Lecture Notes for

# Neural Networks
# and Machine Learning

Fully Convolutional Learning I:
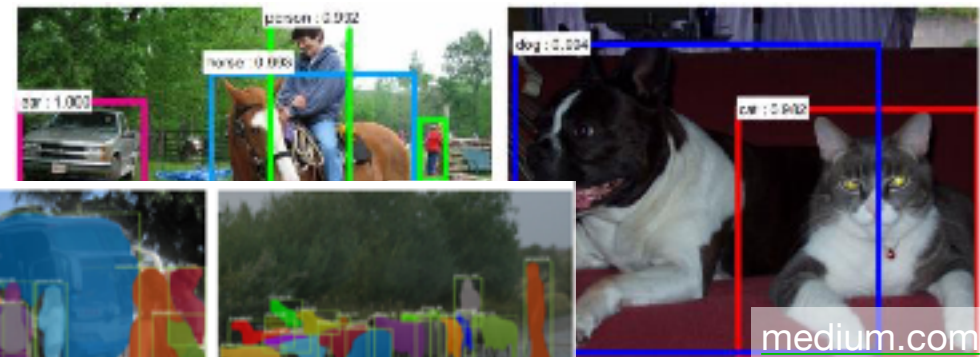Introduction to
Semantic Segmentation

# Logistics and Agenda

- Logistics
  - Lab Grading Update
- Agenda
  - Segmentation
    - Intro to Semantic (this time)
    - Object (partially this time)
    - Instance (next time)

# Types of Fully Convolutional Problems

- Semantic Segmentation

- Object Detection

- Instance Segmentation



medium.com

medium.com

He et al., Mask r-cnn, 2018

# Introduction to Semantic Segmentation

Karandeep Singh @kdpsinghlab · 10h  ···
Statistician: Do you ever use statistics?

ML researcher: Nope. Never.

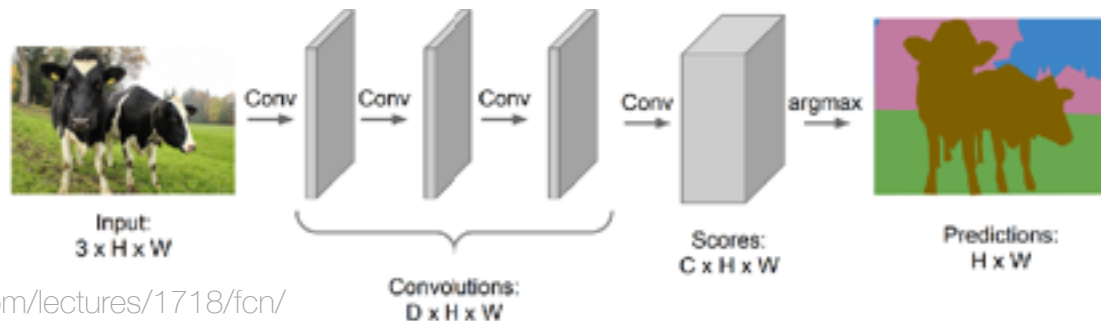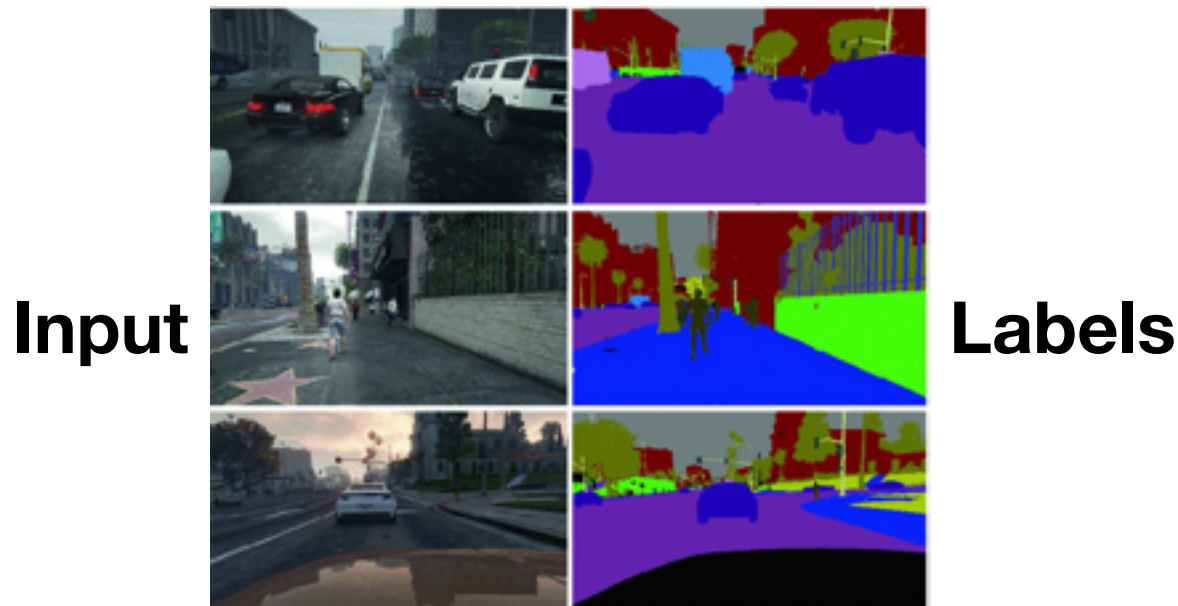Statistician: What about when reading a paper?

