

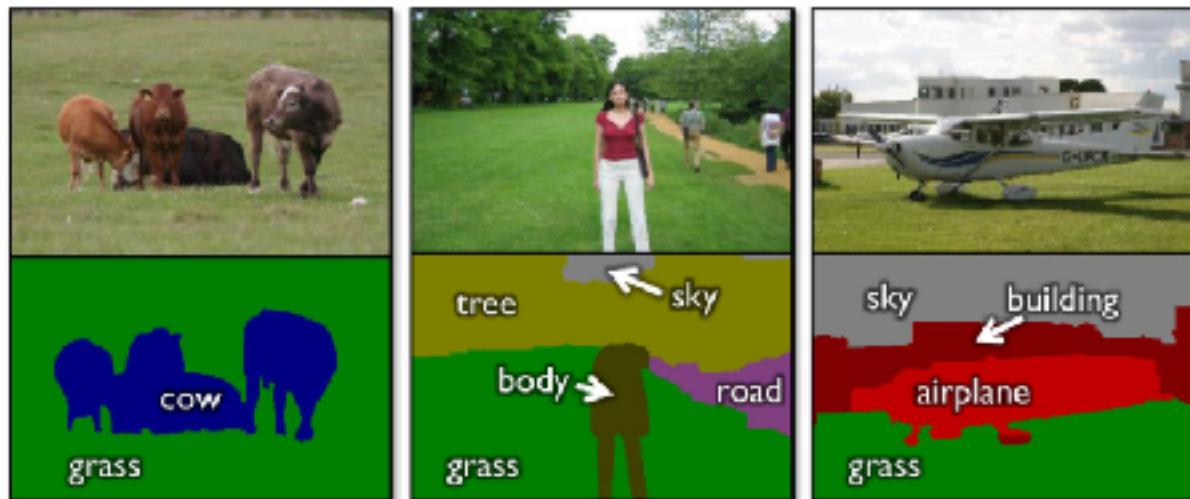
# Semantic Segmentation

# Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)

Classic computer vision problem



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Figure credit: Shelton et al, "TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context", IJCV 2007

# Semantic Segmentation

Detect instances,  
give category, label  
pixels

"simultaneous  
detection and  
segmentation" (SDS)

Lots of recent work  
(MS-COCO)

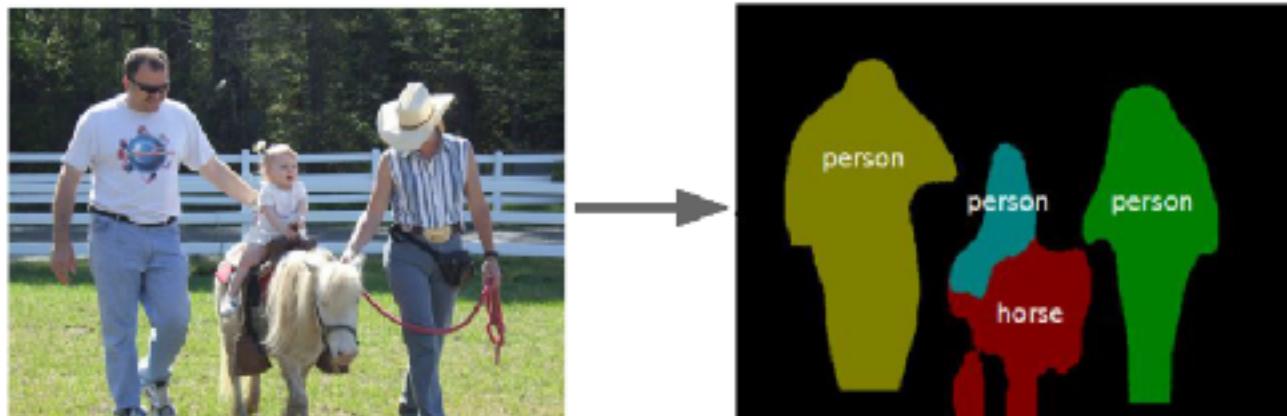


Figure credit: Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

# Semantic Segmentation

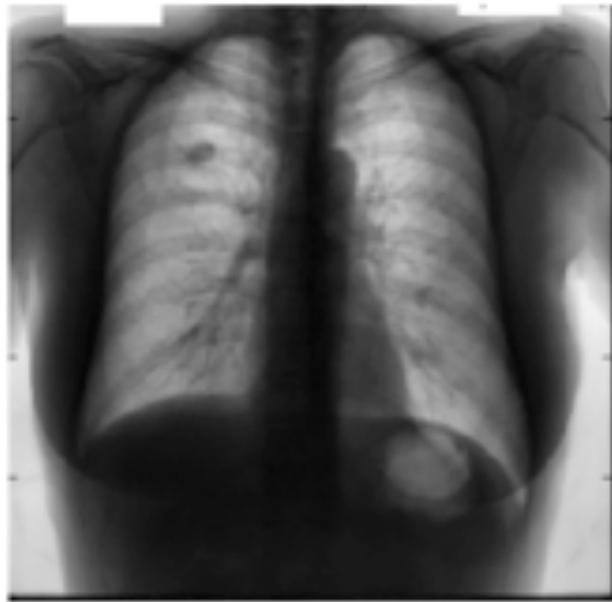
DeepLab V3 xception\_cityscapes\_trainfine (GTX980M) INPUT\_SIZE=1539  
Prediction time: 403ms (2.5 fps) AVG: 356ms (2.8 fps)



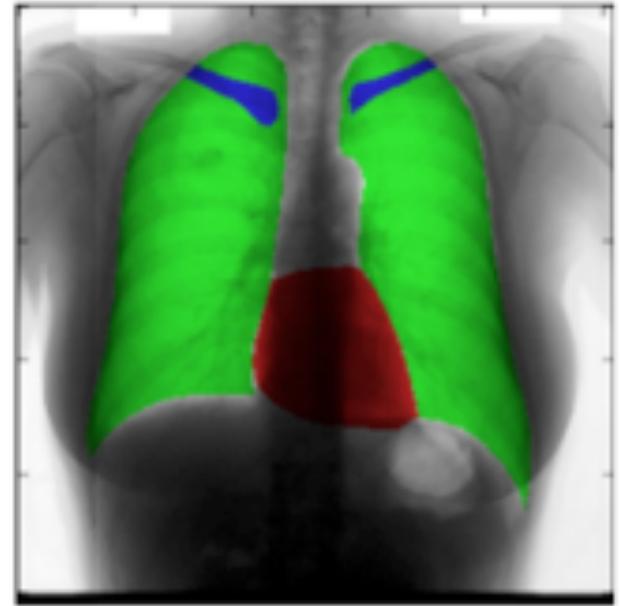
A real-time segmented road scene for autonomous driving.  
<https://www.youtube.com/watch?v=ATlcEDSPWXY>)

# Semantic Segmentation

---



Input Image



Segmented Image

A chest x-ray with the heart (red), lungs (green), and clavicles (blue) are segmented.

<https://arxiv.org/abs/1701.08816>

# Goal of Semantic Segmentation



Input

segmented

1: Person  
2: Purse  
3: Plants/Grass  
4: Sidewalk  
5: Building/Structures

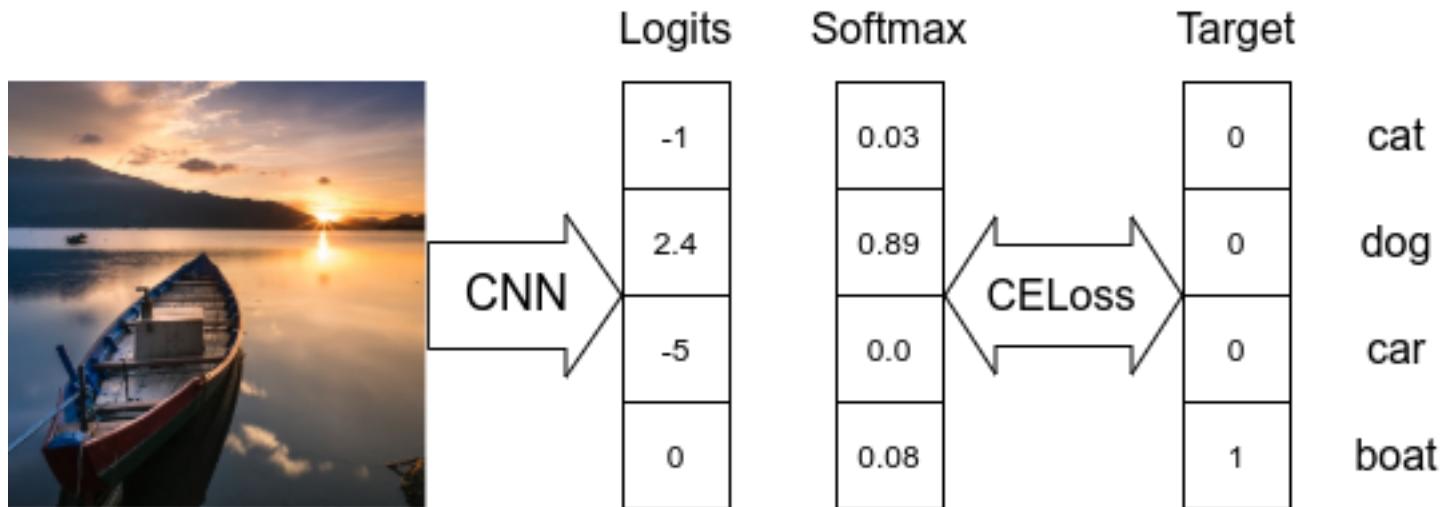
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	3	1	1	1	3	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	1	1	1	1	1	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	1	1	1	1	3	3	3	3	5	5	5	5	5	5
5	5	3	3	3	3	3	3	1	1	1	3	3	5	5	5	5	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	4	4	4	4	5	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	4	4	4	4	4	5	5	5	5	5
4	4	4	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4
3	3	3	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4

Semantic Labels

low-resolution prediction map only for clarity,  
In reality label resolution = input image resolution

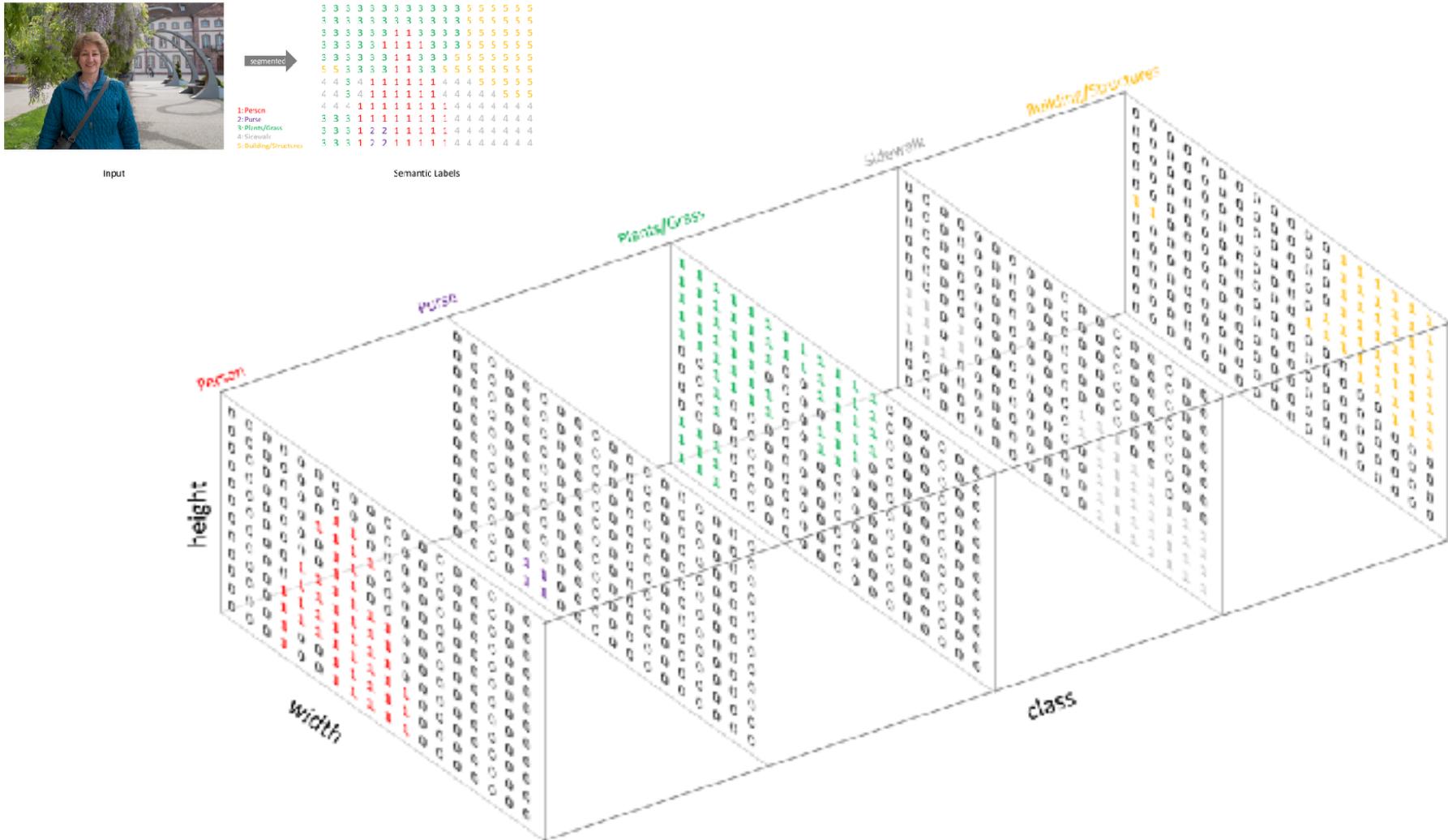
# Recall Image Classification

- ▶ Target by one-hot encoding class labels



# Targets in Semantic Segmentation

- ▶ Target by one-hot encoding class labels



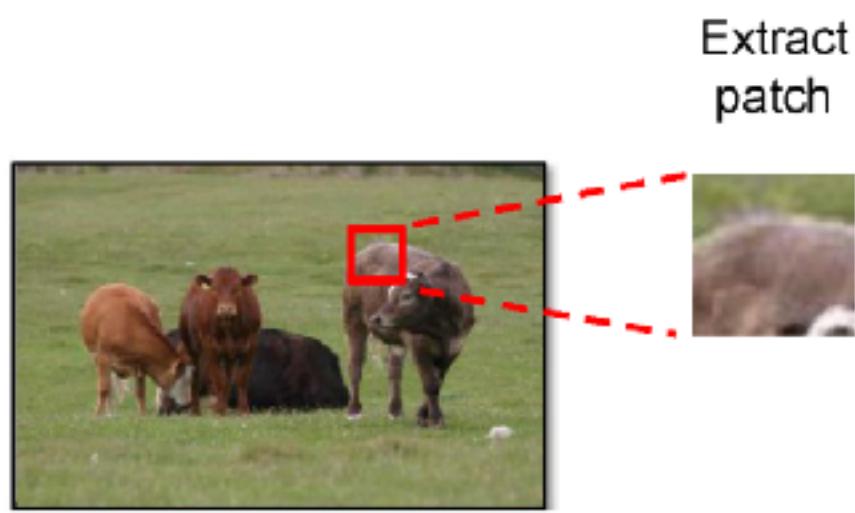
# Semantic Segmentation via Classification

