$$\begin{pmatrix}
 w_{11} & 0 & 0 & 0 \\
 w_{12} & w_{11} & 0 & 0 \\
 w_{13} & w_{12} & 0 & 0 \\
 0 & w_{13} & 0 & 0 \\
 w_{21} & 0 & w_{11} & 0 \\
 w_{22} & w_{21} & w_{12} & w_{11} \\
 w_{23} & w_{22} & w_{13} & w_{12} \\
 0 & w_{23} & 0 & w_{13} \\
 w_{31} & 0 & w_{21} & 0 \\
 w_{32} & w_{31} & w_{22} & w_{21} \\
 w_{33} & w_{32} & w_{23} & w_{22} \\
 0 & w_{33} & 0 & w_{23} \\
 0 & 0 & w_{31} & 0 \\
 0 & 0 & w_{32} & w_{31} \\
 0 & 0 & w_{33} & w_{32} \\
 0 & 0 & 0 & w_{33}
 \end{pmatrix}
 \begin{pmatrix}
 z_{11} \\
 z_{12} \\
 z_{21} \\
 z_{22}
 \end{pmatrix}
 =
 \begin{pmatrix}
 x_{11} \\
 x_{12} \\
 x_{13} \\
 x_{14} \\
 x_{21} \\
 x_{22} \\
 x_{23} \\
 x_{24} \\
 x_{31} \\
 x_{32} \\
 x_{33} \\
 x_{34} \\
 x_{41} \\
 x_{42} \\
 x_{43} \\
 x_{44}
 \end{pmatrix}$$

2x2 input, 4x4 output

Not an inverse of the original convolution operation, simply reverses dimension change!

[Source](#)



Trans Convolution

w_{33}	w_{32}	w_{31}
w_{23}	w_{22}	w_{21}
w_{13}	w_{12}	w_{11}

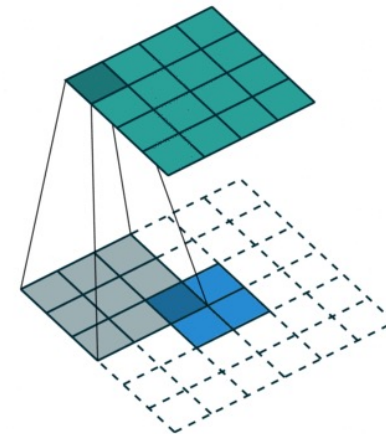
z_{11}	z_{12}
z_{21}	z_{22}

$*^T$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

=

x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}



w_{11}	0	0	0
w_{12}	w_{11}	0	0
w_{13}	w_{12}	0	0
0	w_{13}	0	0
w_{21}	0	w_{11}	0
w_{22}	w_{21}	w_{12}	w_{11}
w_{23}	w_{22}	w_{13}	w_{12}
0	w_{23}	0	w_{13}
w_{31}	0	w_{21}	0
w_{32}	w_{31}	w_{22}	w_{21}
w_{33}	w_{32}	w_{23}	w_{22}
0	w_{33}	0	w_{23}
0	0	w_{31}	0
0	0	w_{32}	w_{31}
0	0	w_{33}	w_{32}
0	0	0	w_{33}

z_{11}
z_{12}
z_{21}
z_{22}

=

x_{11}
x_{12}
x_{13}
x_{14}
x_{21}
x_{22}
x_{23}
x_{24}
x_{31}
x_{32}
x_{33}
x_{34}
x_{41}
x_{42}
x_{43}
x_{44}

$x_{11} = w_{11}z_{11}$

Transposed Convolution

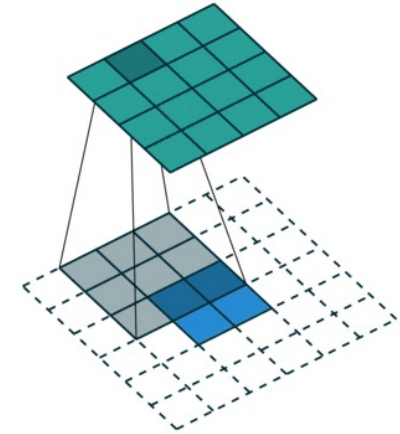
w_{33}	w_{32}	w_{31}
w_{23}	w_{22}	w_{21}
w_{13}	w_{12}	w_{11}
z_{11}	z_{12}	
z_{21}	z_{22}	

$*T$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

=

x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}

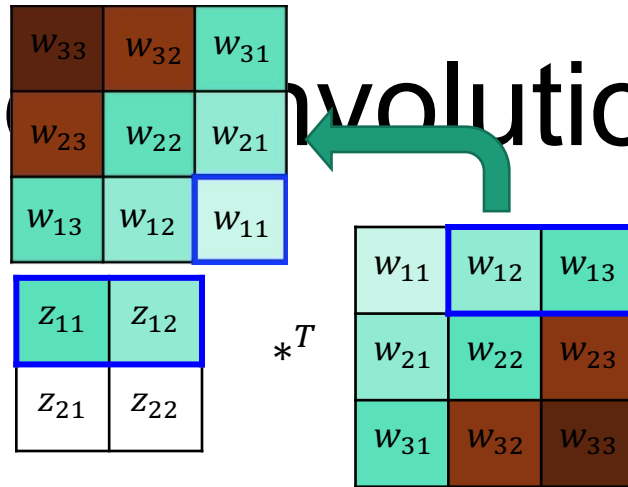


Convolve input with *flipped* filter

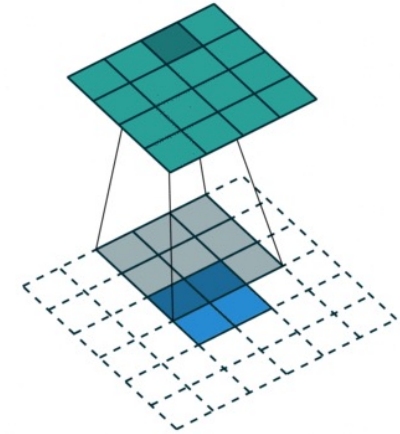
$$x_{12} = w_{12}z_{11} + w_{11}z_{12}$$

w_{11}	0	0	0	$\left(\begin{array}{c} z_{11} \\ z_{12} \\ z_{21} \\ z_{22} \end{array} \right) =$
w_{12}	w_{11}	0	0	
w_{13}	w_{12}	0	0	
0	w_{13}	0	0	
w_{21}	0	w_{11}	0	
w_{22}	w_{21}	w_{12}	w_{11}	
w_{23}	w_{22}	w_{13}	w_{12}	
0	w_{23}	0	w_{13}	
w_{31}	0	w_{21}	0	
w_{32}	w_{31}	w_{22}	w_{21}	
w_{33}	w_{32}	w_{23}	w_{22}	
0	w_{33}	0	w_{23}	
0	0	w_{31}	0	
0	0	w_{32}	w_{31}	
0	0	w_{33}	w_{32}	
0	0	0	w_{33}	

