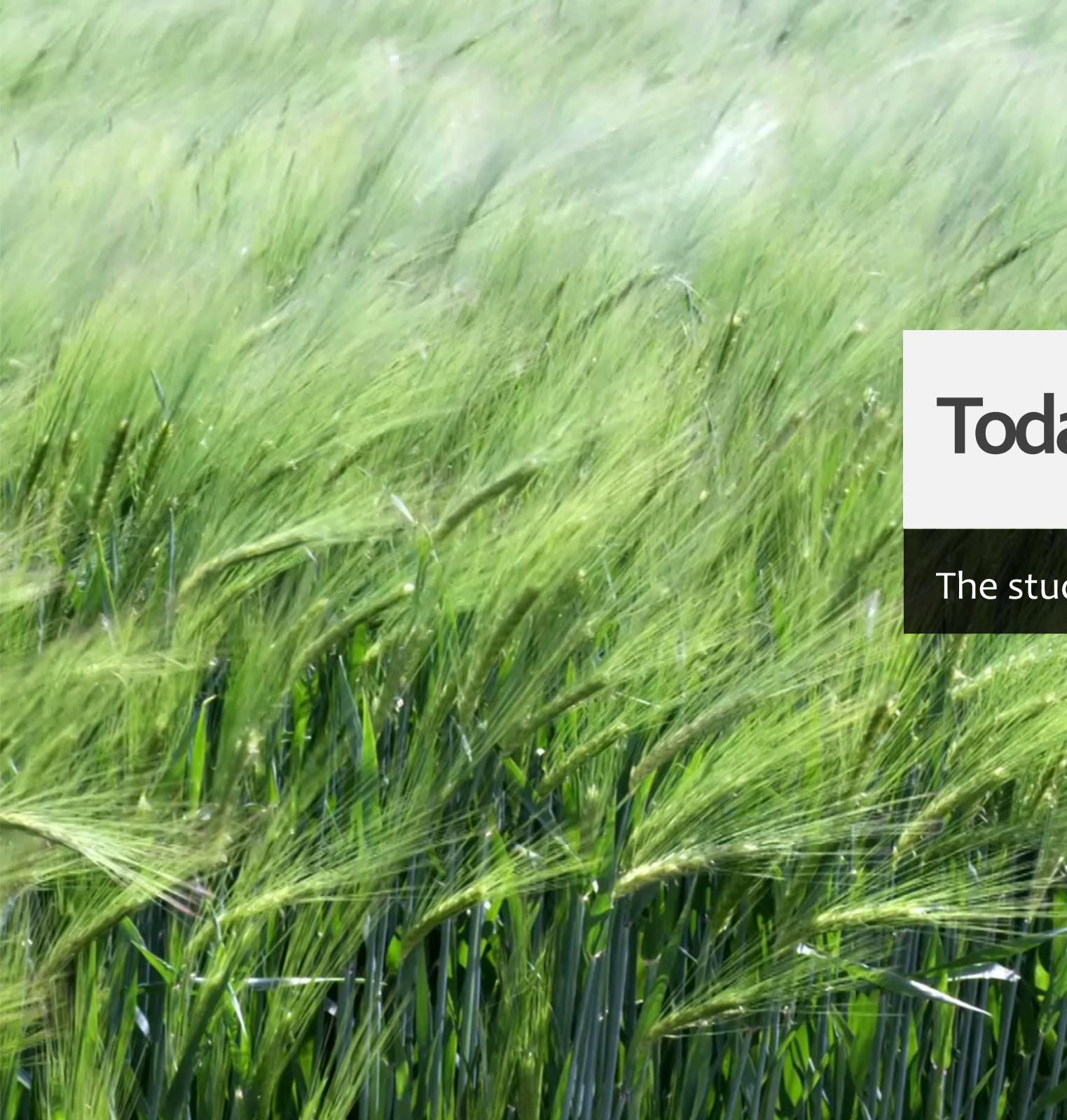


Dr. Muhammad Moazam Fraz
Department of Computing, NUST SEECS
<http://seeecs.nust.edu.pk/faculty/mmfraz.html>
<http://vision.seeecs.edu.pk>



Today's Class : Learning Outcomes

The students will learn about

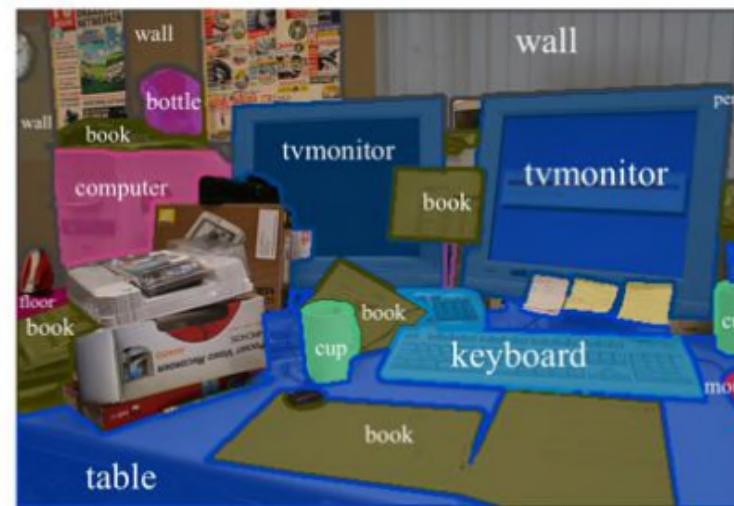
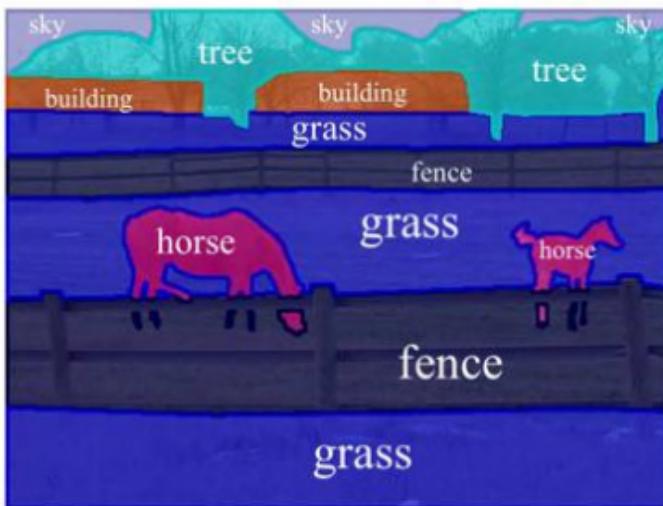
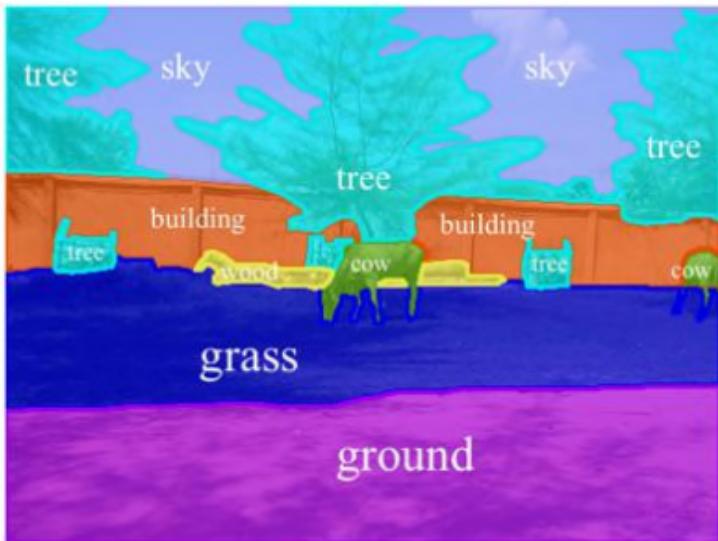
- Semantic Segmentation
- Fully Convolutional Network
- Transposed Convolutions
- U-Net Architecture

Today's Class



What is semantic segmentation?

1. Idea: recognizing, understanding what's in the image in pixel level.
2. A lot more difficult
(Most of the traditional methods cannot tell different objects.)



"Two men riding on a bike in front of a building on the road. And there is a car."

What is semantic segmentation?

1. Idea: recognizing, understanding what's in the image in pixel level.
2. A lot more difficult
(Most of the traditional methods cannot tell different objects.)

No worries, even the best ML researchers find it very challenging.

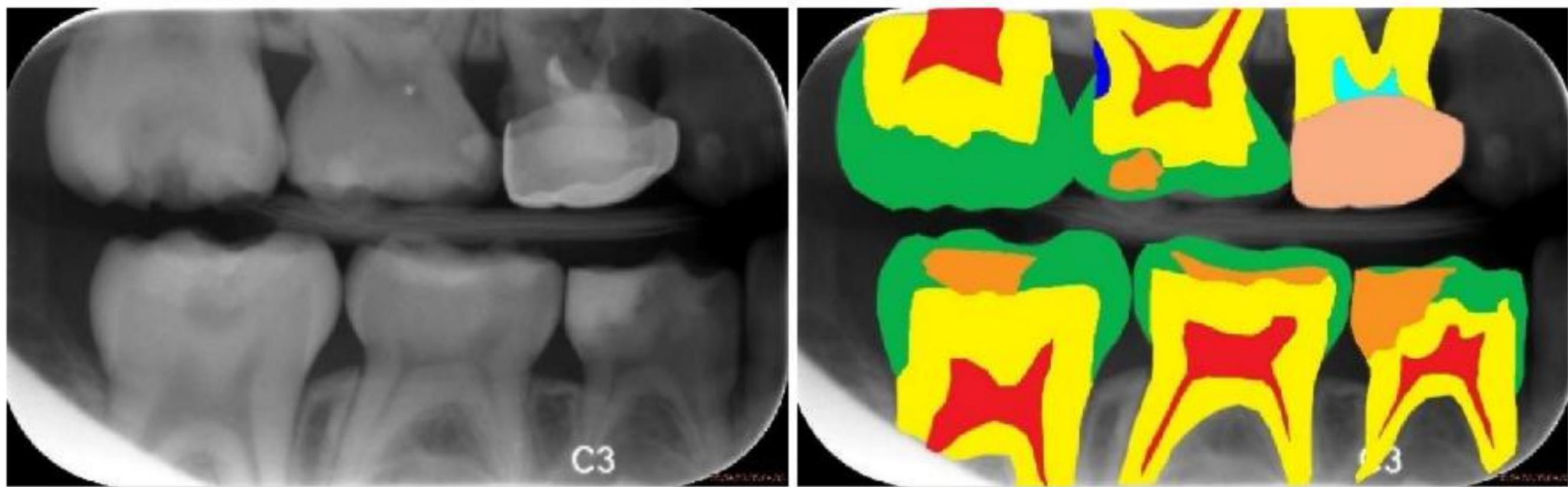
3. Output: regions with different (and limited number of) classes
 1. COCO detection challenge: 80 classes.
 2. PASCAL VOC challenge: 21 classes

Why semantic segmentation?

1. robot vision and understanding
2. autonomous driving



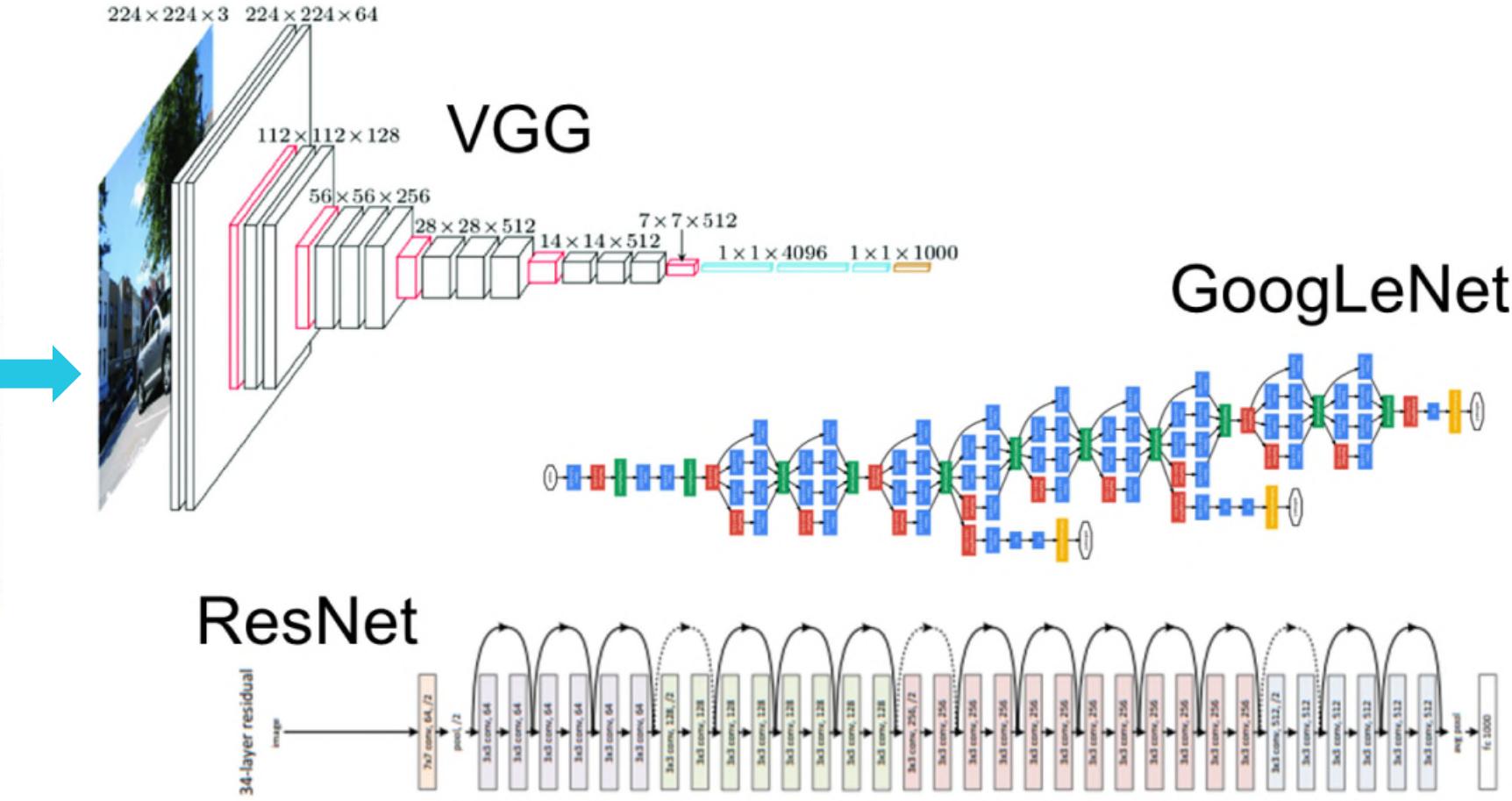
1. Why semantic segmentation?
3. medical purposes (ISBI Challenge)



Semantic segmentation before deep learning

1. Clustering methods
2. Motion & Interactive Segmentation
3. Partial differential equation-based methods
 - Level-set methods
 - Parametric methods (Active Contours, Snakes)
 - Fast marching methods
4. Graph partitioning methods
 - Markov random fields
5. Watershed transformation

So far: Image Classification



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

Computer Vision Tasks: Semantic Segmentation

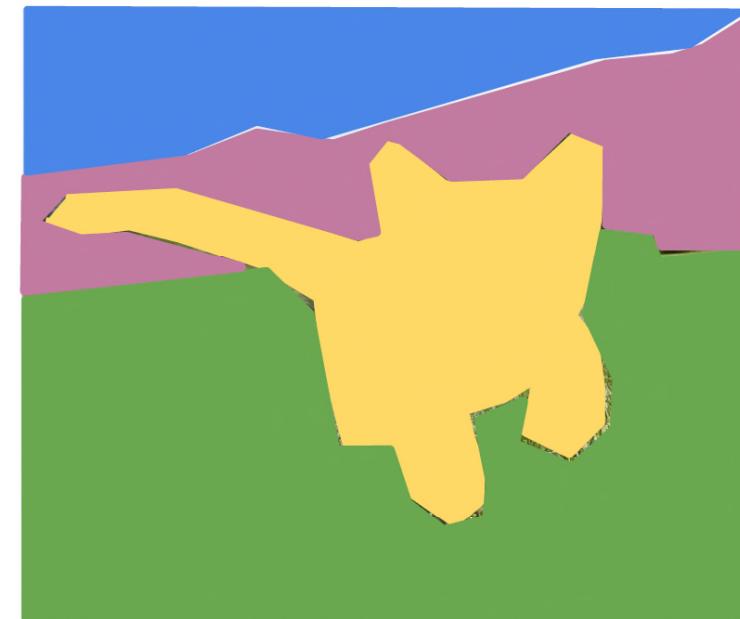
Classification



CAT

No spatial extent

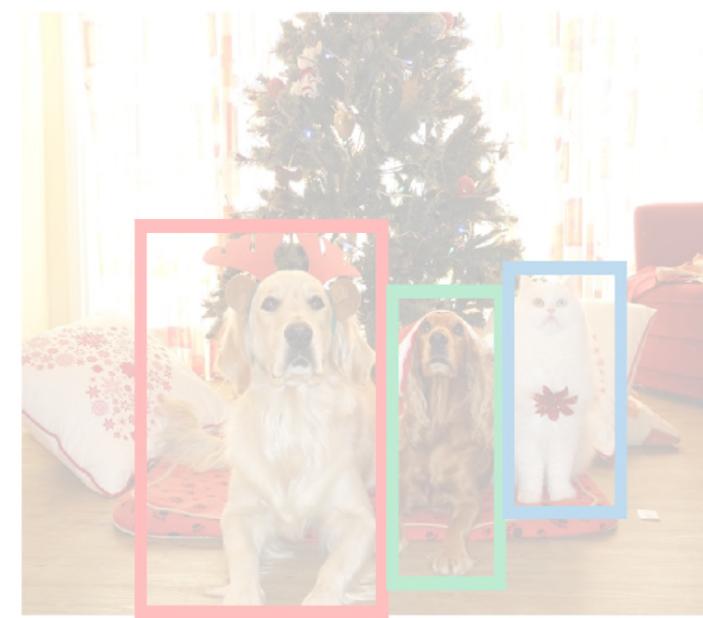
Semantic
Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Objects

