

Computer Vision

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للغعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

February 7, 2025

Table of Contents

1. Introduction
2. Image Segmentation
3. Adapting CNNs to Segmentation Tasks
4. Upsampling Operations
5. Residual Connections and U-Net
6. Instance and Panoptic Segmentation

Learning Outcomes

- ▶ Understand the fundamentals of image segmentation and its importance.
- ▶ Understand how Convolutional Neural Networks (CNNs) are adapted for segmentation tasks.
- ▶ Understand different upsampling techniques used in segmentation models.
- ▶ Understand the role of residual connections and the U-Net architecture in segmentation.
- ▶ Differentiate between instance segmentation and panoptic segmentation.

Image Classification

- ▶ Previously, we discussed Image Classification
- ▶ A core task in Computer Vision



This image by Nikita is
licensed under CC-BY 2.0

(assume given a set of possible labels)
{dog, cat, truck, plane, ...}



cat

Computer Vision Tasks



Classification



CAT

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No spatial extent

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

Mask
is Binary
Segmentation

[This image is CC0 public domain](#)

* Localization is finding coordinate in space

uses argmax after softmax
p

* Granularity is

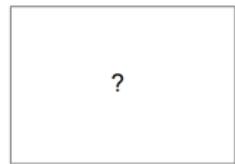
bind boundaries

Semantic Segmentation



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.

*class level
for each
pixel*

Semantic Segmentation

[Full image](#)



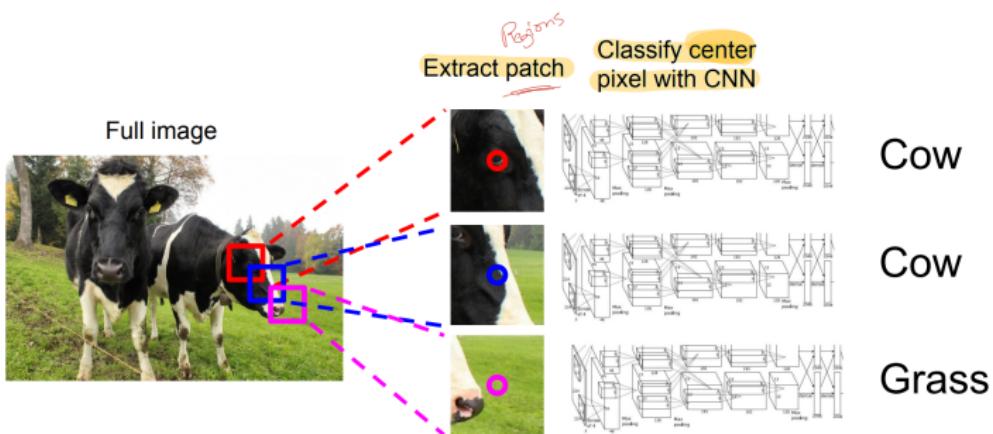
Semantic Segmentation

Full image



- ▶ Impossible to classify without context
- ▶ How do we include context?

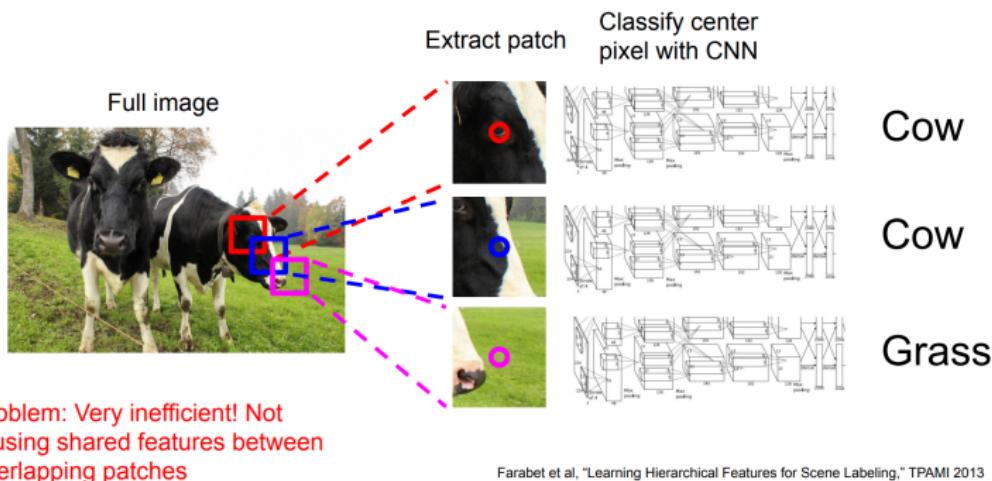
Semantic Segmentation Idea: Sliding Window



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 (cont.)