Transposed Convolution



x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}

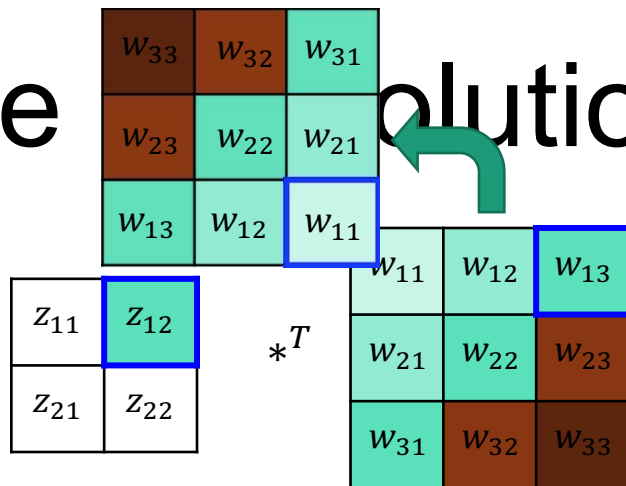


Convolve input with *flipped* filter

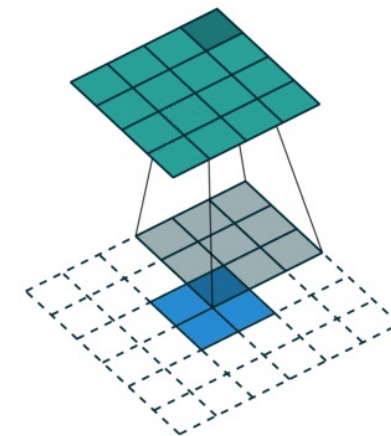
$$\begin{pmatrix}
 w_{11} & 0 & 0 & 0 \\
 w_{12} & w_{11} & 0 & 0 \\
 w_{13} & w_{12} & 0 & 0 \\
 0 & w_{13} & 0 & 0 \\
 w_{21} & 0 & w_{11} & 0 \\
 w_{22} & w_{21} & w_{12} & w_{11} \\
 w_{23} & w_{22} & w_{13} & w_{12} \\
 0 & w_{23} & 0 & w_{13} \\
 w_{31} & 0 & w_{21} & 0 \\
 w_{32} & w_{31} & w_{22} & w_{21} \\
 w_{33} & w_{32} & w_{23} & w_{22} \\
 0 & w_{33} & 0 & w_{23} \\
 0 & 0 & w_{31} & 0 \\
 0 & 0 & w_{32} & w_{31} \\
 0 & 0 & w_{33} & w_{32} \\
 0 & 0 & 0 & w_{33}
 \end{pmatrix}
 \begin{pmatrix}
 z_{11} \\
 z_{12} \\
 z_{21} \\
 z_{22}
 \end{pmatrix}
 =
 \begin{pmatrix}
 x_{11} \\
 x_{12} \\
 x_{13} \\
 x_{14} \\
 x_{21} \\
 x_{22} \\
 x_{23} \\
 x_{24} \\
 x_{31} \\
 x_{32} \\
 x_{33} \\
 x_{34} \\
 x_{41} \\
 x_{42} \\
 x_{43} \\
 x_{44}
 \end{pmatrix}$$

$x_{13} = w_{13}z_{11} + w_{12}z_{12}$

Transpose Convolution



x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}



Convolve input with *flipped* filter

$$\begin{pmatrix}
 w_{11} & 0 & 0 & 0 \\
 w_{12} & w_{11} & 0 & 0 \\
 w_{13} & w_{12} & 0 & 0 \\
 0 & w_{13} & 0 & 0 \\
 w_{21} & 0 & w_{11} & 0 \\
 w_{22} & w_{21} & w_{12} & w_{11} \\
 w_{23} & w_{22} & w_{13} & w_{12} \\
 0 & w_{23} & 0 & w_{13} \\
 w_{31} & 0 & w_{21} & 0 \\
 w_{32} & w_{31} & w_{22} & w_{21} \\
 w_{33} & w_{32} & w_{23} & w_{22} \\
 0 & w_{33} & 0 & w_{23} \\
 0 & 0 & w_{31} & 0 \\
 0 & 0 & w_{32} & w_{31} \\
 0 & 0 & w_{33} & w_{32} \\
 0 & 0 & 0 & w_{33}
 \end{pmatrix}
 \begin{pmatrix}
 z_{11} \\
 z_{12} \\
 z_{21} \\
 z_{22}
 \end{pmatrix}
 =
 \begin{pmatrix}
 x_{11} \\
 x_{12} \\
 x_{13} \\
 x_{14} \\
 x_{21} \\
 x_{22} \\
 x_{23} \\
 x_{24} \\
 x_{31} \\
 x_{32} \\
 x_{33} \\
 x_{34} \\
 x_{41} \\
 x_{42} \\
 x_{43} \\
 x_{44}
 \end{pmatrix}$$

$x_{14} = w_{13}z_{12}$

Trans Convolution

w_{33}	w_{32}	w_{31}
w_{23}	w_{22}	w_{21}
w_{13}	w_{12}	w_{11}

z_{11}	z_{12}
z_{21}	z_{22}

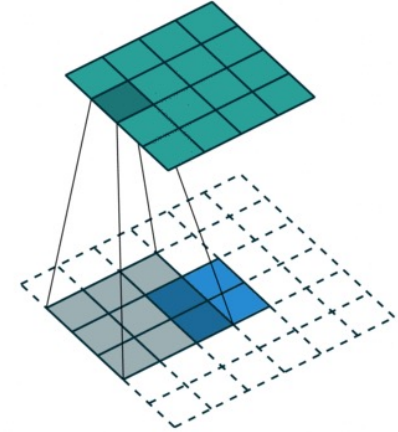
$*^T$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

=

x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}

Convolve input with *flipped* filter



w_{11}	0	0	0
w_{12}	w_{11}	0	0
w_{13}	w_{12}	0	0
0	w_{13}	0	0
w_{21}	0	w_{11}	0
w_{22}	w_{21}	w_{12}	w_{11}
w_{23}	w_{22}	w_{13}	w_{12}
0	w_{23}	0	w_{13}
w_{31}	0	w_{21}	0
w_{32}	w_{31}	w_{22}	w_{21}
w_{33}	w_{32}	w_{23}	w_{22}
0	w_{33}	0	w_{23}
0	0	w_{31}	0
0	0	w_{32}	w_{31}
0	0	w_{33}	w_{32}
0	0	0	w_{33}

z_{11}
z_{12}
z_{21}
z_{22}

=

x_{11}
x_{12}
x_{13}
x_{14}
x_{21}
x_{22}
x_{23}
x_{24}
x_{31}
x_{32}
x_{33}
x_{34}
x_{41}
x_{42}
x_{43}
x_{44}