Instance
Segmentation



DOG, DOG, CAT

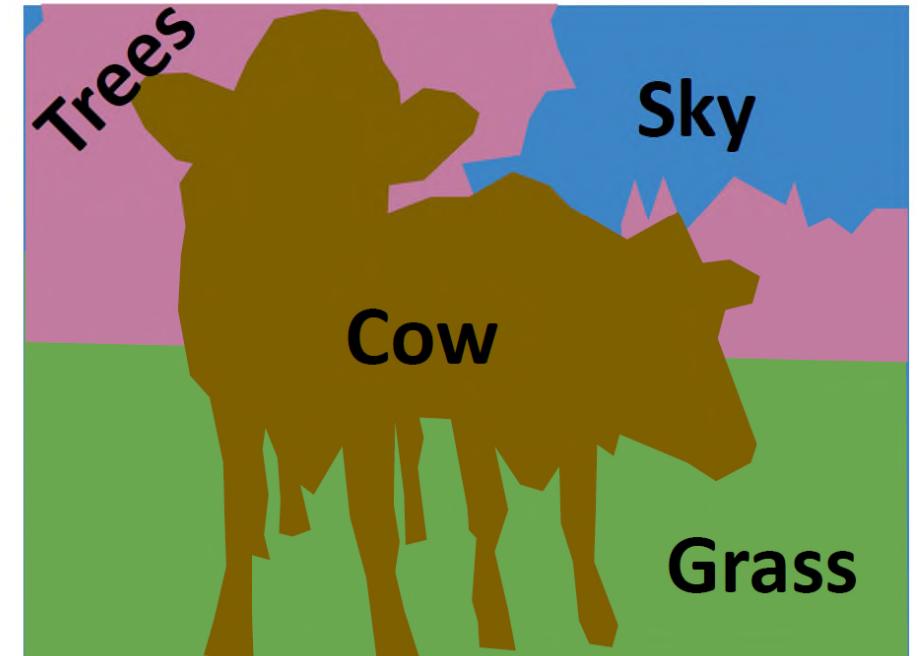
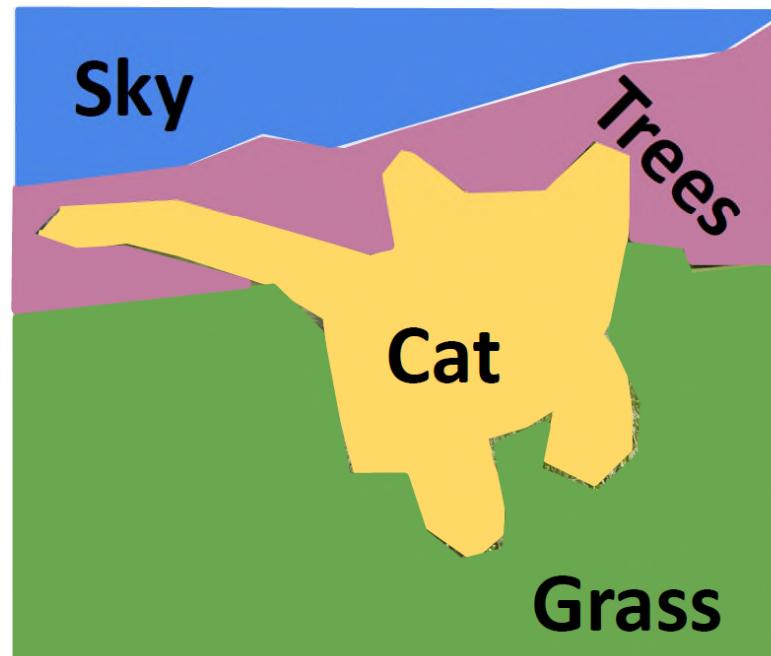
Semantic Segmentation

Label each pixel in the image with a category label

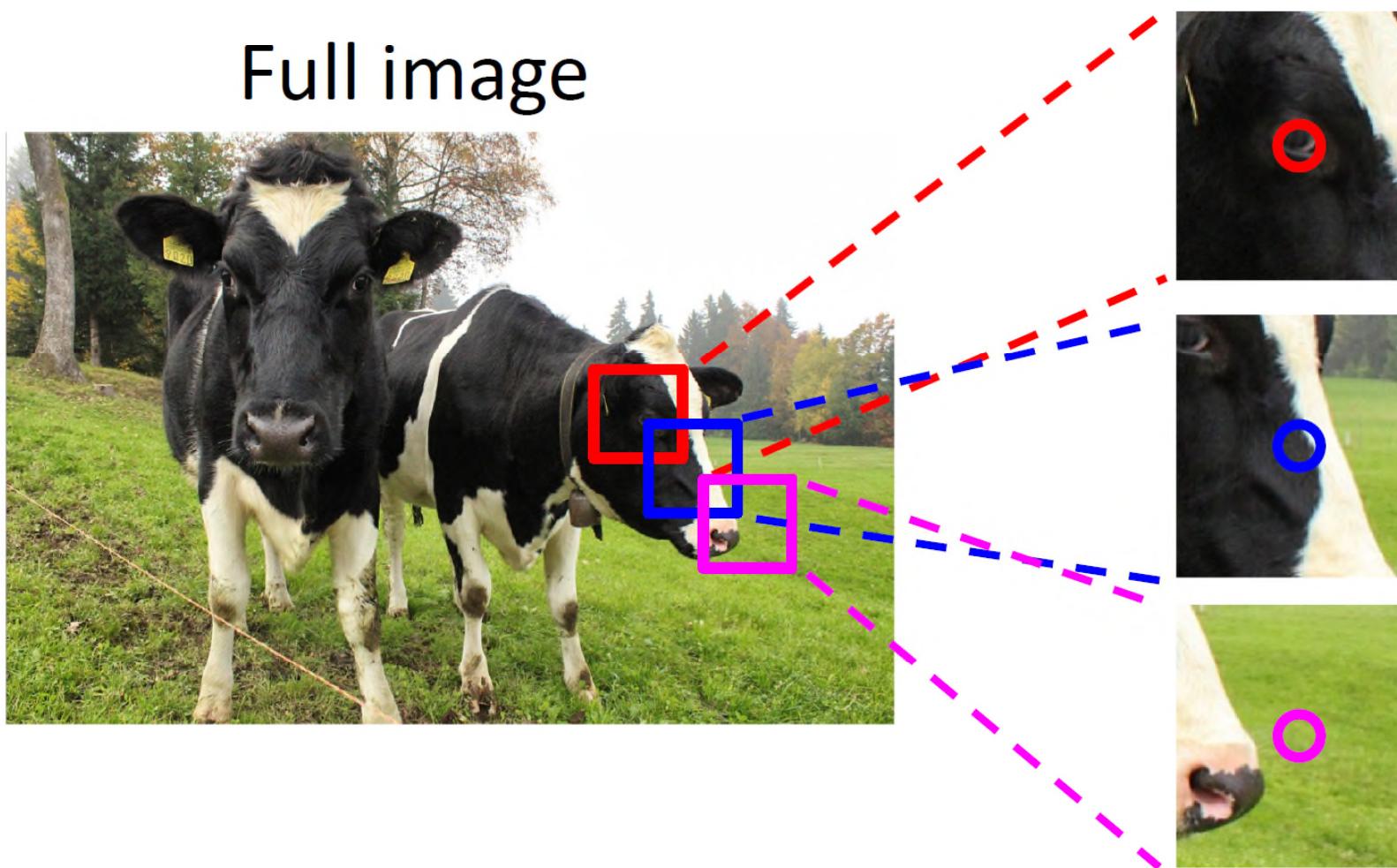
Don't differentiate instances,
only care about pixels



[This image is CC0 public domain](#)

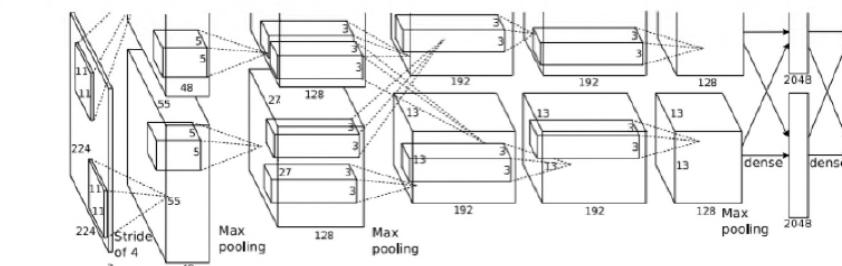
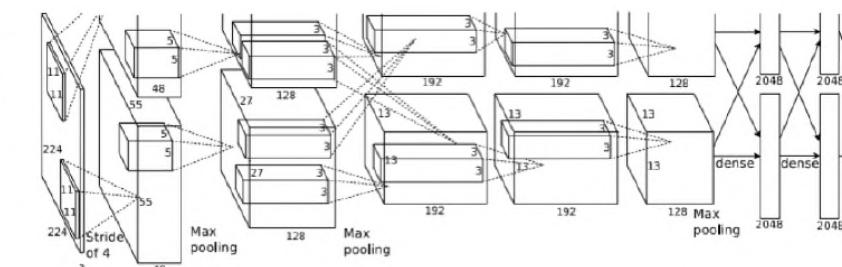
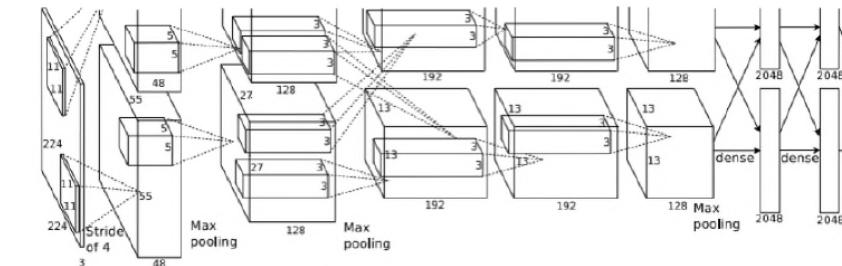


Semantic Segmentation Idea: Sliding Window



Extract
patch

Classify center
pixel with CNN



Cow

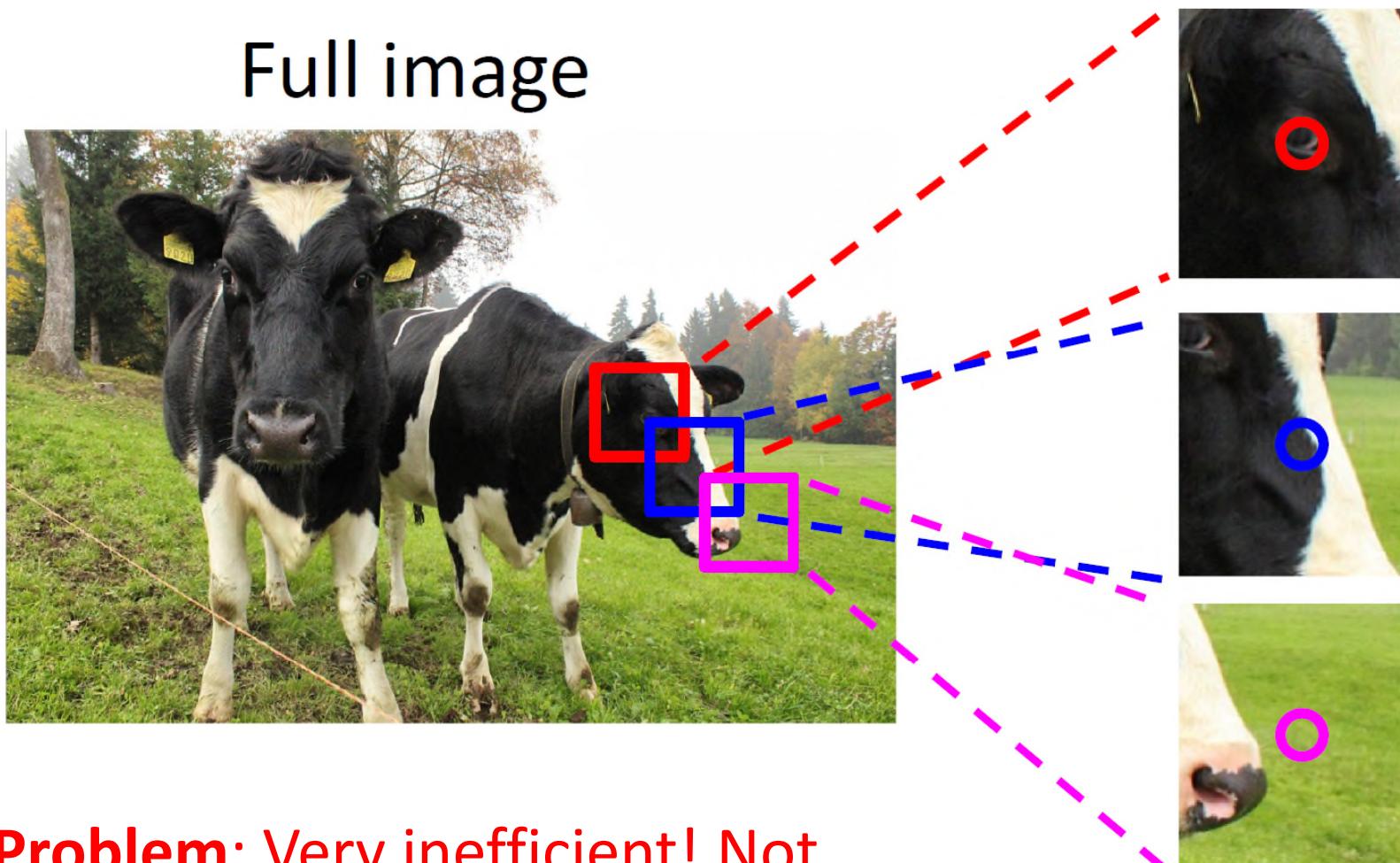
Cow

Grass

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

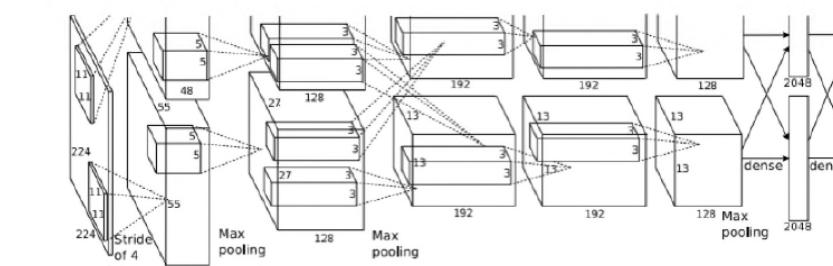
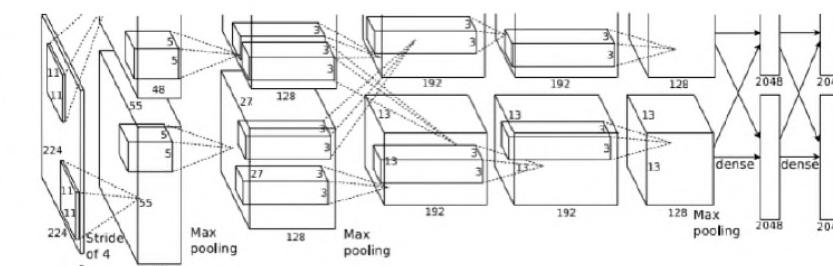
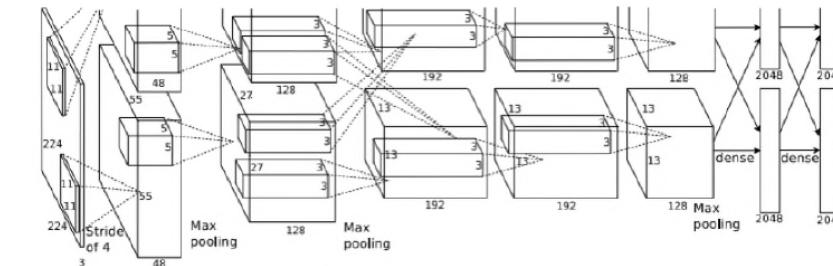
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window



Extract
patch

Classify center
pixel with CNN



Cow

Cow

Grass

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

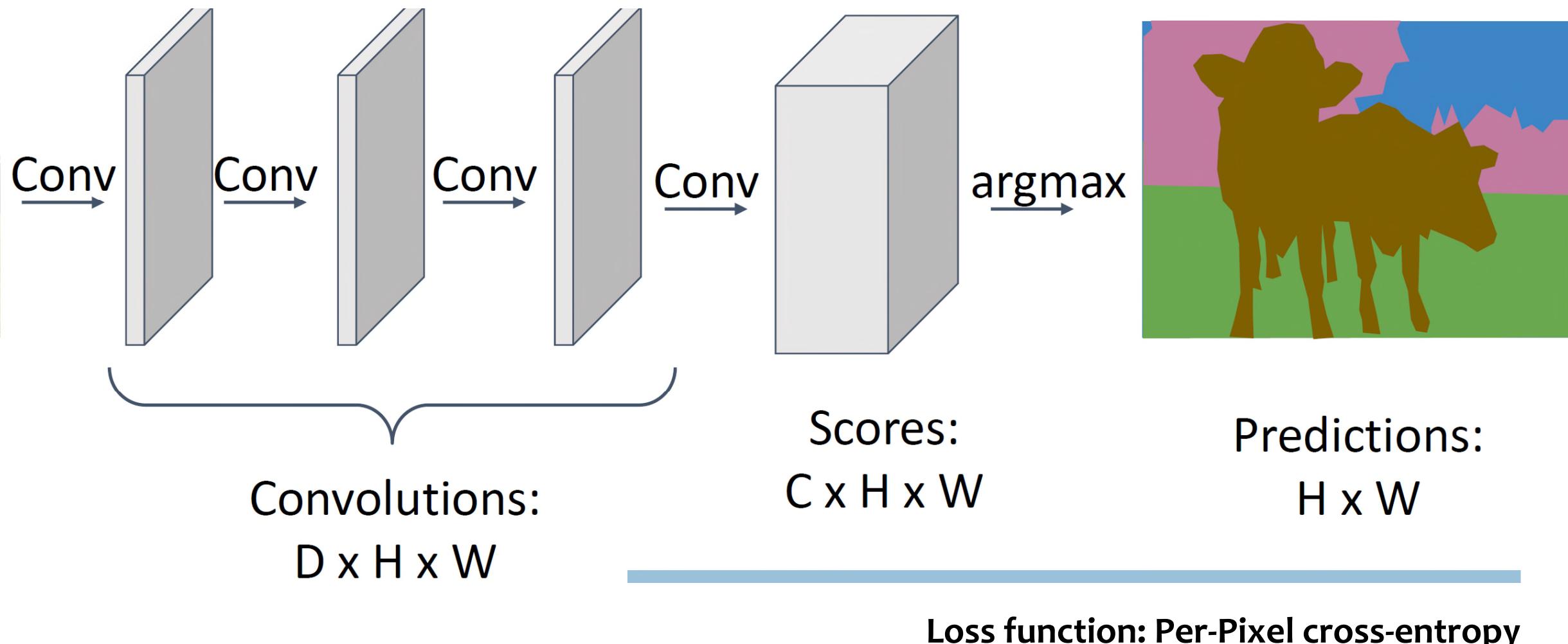
Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation: Fully Convolutional Network