---

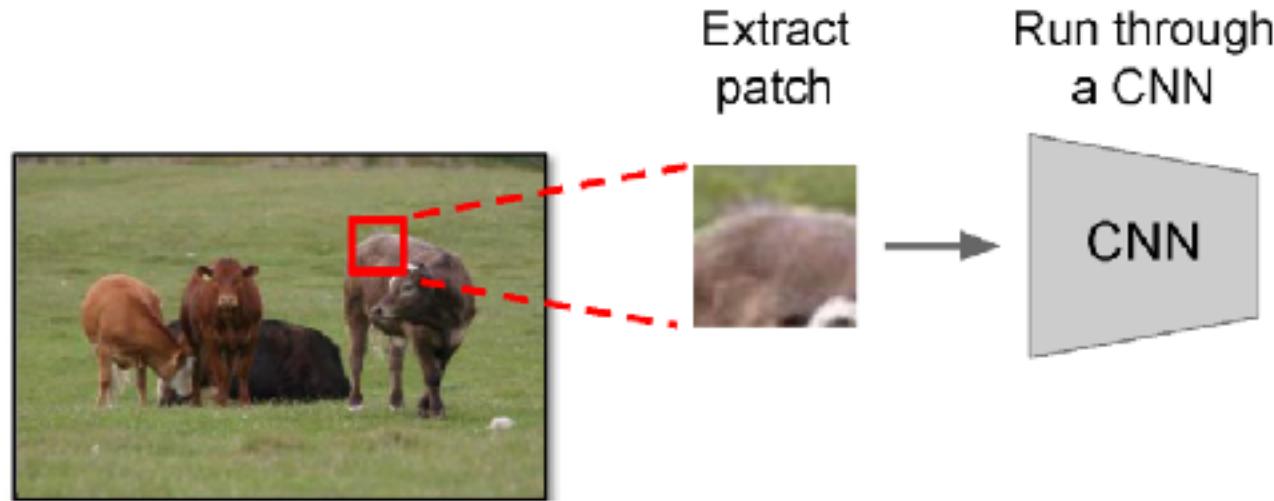


# Semantic Segmentation via Classification

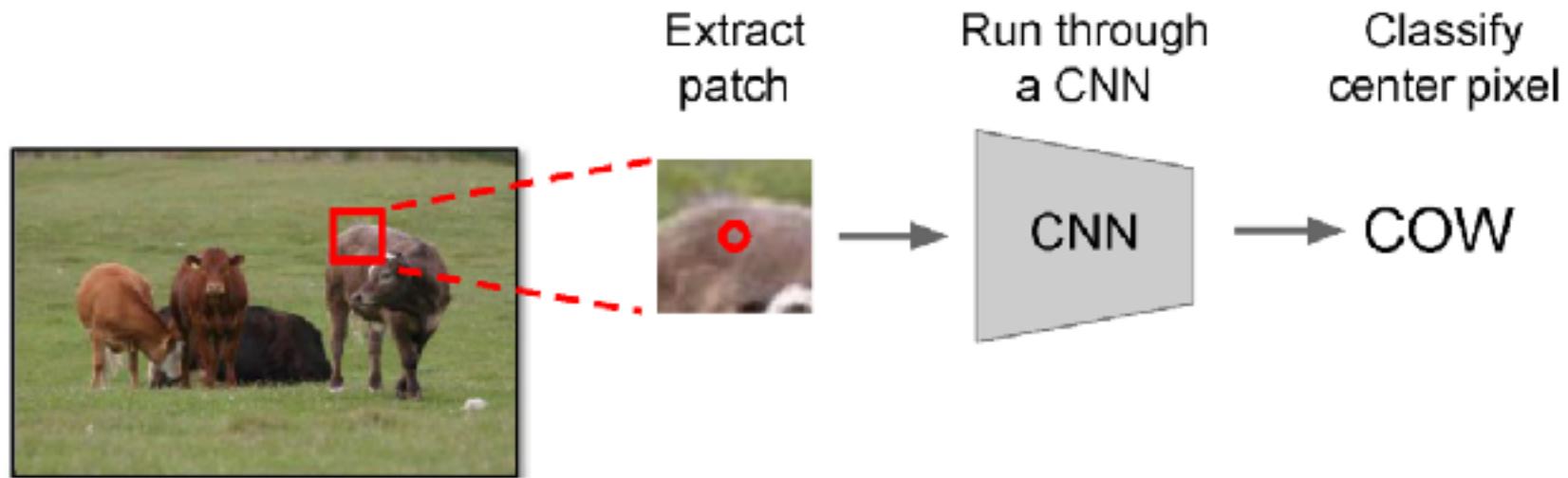
---



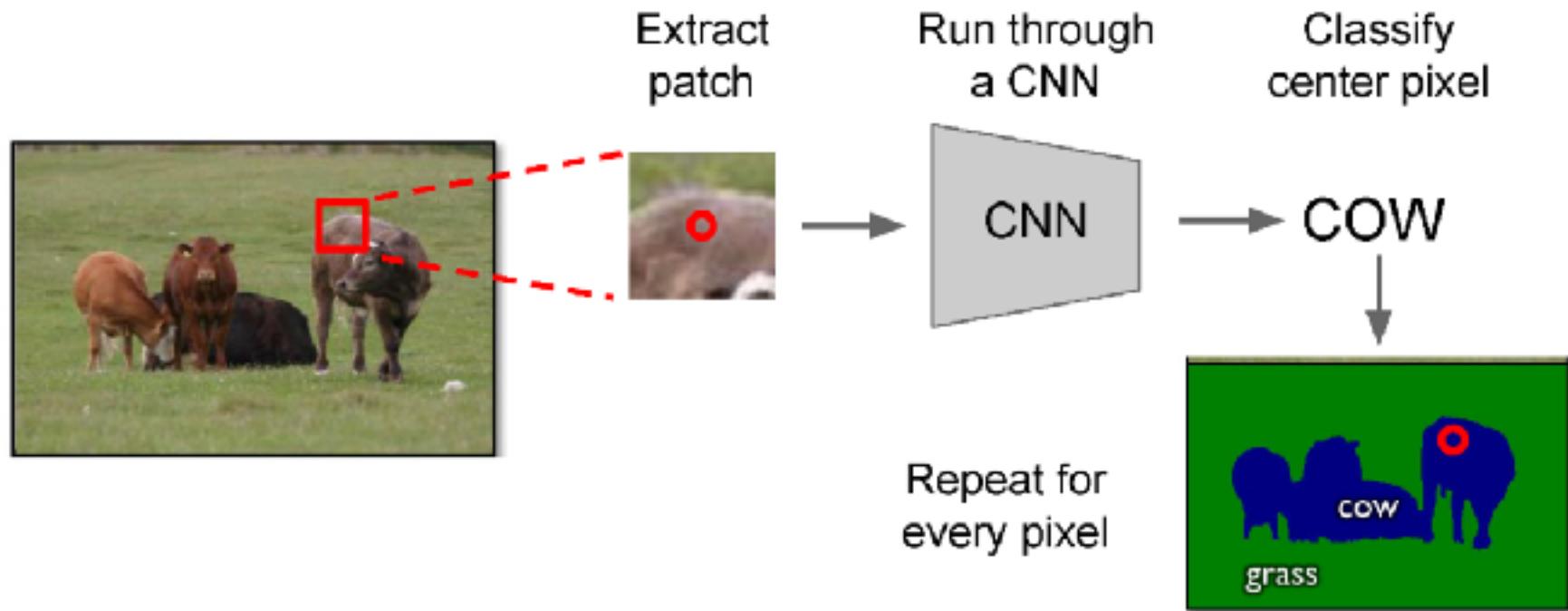
# Semantic Segmentation via Classification



# Semantic Segmentation via Classification

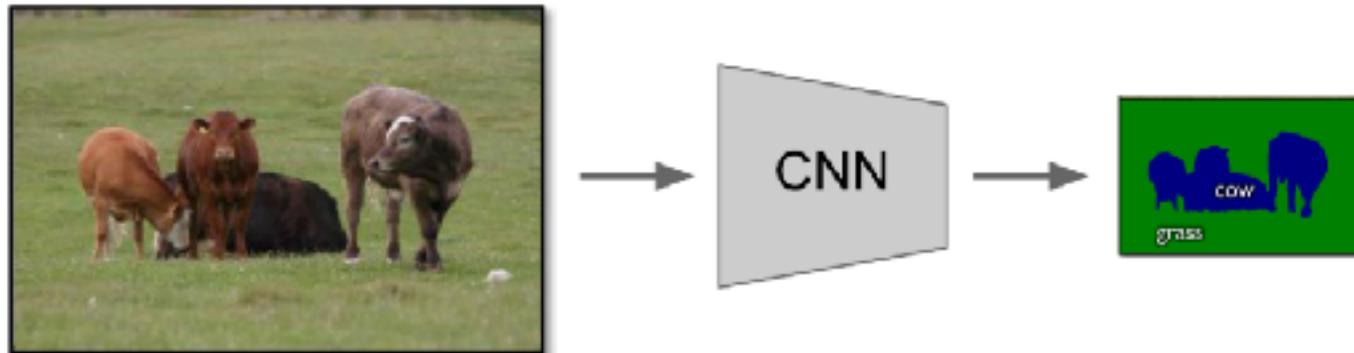


# Semantic Segmentation via Classification



# Semantic Segmentation via Classification

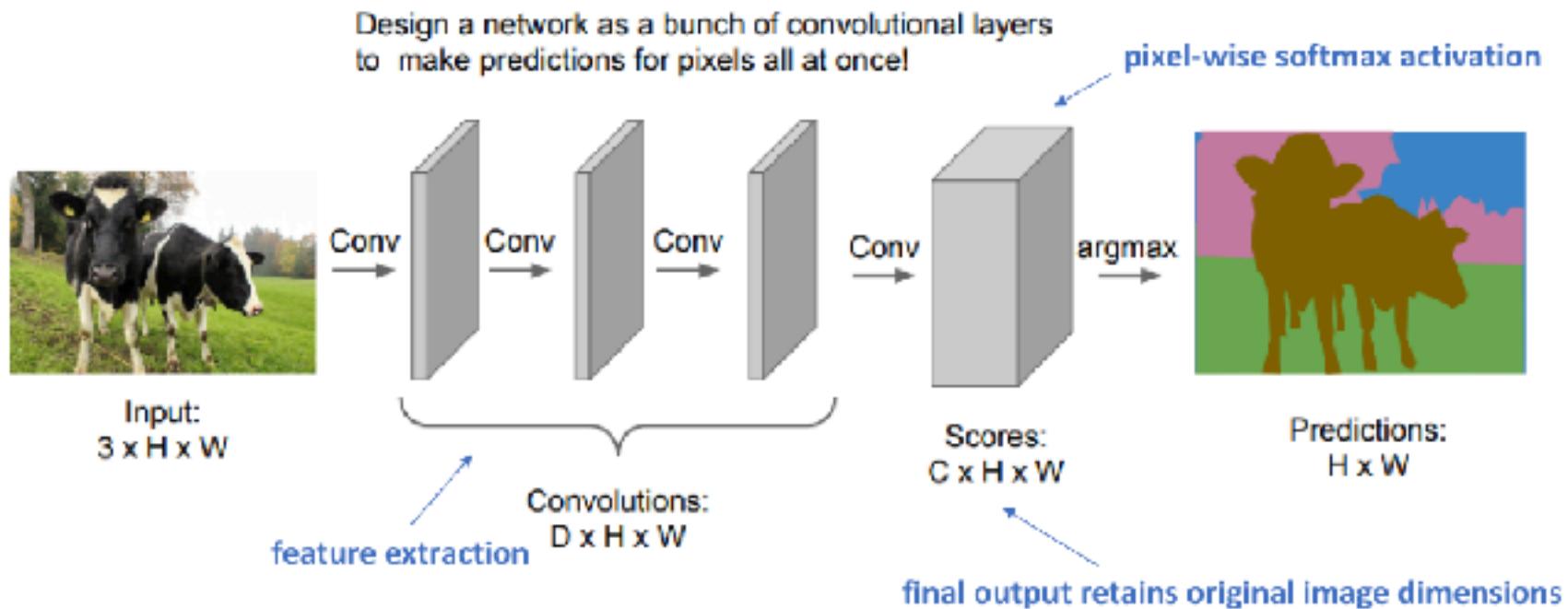
Run “fully convolutional” network  
to get all pixels at once



Smaller output  
due to pooling

- ▶ computationally expensive and extremely slow to perform forward pass independently for each pixel

# Semantic Segmentation: Naive Approach

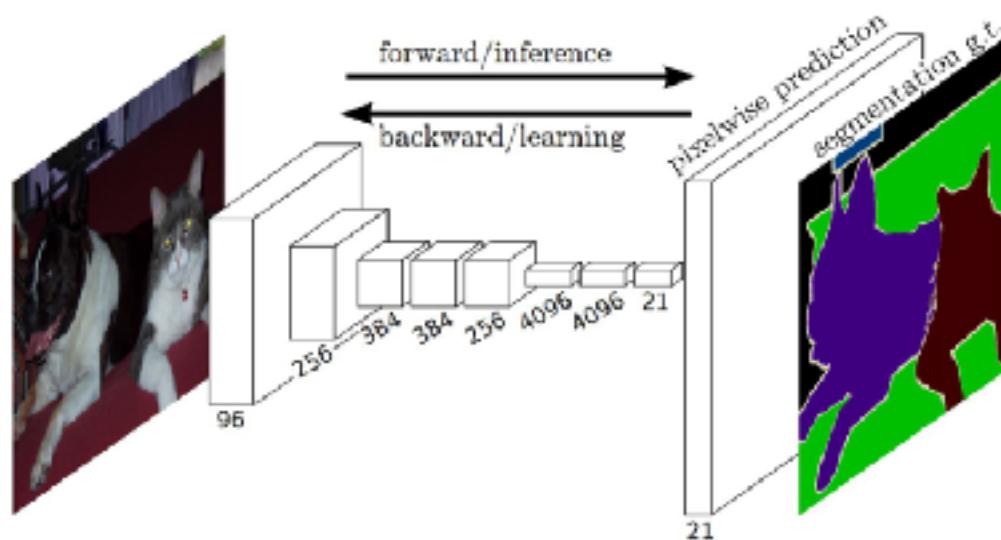


**Downside:** Preserving image dimensions throughout entire network will be computationally expensive.

- ▶ Computationally expensive and memory intensive to preserve the full resolution throughout the network.

# Semantic Segmentation: Upsampling

- In order to maintain expressiveness, we typically need to increase the number of feature maps (channels) as we get deeper in the network.

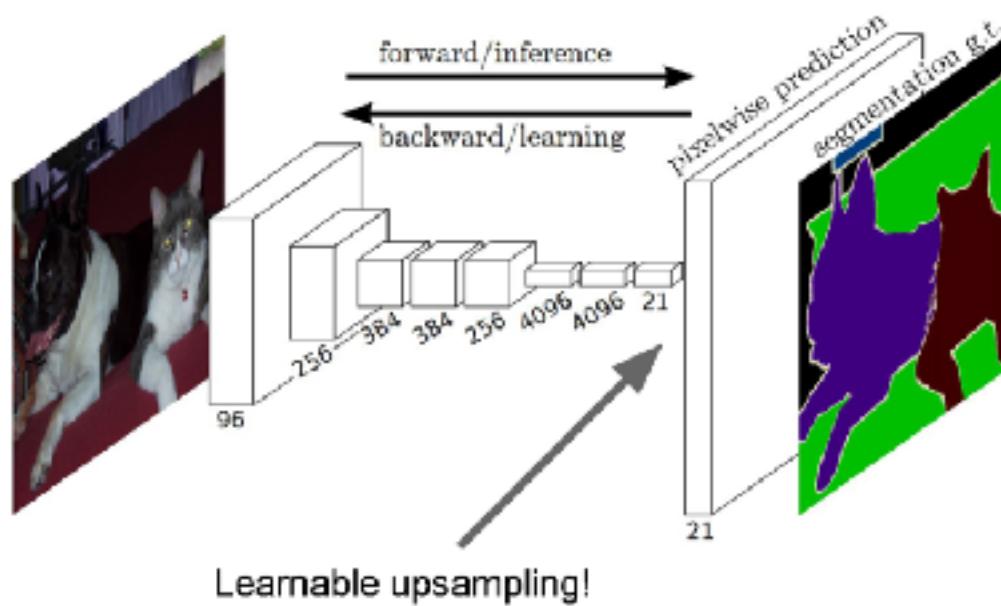


Long, Shelhamer, and Darrell, 'Fully Convolutional Networks for Semantic Segmentation', CVPR 2015

<https://arxiv.org/pdf/1411.4038.pdf>

# Semantic Segmentation

- In order to maintain expressiveness, we typically need to increase the number of feature maps (channels) as we get deeper in the network.