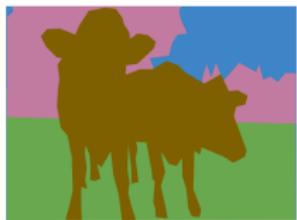
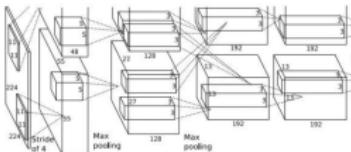


Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

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

Semantic Segmentation Idea: Convolution

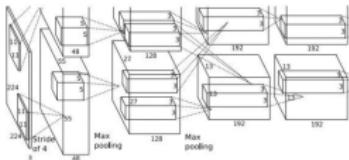
Full image



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

Semantic Segmentation Idea: Convolution (cont.)

Full image

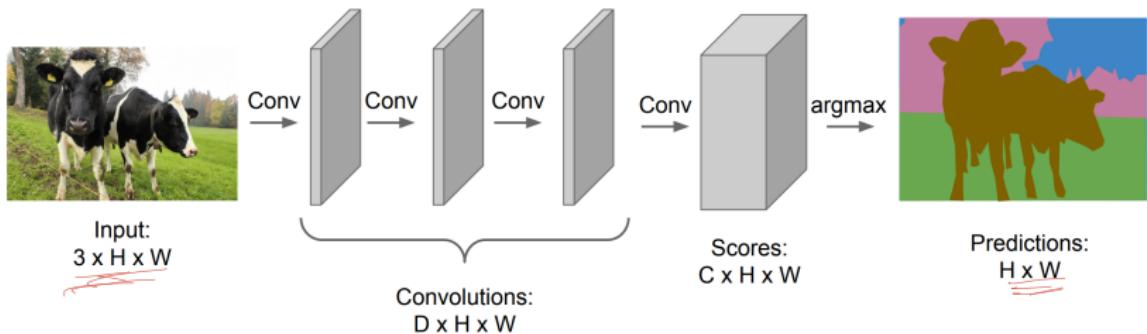


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

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.

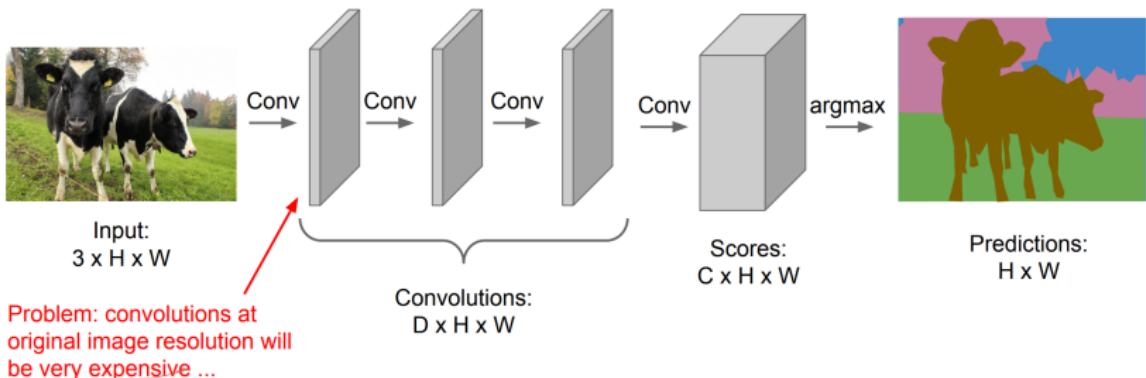
Semantic Segmentation Idea: Fully Convolutional

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



Semantic Segmentation Idea: Fully Convolutional (cont.)

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



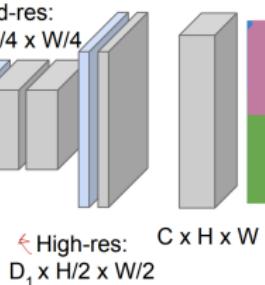
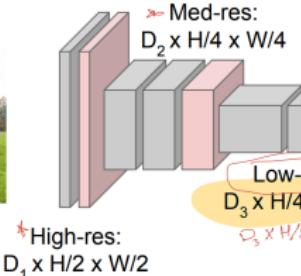
Semantic Segmentation Idea: Fully Convolutional (cont.)



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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

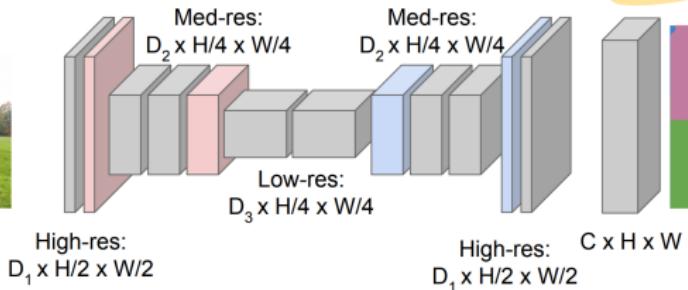
Semantic Segmentation Idea: Fully Convolutional (cont.)