ML: Nope. Never.

Statistician: Ok. So if you're reading an ML paper comparing lots of models, how do you know which one is the best?

ML: Bold font.

# Semantic Segmentation

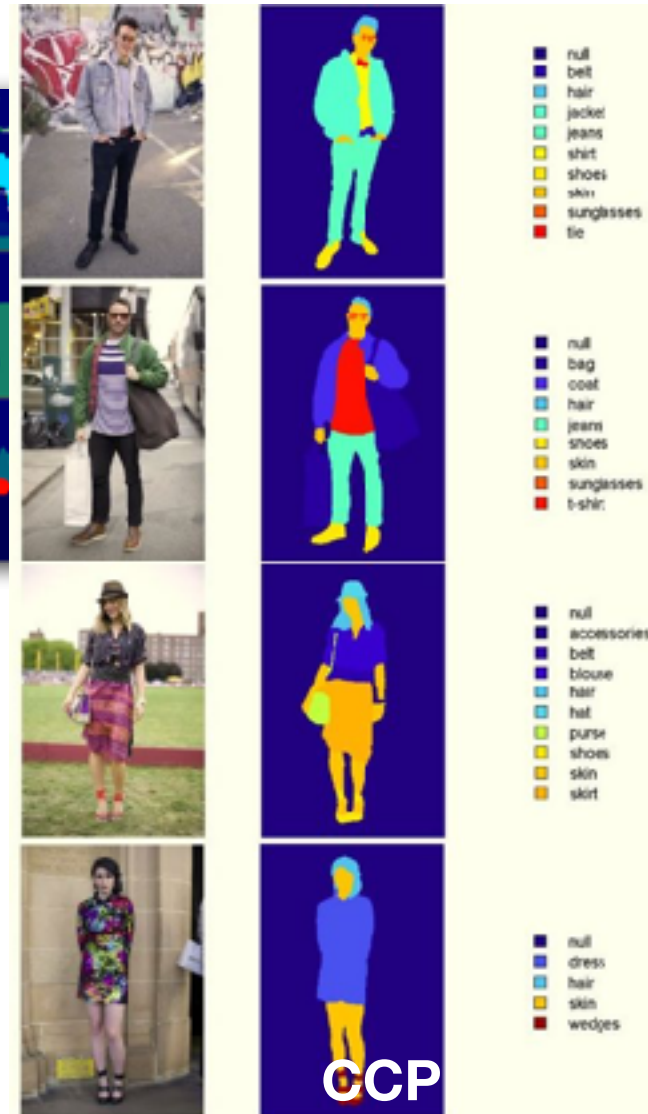- Given a set of pixels, classify each pixel according to what instance it belongs