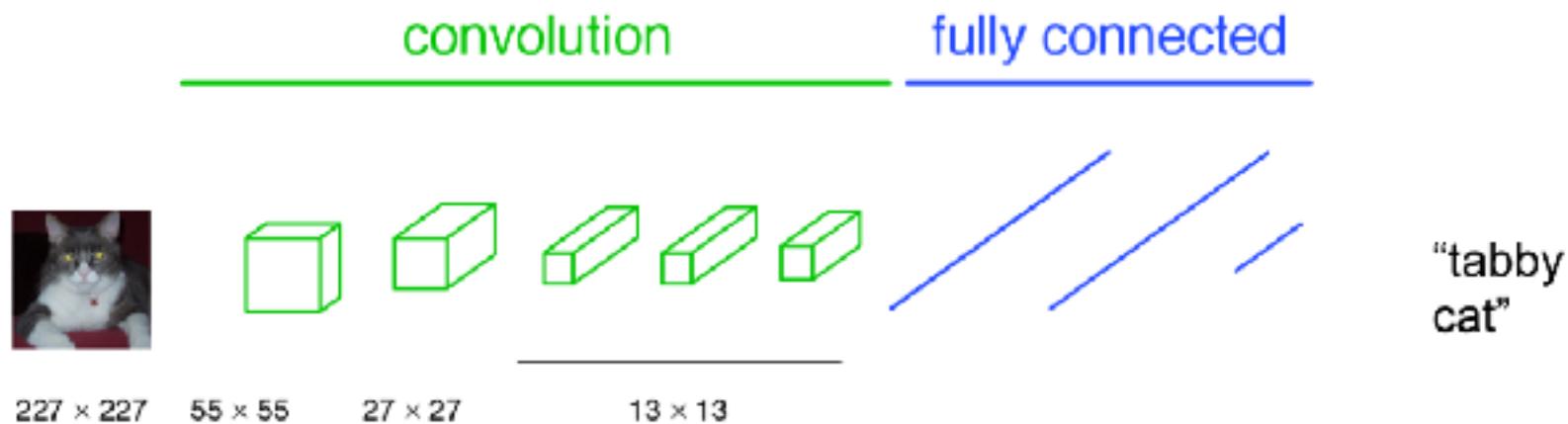


Long, Shelhamer, and Darrell, 'Fully Convolutional Networks for Semantic Segmentation', CVPR 2015