$$x_{21} = w_{21}z_{11} + w_{11}z_{21}$$

Transposed convolution

w_{33}	w_{32}	w_{31}
w_{23}	w_{22}	w_{21}
w_{13}	w_{12}	w_{11}

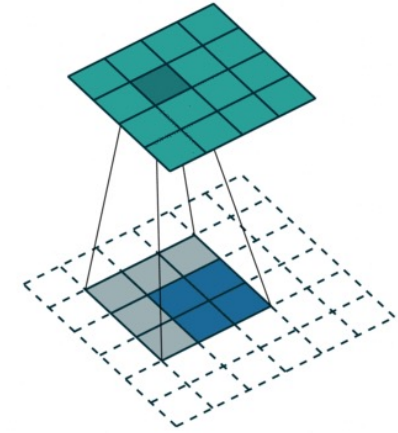
z_{11}	z_{12}
z_{21}	z_{22}

$*T$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

=

x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}



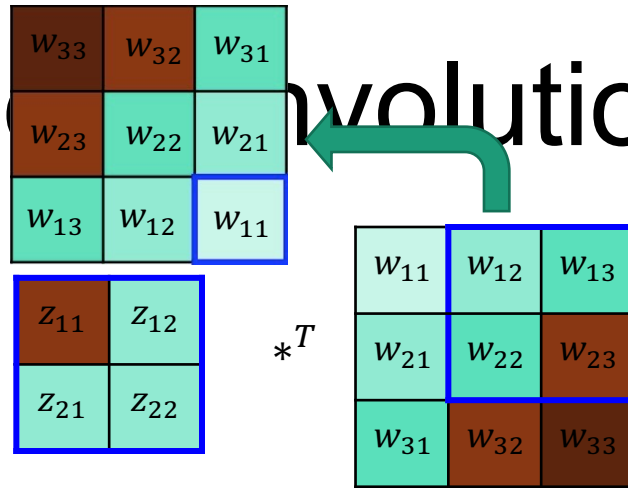
Convolve input with *flipped* filter

$$\begin{pmatrix}
 w_{11} & 0 & 0 & 0 \\
 w_{12} & w_{11} & 0 & 0 \\
 w_{13} & w_{12} & 0 & 0 \\
 0 & w_{13} & 0 & 0 \\
 w_{21} & 0 & w_{11} & 0 \\
 \mathbf{w_{22} & w_{21} & w_{12} & w_{11}} \\
 w_{23} & w_{22} & w_{13} & w_{12} \\
 0 & w_{23} & 0 & w_{13} \\
 w_{31} & 0 & w_{21} & 0 \\
 w_{32} & w_{31} & w_{22} & w_{21} \\
 w_{33} & w_{32} & w_{23} & w_{22} \\
 0 & w_{33} & 0 & w_{23} \\
 0 & 0 & w_{31} & 0 \\
 0 & 0 & w_{32} & w_{31} \\
 0 & 0 & w_{33} & w_{32} \\
 0 & 0 & 0 & w_{33}
 \end{pmatrix}
 \begin{pmatrix}
 z_{11} \\
 z_{12} \\
 z_{21} \\
 z_{22}
 \end{pmatrix}
 =$$

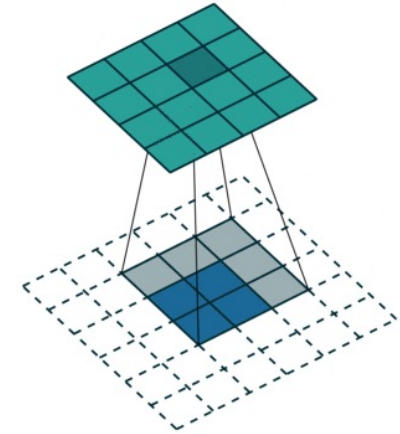
$$\begin{pmatrix}
 x_{11} \\
 x_{12} \\
 x_{13} \\
 x_{14} \\
 x_{21} \\
 \mathbf{x_{22}} \\
 x_{23} \\
 x_{24} \\
 x_{31} \\
 x_{32} \\
 x_{33} \\
 x_{34} \\
 x_{41} \\
 x_{42} \\
 x_{43} \\
 x_{44}
 \end{pmatrix}$$

$$x_{22} = w_{22}z_{11} + w_{21}z_{12} + w_{12}z_{21} + w_{11}z_{22}$$

Transposed Convolution



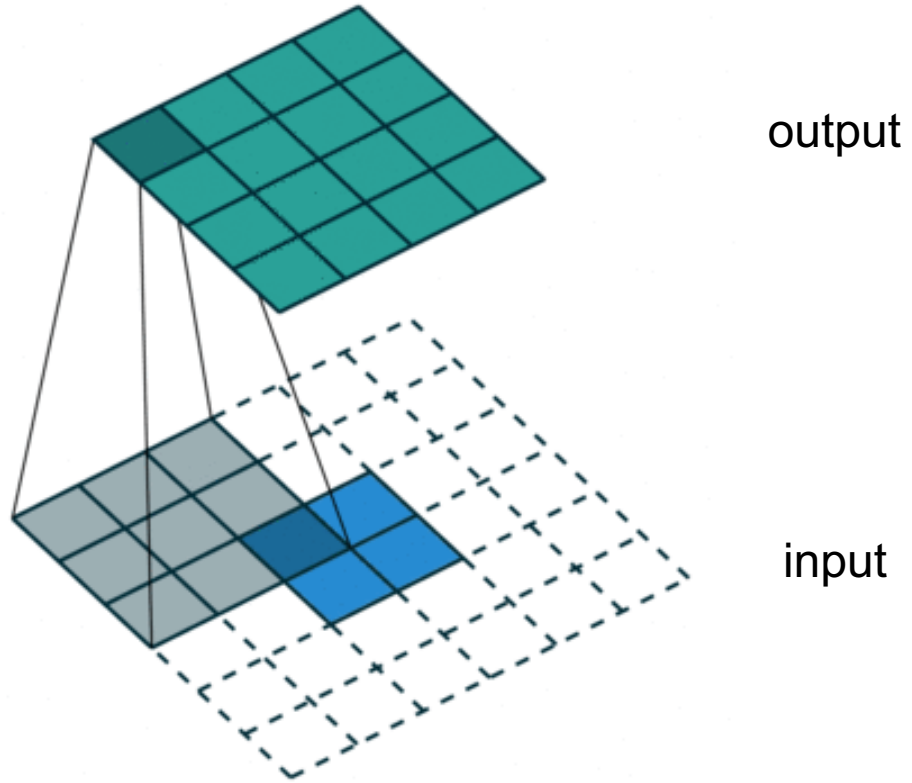
x_{11}	x_{12}	x_{13}	x_{14}
x_{21}	x_{22}	x_{23}	x_{24}
x_{31}	x_{32}	x_{33}	x_{34}
x_{41}	x_{42}	x_{43}	x_{44}