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
max
Truncation



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network Upsampling: Unpooling



Done (transposed)
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

Repeated pixels

② “Bed of Nails”

1	2
3	4

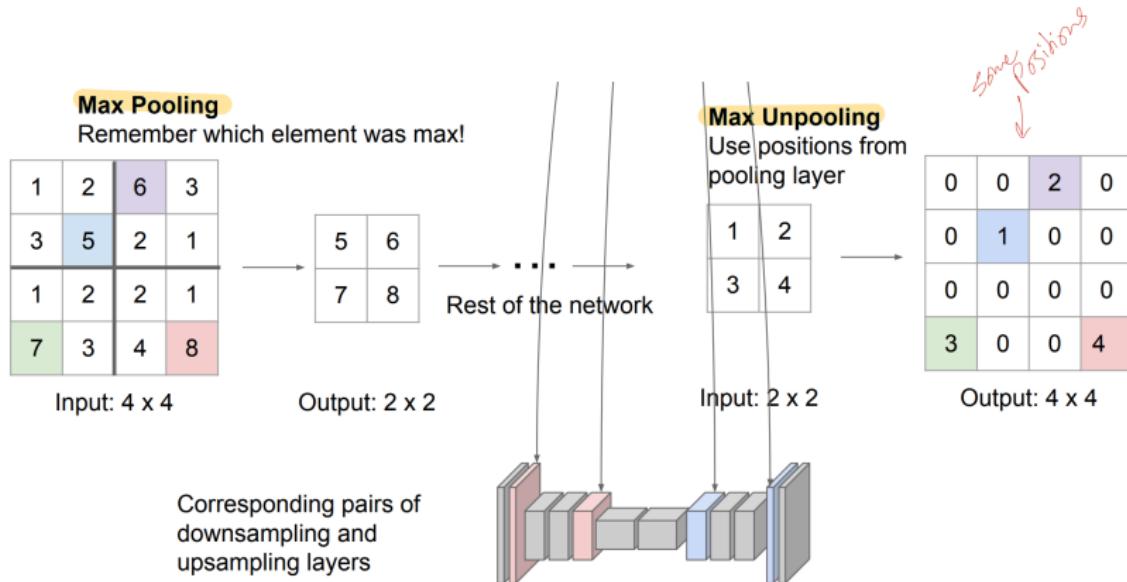
Input: 2 x 2

Zeros pixels

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

Output: 4 x 4

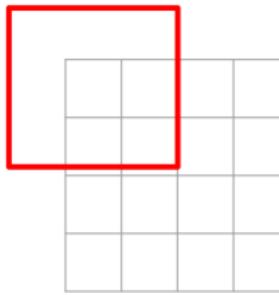
In-Network Upsampling: Max Unpooling



Learnable Upsampling: Transposed Convolution

Reverse
Stride

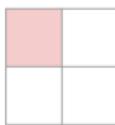
Recall: Normal 3×3 convolution, stride 2 pad 1



Input: 4×4

① Padding
② Filter

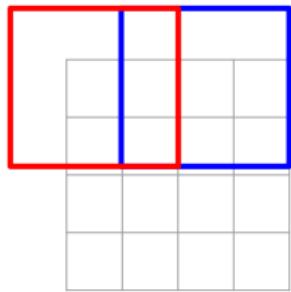
Dot product
between filter
and input



Output: 2×2

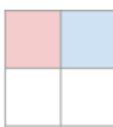
Learnable Upsampling: Transposed Convolution (cont.)

Recall: Normal 3×3 convolution, stride 2 pad 1



Input: 4×4

Dot product
between filter
and input



Output: 2×2

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

Stride gives ratio between
movement in input and
output

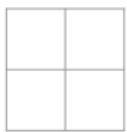
We can interpret strided
convolution as “learnable
downsampling”.

Learnable Upsampling: Transposed Convolution (cont.)