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

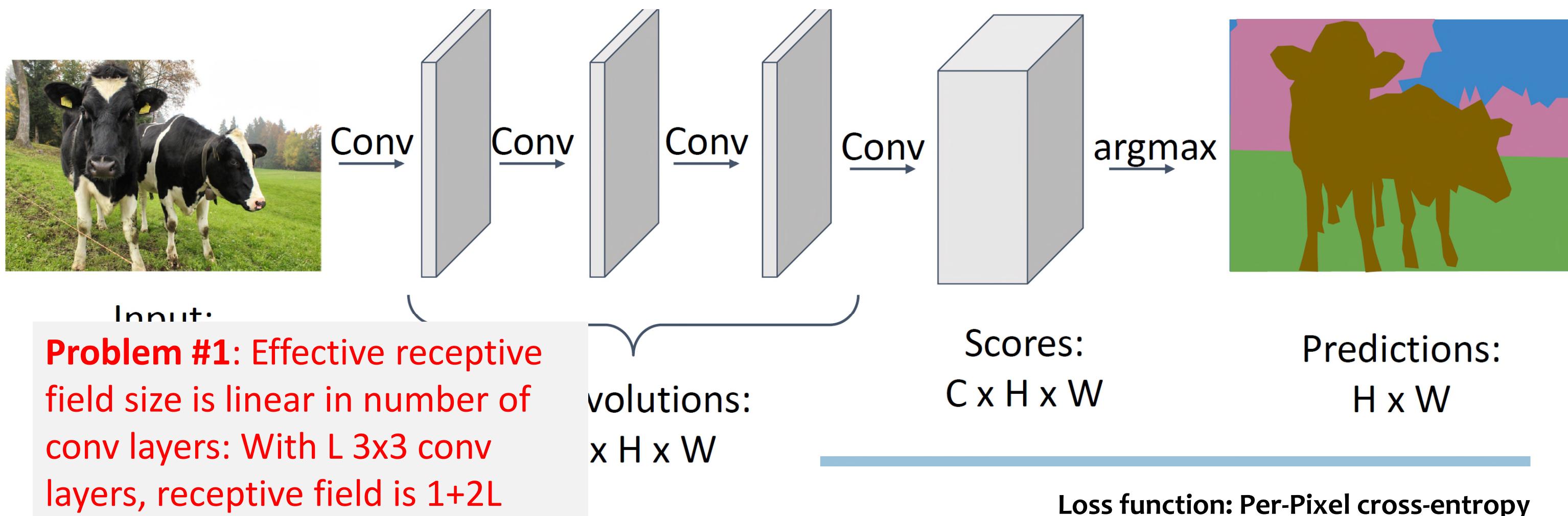


Input:
 $3 \times H \times W$



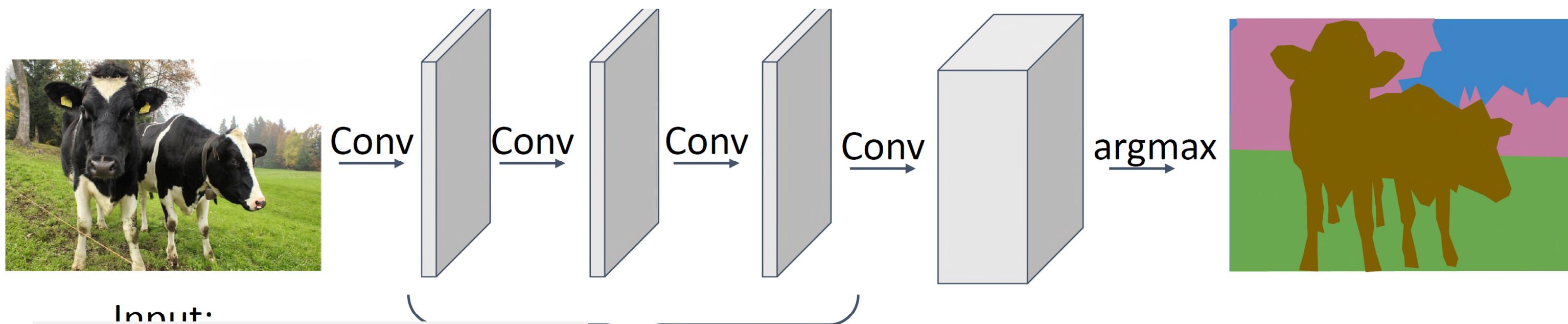
Semantic Segmentation: Fully Convolutional Network

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Semantic Segmentation: Fully Convolutional Network

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Input:

Problem #1: Effective receptive field size is linear in number of conv layers: With L 3×3 conv layers, receptive field is $1+2L$

volutions:
 $x H x W$

Problem #2: Convolution on high res images is expensive! Recall ResNet stem aggressively downsamples

ons:
✓
entropy

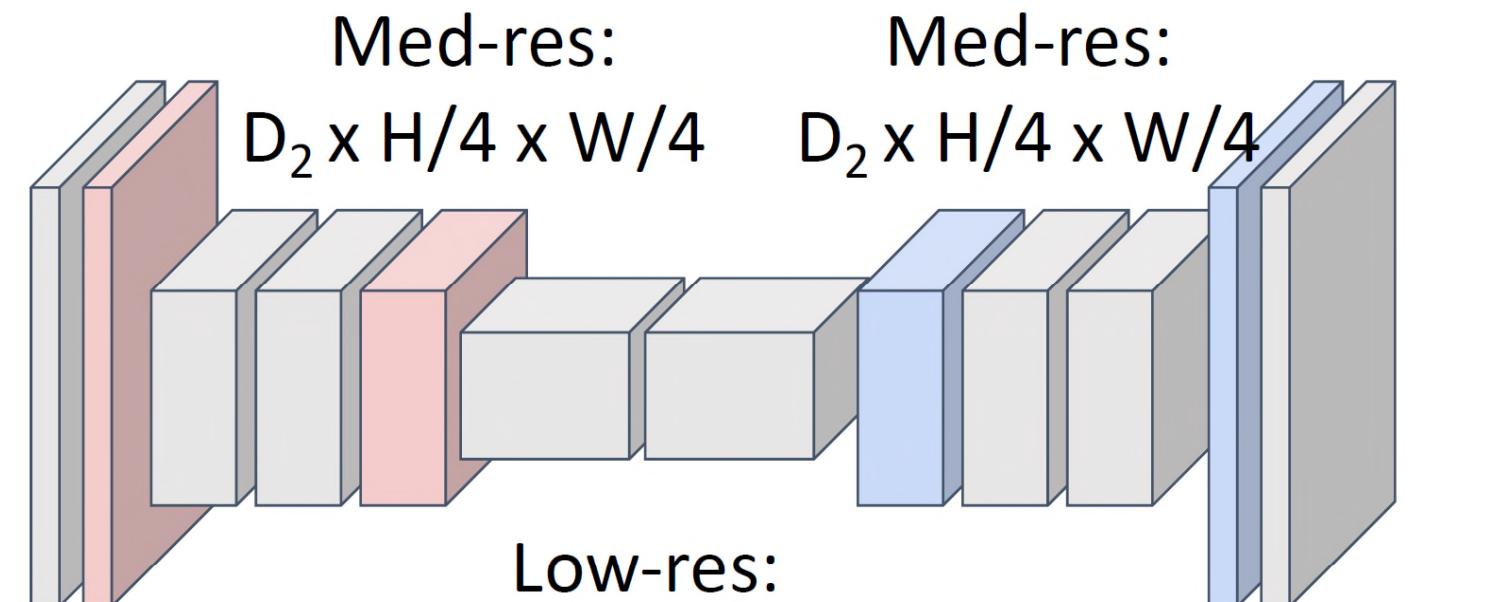
Semantic Segmentation: Fully Convolutional Network

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



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:
 $H \times W$



Semantic Segmentation: Fully Convolutional Network

Downsampling:
Pooling, Strided convolution

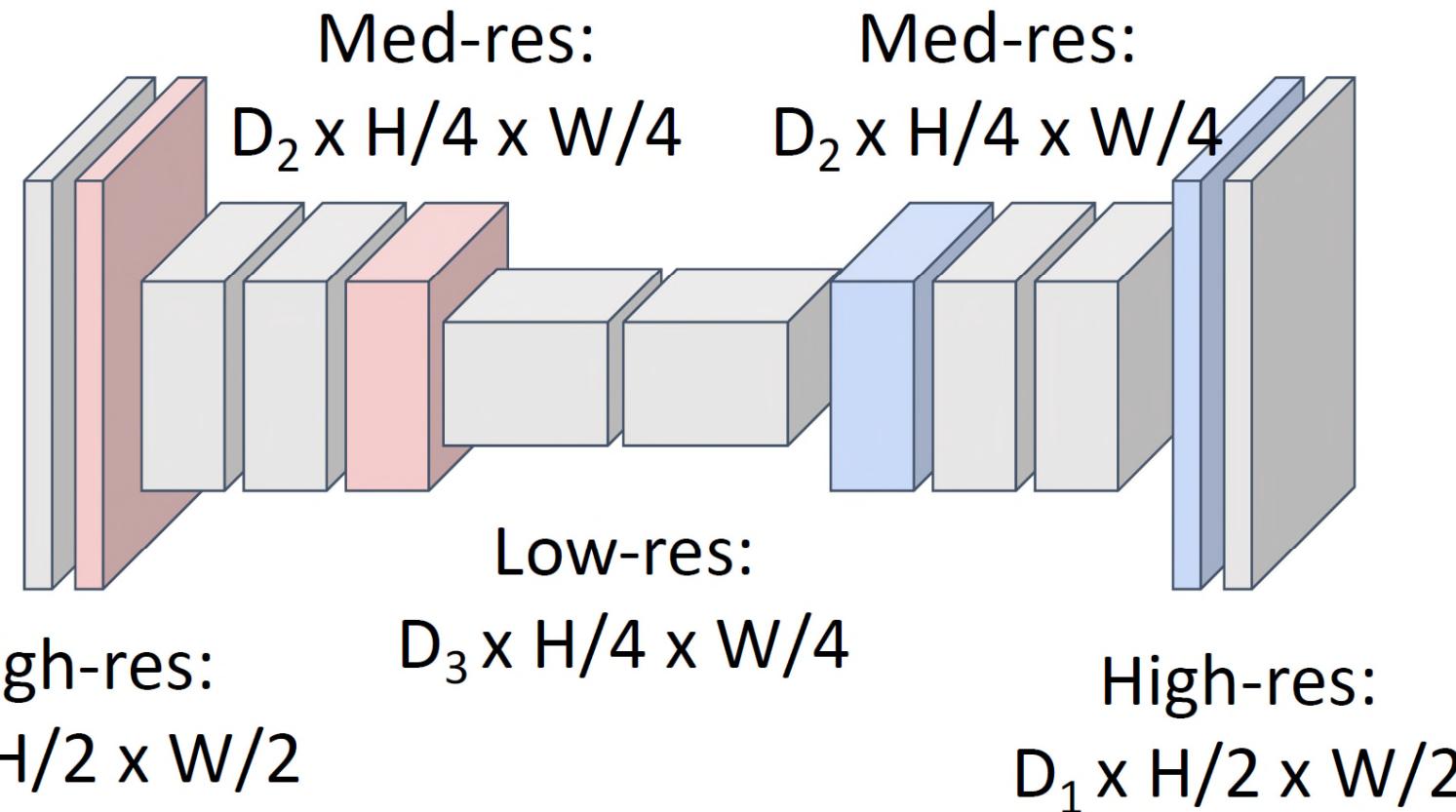
Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!

Upsampling:
???



Input:
 $3 \times H \times W$

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

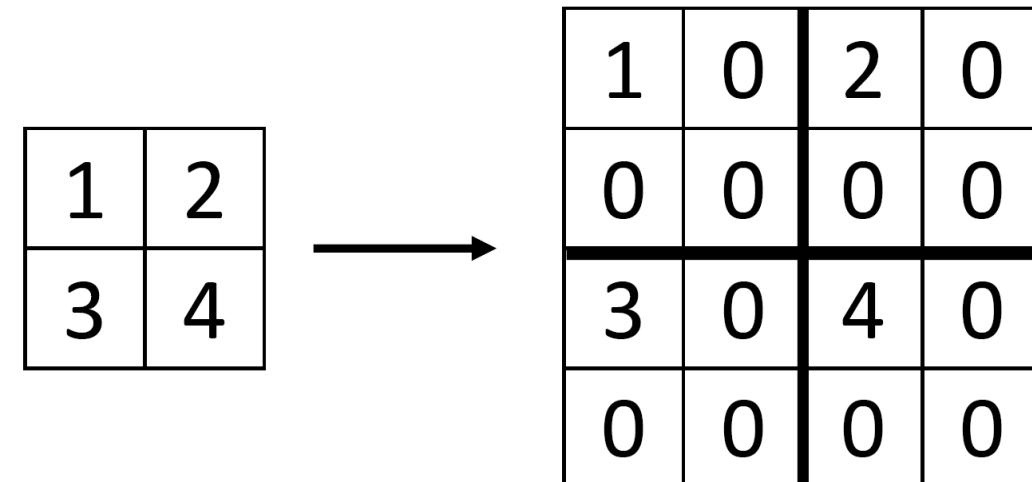


Predictions:
 $H \times W$

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Unpooling”

Bed of Nails



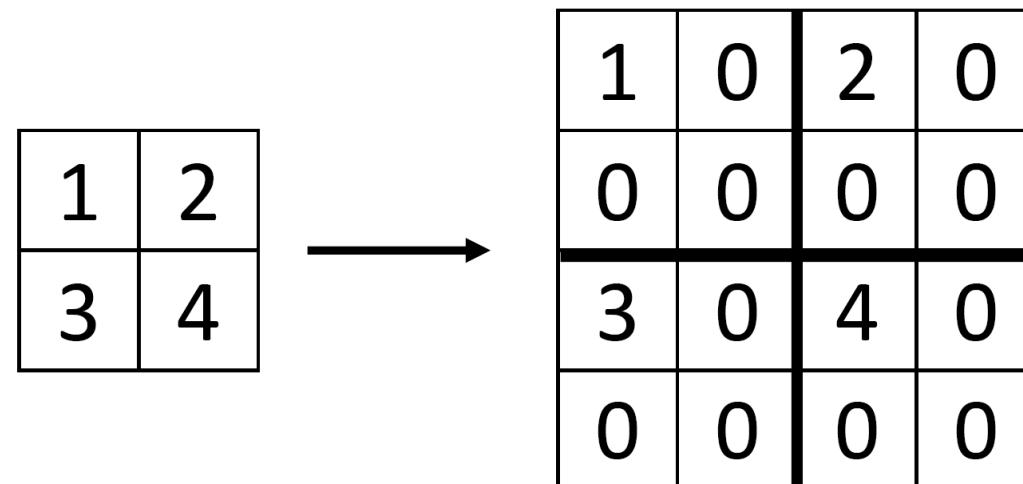
Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Unpooling”

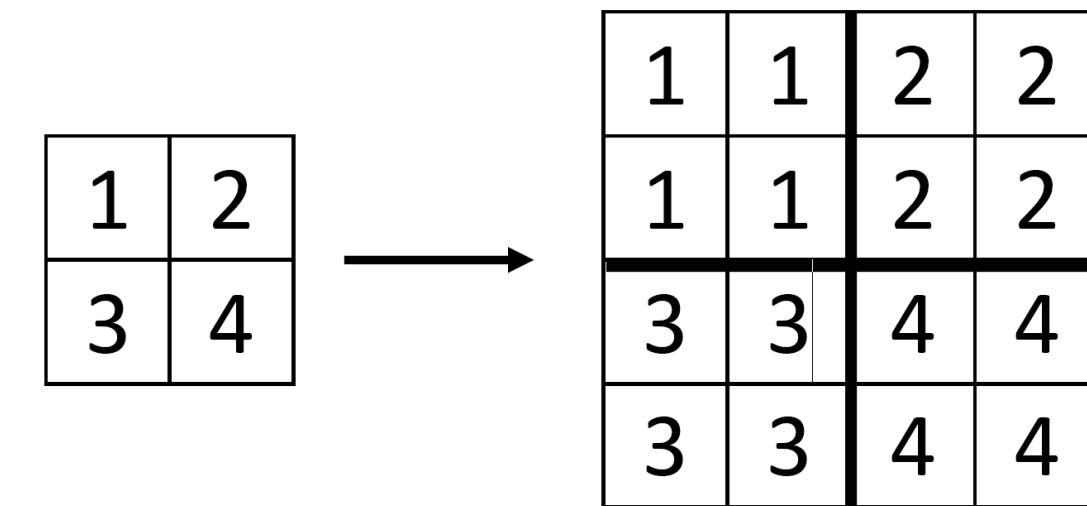
Bed of Nails



Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

Nearest Neighbour

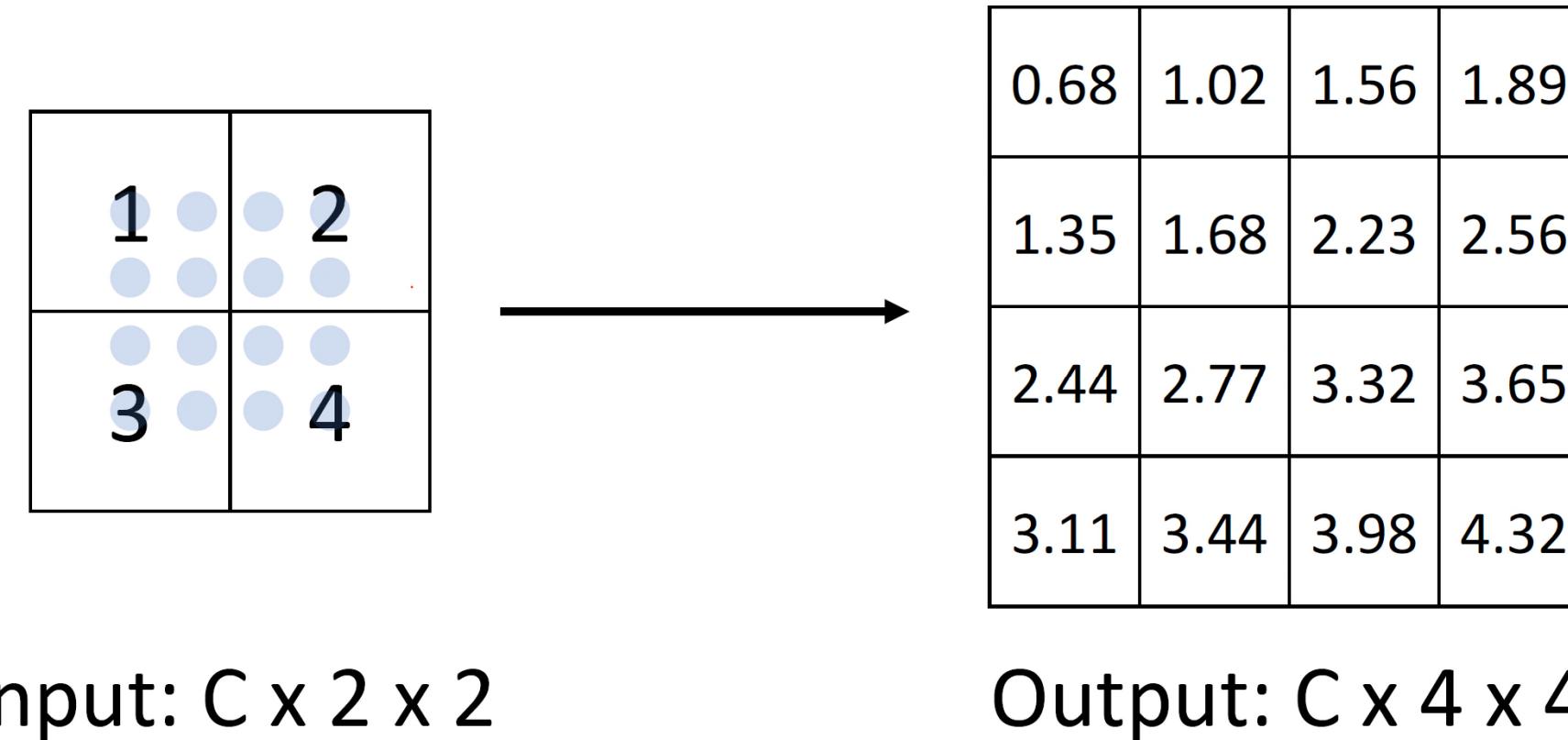


Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Bicubic Interpolation”



Use **three** closest neighbors in x and y to construct **cubic** approximations
(This is how we normally resize images!)