Convolve input with *flipped* filter

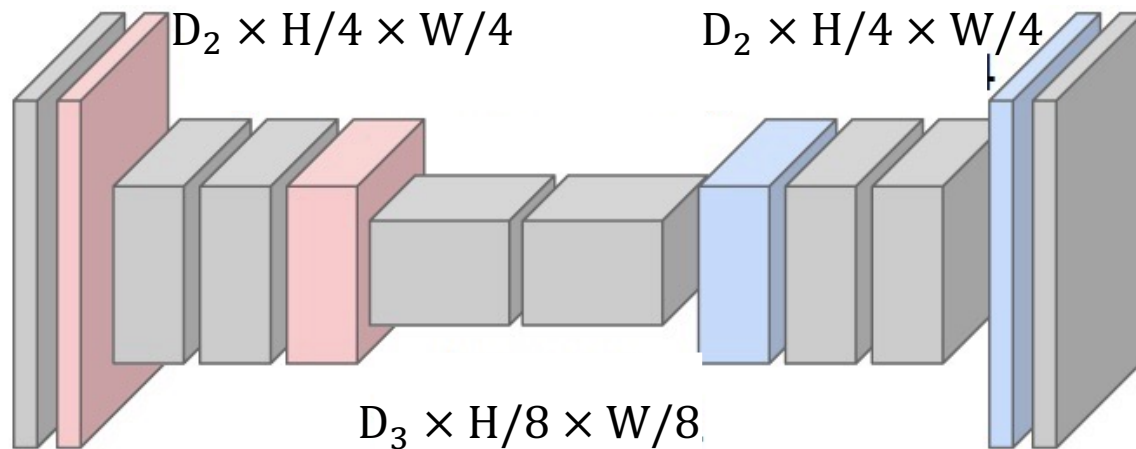
$$\begin{pmatrix}
 w_{11} & 0 & 0 & 0 \\
 w_{12} & w_{11} & 0 & 0 \\
 w_{13} & w_{12} & 0 & 0 \\
 0 & w_{13} & 0 & 0 \\
 w_{21} & 0 & w_{11} & 0 \\
 w_{22} & w_{21} & w_{12} & w_{11} \\
 \mathbf{w_{23}} & \mathbf{w_{22}} & \mathbf{w_{13}} & \mathbf{w_{12}} \\
 0 & w_{23} & 0 & w_{13} \\
 w_{31} & 0 & w_{21} & 0 \\
 w_{32} & w_{31} & w_{22} & w_{21} \\
 w_{33} & w_{32} & w_{23} & w_{22} \\
 0 & w_{33} & 0 & w_{23} \\
 0 & 0 & w_{31} & 0 \\
 0 & 0 & w_{32} & w_{31} \\
 0 & 0 & w_{33} & w_{32} \\
 0 & 0 & 0 & w_{33}
 \end{pmatrix}
 \begin{pmatrix}
 z_{11} \\
 z_{12} \\
 z_{21} \\
 z_{22}
 \end{pmatrix}
 =
 \begin{pmatrix}
 x_{11} \\
 x_{12} \\
 \mathbf{x_{23}} \\
 x_{24} \\
 x_{31} \\
 x_{32} \\
 x_{33} \\
 x_{34} \\
 x_{41} \\
 x_{42} \\
 x_{43} \\
 x_{44}
 \end{pmatrix}
 \quad
 x_{23} = w_{23}z_{11} + w_{22}z_{12} + w_{13}z_{21} + w_{12}z_{22}$$