Stride + Pad
 $2+1=3$

① Filter Padding
② Padding

3 x 3 transposed convolution, stride 2 pad 1



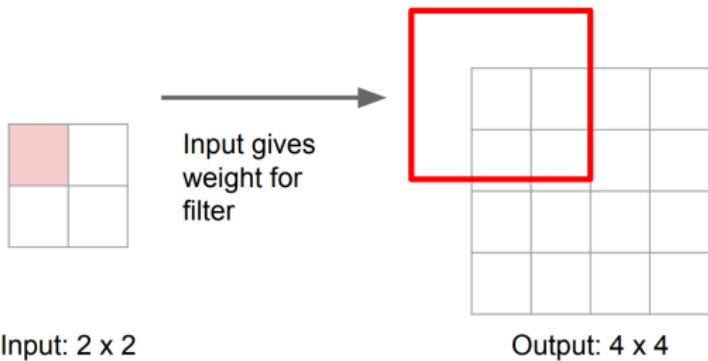
Input: 2 x 2



Output: 4 x 4

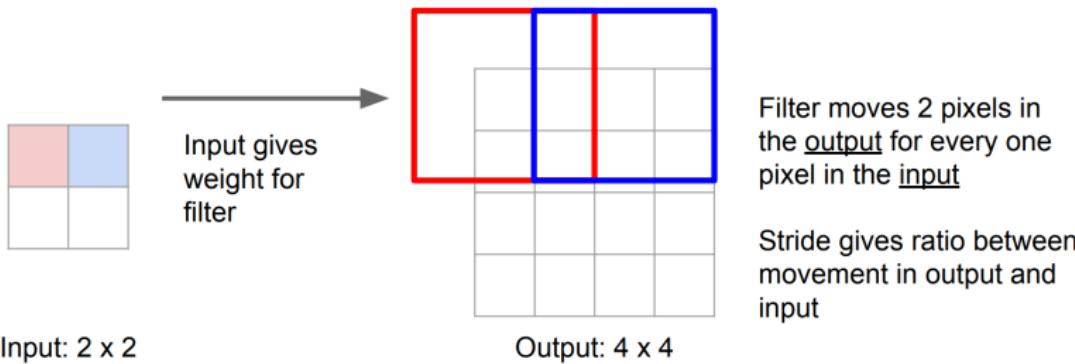
Learnable Upsampling: Transposed Convolution (cont.)

3 x 3 transposed convolution, stride 2 pad 1

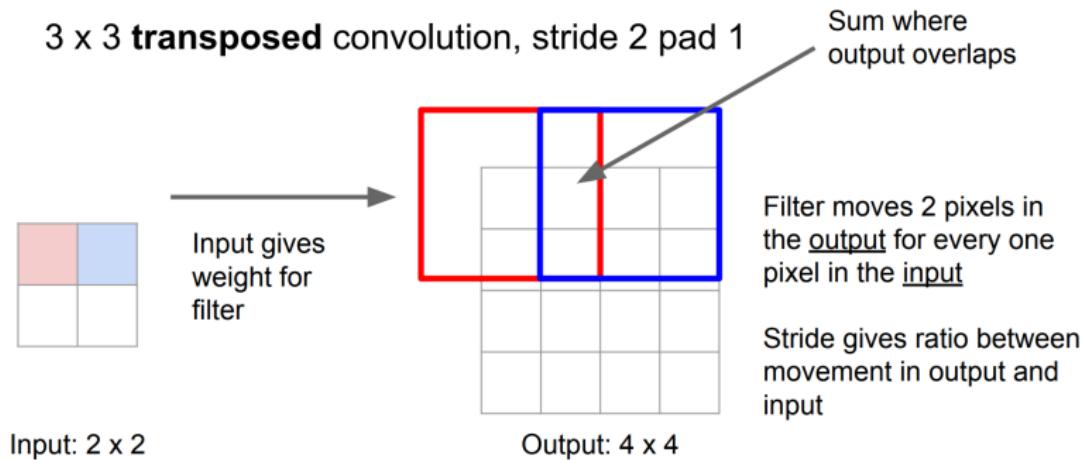


Learnable Upsampling: Transposed Convolution (cont.)

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



Learnable Upsampling: Transposed Convolution (cont.)



~~ext
upsampling
(Transpose)~~

5	1
3	7

input = 2×2

Kernel = 3×3

Stride = 1



1	2	3	4
5	$5+1=6$ $=6$	$5+1=6$ $=6$	1
$5+3=8$	$5+1+3=9$ $=9$ $+7=16$	$5+1+3=9$ $=9$ $+7=16$	$1+7=8$ $=8$
$5+3=8$	$5+1+3=9$ $=9$ $+7=16$	$5+1+3=9$ $=9$ $+7=16$	$1+7=8$ $=8$
3	$3+7=10$ $=10$	$3+7=10$ $=10$	7

Transposed Convolution Formula :