**Input**

**Labels**



https://tjmachinelearning.com/lectures/1718/fcn/

# Popular Semantic Segmentation Datasets

**COCO** http://cocodataset.org/



**Cityscapes**

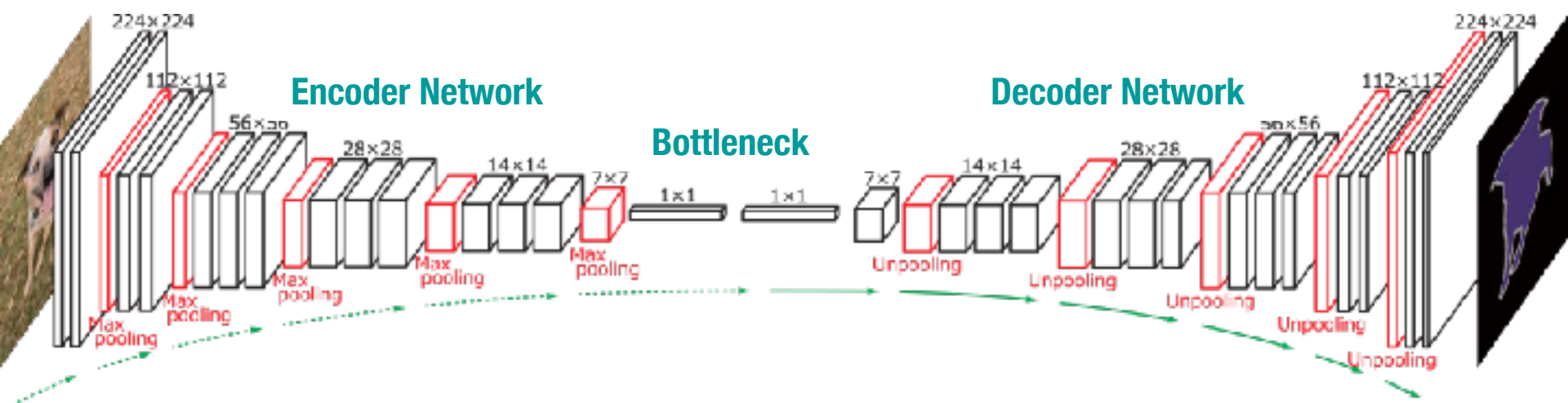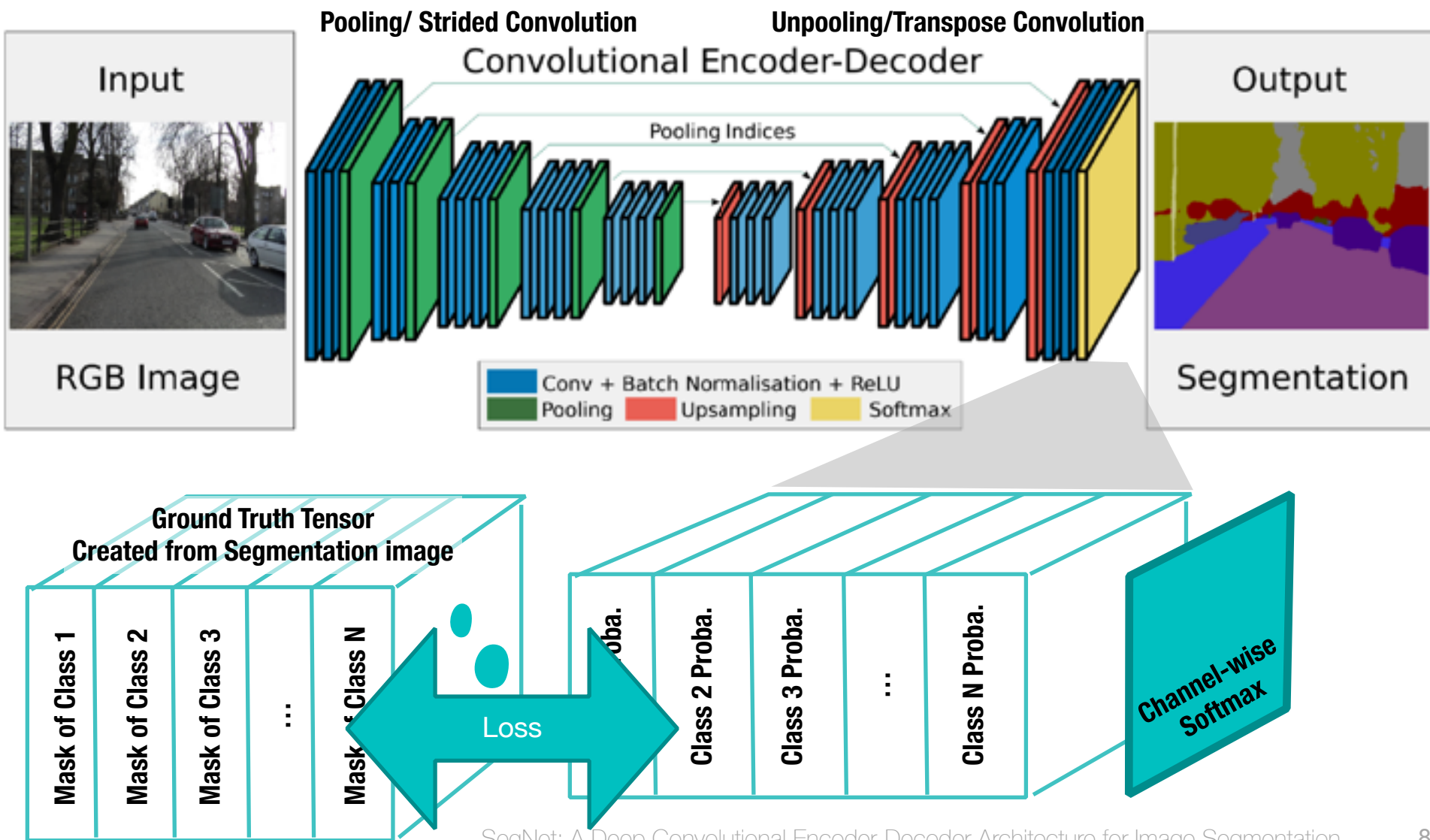**CCP**