<https://arxiv.org/pdf/1411.4038.pdf>

# a classification network

---

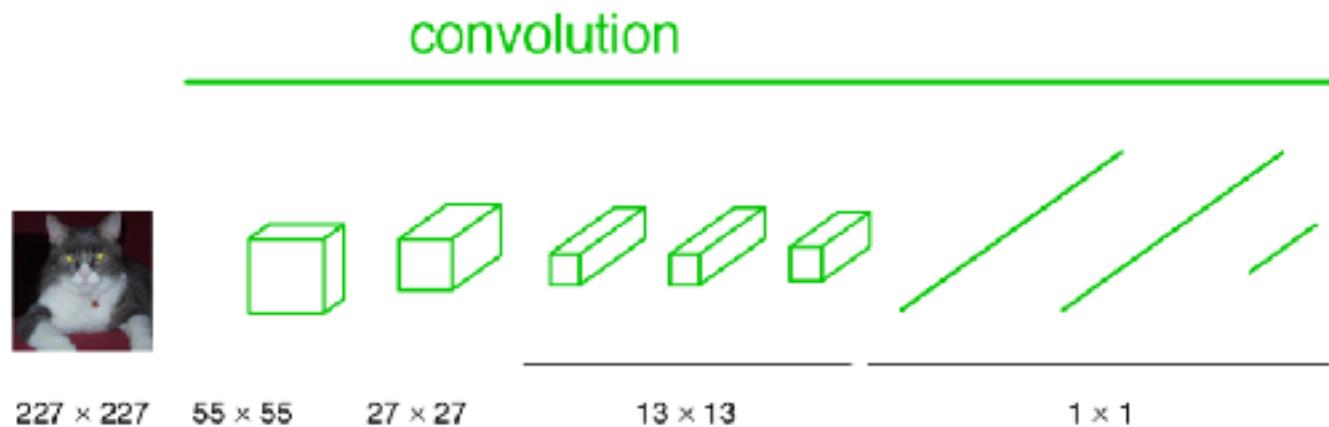


Slide credit: Jonathan Long



# becoming fully convolutional

---



Slide credit: Jonathan Long



# becoming fully convolutional

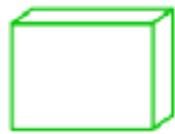
---

convolution

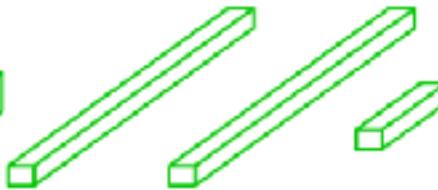
---



$H \times W$



$H/16 \times W/16$



$H/32 \times W/32$

Slide credit: Jonathan Long

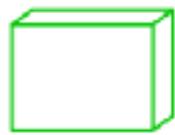


# upsampling output

convolution



$H \times W$



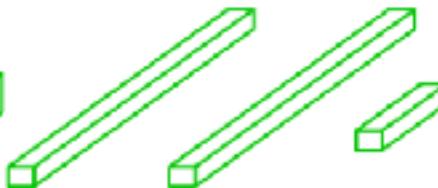
$H/4 \times W/4$



$H/8 \times W/8$



$H/16 \times W/16$



$H/32 \times W/32$

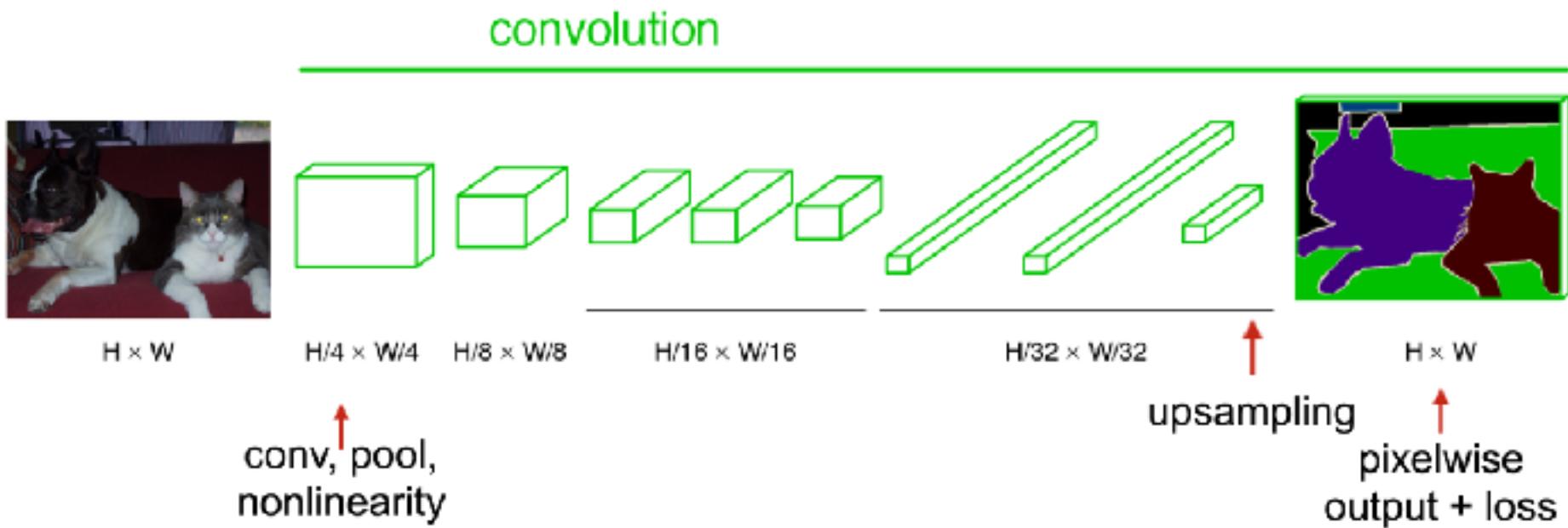


$H \times W$

Slide credit: Jonathan Long



# end-to-end, pixels-to-pixels network

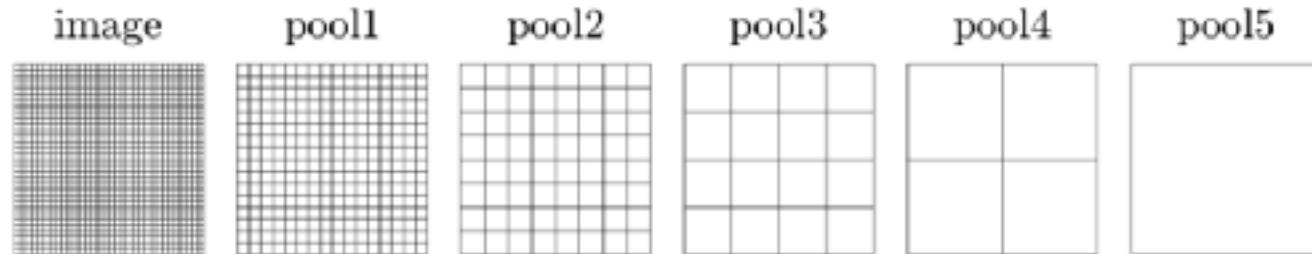


Slide credit: Jonathan Long

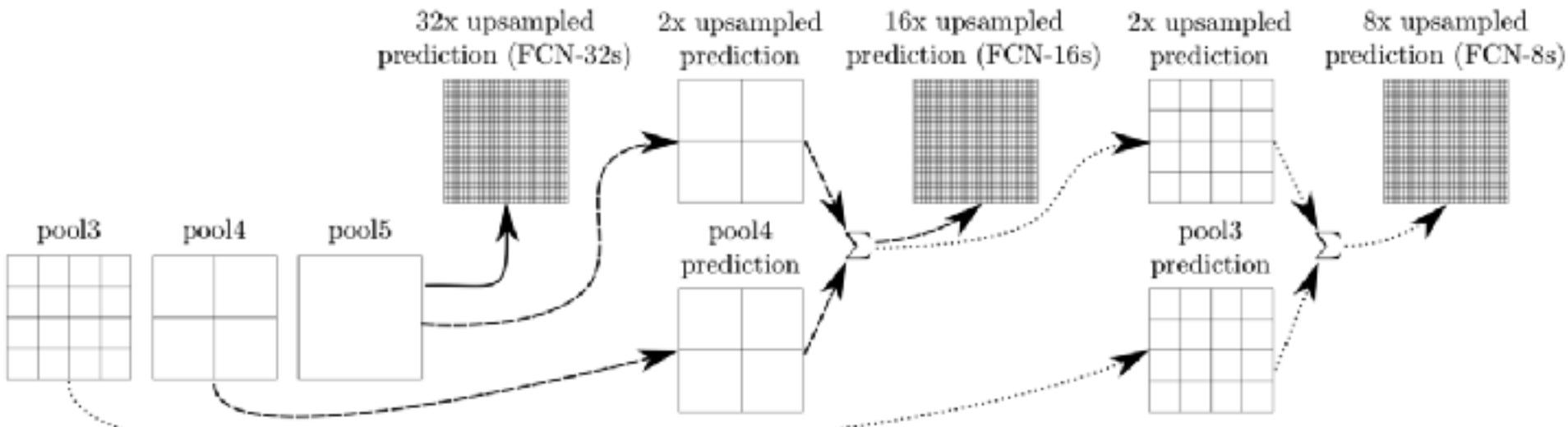


# skip layers

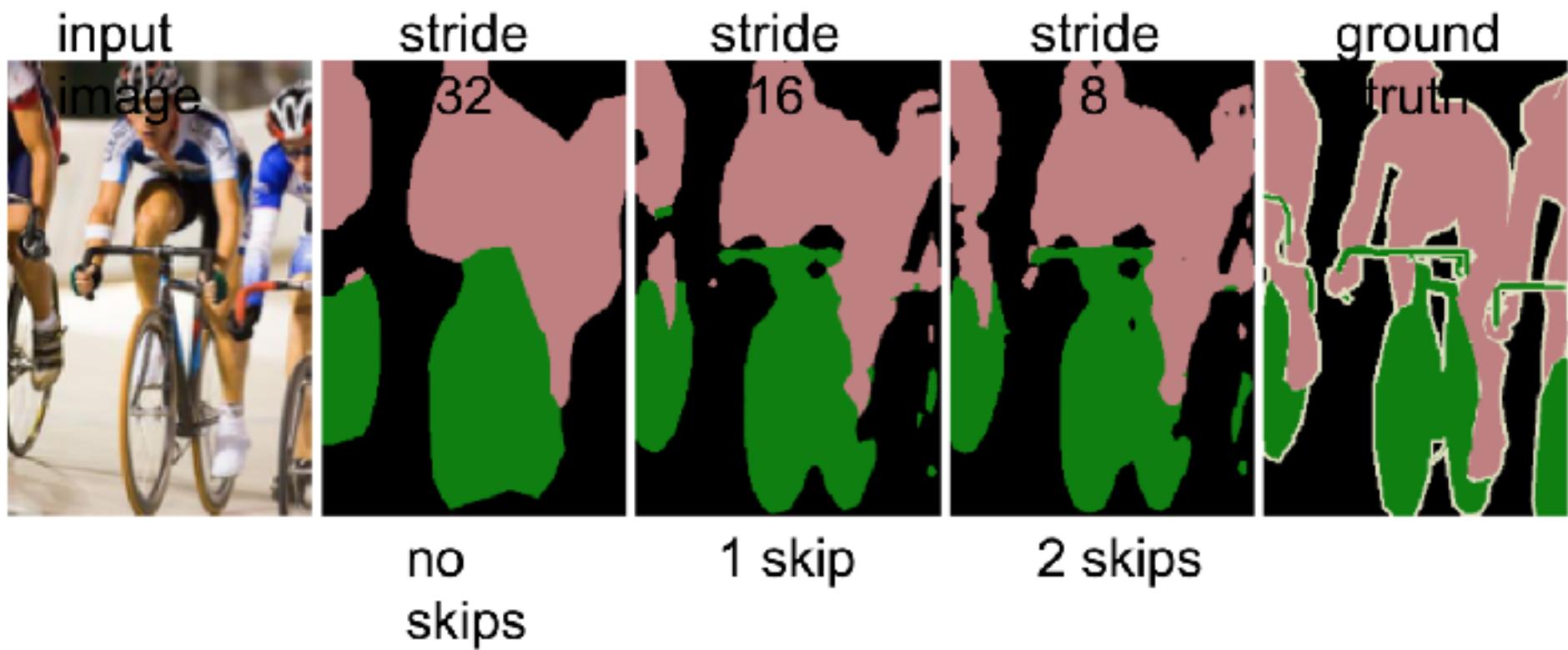
## ▶ Encoder/Downsampling



## ▶ Learnable Upsampling

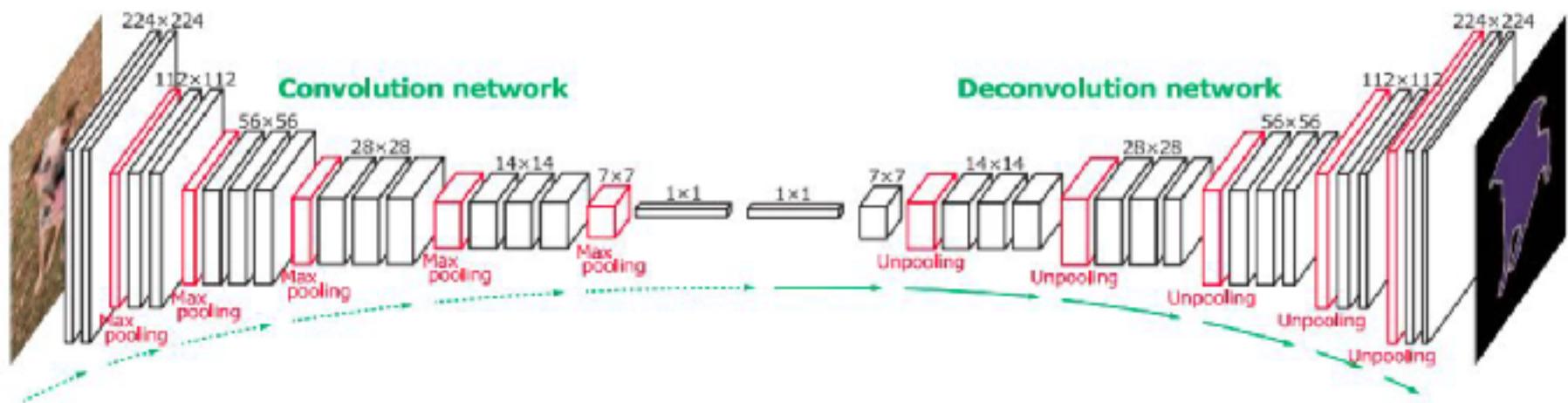


# skip layer refinement



▶ Slide credit: Jonathan Long

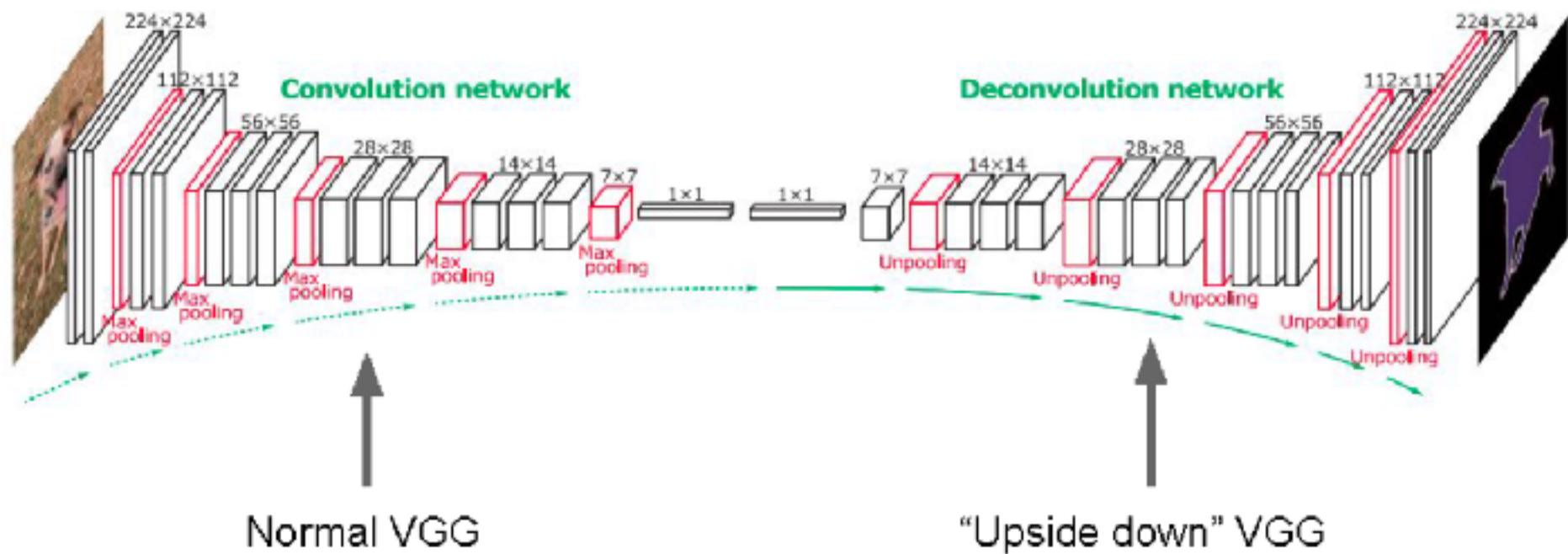
# Semantic Segmentation: Encoder/Decoder



| Deconvolution Network for Semantic Segmentation\*, ICCV 2015

Fei-Fei Li, Andrej Karpathy & Justin Johnson 2016

# Semantic Segmentation: Encoder/Decoder



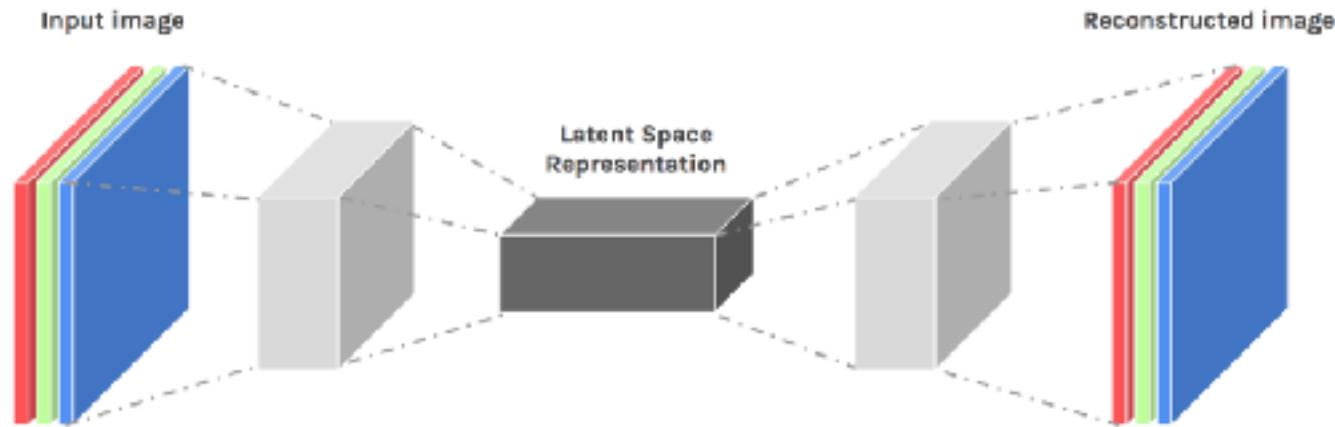
Normal VGG

“Upside down” VGG

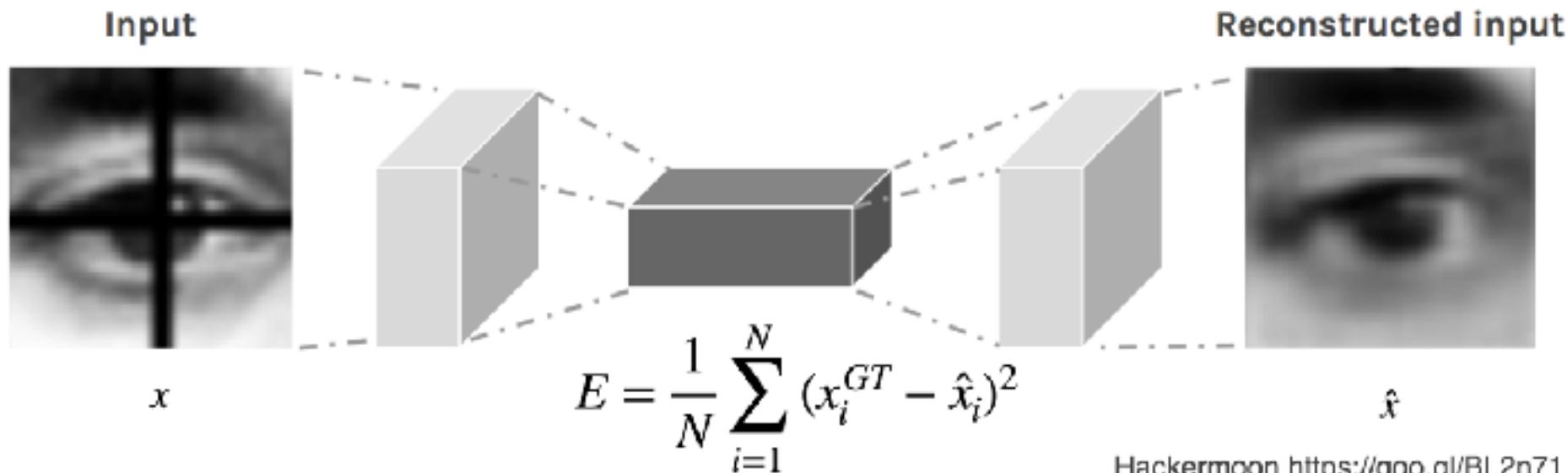
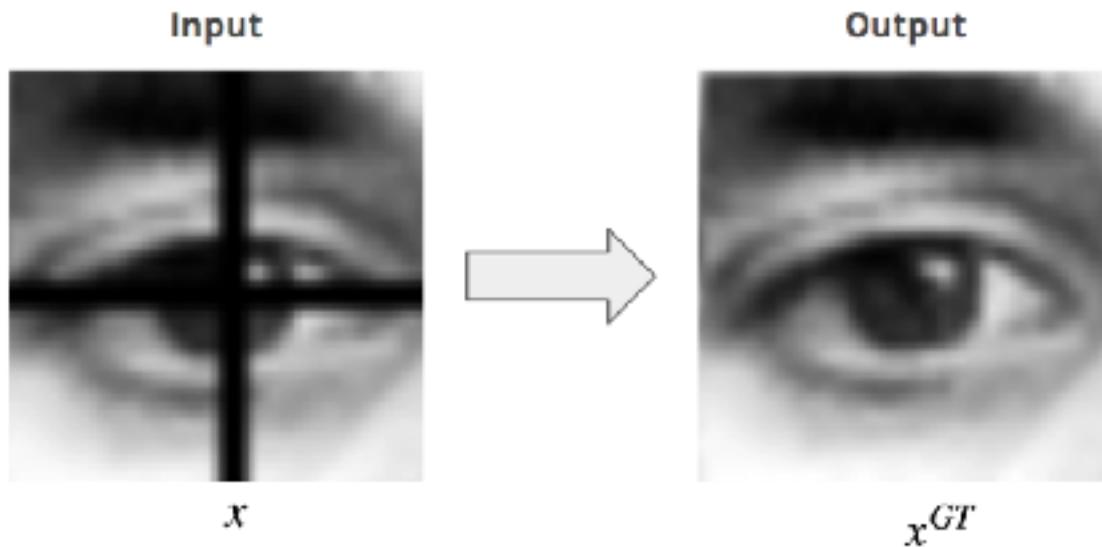
# Recap: Convolutional Autoencoders

---

- ▶ Replace *fully connected layers by convolutional layers*. This helps the network extract **visual features** from the images, and therefore obtain a much more accurate latent space representation. The reconstruction process uses **upsampling** and convolutions.

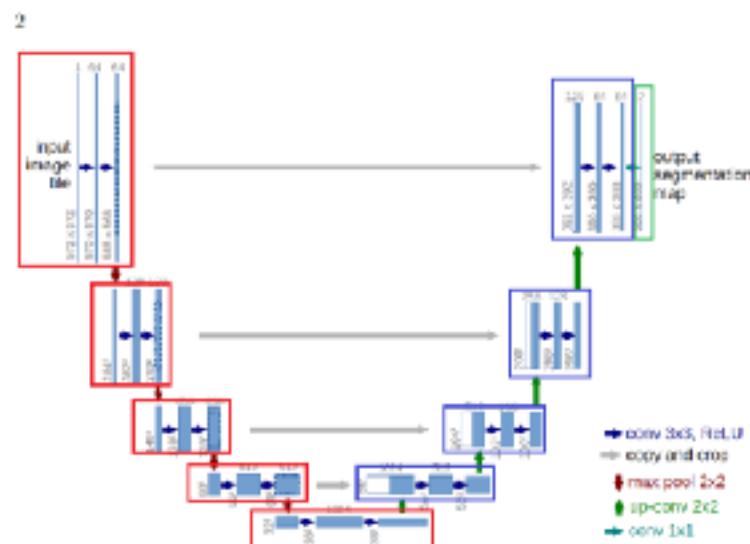
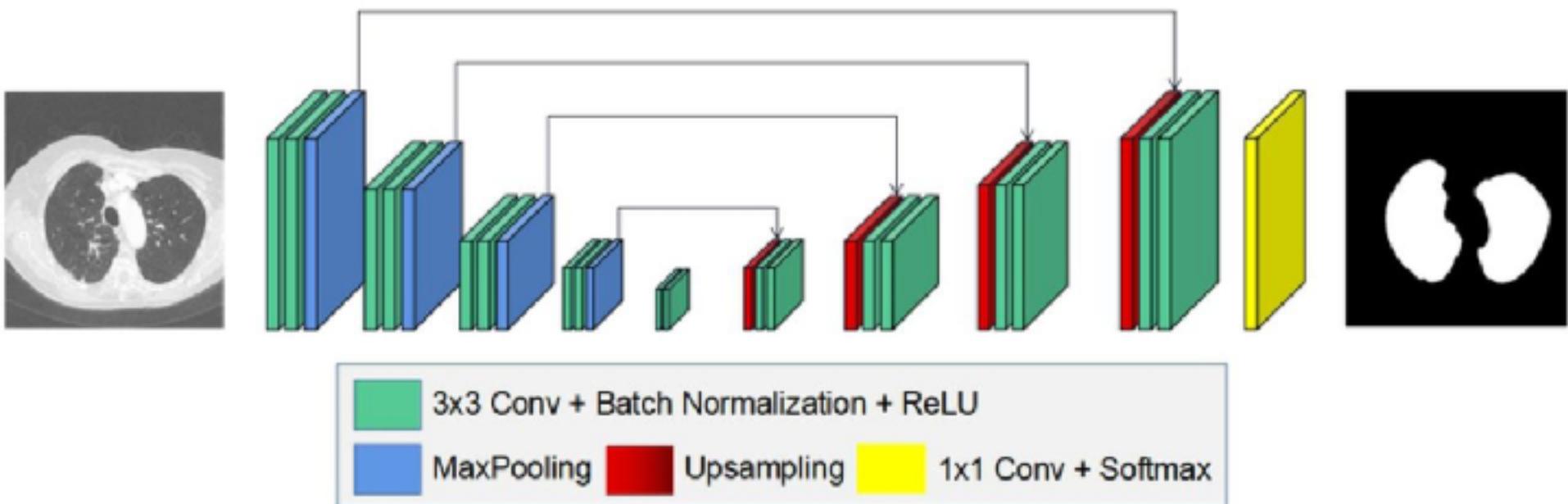