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Max Unpooling”

Max Pooling: Remember which position had the max

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



5	6
7	8



Rest
of
net



1	2
3	4



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

Max Unpooling: Place into remembered positions

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Bilinear Interpolation”

1	2
3	4
3	4



1.00	1.25	1.75	2.00
1.50	1.75	2.25	2.50
2.50	2.75	3.25	3.50
3.00	3.25	3.75	4.00

Input: $C \times 2 \times 2$

Output: $C \times 4 \times 4$

$$f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$$
$$j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$$

Use two closest neighbors in x and y
to construct linear approximations

Semantic Segmentation: Fully Convolutional Network

In-Network Upsampling: “Max Unpooling”

Max Pooling: Remember which position had the max

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

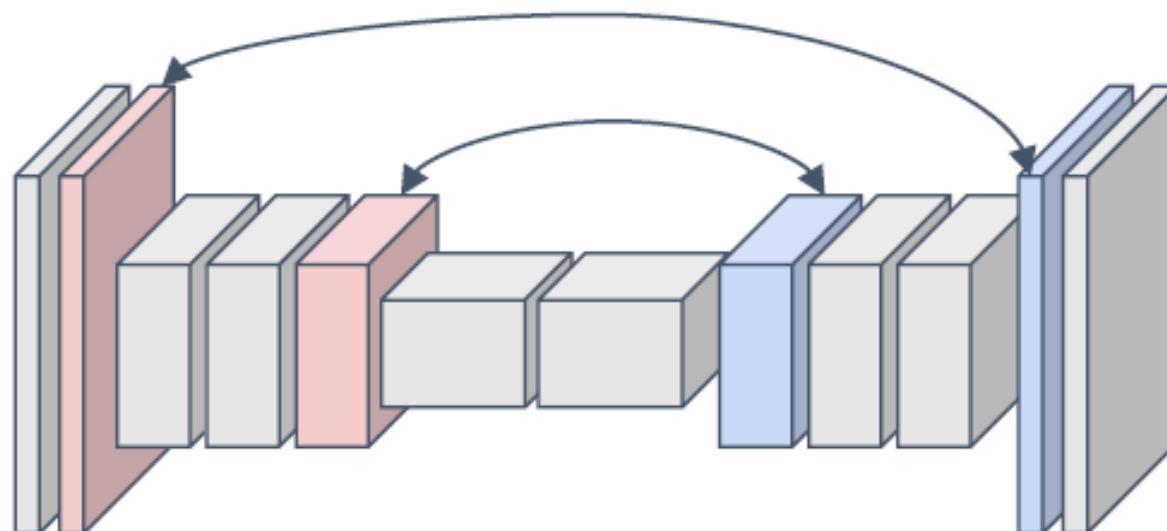
5	6
7	8

Rest
of
net

1	2
3	4

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

Max Unpooling: Place into remembered positions

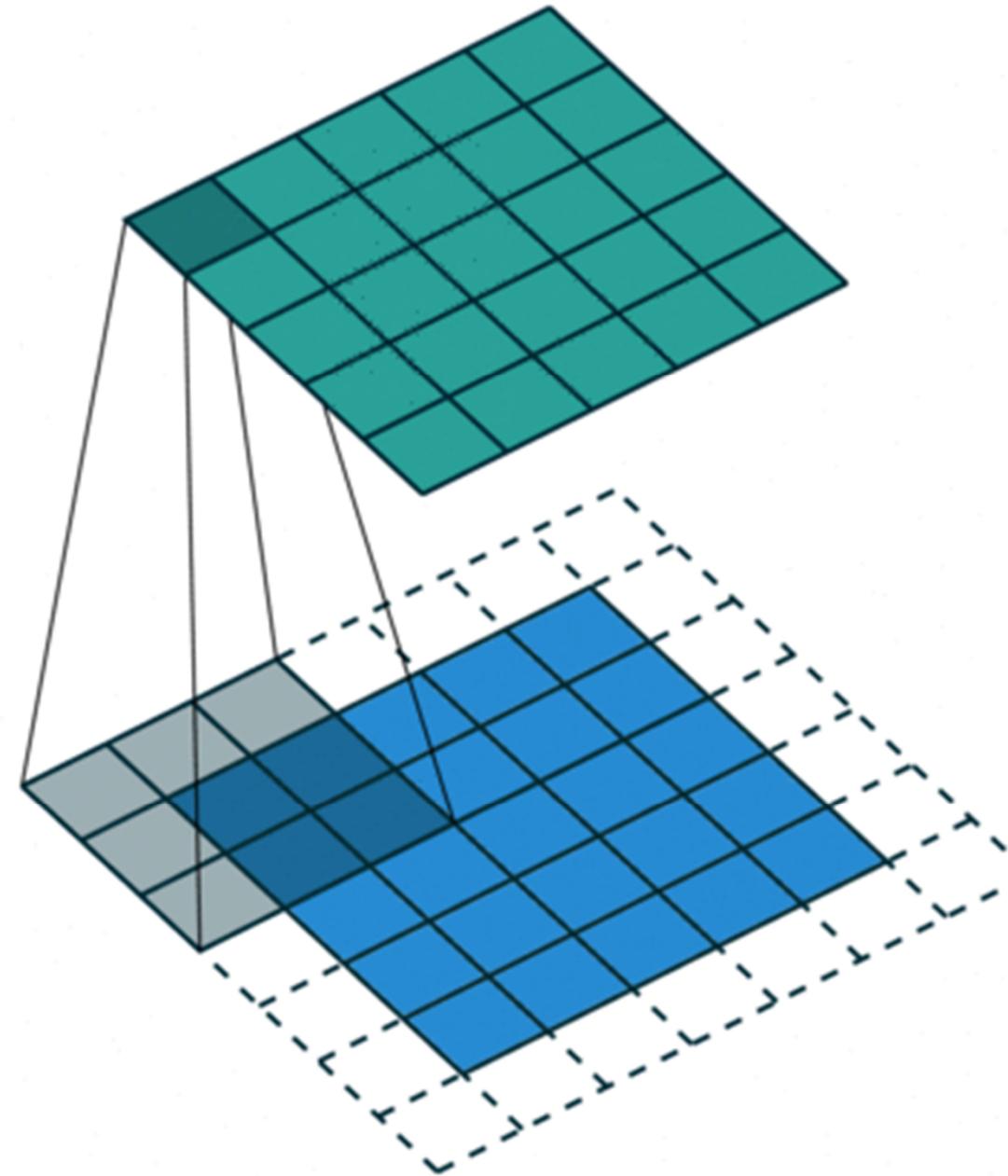


Pair each **downsampling** layer
with an **upsampling** layer

Up-sampling

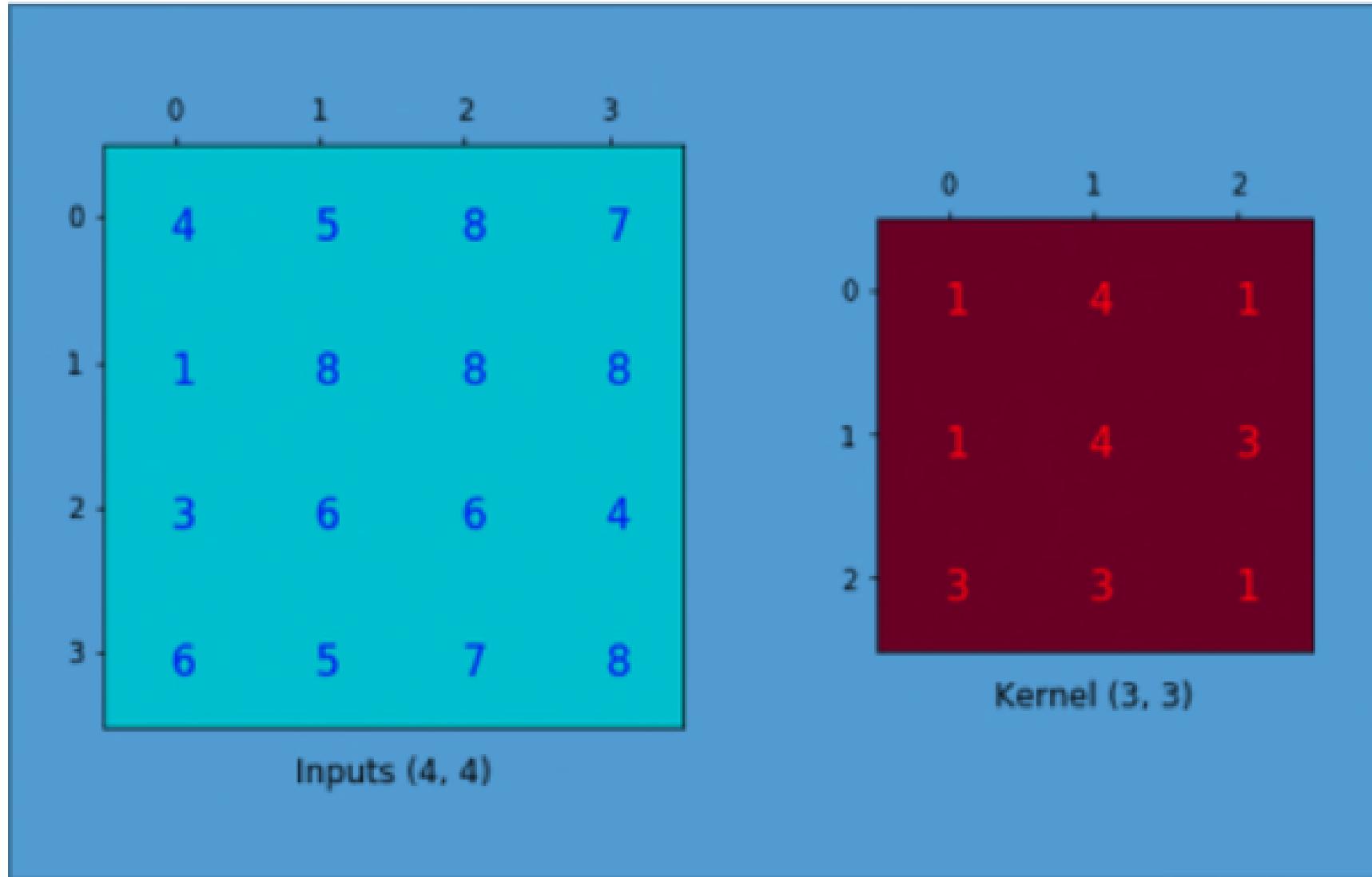
- There are various methods to conduct up-sampling operation:
- Interpolation
 - Nearest neighbor interpolation
 - Bi-linear interpolation
 - Bi-cubic interpolation
- All these methods involve some interpolation method which we need to chose when deciding a network architecture.
- It is like a **manual feature engineering** and there is **nothing that the network can learn about.**
- Convolutions
 - Transposed Convolutions
 - Dilated Convolutions

Normal convolutions



2D convolution using a kernel size of 3, stride of 1 and padding

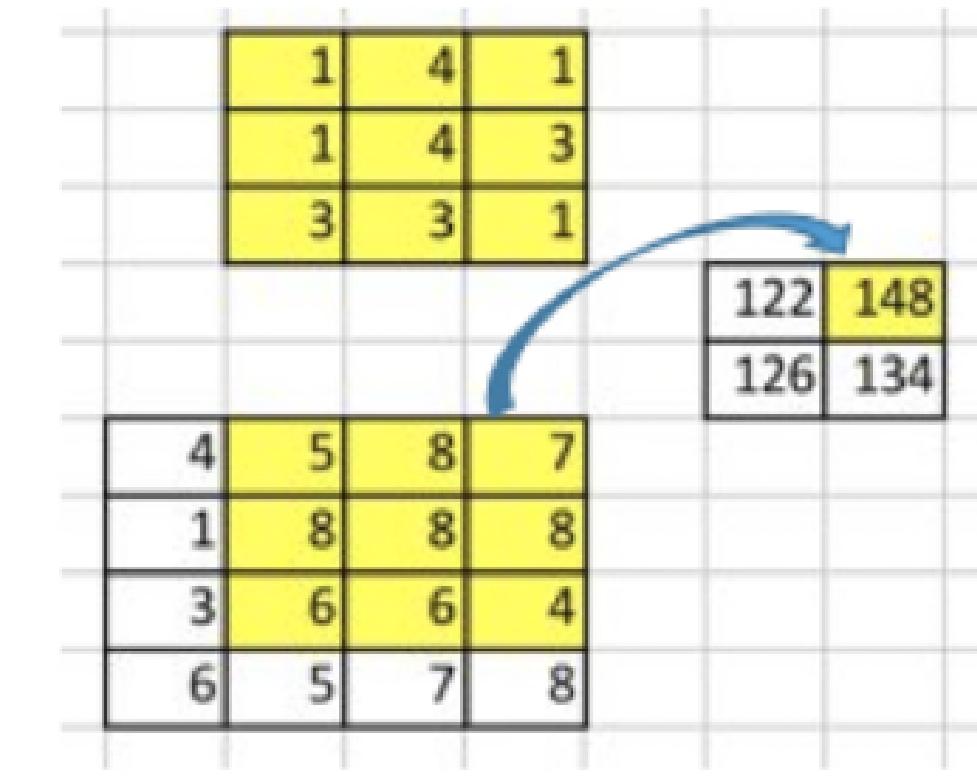
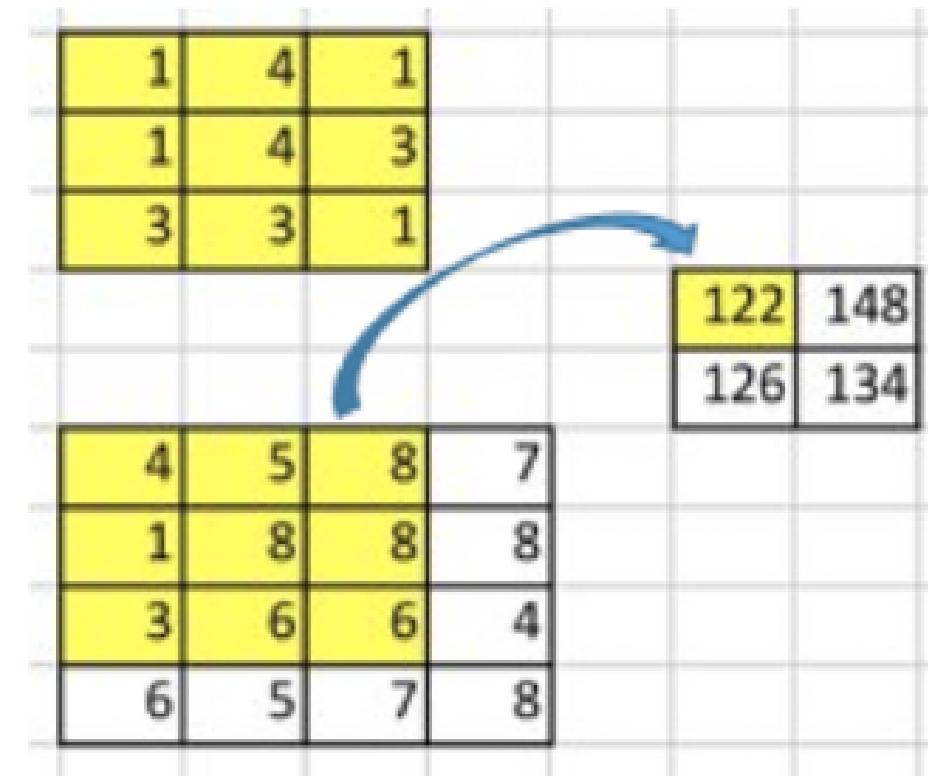
Normal convolutions