# Early Training Methods (Pre 2018)



- Init Encoder with traditional CNN (like VGG or DarkNet)
- Freeze encoder and train decoder with segmented image maps
- Unfreeze encoder and fine tune
  - Repeat tuning as needed

# Putting it all together



**Pooling/ Strided Convolution**      **Unpooling/Transpose Convolution**

Convolutional Encoder-Decoder

Input — RGB Image

Pooling Indices

- Conv + Batch Normalisation + ReLU
- Pooling
- Upsampling
- Softmax

Output — Segmentation

**Ground Truth Tensor Created from Segmentation image**

Mask of Class 1 | Mask of Class 2 | Mask of Class 3 | ... | Mask of Class N

Loss

Class 1 Proba. | Class 2 Proba. | Class 3 Proba. | ... | Class N Proba.

Channel-wise Softmax

SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation

# Putting it all together



**Pooling/ Strided Convolution**

**Unpooling/Transpose Convolution**

Input — RGB Image

Convolutional Encoder-Decoder

Pooling Indices

Output — Segmentation

- Conv + Batch Normalisation + ReLU
- Pooling
- Upsampling
- Softmax

**Self Test:**
Does it change the architecture if the Image input size changes?

Class 1 Proba. | Class 2 Proba. | Class 3 Proba. | ... | Class N Proba.

Channel-wise Softmax
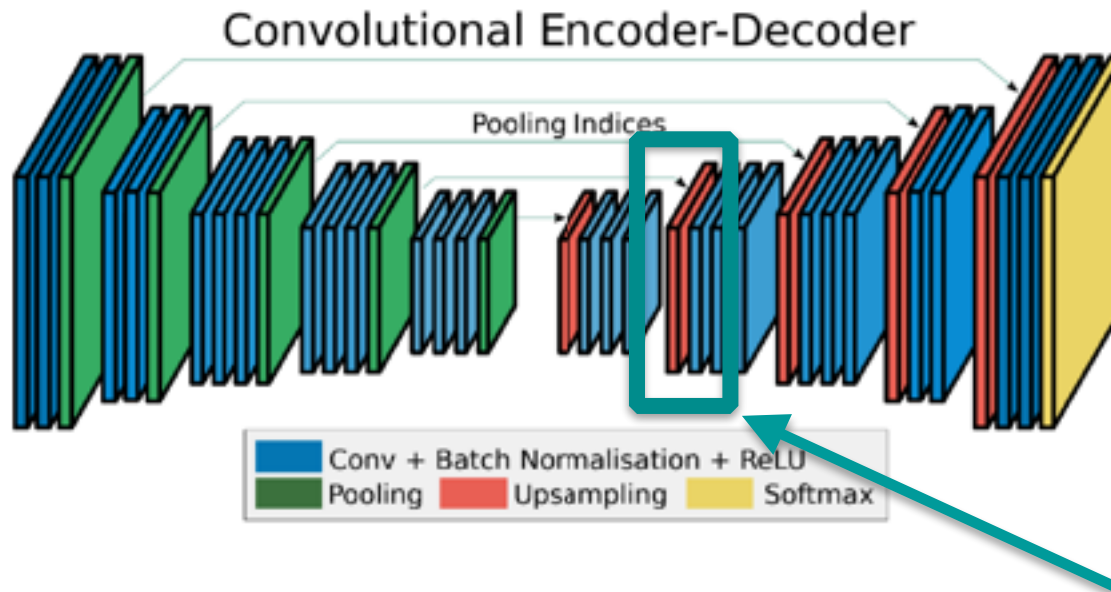
# Upsampling Layers