↑ input feature map
↑ stride
↑ Kernel size
→ padding

$$\text{output size} = (\text{Input } n - 1) * S + F - 2P$$

$$= (2 - 1) * 1 + 3
= 1 + 3 = \boxed{4}$$

Transposed

5	6	6	1
8	16	16	8
8	16	16	8
3	10	10	7

4×4

* There is some slides missing afterward...
check Github

Ex 2

1	2
5	8

input = 2×2

1	1	1
1	1	1
1	1	1

Kernel = 3×3

Stride = 2



1	1	$1+2$	2	2
1	1	$1+2$	2	2
$1+5$	$1+5$	$1+2 + 5+8$	$2+8$	$2+8$
5	5	$5+8$	8	8

Transposed Convolution formula :

↑ input feature map
 ↑ stride
 ↑ Kernel size
 → Padding

$$\text{output size} = (\text{Input n - 1}) * S + F - 2P$$

$$= (2 - 1) * 2 + 3 \\ = 2 + 3 = \boxed{5}$$

1	1	3	2	2
1	1	3	2	2
6	6	16	10	10
5	5	13	8	8
5	5	13	8	8

5×5

ex3

3	5
2	6

input = 2×2

1	2
2	1

Kernal = 2×2



3	$6+5$	10
$6+2$	$3+10$ $+4+6$	$5+12$
4	$2+12$	6



3	11	10
8	23	17
4	14	6

3×3

Stride = 1

⇒ Transposed Convolution formula :

↑ input feature map
↑ stride
↑ Kernel size
→ padding

$$\text{output size} = (\text{Input n} - 1) * S + F - 2P$$

$$(2-1)*1 + 2 = 1 + 2 = \boxed{3}$$

ex 4

3	5
2	6

input = 2×2

1	2
2	1

Kernal = 2×2



3	6	5	10
6	3	10	5
2	4	6	12
4	2	12	6

4×4

Stride = 2

→ Transposed Convolution formula :

↑ input feature map
 ↑ stride
 ↑ Kernel size
 ↓ padding

$$\text{output size} = (\text{Input size} - 1) * S + F - 2P$$

$$\begin{aligned}
 &= (2 - 1) * 2 + 2 \\
 &= 2 + 2 = 4
 \end{aligned}$$

Learnable Upsampling: Transposed Convolution (cont.)



Input Kernel

0	1
2	3

Transposed Conv (Stride 2)

1	4
2	3

$$= \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|} \hline 1 & 4 \\ \hline 2 & 3 \\ \hline \end{array} + \begin{array}{|c|c|} \hline \end{array} + \begin{array}{|c|c|} \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & 4 \\ \hline 0 & 0 & 2 & 3 \\ \hline 2 & 8 & 3 & 12 \\ \hline 4 & 6 & 6 & 9 \\ \hline \end{array}$$

Learnable Upsampling: Transposed Convolution (cont.)



Input Kernel

0	1
2	3

Transposed Conv (Stride 1)

4	1
2	3

$$= \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|} \hline 4 & 1 \\ \hline 2 & 3 \\ \hline \end{array} + \begin{array}{|c|c|} \hline \end{array} + \begin{array}{|c|c|} \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 4 & 1 \\ \hline 8 & 16 & 6 \\ \hline 4 & 12 & 9 \\ \hline \end{array}$$

The output $O(x, y)$ of a transposed convolution is computed as:

$$O(x, y) = \sum_{i,j} I(i, j) \cdot K(x - i \cdot s, y - j \cdot s)$$

where:

- ▶ $O(x, y)$ is the output at position (x, y) ,
- ▶ $I(i, j)$ is the input value at (i, j) ,
- ▶ $K(x', y')$ is the kernel value at (x', y') ,
- ▶ s is the stride.

Transposed Convolution Formula :

$$\rightarrow \text{output size} = \frac{(Input - 1) * S + F - 2P}{\text{stride}} + 1$$

Some
padding

$$(n - 1) * S + F - 2P + C$$

Output padding

Transposed Convolution Output Size Formula

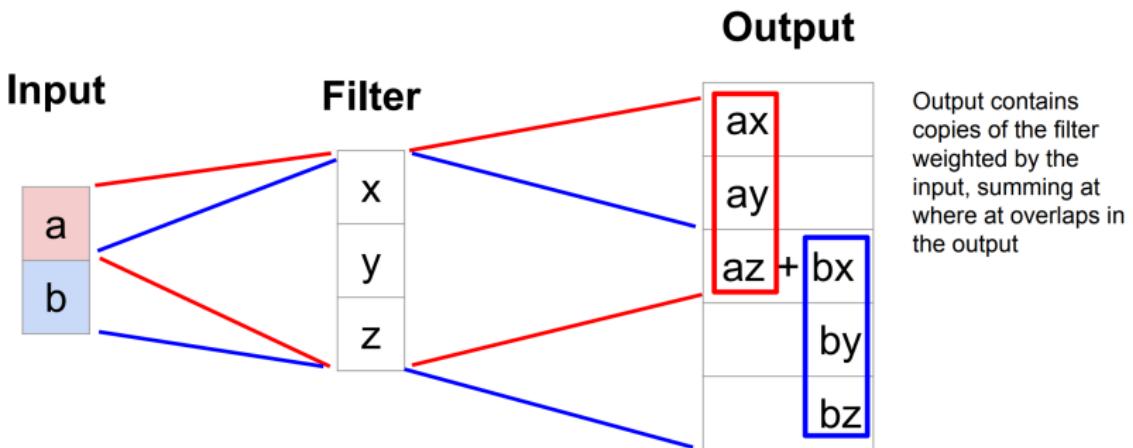
$$H_{out} = (H_{in} - 1) \times \text{stride}[0] - 2 \times \text{padding}[0] \\ + \text{dilation}[0] \times (\text{kernel_size}[0] - 1) \\ + \text{output_padding}[0] + 1 \quad (1)$$