# Example 1: Ultra-basic image reconstruction



# U-Net:

## Convolutional Networks for Biomedical Image Segmentation



# Methods for upsampling

# Methods for upsampling

## Nearest Neighbor

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

## "Bed of Nails"

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

## 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

Output: 2 x 2

## Max Unpooling

Use positions from pooling layer

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

# Transposed Convolution

## ▶ Convolution

$$\begin{array}{|c|c|} \hline O_{11} & O_{12} \\ \hline O_{21} & O_{22} \\ \hline \end{array} = \text{Convolution} \left( \begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21} & X_{22} & X_{23} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array}, \begin{array}{|c|c|} \hline F_{11} & F_{12} \\ \hline F_{21} & F_{22} \\ \hline \end{array} \right)$$

$$\begin{array}{|c|c|c|} \hline X_{11} & X_{12}F_{11} & X_{13} \\ \hline X_{21} & X_{22}F_{21} & X_{23} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|} \hline X_{11} & X_{12}F_{11} & X_{13}F_{12} \\ \hline X_{21} & X_{22}F_{21} & X_{23}F_{22} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21}F_{11} & X_{22}F_{12} & X_{23} \\ \hline X_{31}F_{21} & X_{32}F_{22} & X_{33} \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21} & X_{22}F_{11} & X_{23}F_{12} \\ \hline X_{31} & X_{32}F_{21} & X_{33}F_{22} \\ \hline \end{array}$$

## ▶ Full Convolution

$$\begin{array}{|c|c|c|} \hline \delta E / \delta X_{11} & \delta E / \delta X_{12} & \delta E / \delta X_{13} \\ \hline \delta E / \delta X_{21} & \delta E / \delta X_{22} & \delta E / \delta X_{23} \\ \hline \delta E / \delta X_{31} & \delta E / \delta X_{32} & \delta E / \delta X_{33} \\ \hline \end{array} = \text{Full_Convolution} \left( \begin{array}{|c|c|c|} \hline \delta E / \delta O_{11} & \delta E / \delta O_{12} \\ \hline \delta E / \delta O_{21} & \delta E / \delta O_{22} \\ \hline \delta E / \delta O_{31} & \delta E / \delta O_{32} \\ \hline \end{array}, \begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & F_{11} \\ \hline \end{array} \right)$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & r_1\delta O_{11} + \delta O_{12} \\ \hline & \delta O_{21} + \delta O_{22} \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & r_2\delta O_{21} + \delta O_{22} \\ \hline & \delta O_{11} + \delta O_{12} \\ \hline \end{array}$$

$$\delta O_{11} = F_{11}\delta O_{11} + F_{12}\delta O_{12}$$

$$\begin{array}{|c|c|} \hline \delta O_{11} & \delta O_{12} \\ \hline F_{22} & r_1\delta O_{21} + \delta O_{22} \\ \hline F_{12} & F_{11} \\ \hline \end{array}$$

$$\delta O_{21} = F_{21}\delta O_{21} + F_{22}\delta O_{22}$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & r_1\delta O_{11} + \delta O_{12} \\ \hline & \delta O_{21} + \delta O_{22} \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & r_2\delta O_{21} + \delta O_{22} \\ \hline & \delta O_{11} + \delta O_{12} \\ \hline \end{array}$$

$$\delta O_{11} = F_{11}\delta O_{11} + F_{12}\delta O_{12}$$

$$\begin{array}{|c|c|} \hline \delta O_{11} & \delta O_{12} \\ \hline F_{22} & r_1\delta O_{21} + \delta O_{22} \\ \hline F_{12} & F_{11} \\ \hline \end{array}$$

$$\delta O_{21} = F_{21}\delta O_{21} + F_{22}\delta O_{22}$$

$$\begin{array}{|c|c|c|} \hline F_{22} & F_{21} \\ \hline F_{12} & F_{11} \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline \delta O_{11} & r_1\delta O_{12} + F_{11} \\ \hline \delta O_{21} & \delta O_{22} \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline F_{22} & F_{21} \\ \hline \delta O_{11} & r_2\delta O_{12} + F_{21} \\ \hline \delta O_{21} & \delta O_{22} \\ \hline \end{array}$$

$$\delta O_{11} = F_{11}\delta O_{11} + F_{12}\delta O_{12}$$

$$\begin{array}{|c|c|} \hline \delta O_{11} & \delta O_{12} \\ \hline F_{22} & r_1\delta O_{21} + \delta O_{22} \\ \hline F_{12} & F_{11} \\ \hline \end{array}$$

$$\delta O_{21} = F_{21}\delta O_{21} + F_{22}\delta O_{22}$$

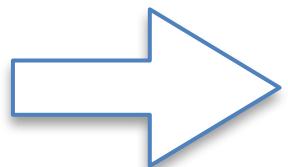
## Summary

---

- ▶ The transposed convolution operation forms the **same connectivity** as the normal convolution but in the **backward direction**.
- ▶ We can use it to conduct **up-sampling**. Moreover, the **weights in the transposed convolution are learnable**. So we do not need a predefined interpolation method.
- ▶ Even though it is called the transposed convolution, it does not mean that we take some existing convolution matrix and use the transposed version.

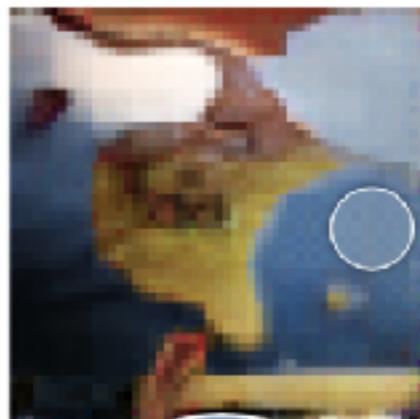
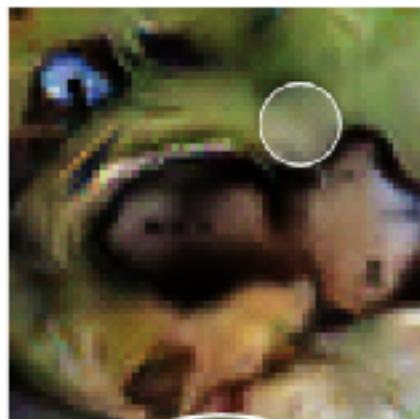
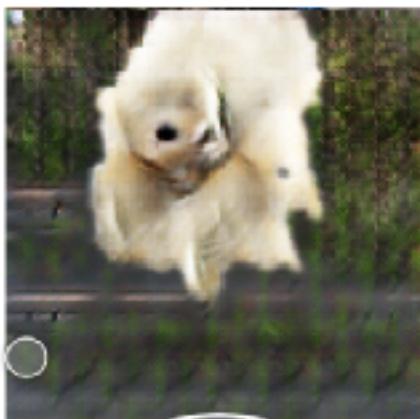
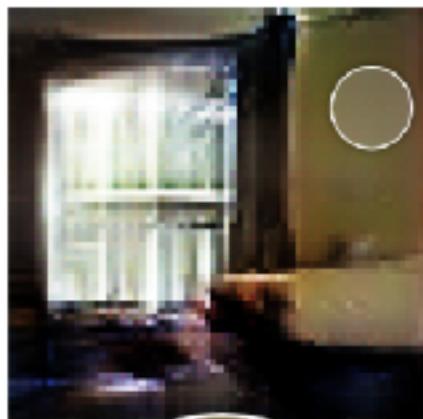
# Transposed Convolution and Checkerboard Artifacts

Supplementary Slides

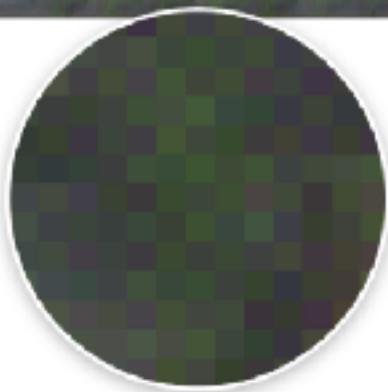


# Checkerboard Artifacts

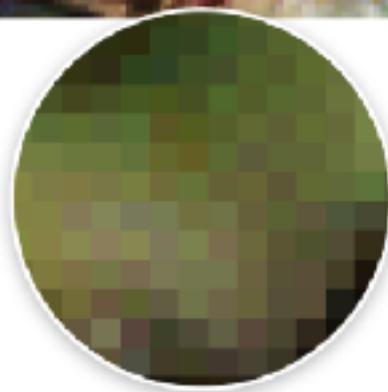
---



Radford, et al., 2015 [1]



Salimans et al., 2016 [2]



Donahue, et al., 2016 [3]

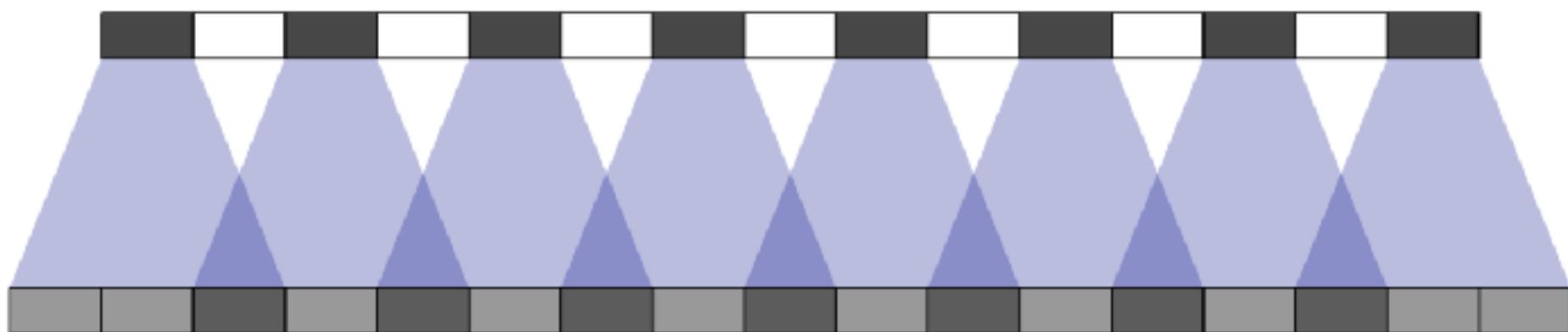


Dumoulin, et al., 2016 [4]