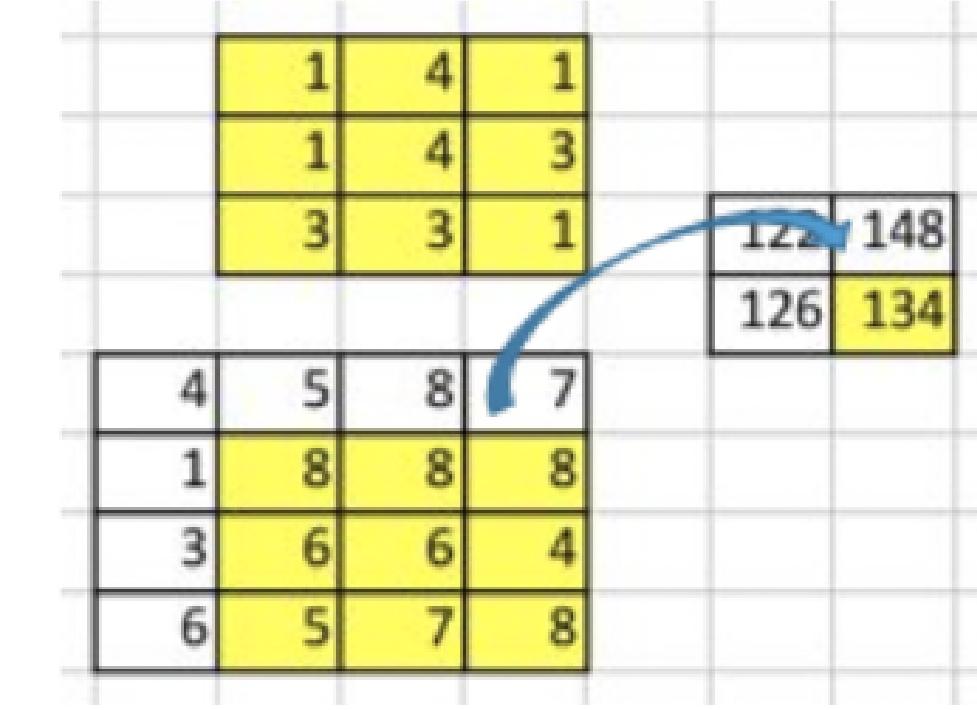
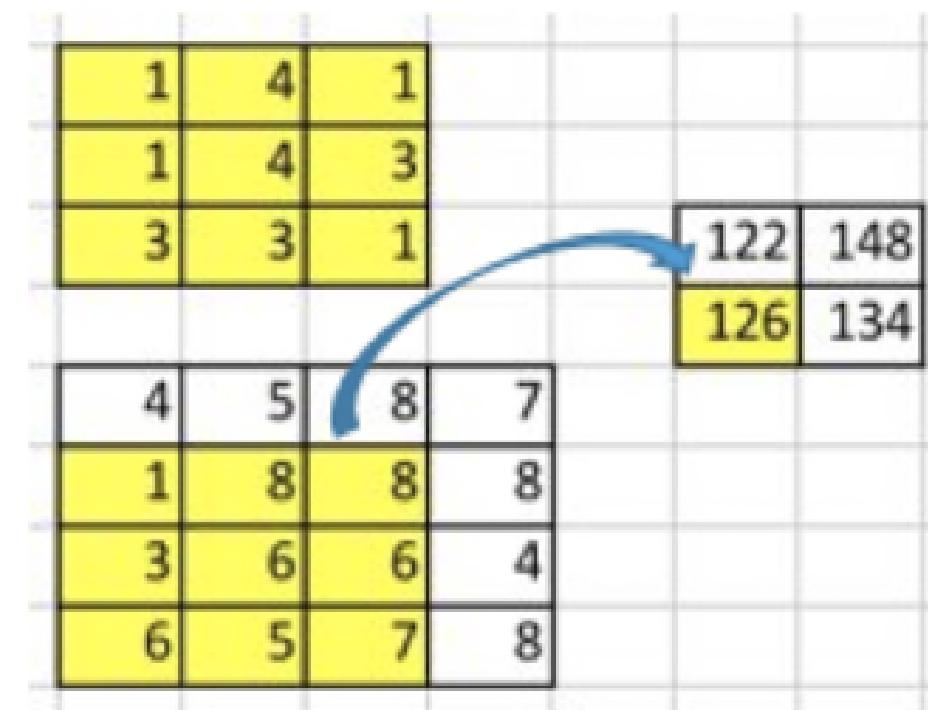
the output is a 2x2 matrix

we have a 4x4 matrix and apply a convolution operation on it with a 3x3 kernel, with no padding, and with a stride of 1

The convolution operation calculates the sum of the element-wise multiplication between the input matrix and kernel matrix.



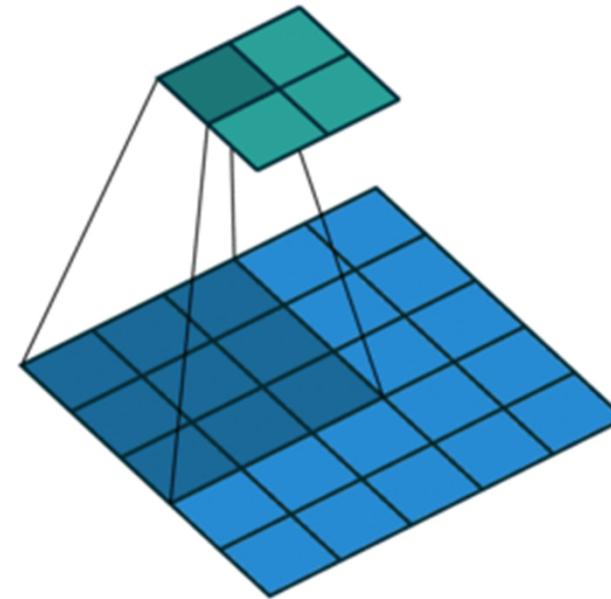
Since we have no padding and the stride of 1, we can do this only 4 times. Hence, the output matrix is 2x2.



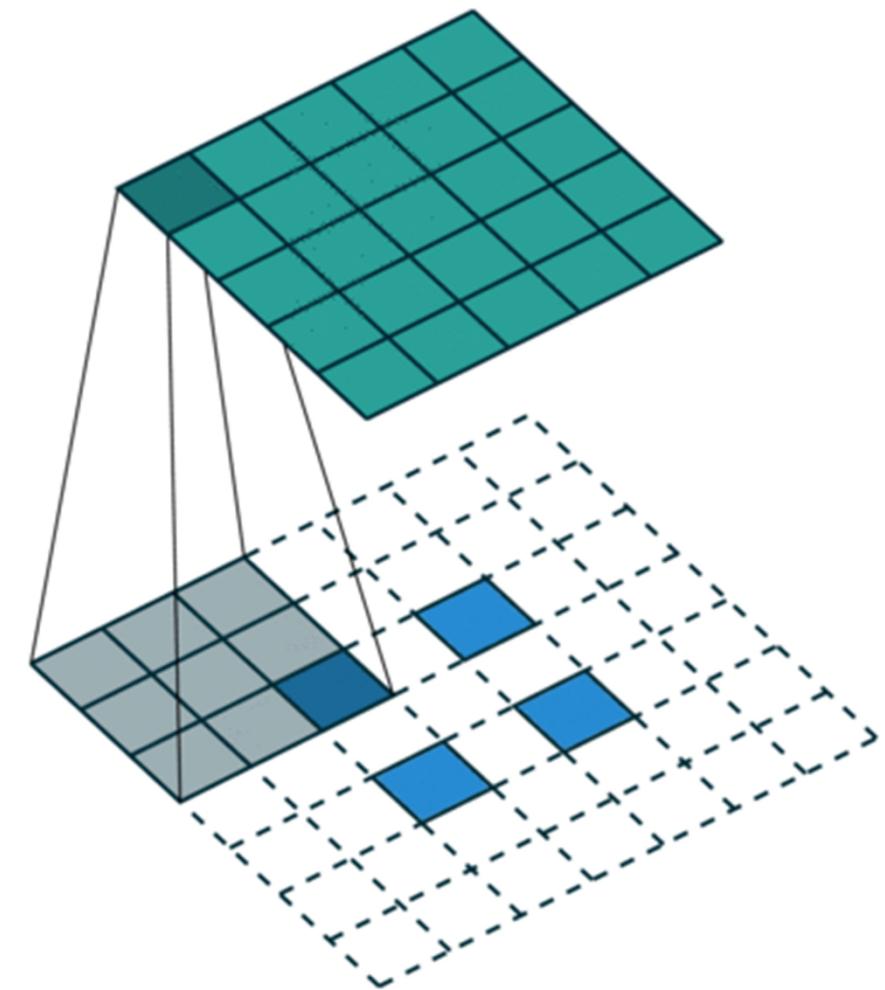
Normal convolutions : One important point

- One important point of such convolution operation is that the positional connectivity exists between the input values and the output values.
- For example, the top left values in the input matrix affect the top left value of the output matrix.
- More concretely, the 3×3 kernel is used to connect the 9 values in the input matrix to 1 value in the output matrix.
- **A convolution operation forms a many-to-one relationship.** Let's keep this in mind as we need it later on.

2D convolution with no padding,
stride of 2 and kernel of 3



Transposed 2D convolution with no padding, stride of 2 and kernel of 3



It guarantees that the output will be a 5×5 image as well, while still performing a normal convolution operation.

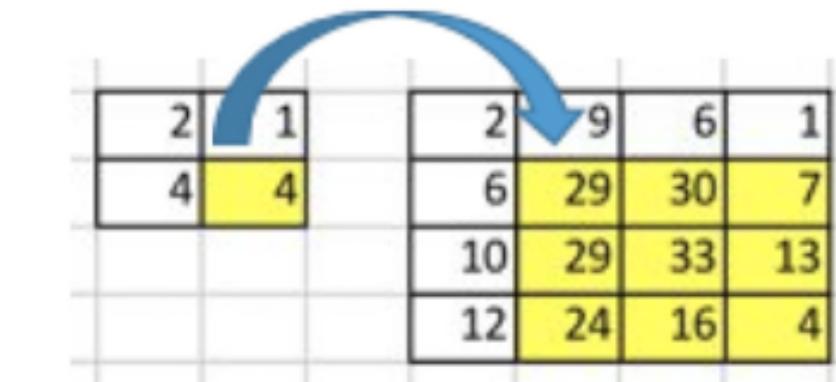
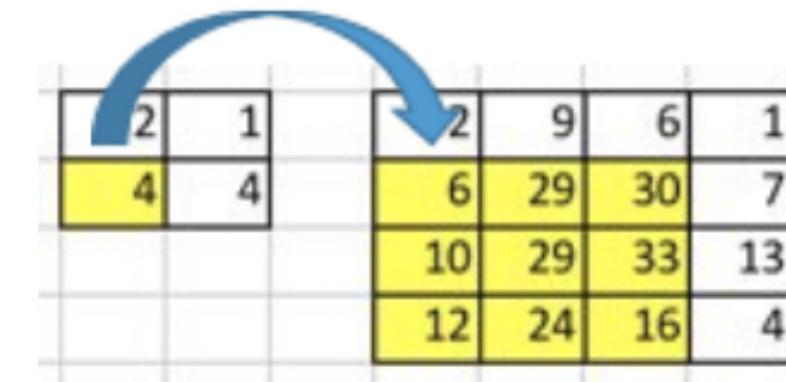
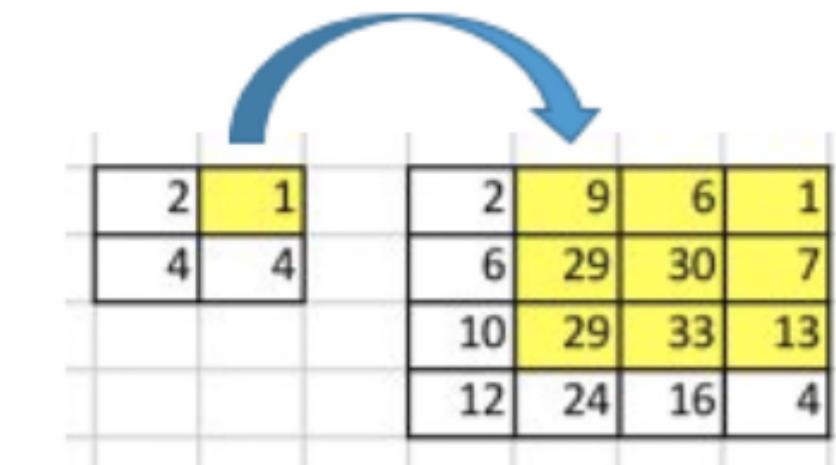
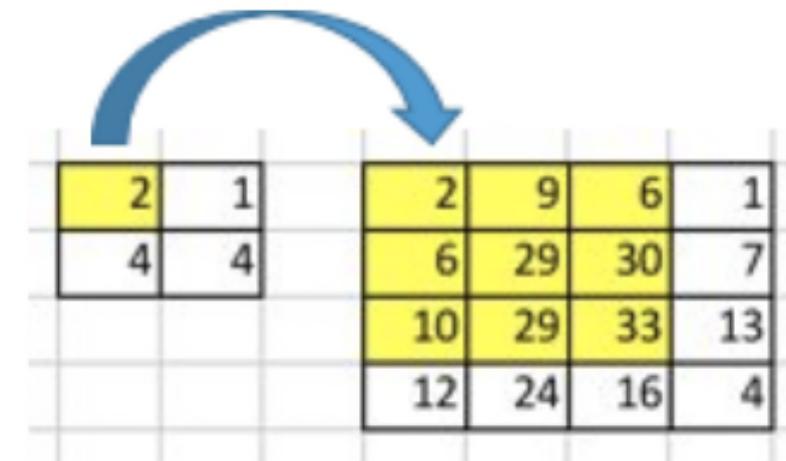
To achieve this, we need to perform some fancy padding on the input.

We want to associate 1 value in a matrix to 9 values in another matrix. It's a **one-to-many relationship**.

This is like going backward of convolution operation, and it is the core idea of **transposed convolution**

For example, we up-sample a 2x2 matrix to a 4x4 matrix.

The operation maintains the 1-to-9 relationship



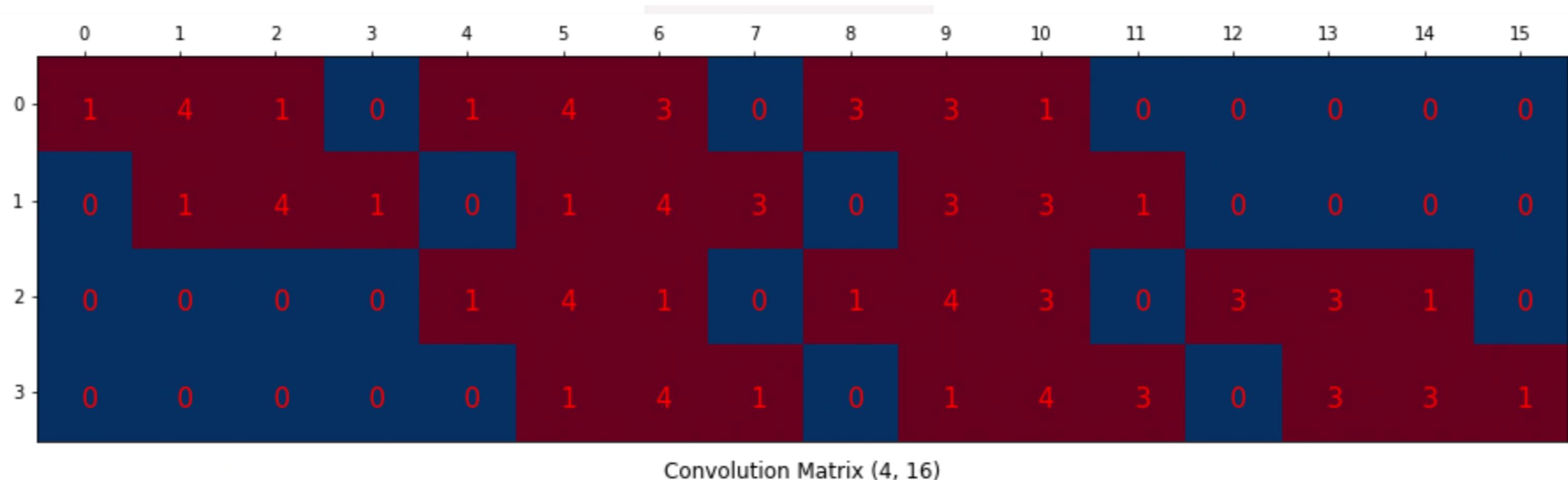
Transposed Convolutions

we need to define the **convolution matrix** and the **transposed convolution matrix**