$$W_{out} = (W_{in} - 1) \times \text{stride}[1] - 2 \times \text{padding}[1] \\ + \text{dilation}[1] \times (\text{kernel_size}[1] - 1) \\ + \text{output_padding}[1] + 1 \quad (2)$$

where:

- ▶ H_{out}, W_{out} - Output height and width.
- ▶ H_{in}, W_{in} - Input height and width.
- ▶ Stride - Step size of the filter movement.
- ▶ Padding - Number of pixels added around the input.
- ▶ Dilation - Spacing between kernel elements.
- ▶ Kernel size - Size of the convolution filter.
- ▶ Output padding - Additional padding applied to the output.

Transposed Convolution: 1D Example



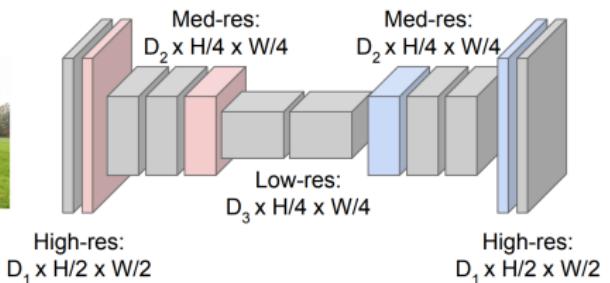
Semantic Segmentation Idea: Fully Convolutional

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



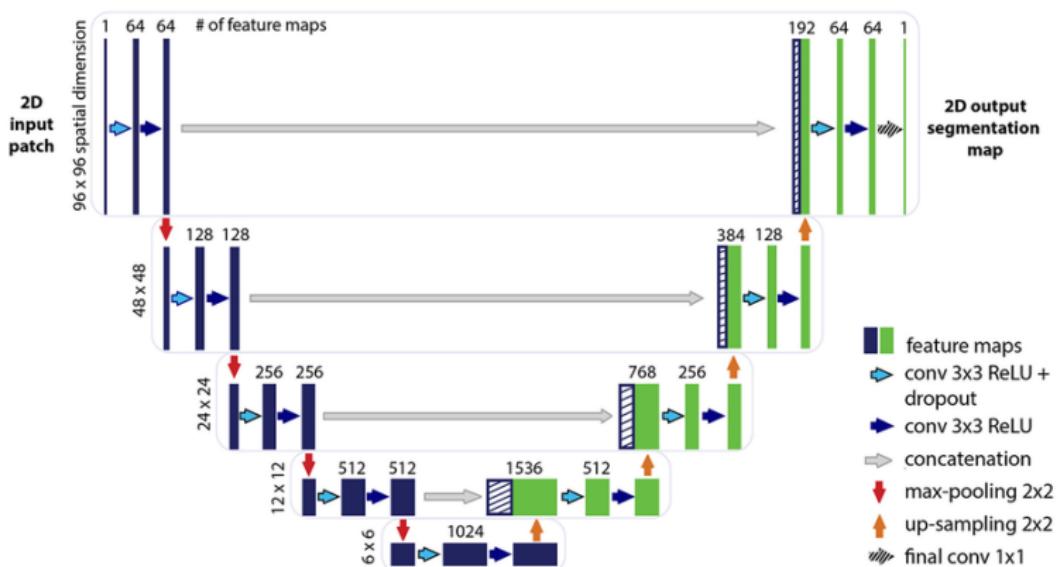
Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

- ▶ The downsampling-then-upsampling approach works well for semantic segmentation.
- ▶ **But... can we do better?**
- ▶ **Problem:** Important details and spatial information may be lost during downsampling.
- ▶ **Solution:** Introduce **residual connections** to preserve spatial information.

Residual Connections in Segmentation

- ▶ Directly connect features from downsampling layers to upsampling layers.
- ▶ Help recover lost spatial details and improve segmentation accuracy.