## Checkerboard Artifacts - 1D Example

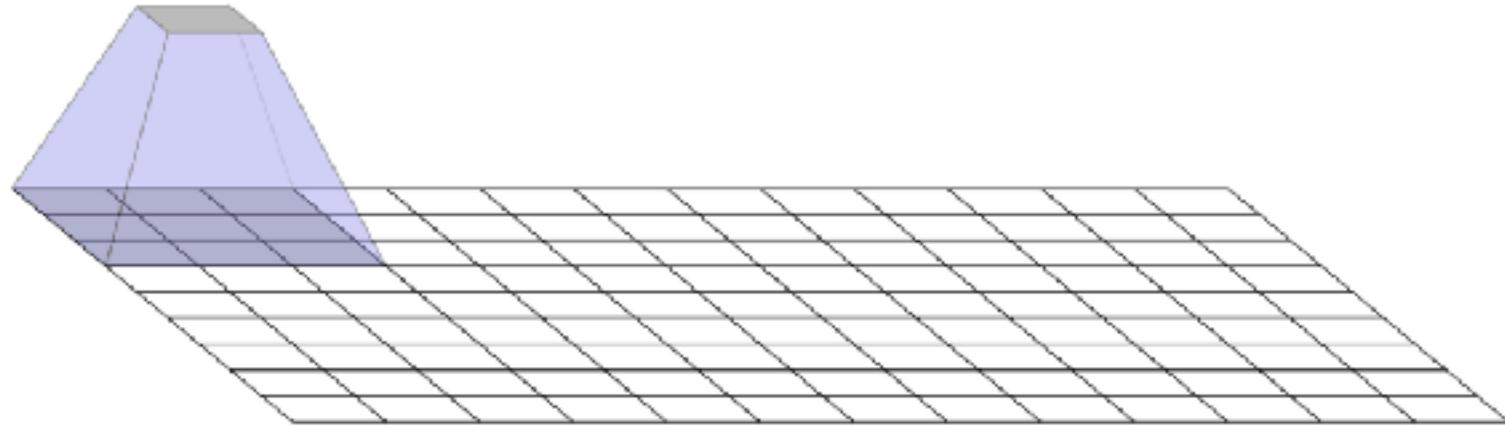
---

- ▶ Stride 2
- ▶ Size 3



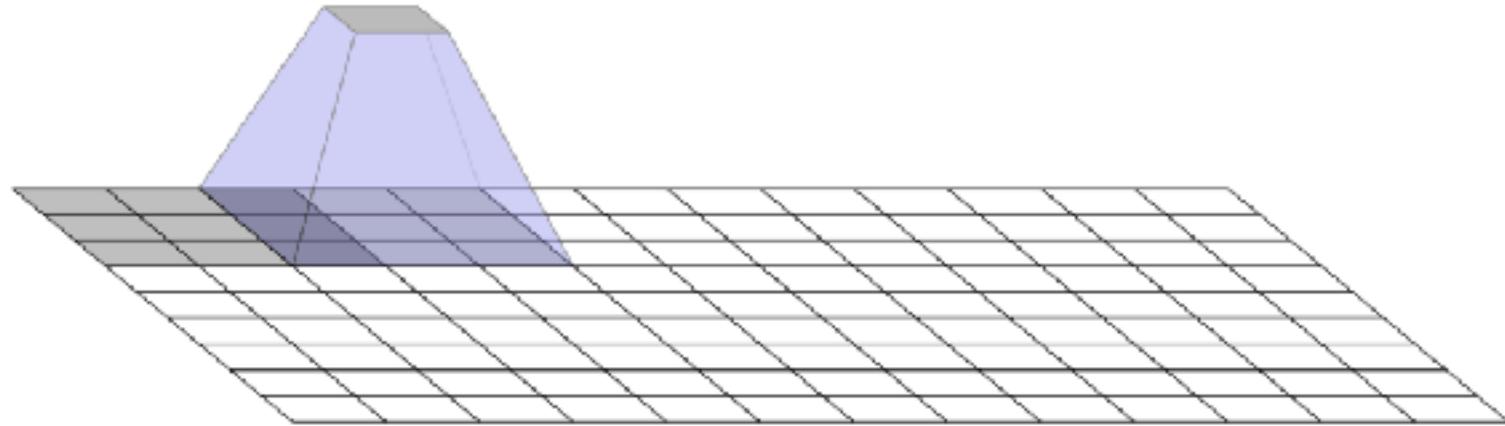
# Transposed Convolution & Overlap

---



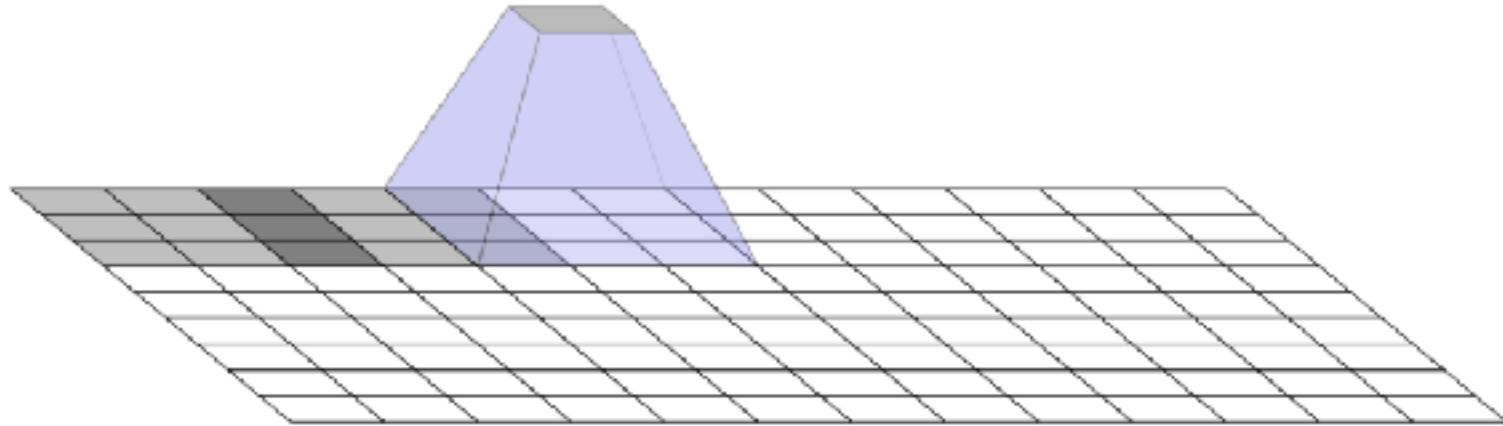
# Transposed Convolution & Overlap

---



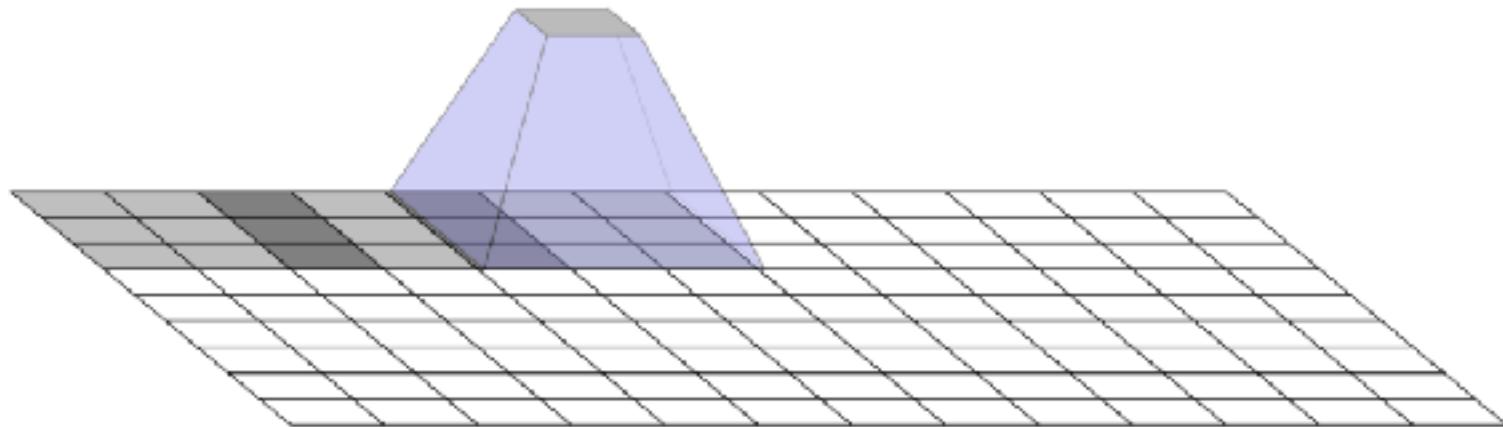
# Transposed Convolution & Overlap

---



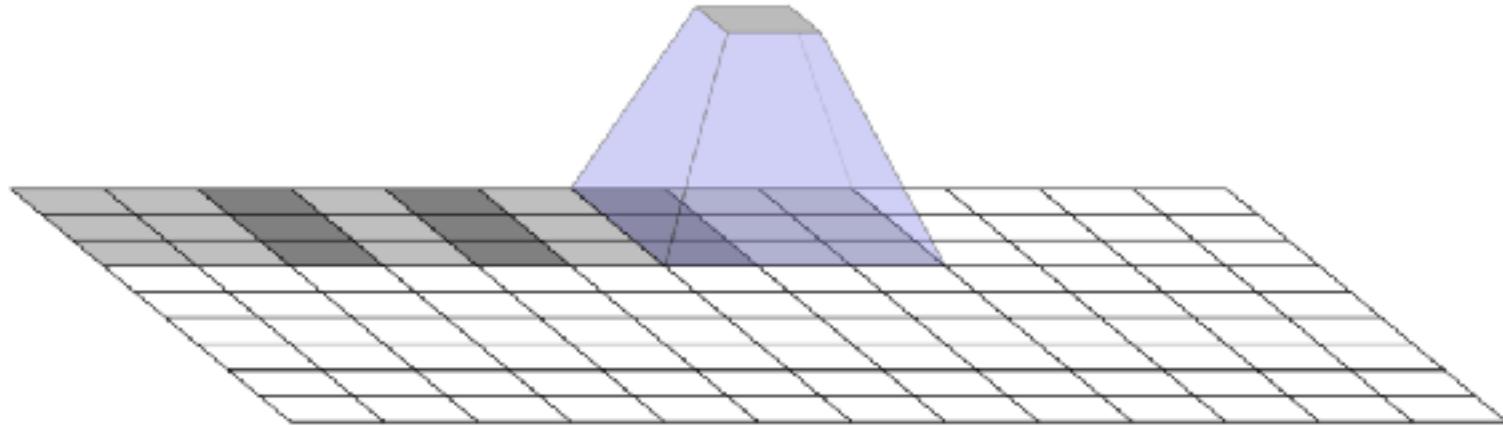
# Transposed Convolution & Overlap

---



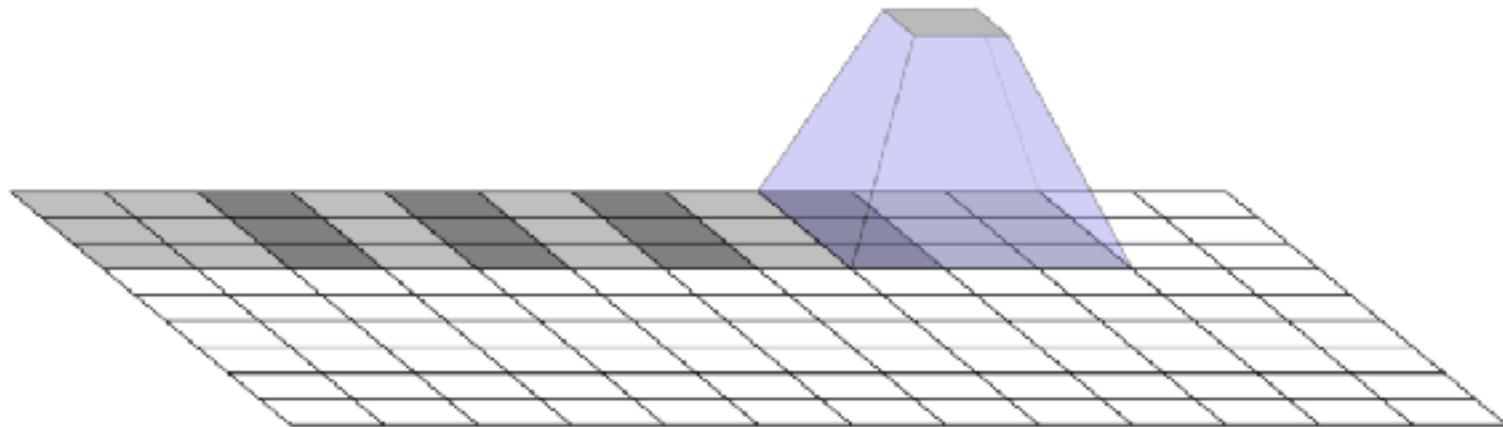
# Transposed Convolution & Overlap

---



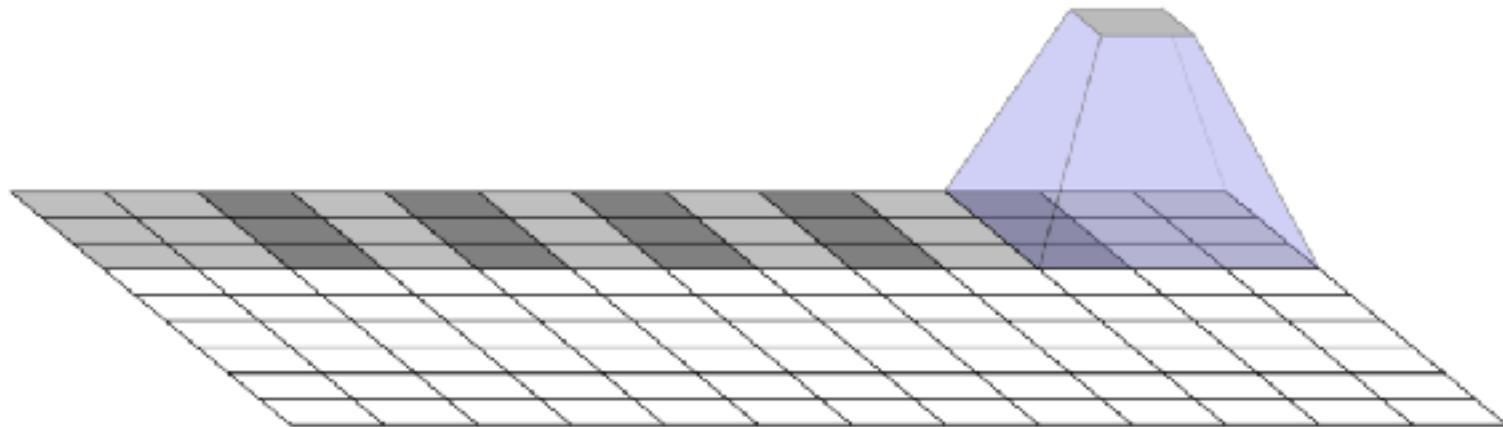
# Transposed Convolution & Overlap

---



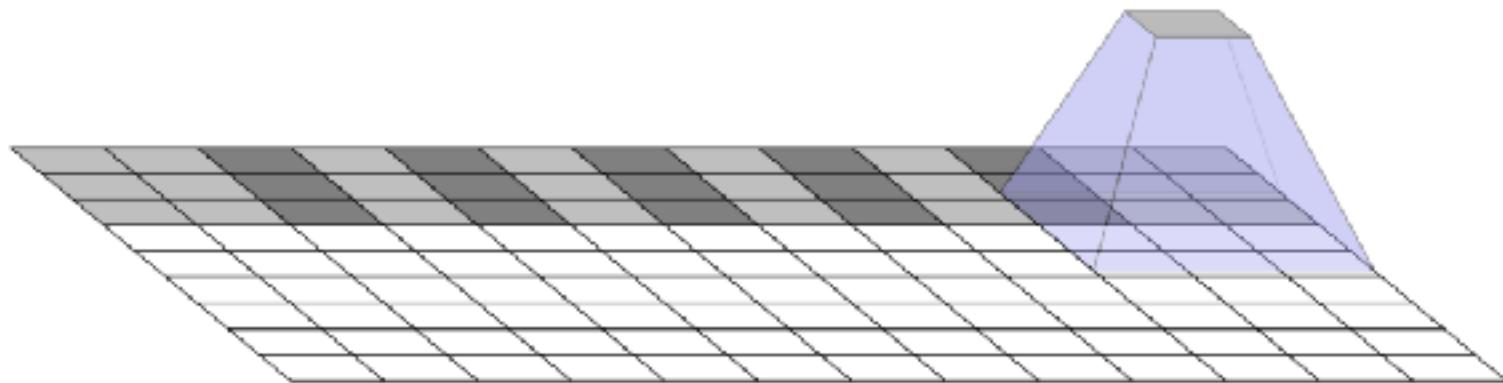
# Transposed Convolution & Overlap

---



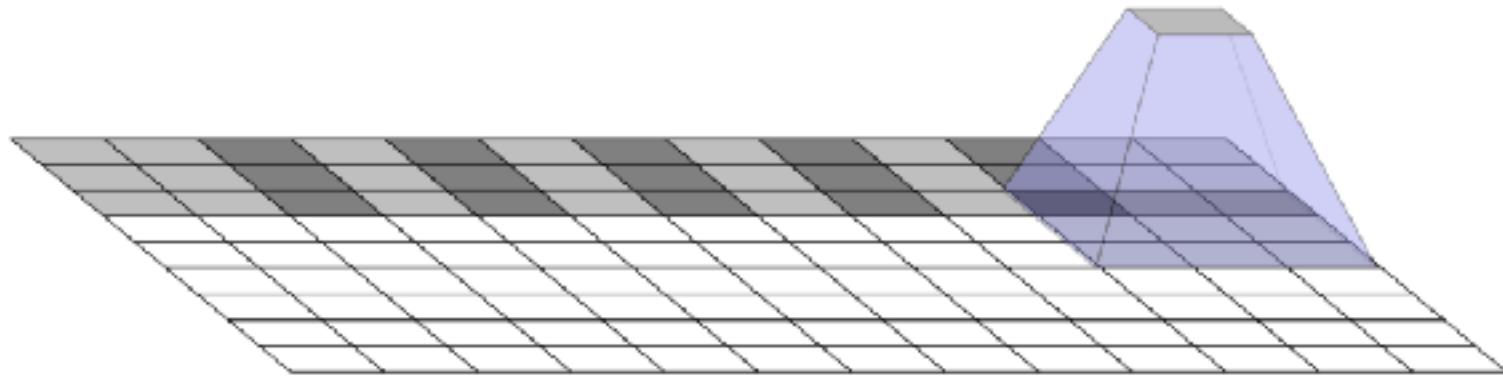
# Transposed Convolution & Overlap

---



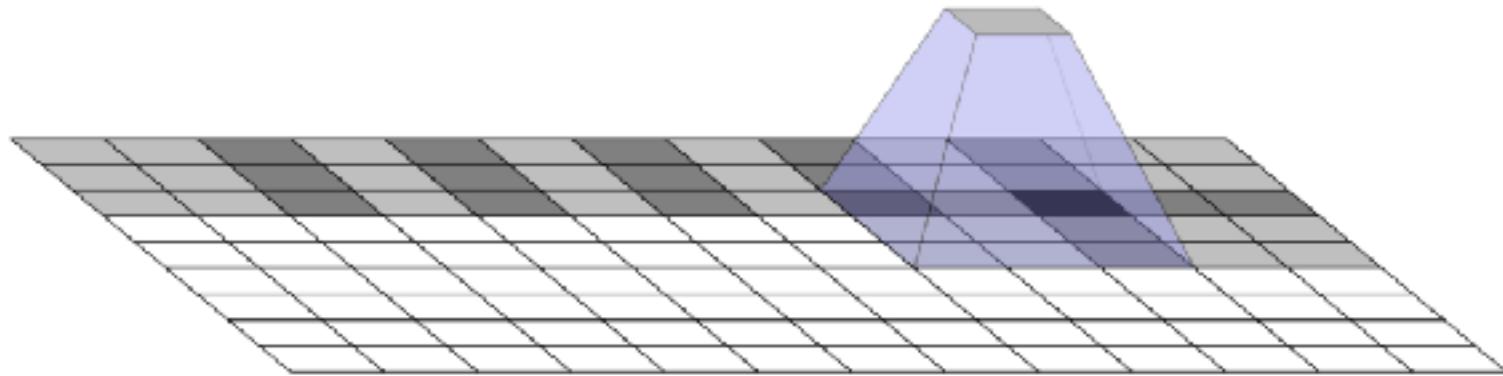
# Transposed Convolution & Overlap

---



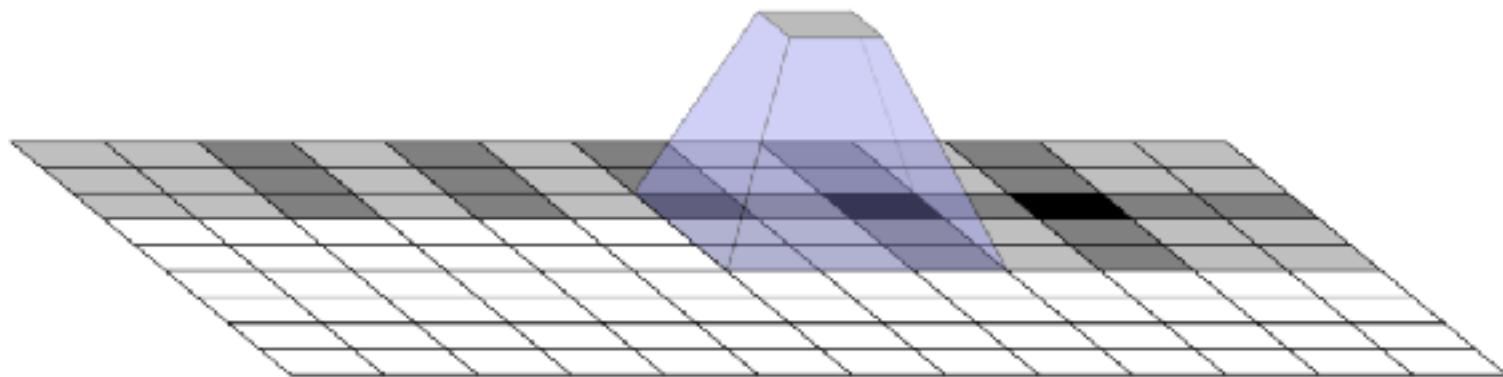
# Transposed Convolution & Overlap

---



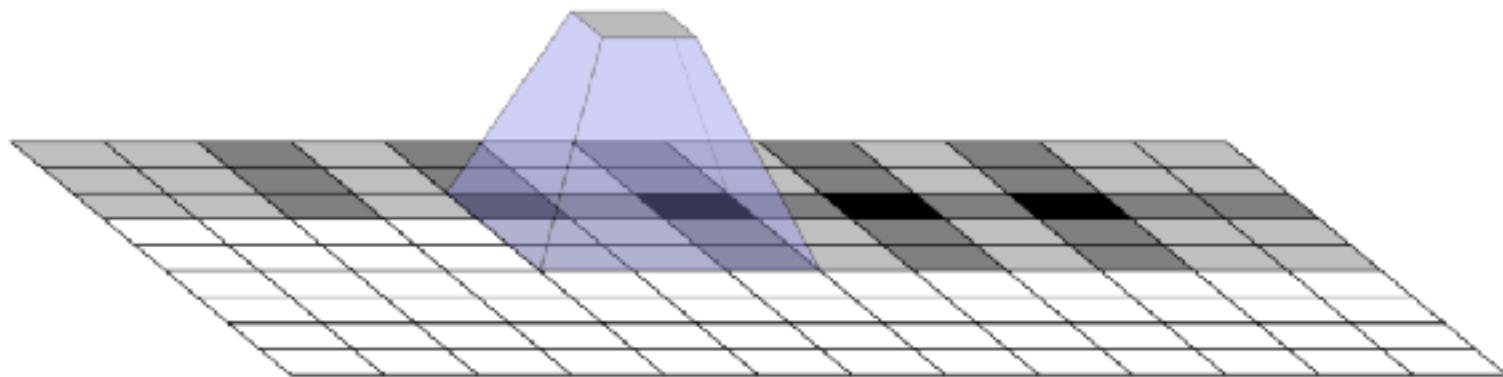
# Transposed Convolution & Overlap

---



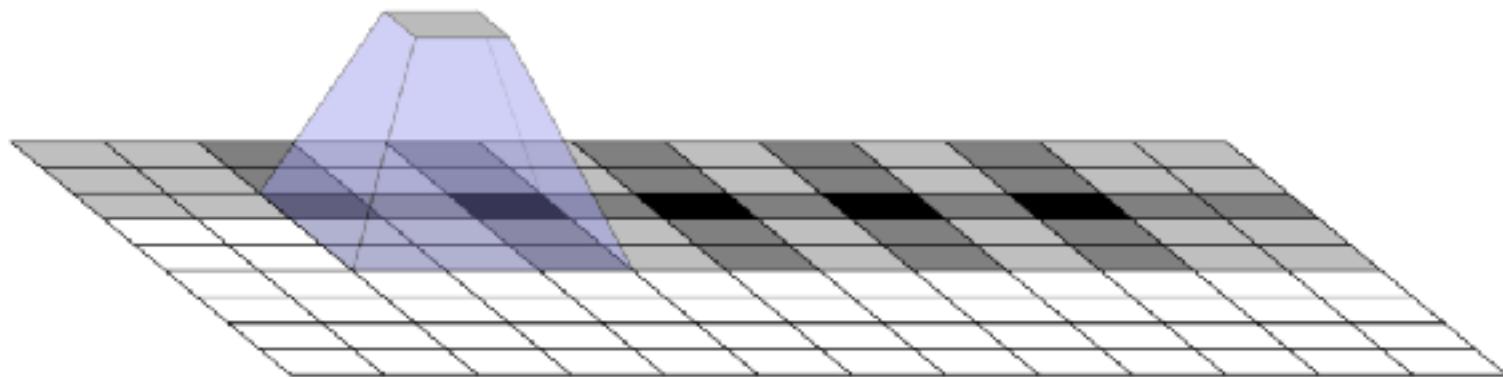
# Transposed Convolution & Overlap

---



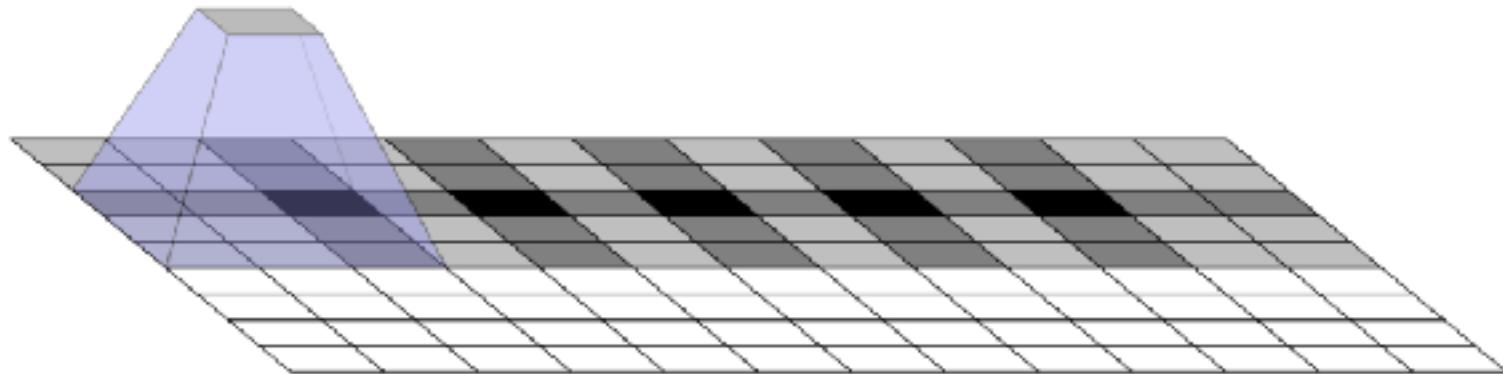
# Transposed Convolution & Overlap

---



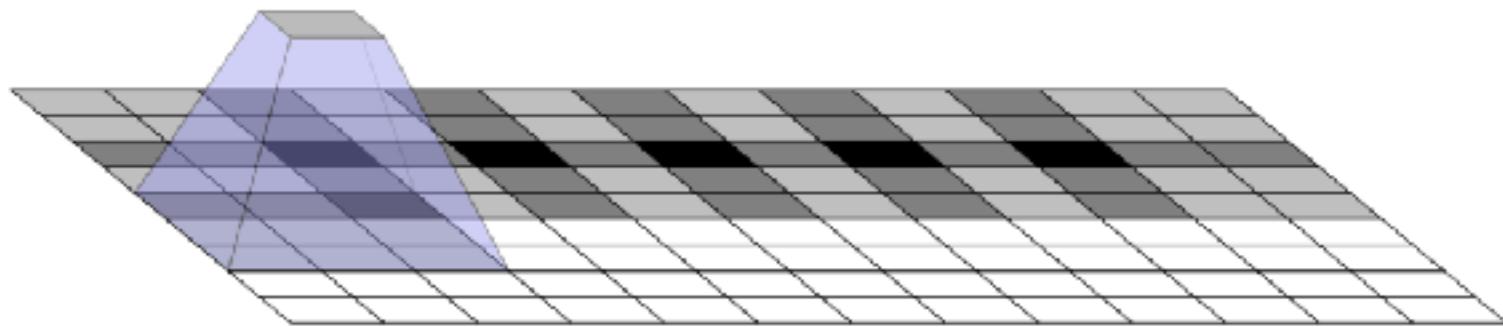
# Transposed Convolution & Overlap

---



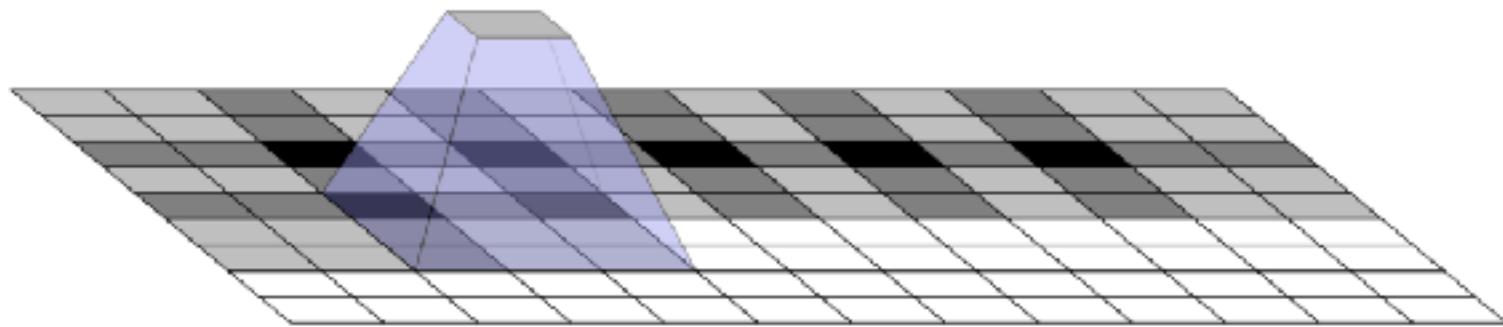