We rearrange the 3×3 kernel into a 4×16 matrix as

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

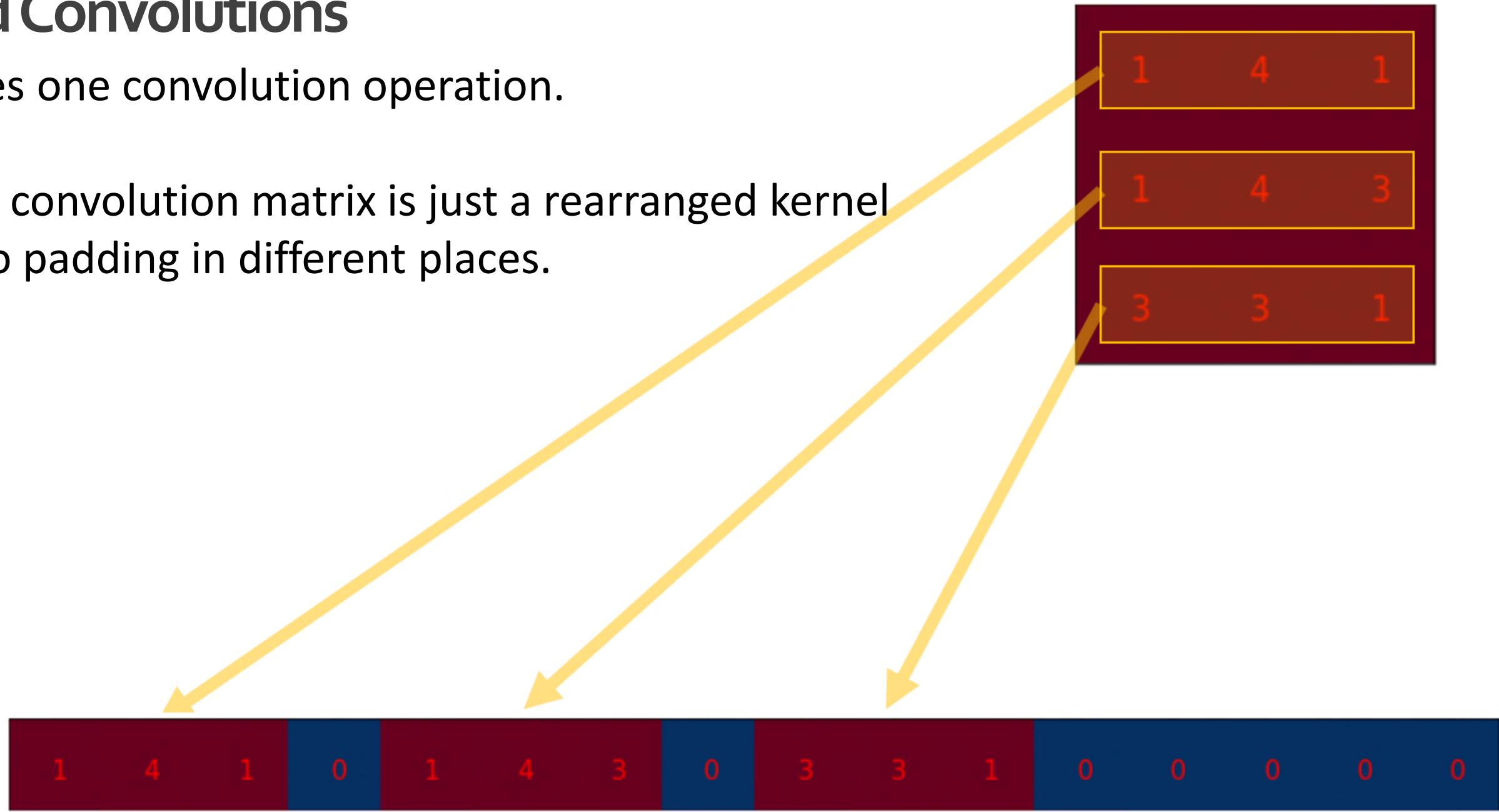
Kernel (3, 3)



Transposed Convolutions

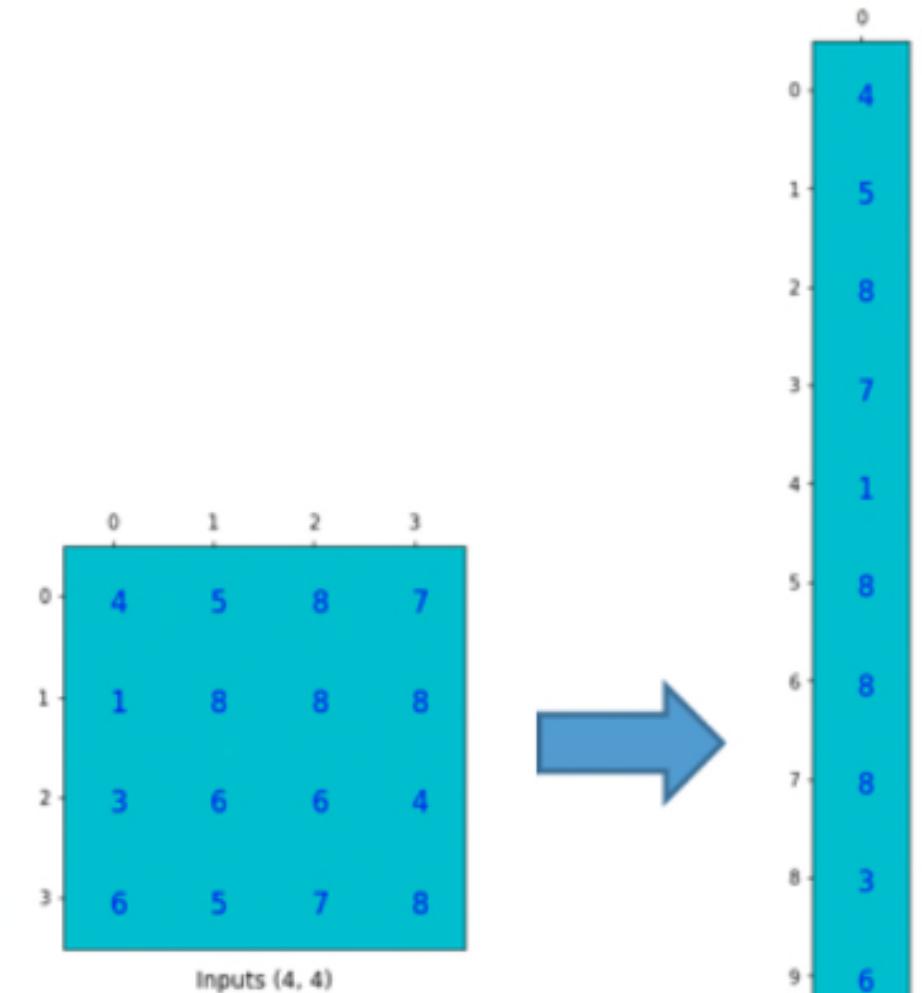
Each row defines one convolution operation.

Each row of the convolution matrix is just a rearranged kernel matrix with zero padding in different places.

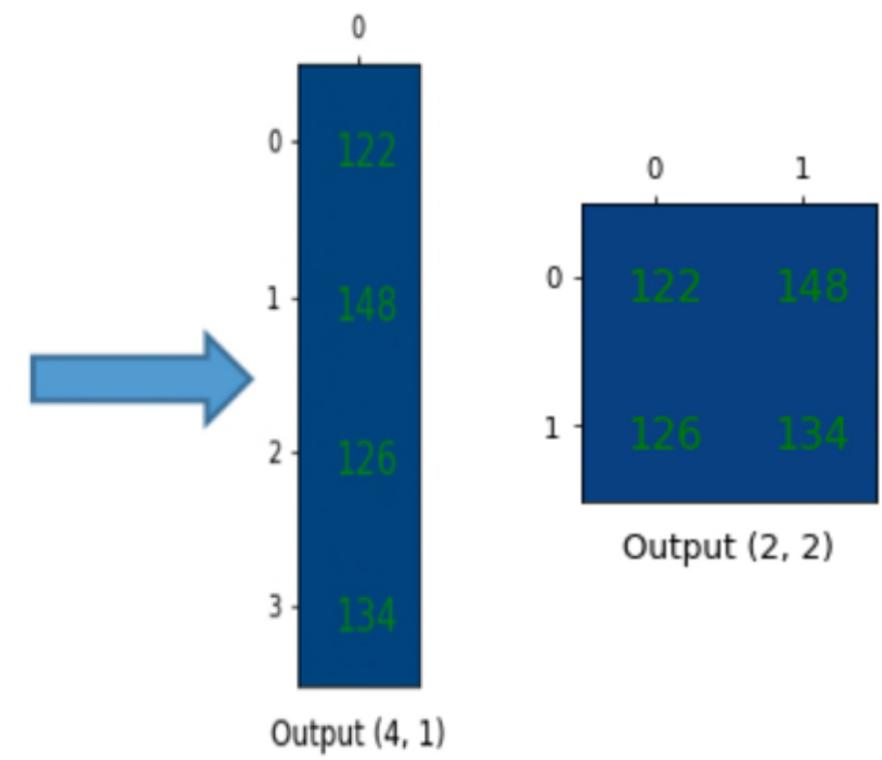
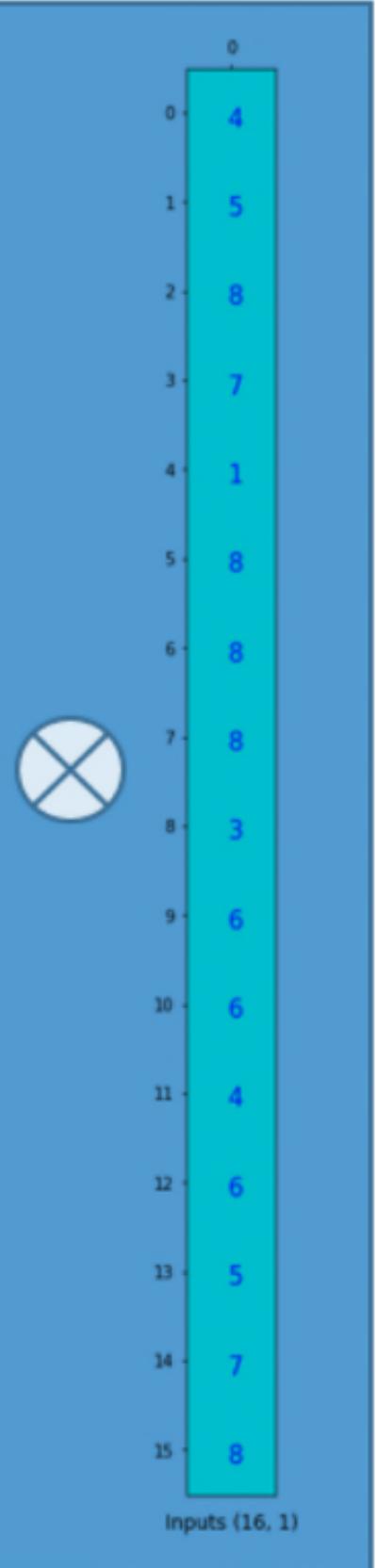
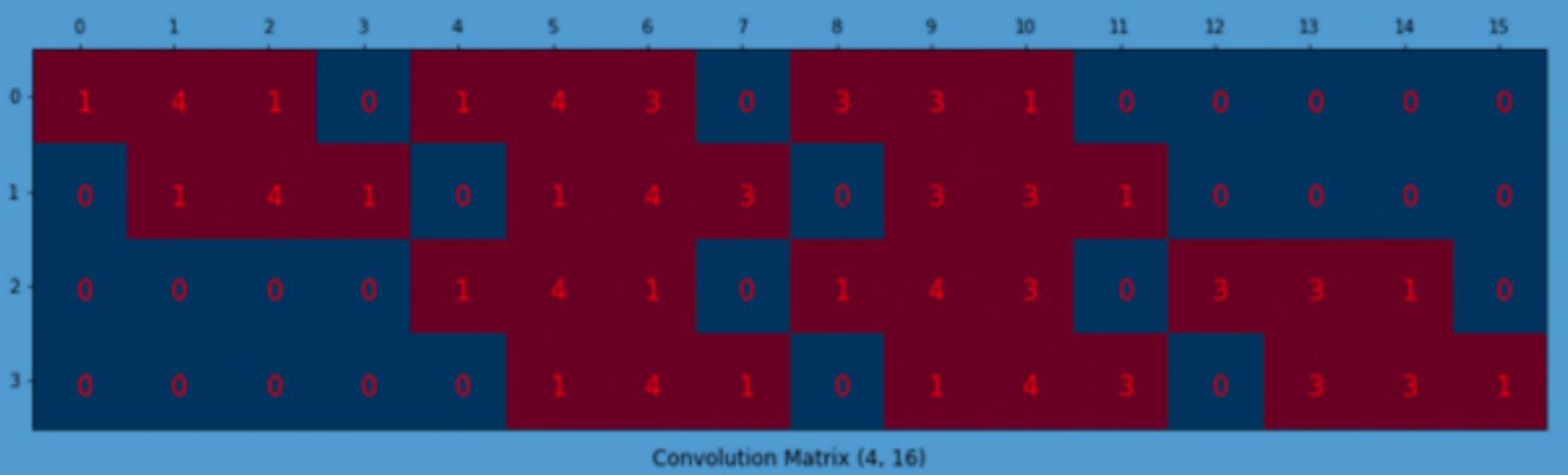


Transposed Convolutions

We flatten the input matrix (4×4) into a column vector (16×1)



We can matrix-multiply the 4×16 convolution matrix with the 16×1 input matrix (16 dimensional column vector)

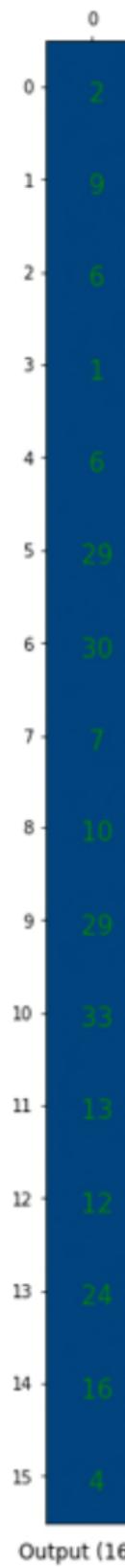
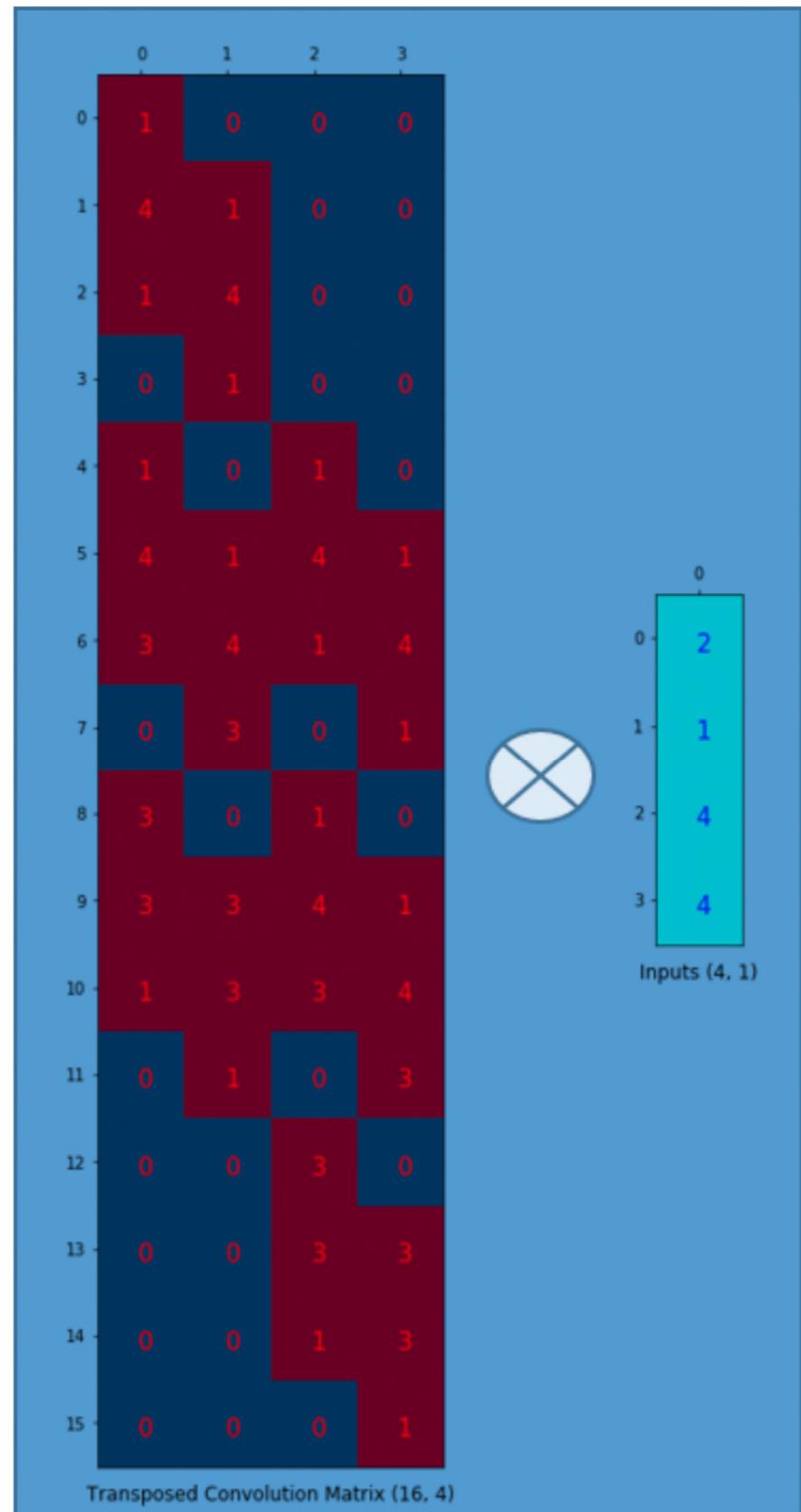


Transposed Convolutions

The point is that with the convolution matrix, you can go from 16 (4x4) to 4 (2x2) because the convolution matrix is 4x16.

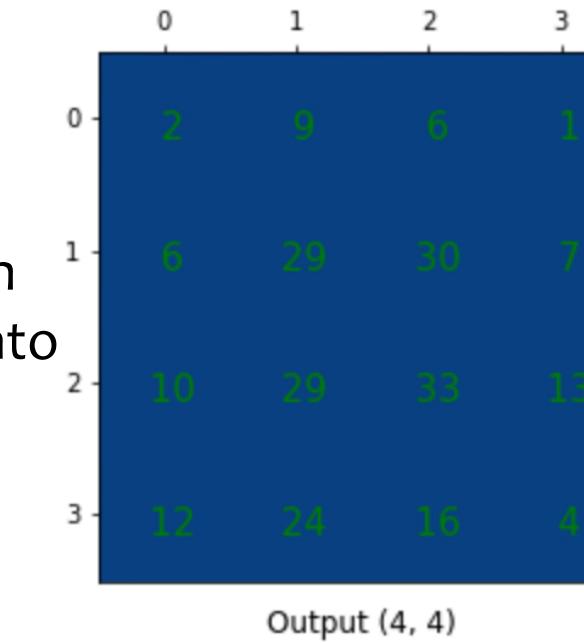
Then, if you have a 16x4 matrix, you can go from 4 (2x2) to 16 (4x4)

- We transpose the convolution matrix C (4x16) to C^T (16x4).
- We can matrix-multiply C^T (16x4) with a column vector (4x1) to generate an output matrix (16x1)
- The transposed matrix connects 1 value to 9 values in the output.



We have just up-sampled a smaller matrix (2x2) into a larger one (4x4).
The transposed convolution maintains the 1 to 9 relationship because of the way it lays out the weights.

The output can
be reshaped into
4x4



The actual weight values in the matrix does not have to come from the original convolution matrix. What's important is that the weight layout is transposed from that of the convolution matrix.

Transposed Convolutions

- Even though it is called the transposed convolution, it does not mean that we take some existing convolution matrix and use the transposed version.
- The main point is that the association between the input and the output is handled in the backward fashion compared with a standard convolution matrix (one-to-many rather than many-to-one association).