Transpose Convolution





Input:
 $3 \times H \times W$

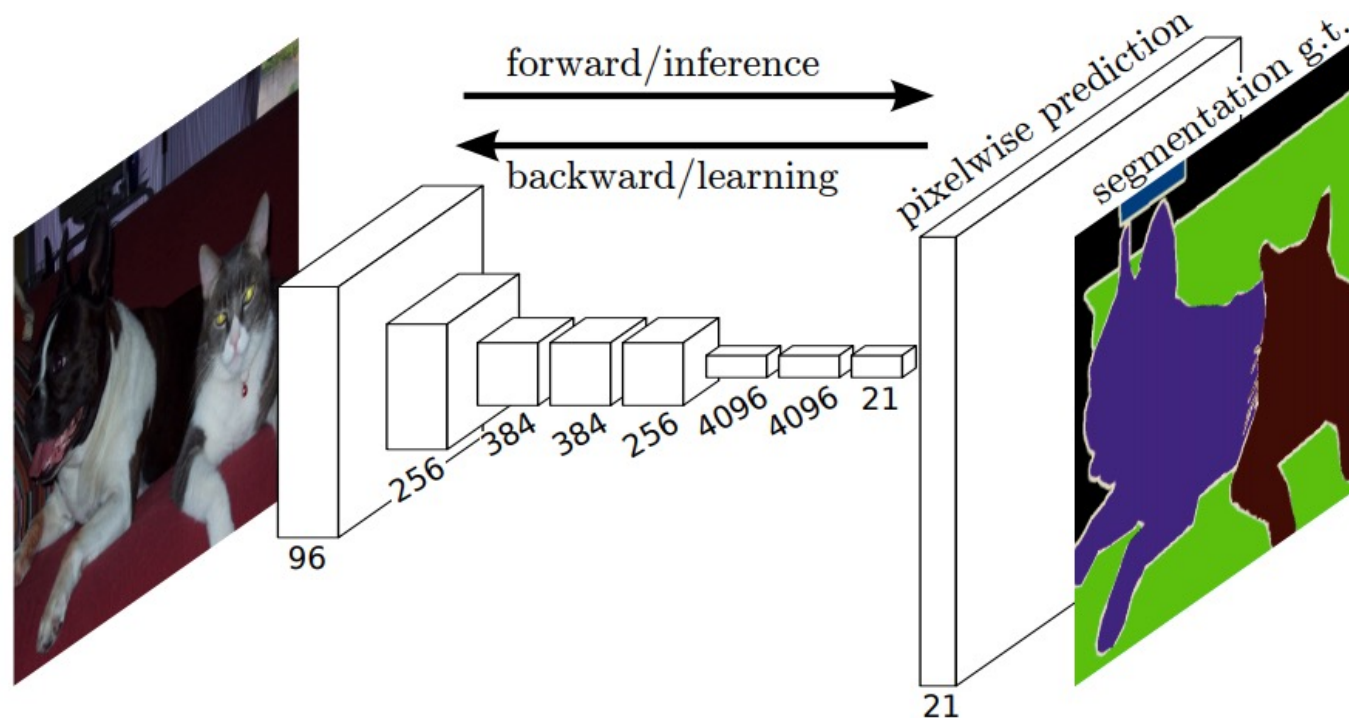


High-res:
 $D_1 \times H/2 \times W/2$

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