# Transposed Convolution & Overlap

---



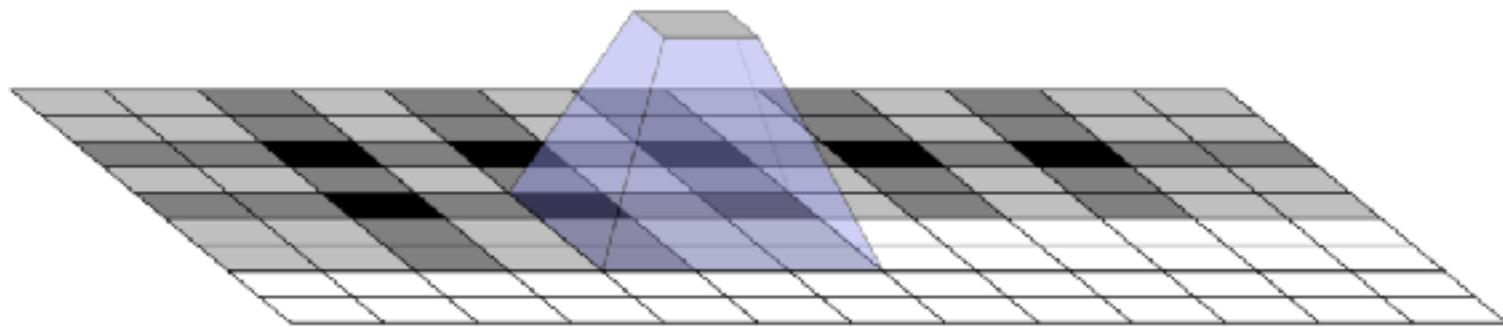
# Transposed Convolution & Overlap

---



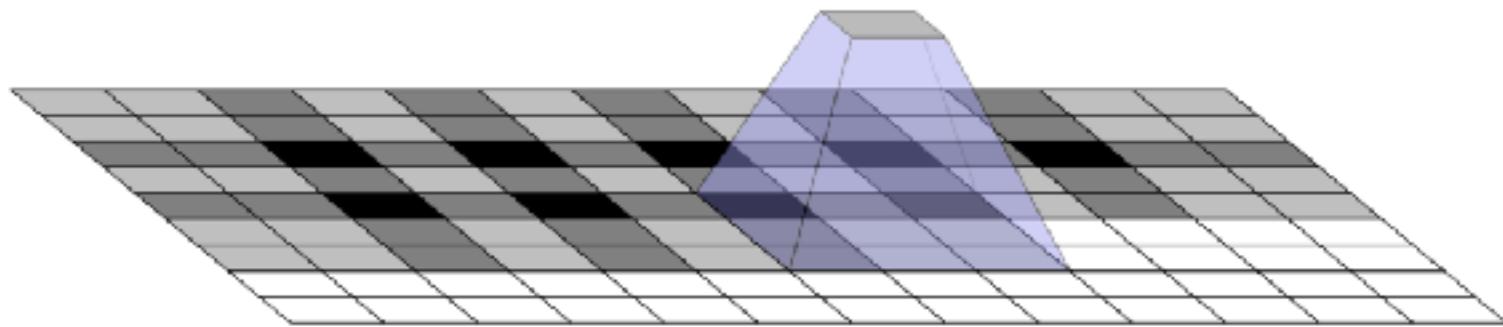
# Transposed Convolution & Overlap

---



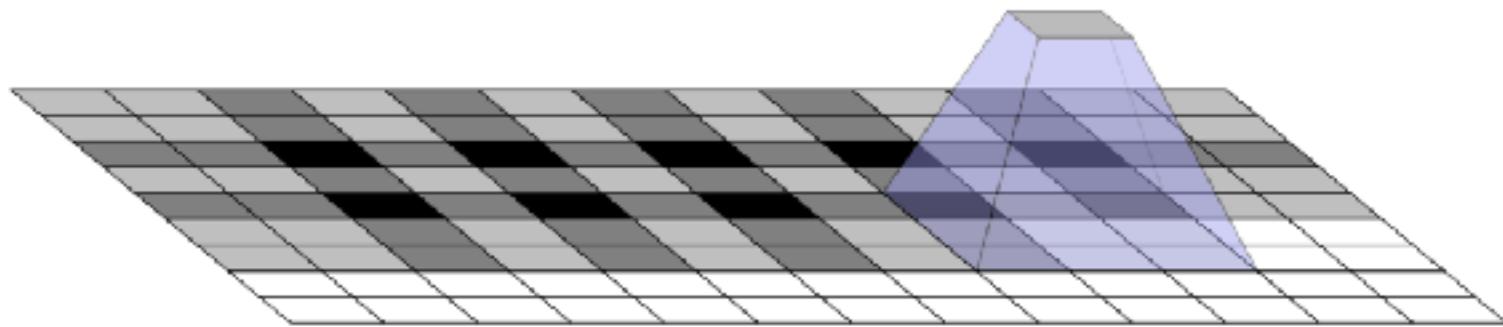
# Transposed Convolution & Overlap

---



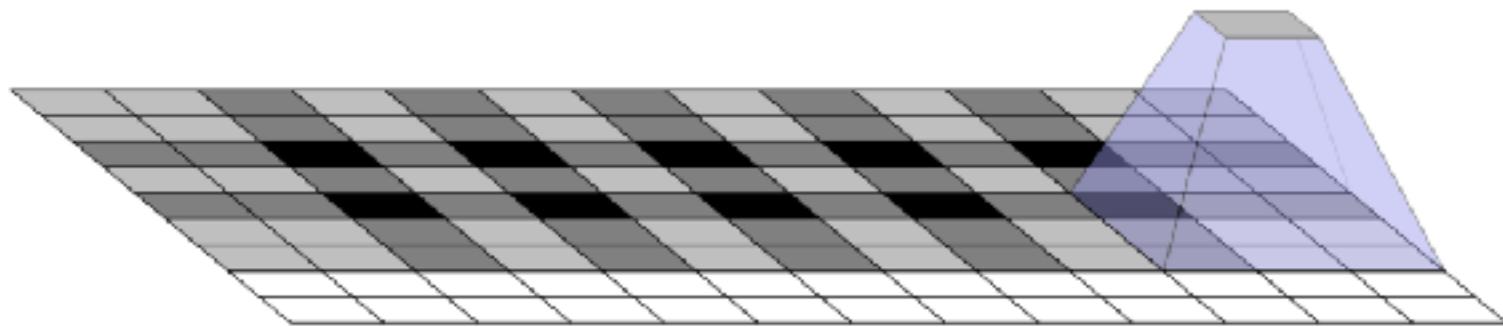
# Transposed Convolution & Overlap

---



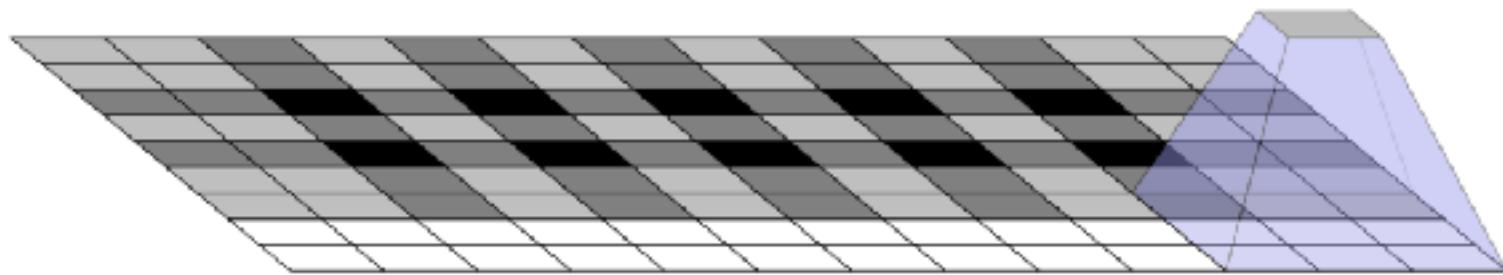
# Transposed Convolution & Overlap

---



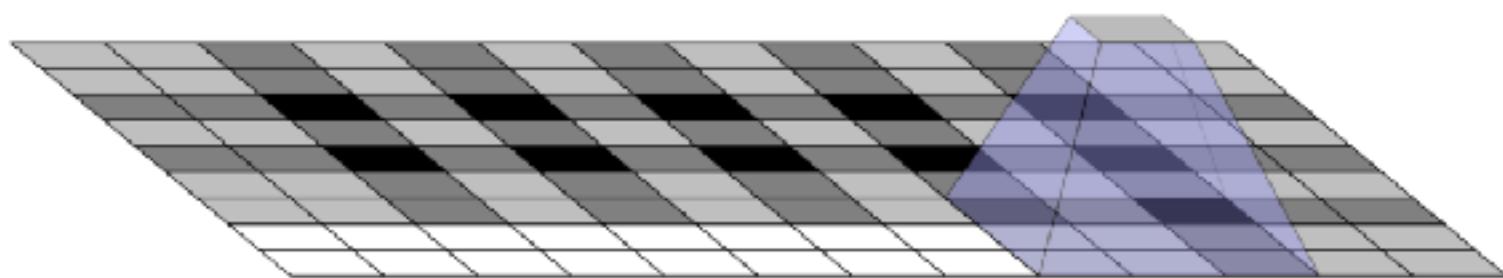
# Transposed Convolution & Overlap

---



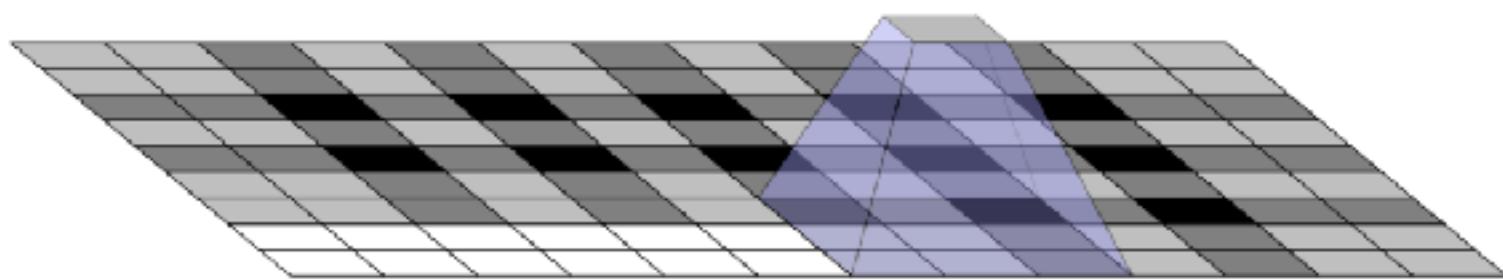
# Transposed Convolution & Overlap

---



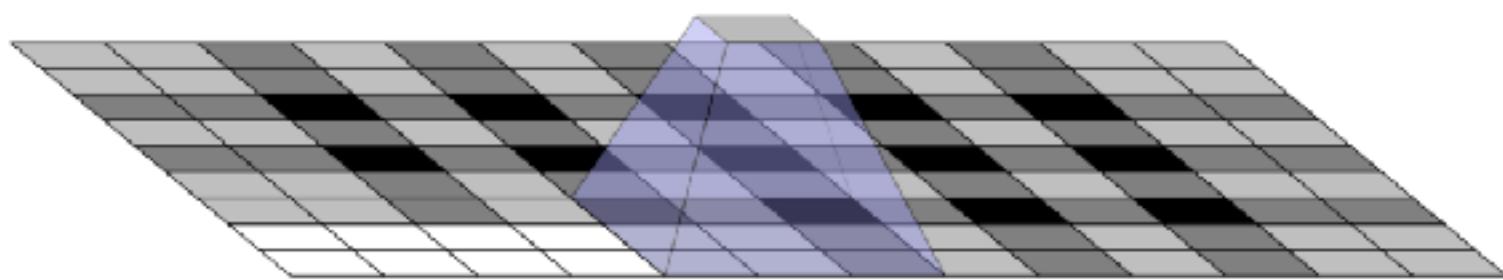
# Transposed Convolution & Overlap

---



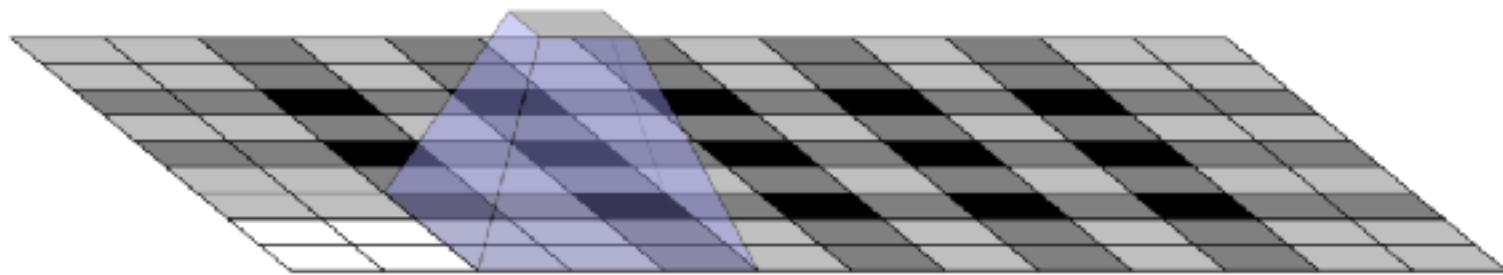
# Transposed Convolution & Overlap

---



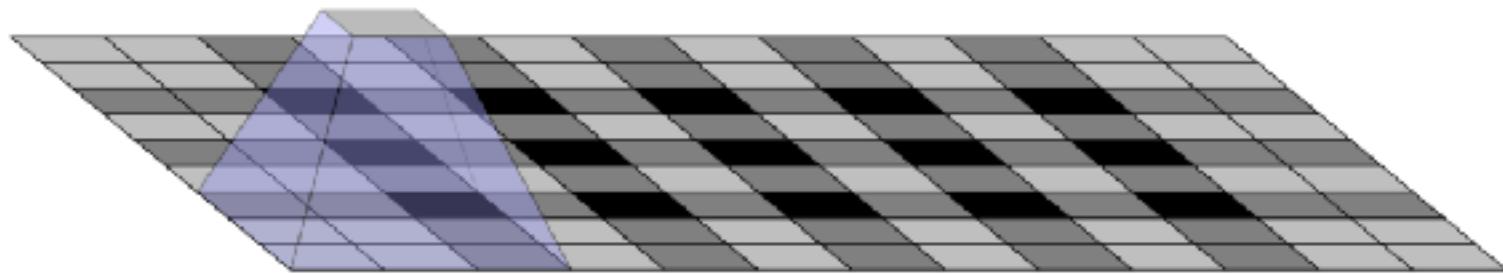
# Transposed Convolution & Overlap

---



# Transposed Convolution & Overlap

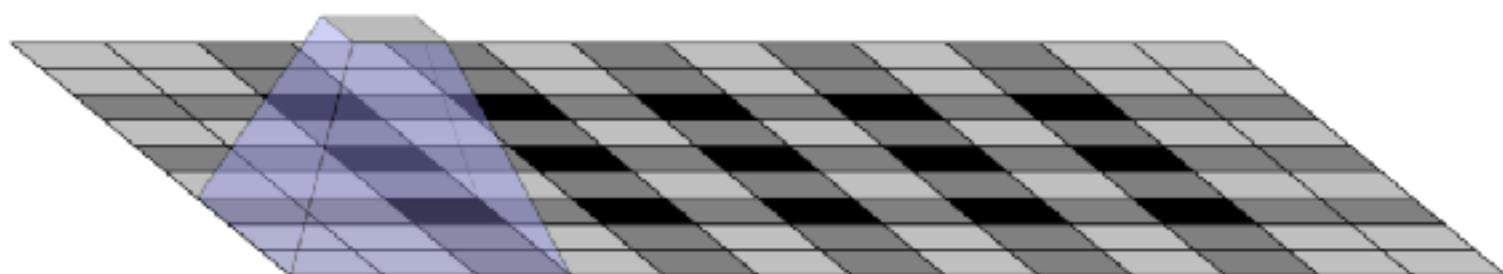
---



## Transposed Convolution & Overlap

---

- In particular, deconvolution has uneven overlap when the **kernel size** (the output window size) is **not divisible** by the **stride** (the spacing between points on the top).



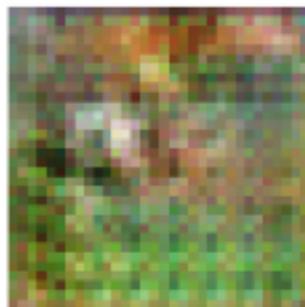
# Resize Convolution

---

- ▶ Idea
  - ▶ Resize the image (using [nearest-neighbor interpolation](#) or [bilinear interpolation](#)) and then do a convolutional layer