Residual Connections in Segmentation

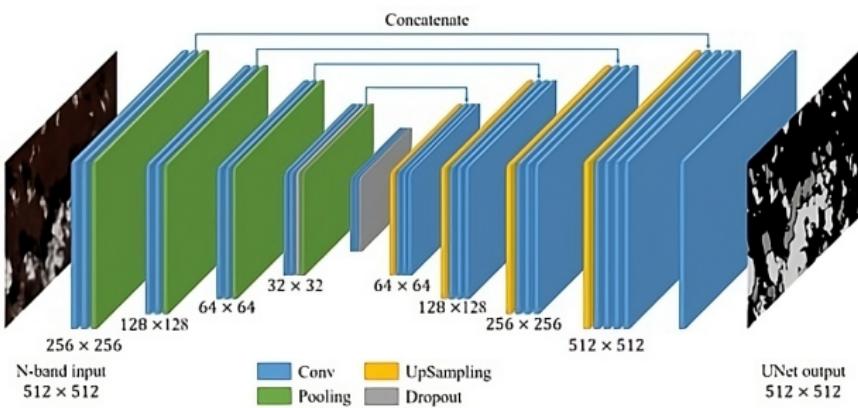
► There are two Types of Residuals:

- **Addition:** Adds features from the encoder to the decoder element-wise.
- **Concatenation:** Concatenates features from the encoder to the decoder along the **channel dimension**.

► Which Is Better?

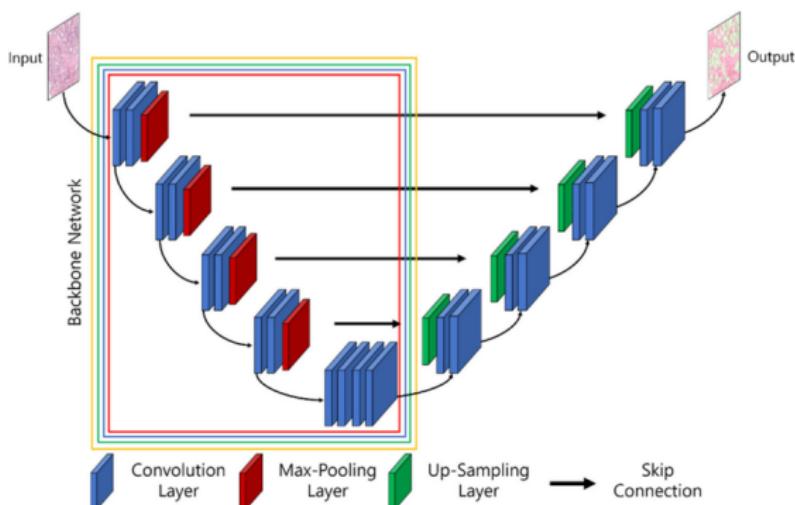
- Concatenation is often better because it retains more feature information from the encoder.
- **Note:** technically, concatenation might be harder to implement because it requires aligning input and output shapes.

- ▶ This architecture, with residual connections, is called **U-Net**.
Difference than AutoEncoder?
- ▶ **Why the name?**
 - The architecture resembles the shape of the letter "U".
** UNet we have Skip-Connections
* Reconstruct
* Autoencoder drop hole
* Generate something (unsupervised)*
 - Features are downsampled in the encoder and upsampled in the decoder, with skip connections in between.
** UNet is supervised*
- ▶ U-Net is widely used for segmentation tasks, especially in biomedical imaging.



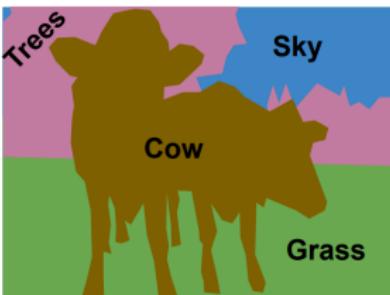
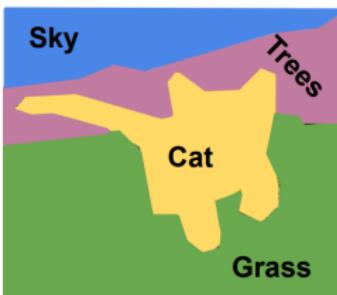
► Pretrained Encoder:

- The encoder can use a pretrained backbone (e.g., ResNet, EfficientNet).
- This helps utilize features learned on large datasets (e.g., ImageNet).
- Only the decoder is trained from scratch for segmentation-specific tasks.