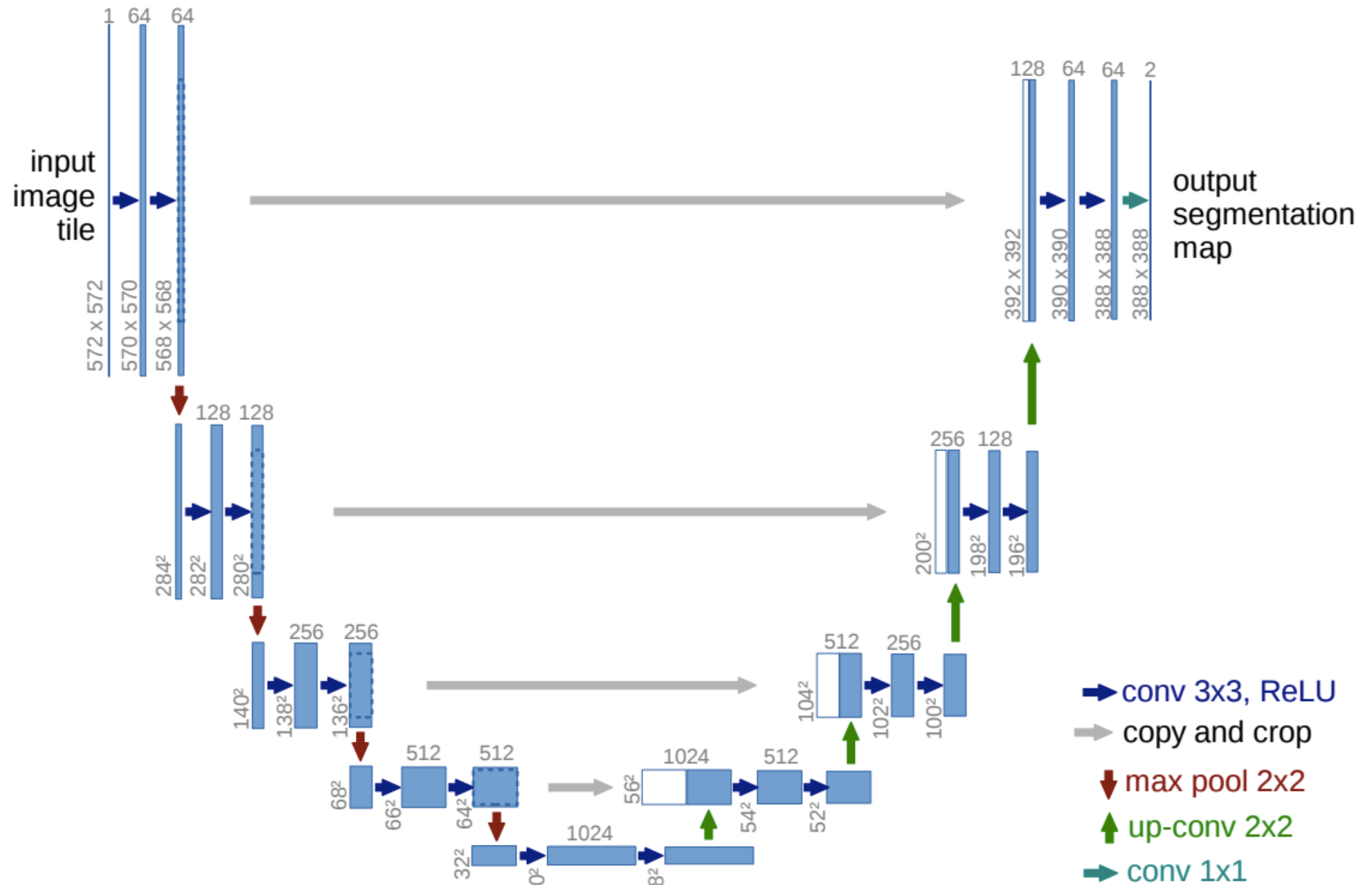


Predictions:
 $H \times W$



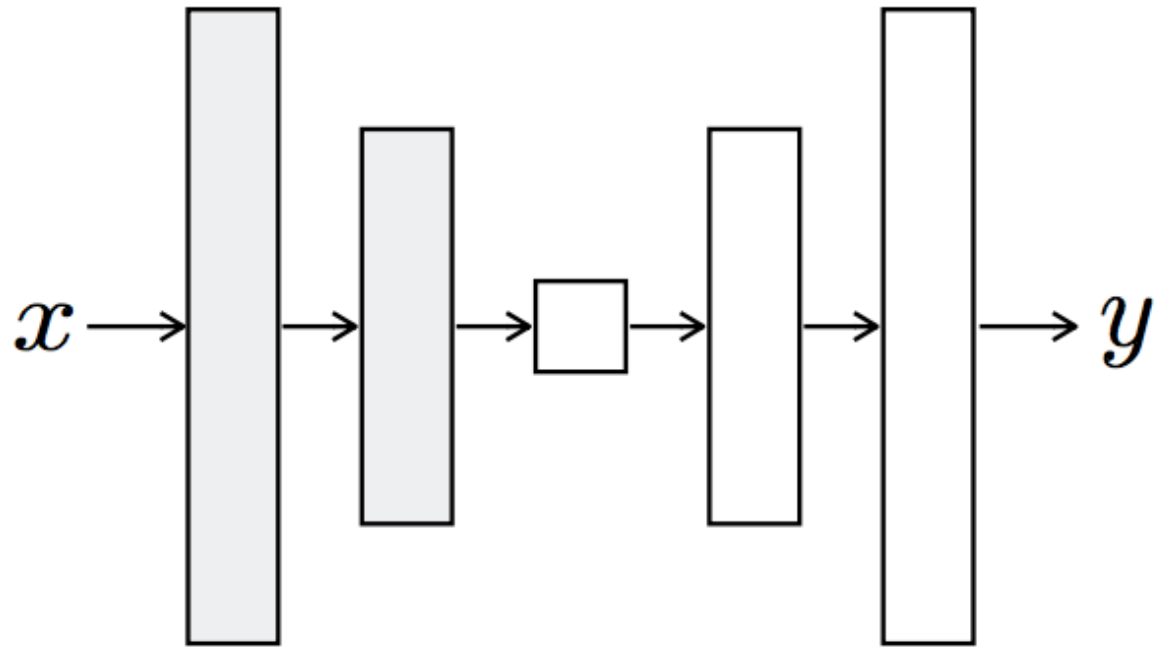
Advanced Techniques in Segmentation

U-Net