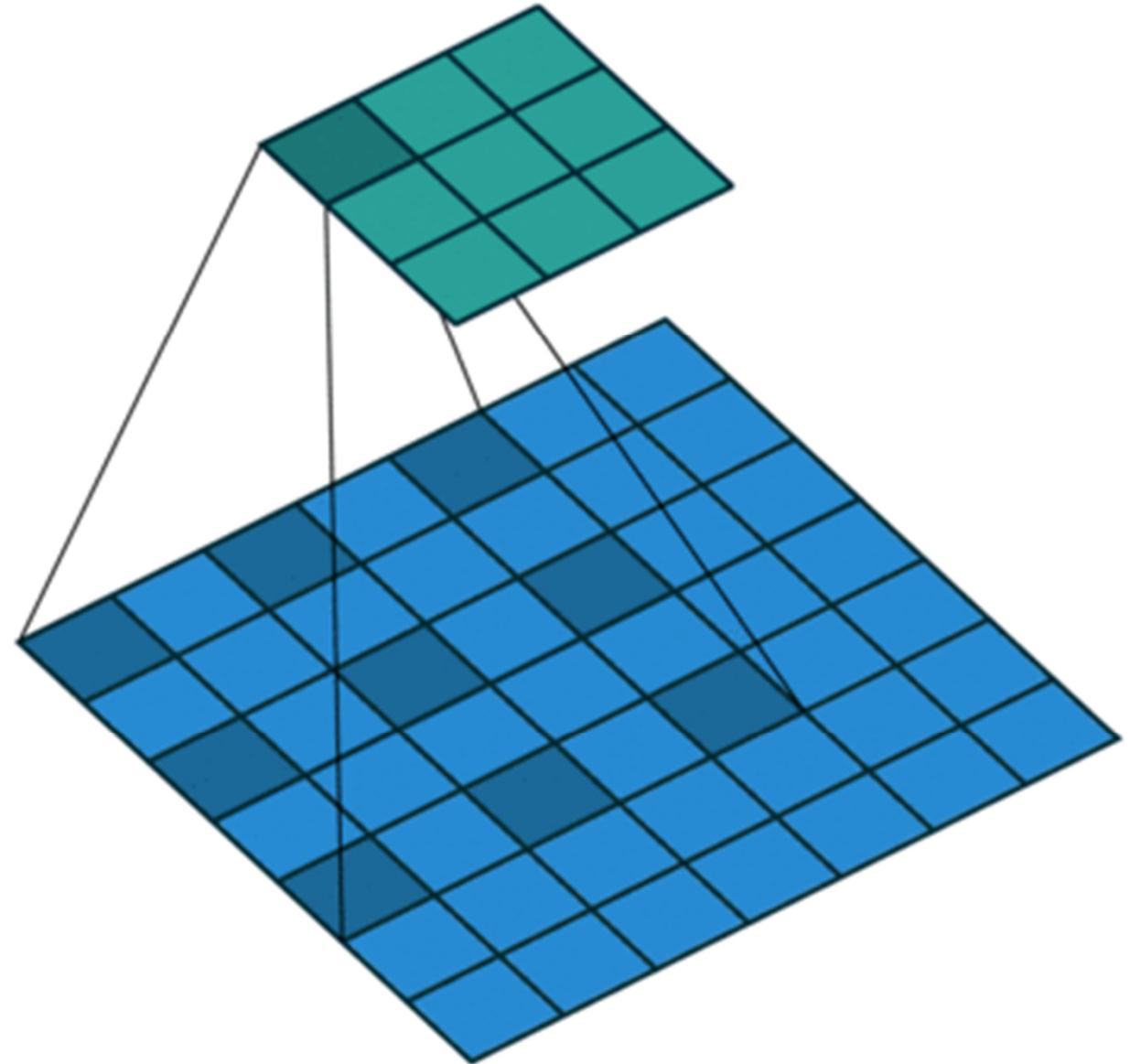
One caution: the transposed convolution is the cause of the checkerboard artifacts in generated images. [This article](#)* recommends an up-sampling operation (i.e., an interpolation method) followed by a convolution operation to reduce such issues. If your main objective is to generate images without such artifacts, it is worth reading the paper to find out more.

*Deconvolution and Checkerboard Artifacts

<https://distill.pub/2016/deconv-checkerboard/>

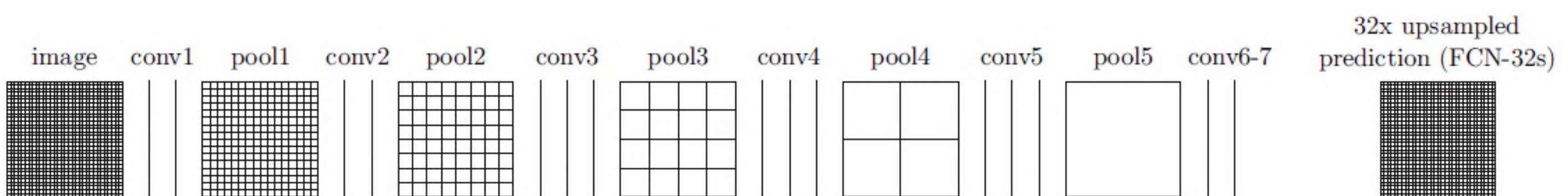
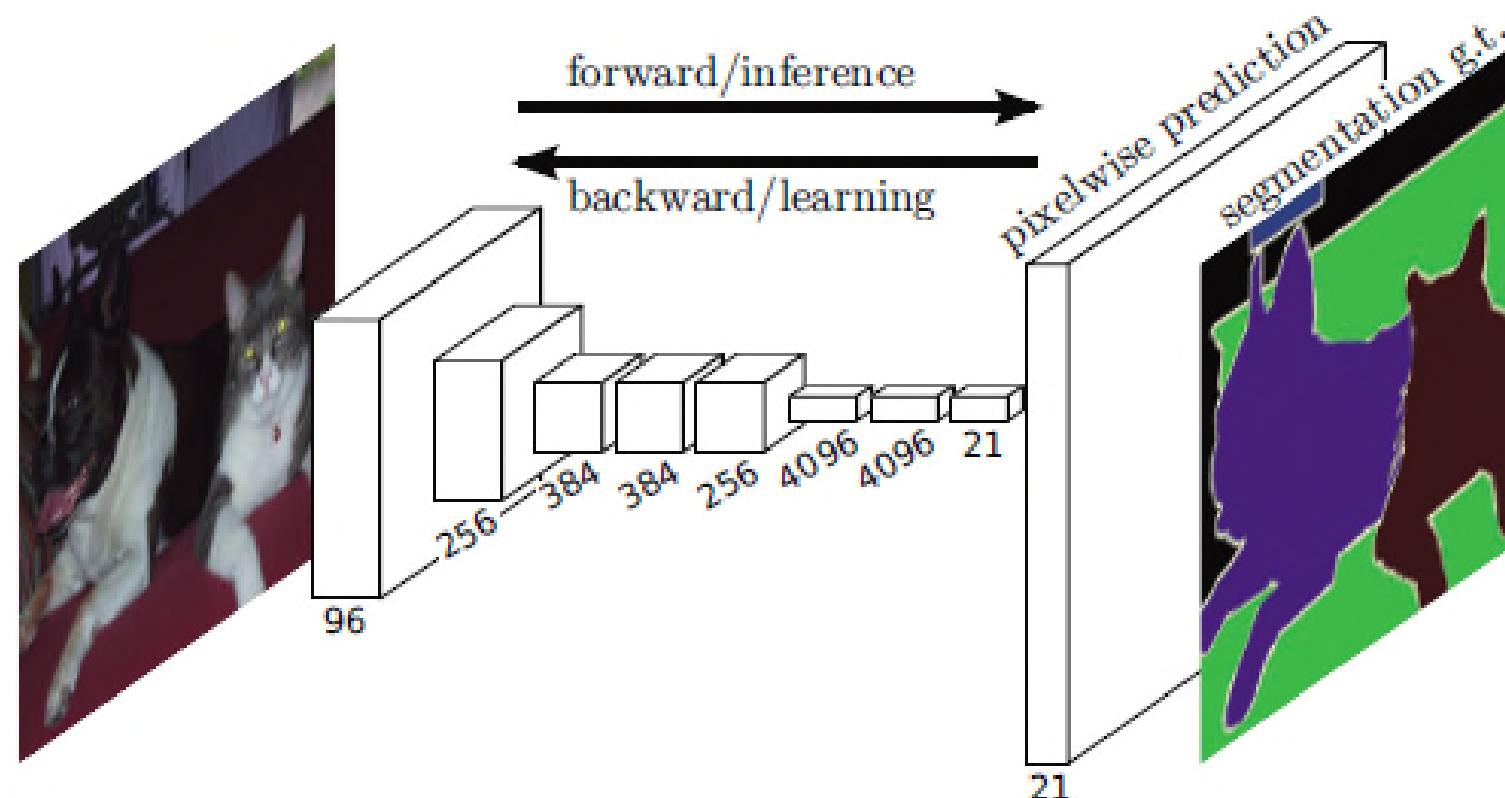
Dilated convolutions

- Dilated convolutions introduce another parameter to convolutional layers called the dilation rate.
- This defines a spacing between the values in a kernel. A 3×3 kernel with a dilation rate of 2 will have the same field of view as a 5×5 kernel, while only using 9 parameters.
(Imagine taking a 5×5 kernel and deleting every second column and row).
- This delivers a wider field of view at the same computational cost. Dilated convolutions are particularly popular in the field of real-time segmentation.
- Use them if you need a wide field of view and cannot afford multiple convolutions or larger kernels.



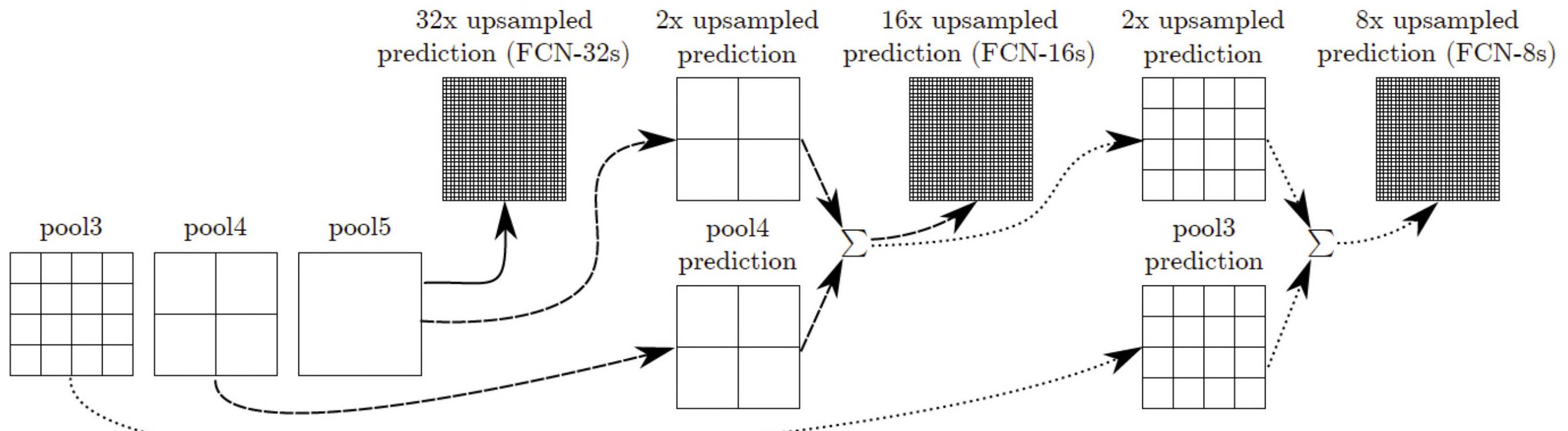
2D Atrous convolution using a 3×3 kernel with a dilation rate of 2 and no padding

Back to FCN

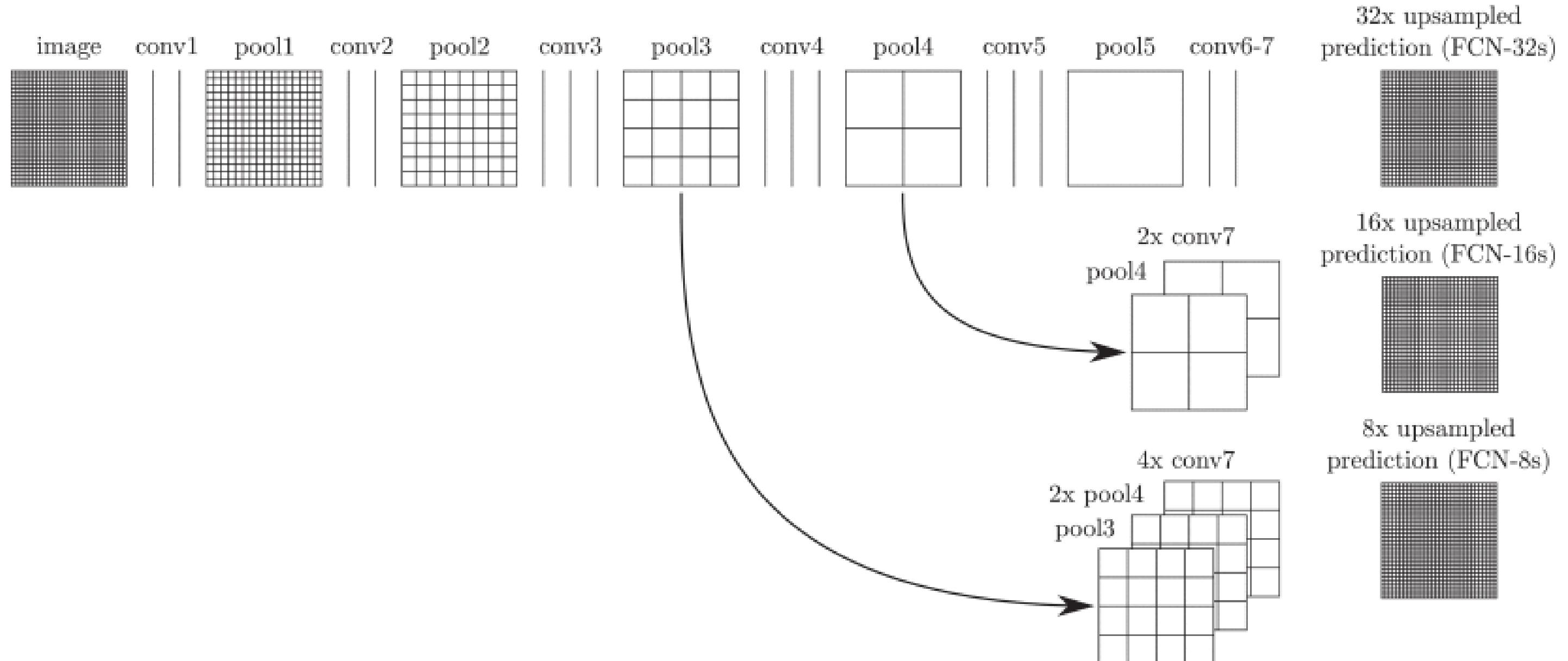


FCNs

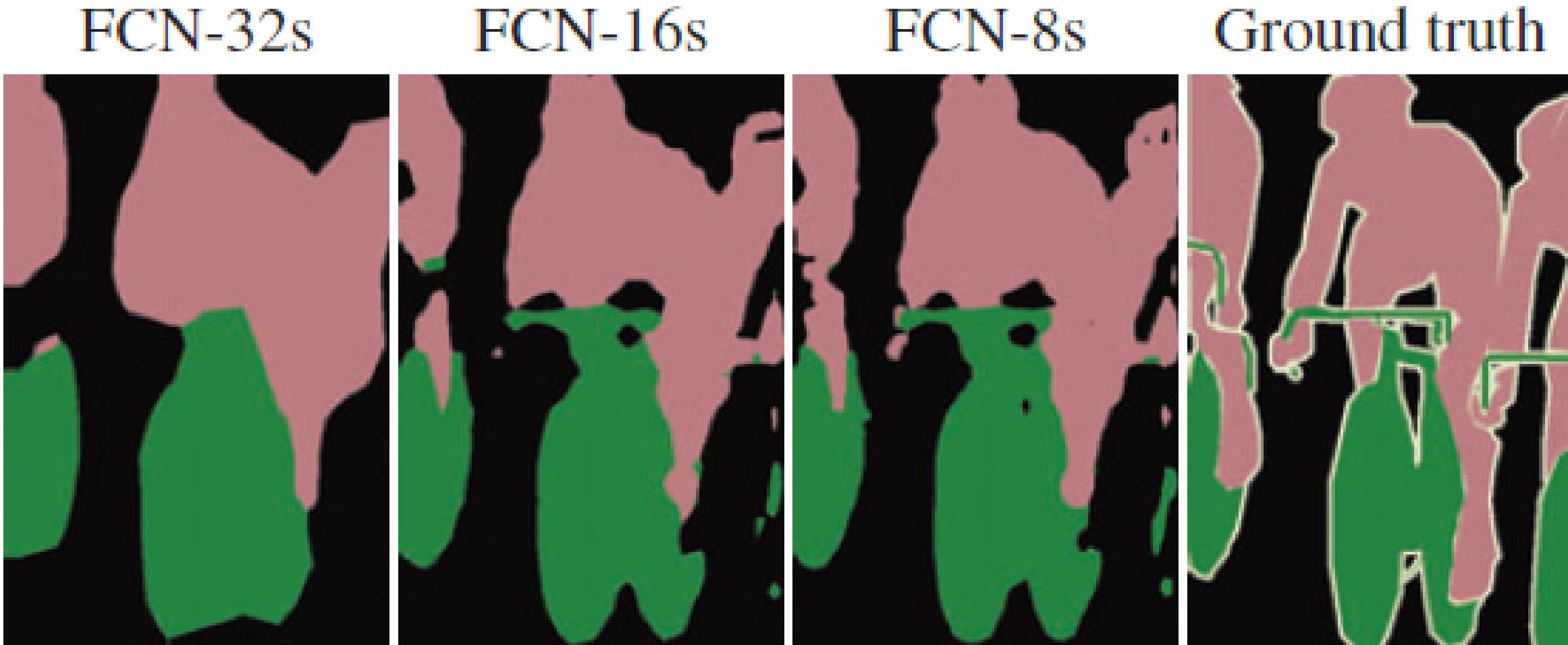
- Deep features can be obtained when going deeper, spatial location information is also lost when going deeper.
- That means output from shallower layers have more location information.
- If we combine both, we can enhance the result.
- To combine, we fuse the output (by element-wise addition):