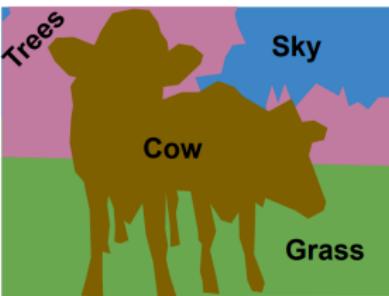
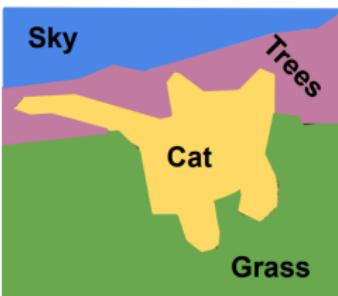
Semantic Segmentation

- ▶ Label each pixel in the image with a category label



Semantic Segmentation

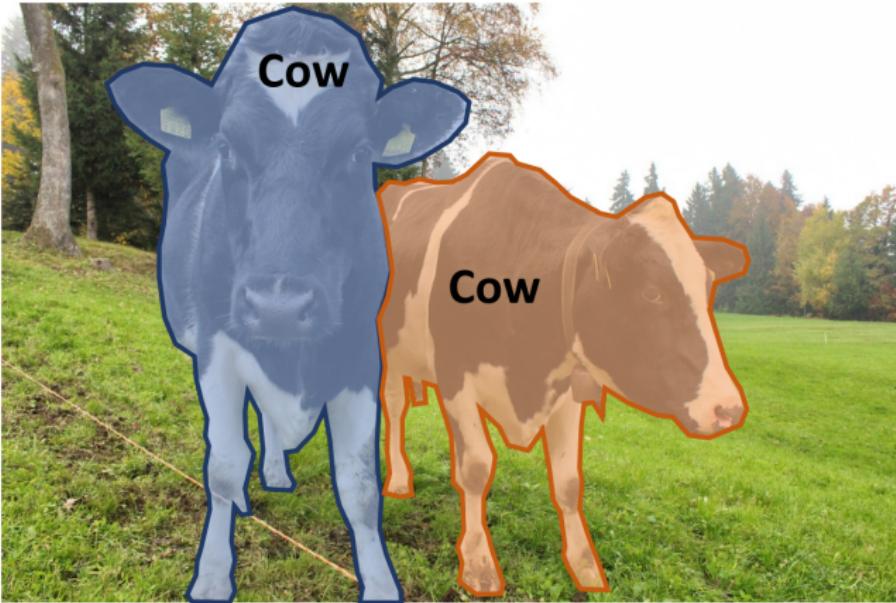
- ▶ Label each pixel in the image with a category label



- ▶ Does not differentiate instances, only care about pixels

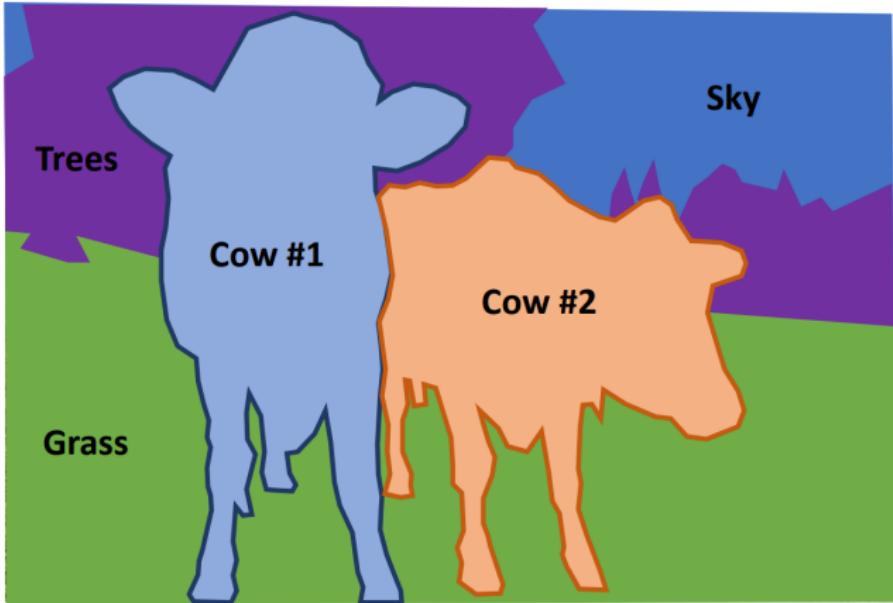
Instance Segmentation

- ▶ Separate object instances, but only things



Panoptic Segmentation

- ▶ Label all pixels in the image (both things and stuff)



These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: Deep Learning for Computer Vision
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV Deep Learning for Computer Vision: Fundamentals and Applications
- ▶ Justin Johnson, UMich EECS 498.008/598.008: Deep Learning for Computer Vision