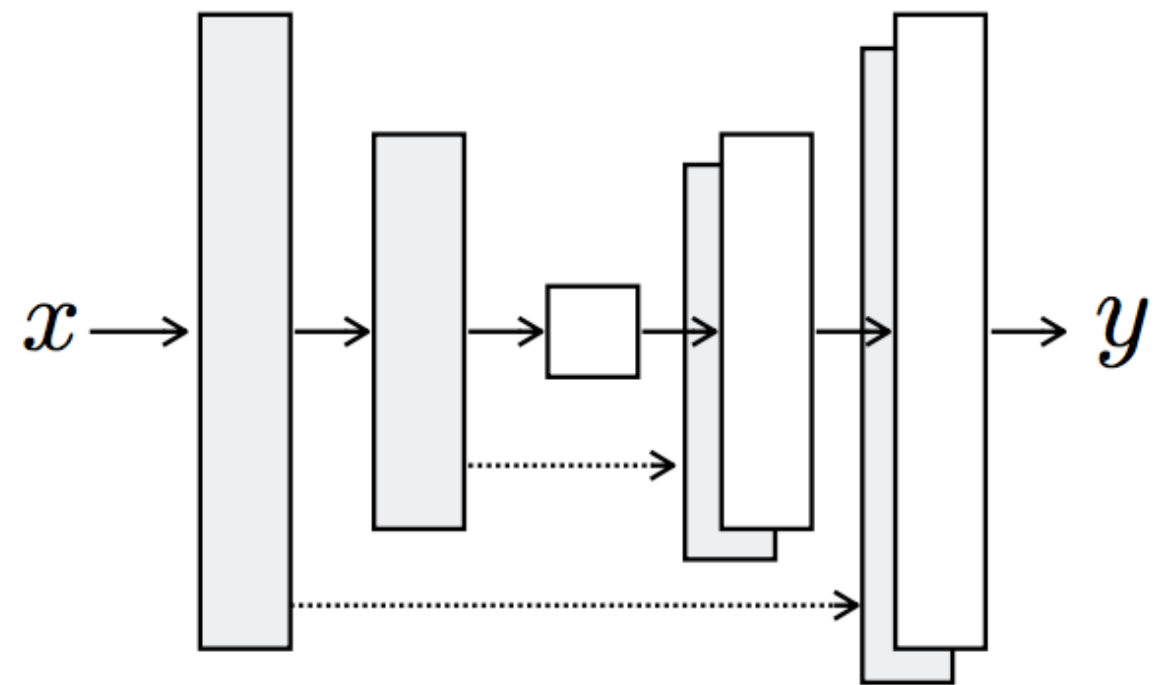


U-Net

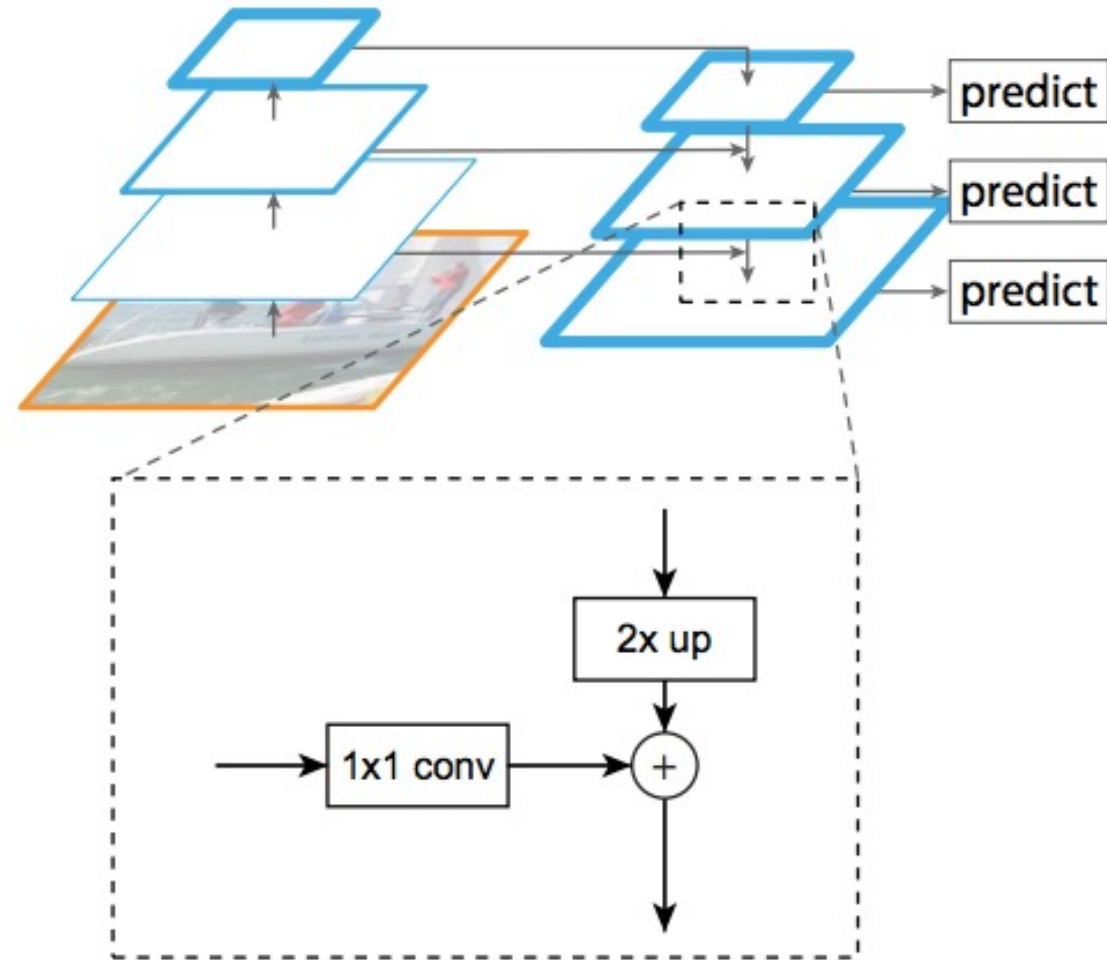
Encoder-decoder



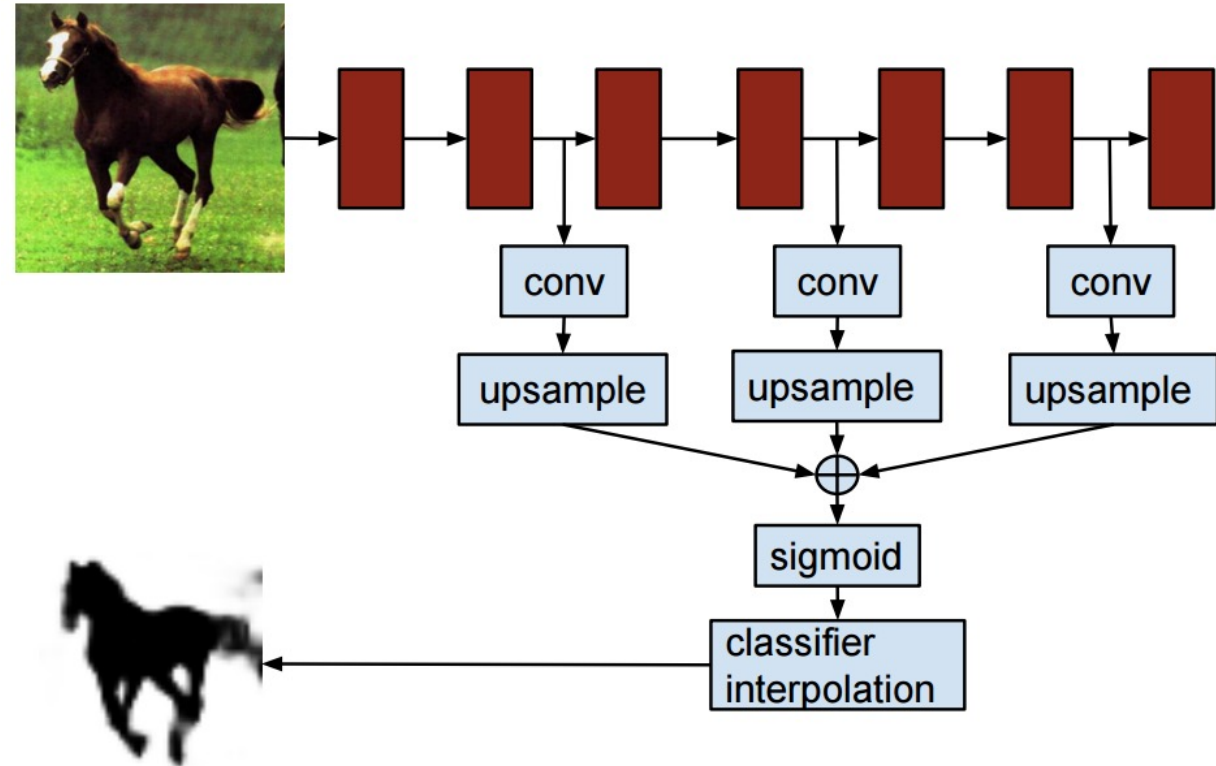
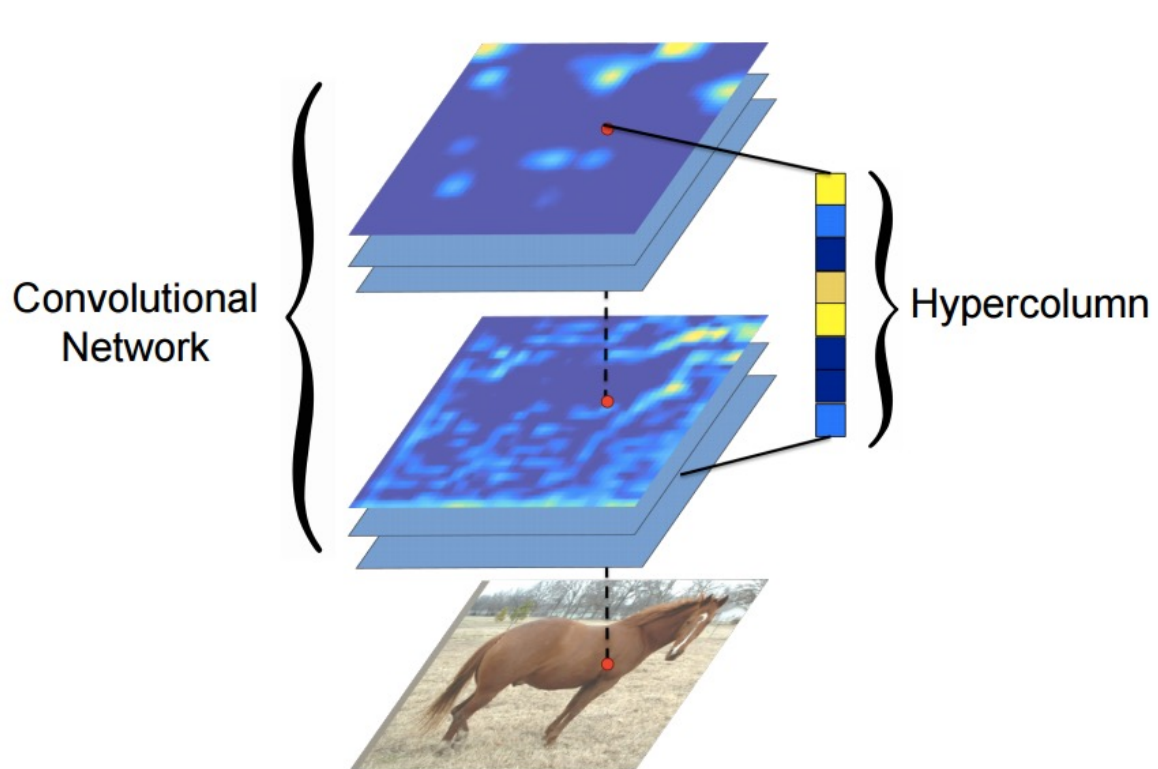
U-Net