FCNs



Visual Demonstration



Layers in Keras

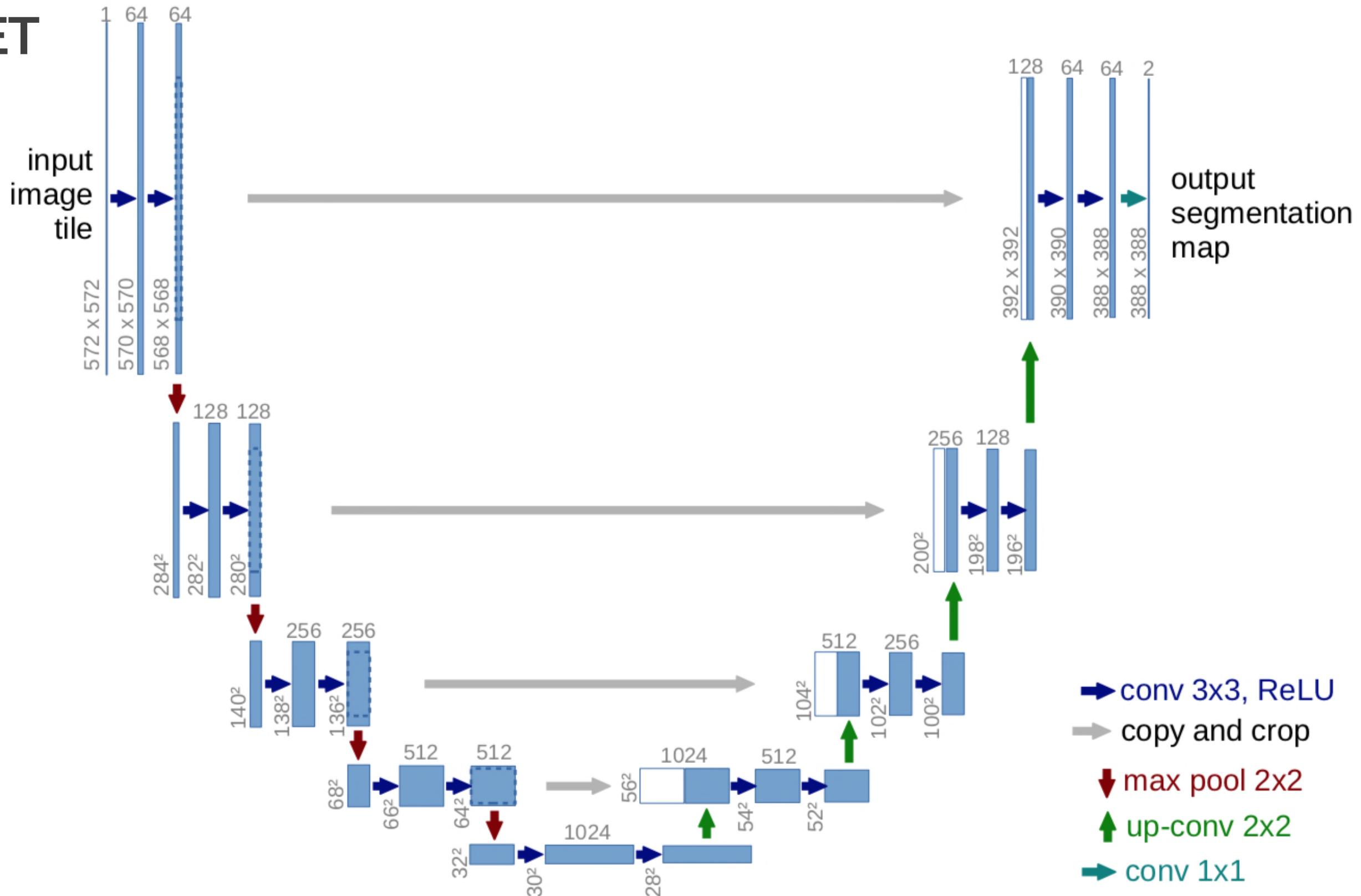
UpSampling2D is just a simple scaling up of the image by using nearest neighbour or bilinear upsampling, so nothing smart.

Advantage is it's cheap.

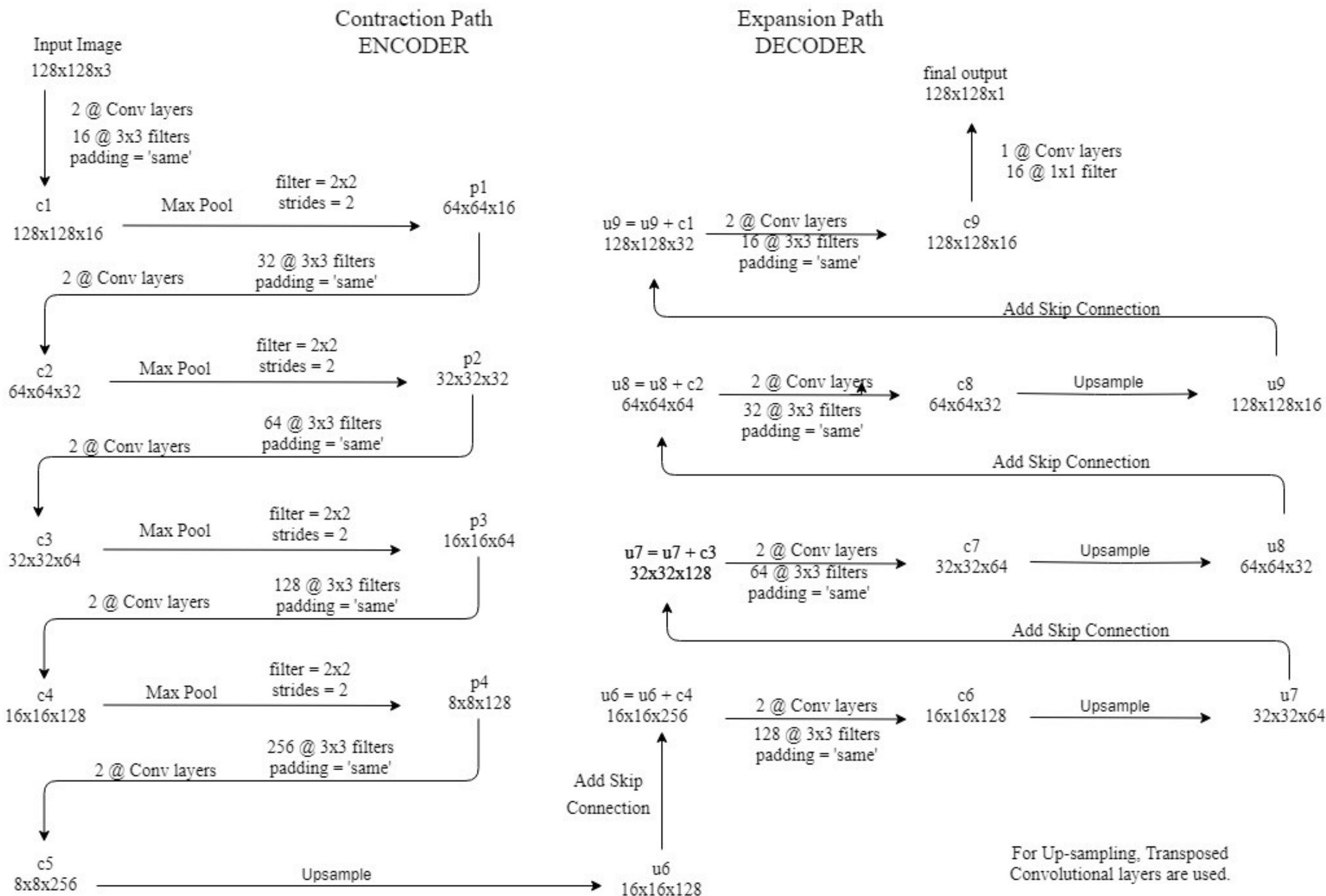
Conv2DTranspose is a convolution operation whose kernel is learnt (just like normal conv2d operation) while training your model

Interesting: <https://www.kaggle.com/mpalermo/remove-grideffect-on-generated-images/notebook>

UNET



UNET



The parameters for each Transpose Convolution are such that, the height and width of the image are doubled while the depth (no. of channels) is halved

UNET

- 2@Conv layers means that two consecutive Convolution Layers are applied
- c_1, c_2, \dots, c_9 are the output tensors of Convolutional Layers
- p_1, p_2, p_3 and p_4 are the output tensors of Max Pooling Layers
- u_6, u_7, u_8 and u_9 are the output tensors of up-sampling (transposed convolutional) layers
- The **left hand side is the contraction path (Encoder)** where we apply regular convolutions and max pooling layers.
- In the Encoder, the size of the image gradually reduces while the depth gradually increases. Starting from $128 \times 128 \times 3$ to $8 \times 8 \times 256$
- This basically means the network learns the “WHAT” information in the image, however it has lost the “WHERE” information
- **The right hand side is the expansion path (Decoder)** where we apply transposed convolutions along with regular convolutions
- In the decoder, the size of the image gradually increases and the depth gradually decreases. Starting from $8 \times 8 \times 256$ to $128 \times 128 \times 1$
- Intuitively, the Decoder recovers the “WHERE” information (precise localization) by gradually applying up-sampling
- To get better precise locations, at every step of the decoder we use **skip connections** by concatenating the output of the transposed convolution layers with the feature maps from the Encoder at the same level:
 $u_6 = u_6 + c_4$
 $u_7 = u_7 + c_3$
 $u_8 = u_8 + c_2$
 $u_9 = u_9 + c_1$
- After every concatenation we again apply two consecutive regular convolutions so that the model can learn to assemble a more precise output

UNet

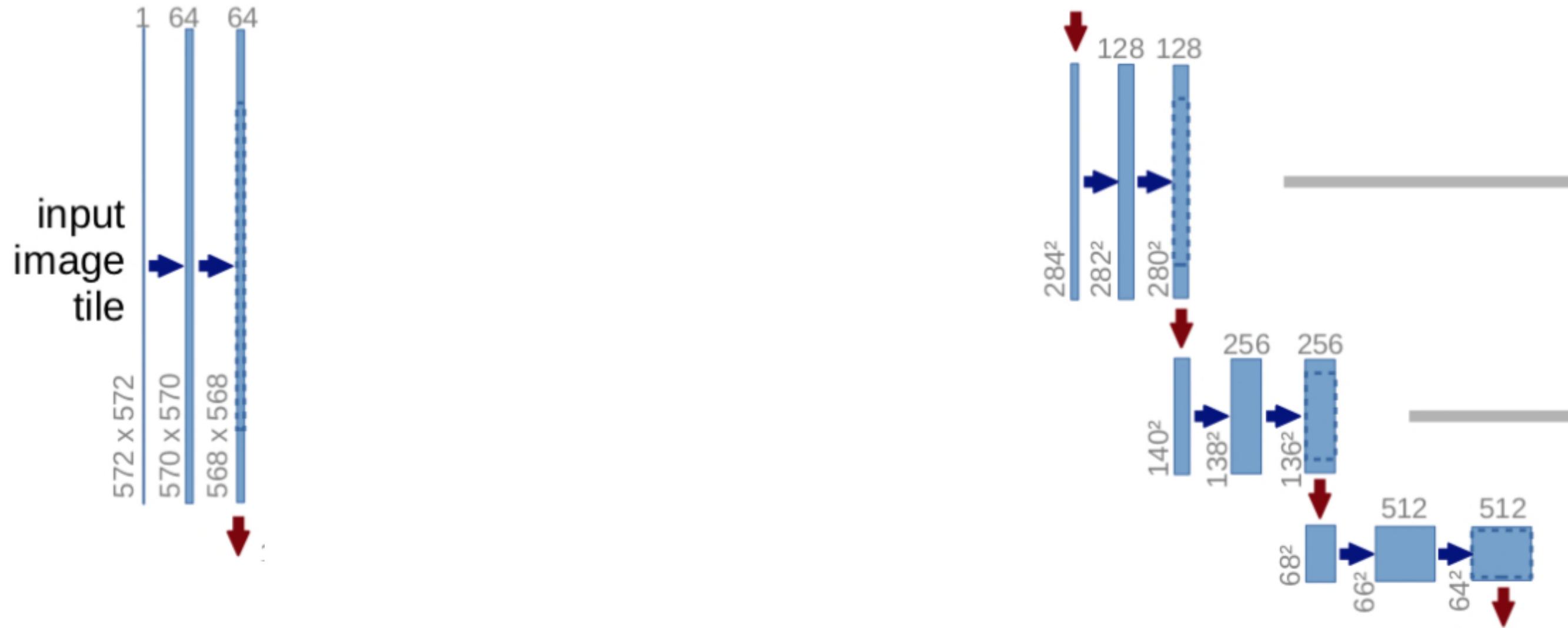
On a high level, we have the following relationship:

Input (128x128x1) => Encoder =>(8x8x256) => Decoder =>Output (128x128x1)

Implementation: <https://github.com/hlamba28/UNET-TGS>

Explanation:

<https://towardsdatascience.com/understanding-semantic-segmentation-with-unet-6be4f42d4b47>



```

1 conv1 = Conv2D(start_neurons * 1, (3, 3), activation="relu", padding="same")(input_layer)
2 conv1 = Conv2D(start_neurons * 1, (3, 3), activation="relu", padding="same")(conv1)
3 pool1 = MaxPooling2D((2, 2))(conv1)
4 pool1 = Dropout(0.25)(pool1)
    
```

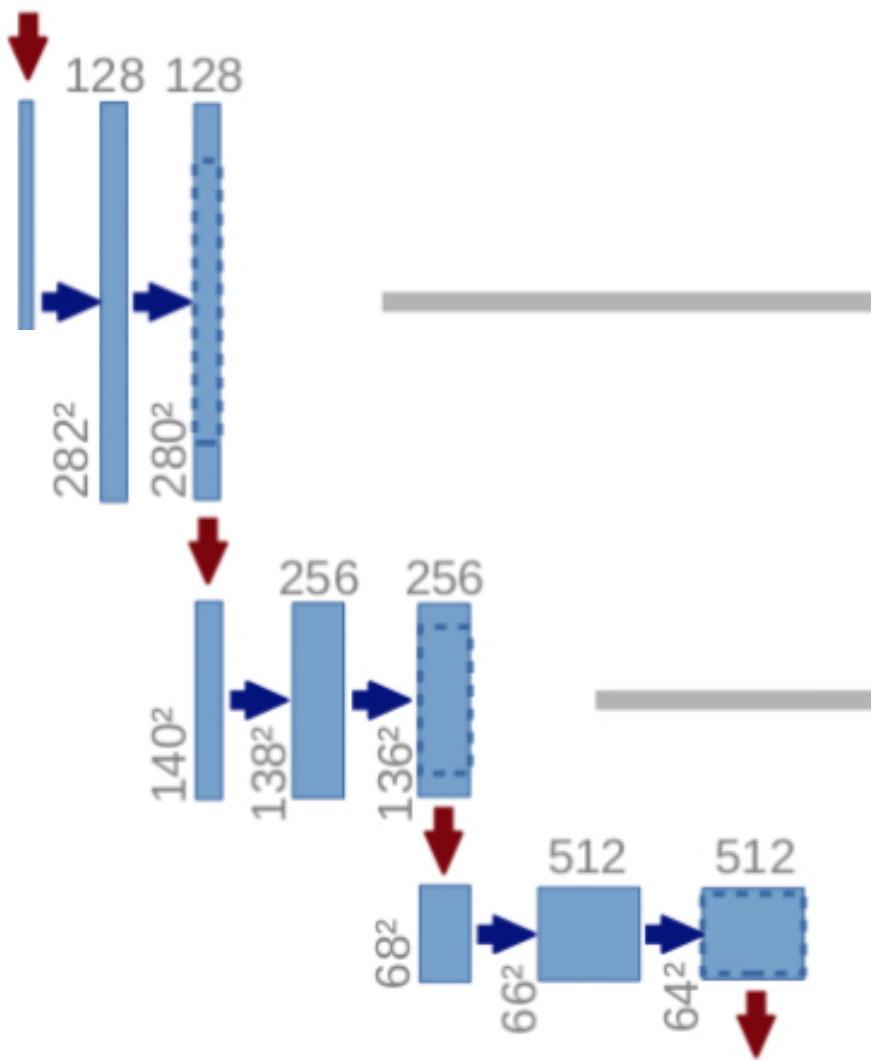
```

conv2 = Conv2D(start_neurons * 2, (3, 3), activation="relu", padding="same")(pool1)
conv2 = Conv2D(start_neurons * 2, (3, 3), activation="relu", padding="same")(conv2)
pool2 = MaxPooling2D((2, 2))(conv2)
pool2 = Dropout(0.5)(pool2)

conv3 = Conv2D(start_neurons * 4, (3, 3), activation="relu", padding="same")(pool2)
conv3 = Conv2D(start_neurons * 4, (3, 3), activation="relu", padding="same")(conv3)
pool3 = MaxPooling2D((2, 2))(conv3)
pool3 = Dropout(0.5)(pool3)

conv4 = Conv2D(start_neurons * 8, (3, 3), activation="relu", padding="same")(pool3)
conv4 = Conv2D(start_neurons * 8, (3, 3), activation="relu", padding="same")(conv4)
pool4 = MaxPooling2D((2, 2))(conv4)
pool4 = Dropout(0.5)(pool4)

```



Reading Material

- **Papers:**
 - A guide to convolution arithmetic for deep learning by Vincent Dumoulin, Francesco Visin at <https://arxiv.org/abs/1603.07285>
 - [2015 CVPR] [FCN]
[Fully Convolutional Networks for Semantic Segmentation](#)
 - [2017 TPAMI] [FCN]
[Fully Convolutional Networks for Semantic Segmentation](#)
 - **Understanding transposed convolutions**
 - <https://www.machinecurve.com/index.php/2019/09/29/understanding-transposed-convolutions/>

Reading Material

- Understanding UNET : Implementation
- <https://towardsdatascience.com/understanding-semantic-segmentation-with-unet-6be4f42d4b47>
- Line by Line Explanation
- <https://towardsdatascience.com/unet-line-by-line-explanation-9b191c76baf5>

Acknowledgments

Various contents in this presentation have been taken from different books, lecture notes, and the web. These solely belong to their owners, and are here used only for clarifying various educational concepts. Any copyright infringement is not intended