Shit Academics Say @Academi... · 22h  ···
not wrong

> monstera adansonii @yourn... · 2d
> everything is peer reviewed if your friends are judgmental enough

# Decoder Network



Convolutional Encoder-Decoder

Pooling Indices

Conv + Batch Normalisation + ReLU
Pooling  Upsampling  Softmax

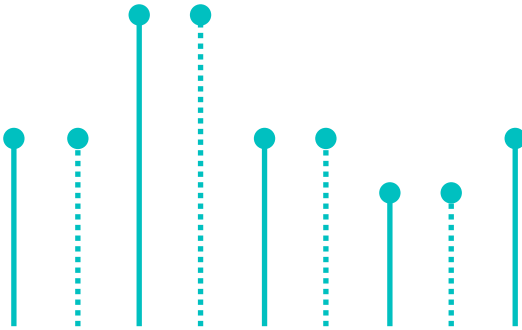Some researcher started calling this **deconvolution**.

If you use that term in this class, **you fail**.

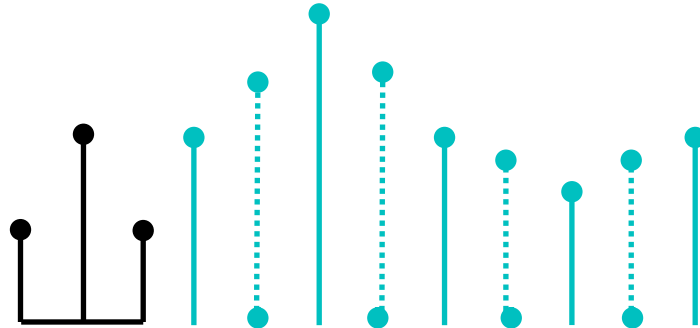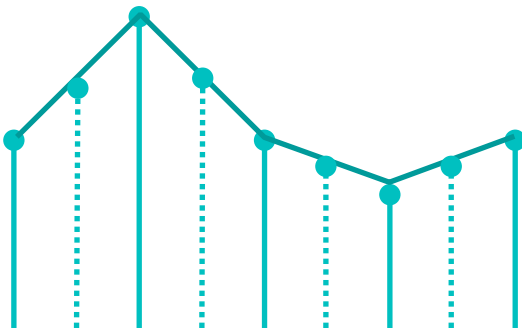This is upsampling and then convolution, but **now the interpolation filters are learned**!!

# Integer Upsampling via Interpolation

**Nearest Neighbor**

**Linear**

**Cubic**

All are equivalent to inserting zeros and applying convolutional filter

# Image Upsampling, Integer Factor

- Insert Zeros
- Convolve

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

| 1 | | 2 | | 3 | | 4 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| 5 | | 6 | | 7 | | 8 | |
| | | | | | | | |
| 9 | | 10 | | 11 | | 12 | |
| | | | | | | | |
| 13 | | 14 | | 15 | | 16 | |
| | | | | | | | |