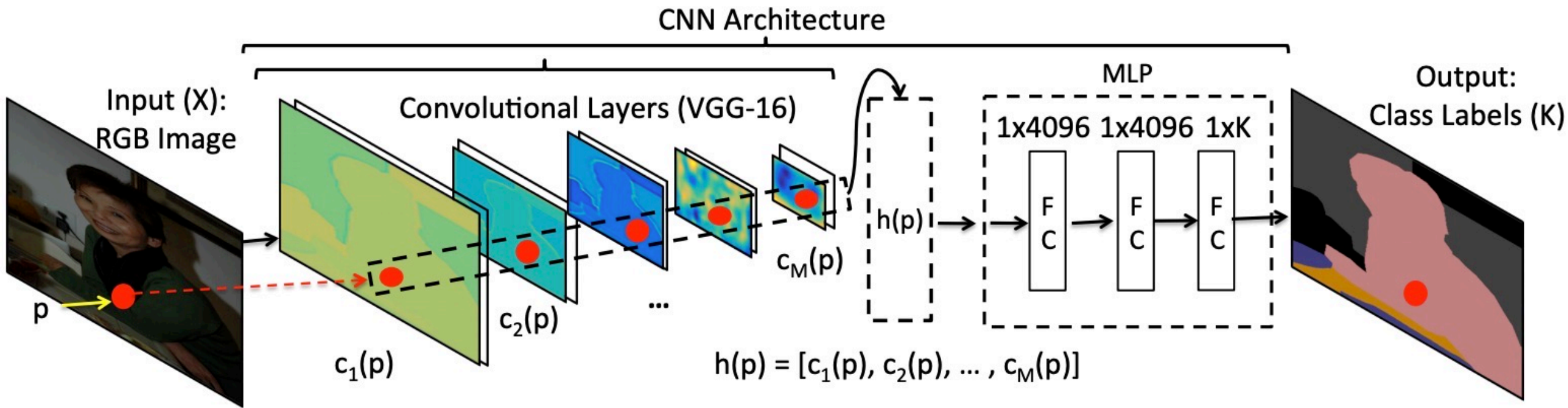
U-Net and FPN



Hypercolumn and Skip Connection



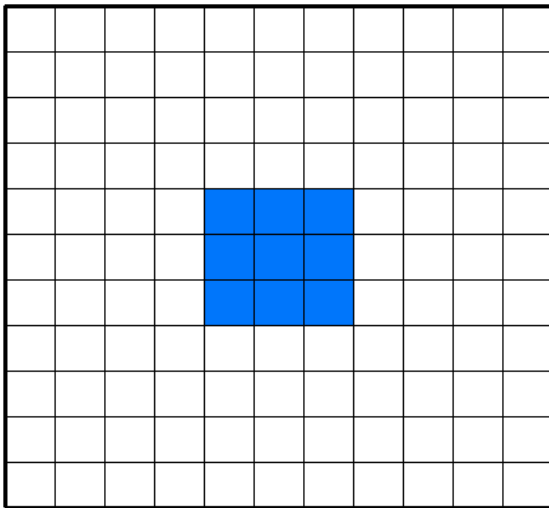
Hypercolumn and Skip Connection



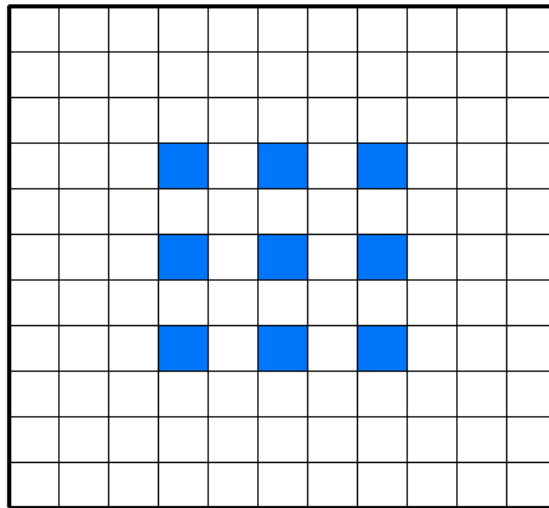
Dilated convolutions

- Idea: instead of reducing spatial resolution of feature maps, use a large sparse filter

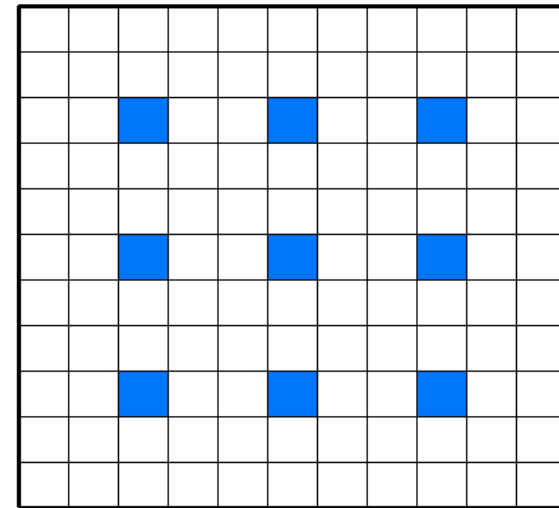
Dilation factor 1



Dilation factor 2

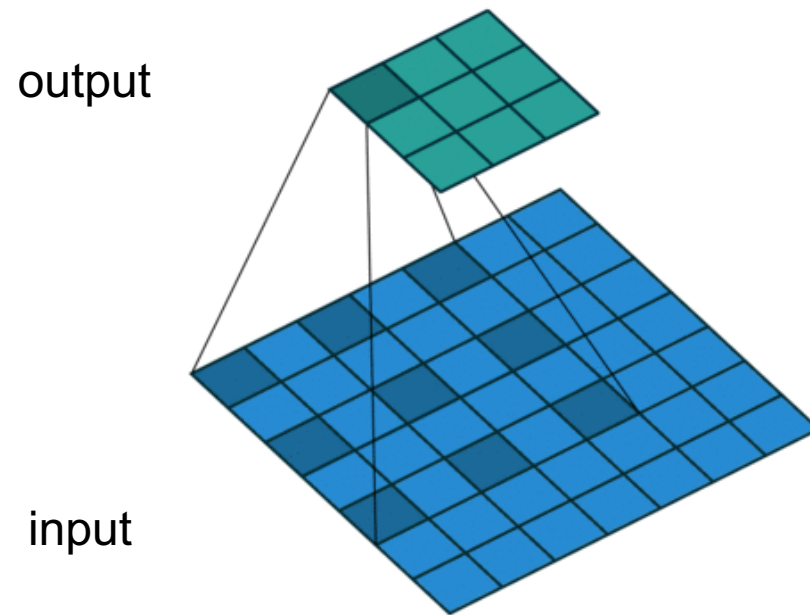


Dilation factor 3

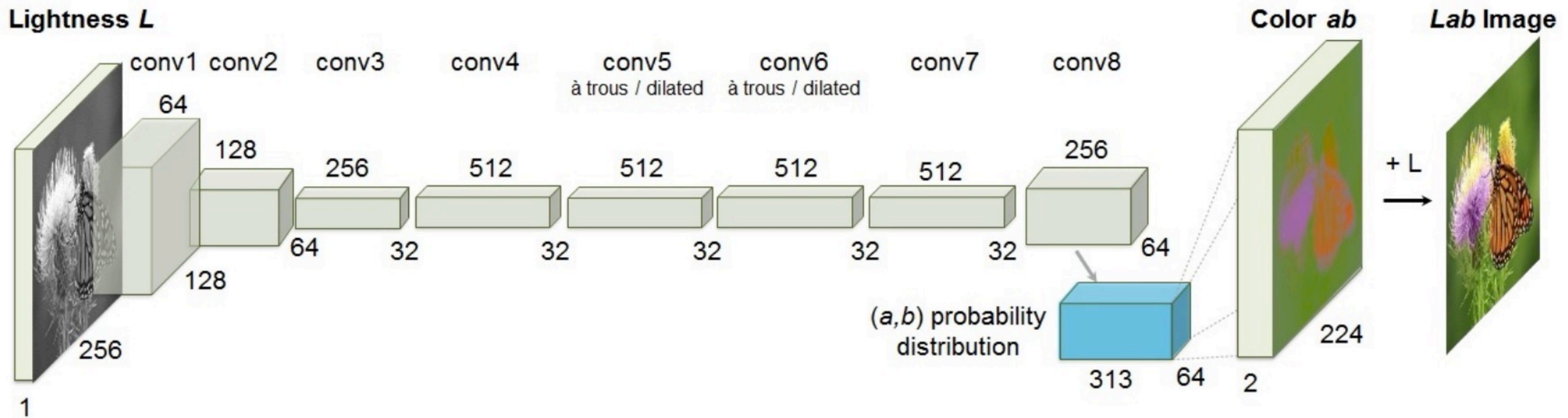


Dilated convolutions

- Increase the receptive field



Dilated convolutions



This Class

- Naïve FCN model for Image Segmentation
- Transpose Convolution
- Advanced Techniques in Segmentation

Next Class

Visualizing Deep Networks