- ▶ NN-Resize Convolution
- ▶ Bilinear-Resize Convolution

# Resize Convolution

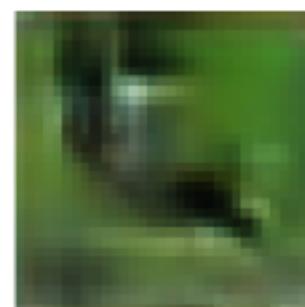
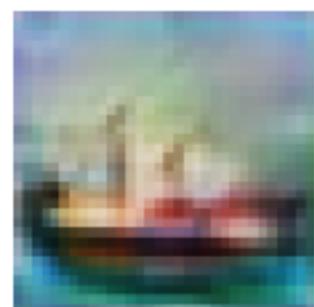
---



Deconv in last two layers.  
Other layers use resize-convolution.  
*Artifacts of frequency 2 and 4.*



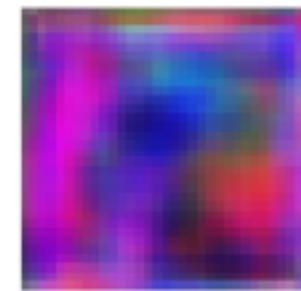
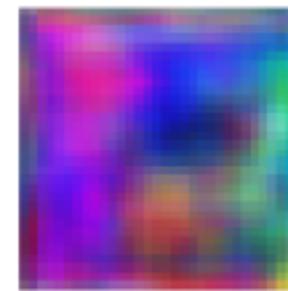
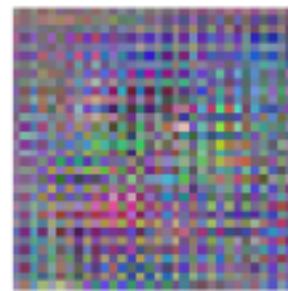
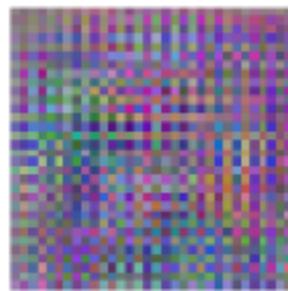
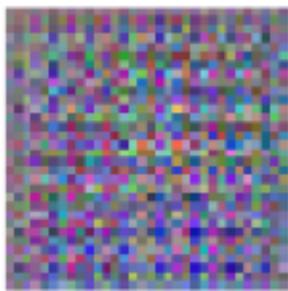
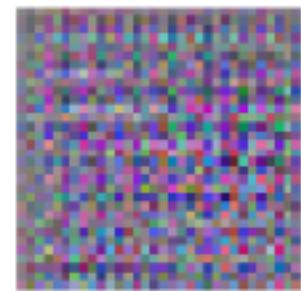
Deconv only in last layer.  
Other layers use resize-convolution.  
*Artifacts of frequency 2.*



All layers use resize-convolution.  
*No artifacts.*

# Resize Convolution

---



Deconvolution in last two layers.  
*Artifacts prior to any training.*

Deconvolution only in last layer.  
*Artifacts prior to any training.*

All layers use resize-convolution.  
*No artifacts before or after training.*

# Resize Convolution

---



Using deconvolution.  
*Heavy checkerboard artifacts.*



Using resize-convolution.  
*No checkerboard artifacts.*

# Artifacts in Gradients - Resize Convolution

---



DeepDream only applying the neural network to a fixed position.  
*Severe artifacts.*



DeepDream applying the network to a different position each step.  
*Reduced artifacts.*