| 0.25 | 0.5 | 0.25 |
|---|---|---|
| 0.5 | 1 | 0.5 |
| 0.25 | 0.5 | 0.25 |

Bilinear Filtering

Bicubic Filter

# Image Upsampling, Integer Factor



**Nearest Neighbor**

`UpSampling2D()`

**Bilinear**

`UpSampling2D(interpolation='bilinear')`

**Bicubic**

**Many Types of Upsampling, with varying computational cost:**

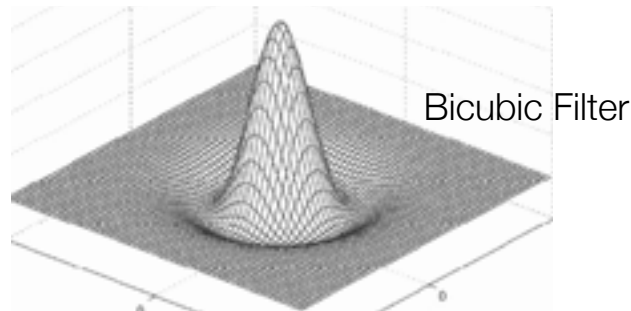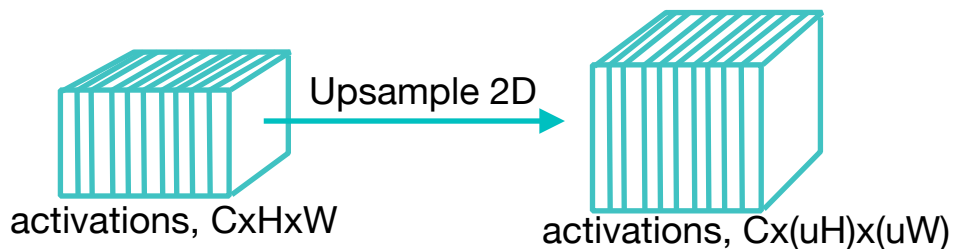area, bicubic, gaussian, lanczos3, lanczos5, mitchellcubic

Upsample 2D

activations, CxHxW

activations, Cx(uH)x(uW)

https://www.cs.toronto.edu/~guerzhoy/320/lec/upsampling.pdf

14

# What about transpose convolution?

Convolution as Matrix Multiplication

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

Transpose

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

like convolving with "reversed coefficients"

**Regular Convolution**

$b$
$a$
$c$

$z$  $y$  $x$

$ax$
$ay+bx$
$az+by+cx$
$bz+cy$
$cz$

**Transpose Convolution**

$b$
$a$
$c$

$x$  $y$  $z$

$az$
$ay+bz$
$ax+by+cz$
$bx+cy$
$cx$

# Transpose Convolution: Strides

Strided Convolution as Matrix Multiplication

$$\begin{bmatrix} y & x & 0 & 0 & 0 \\ 0 & z & y & x & 0 \\ 0 & 0 & 0 & z & y \end{bmatrix} \times \begin{bmatrix} 0 \\ a \\ b \\ c \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ az+by+cx \\ cz \end{bmatrix}$$
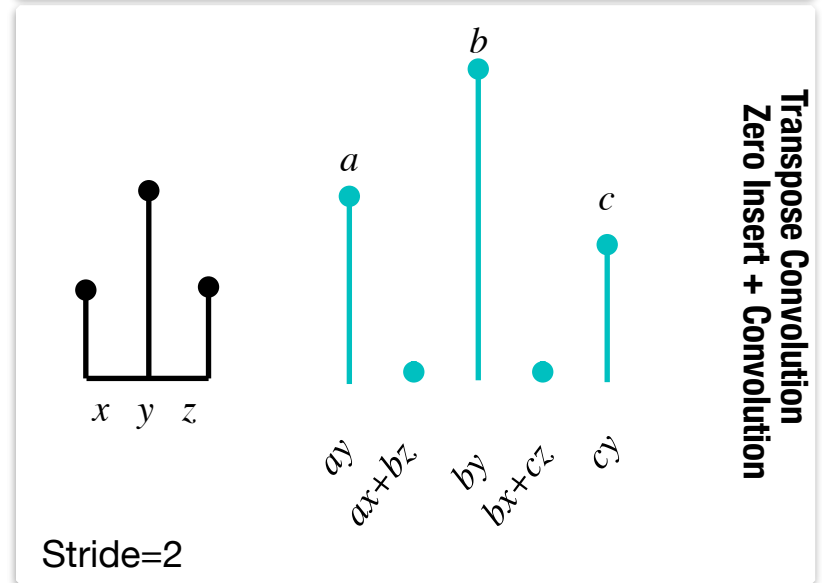
Transpose

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



**Regular Convolution**

Stride=2

**Transpose Convolution
Zero Insert + Convolution**

Stride=2

# Convolution after zero insertion

- Kernel size should be a symmetric multiple of the stride



four pixels + bias

two pixels + bias

one pixel + bias

four pixels + bias

four pixels + bias

3x3 not multiple of stride

4x4 multiple of stride

Bias needs to account for both when different numbers of pixels overlap with the kernel

Multiple of stride ensures that same number of active pixels overlap the kernel.

Stride = 2

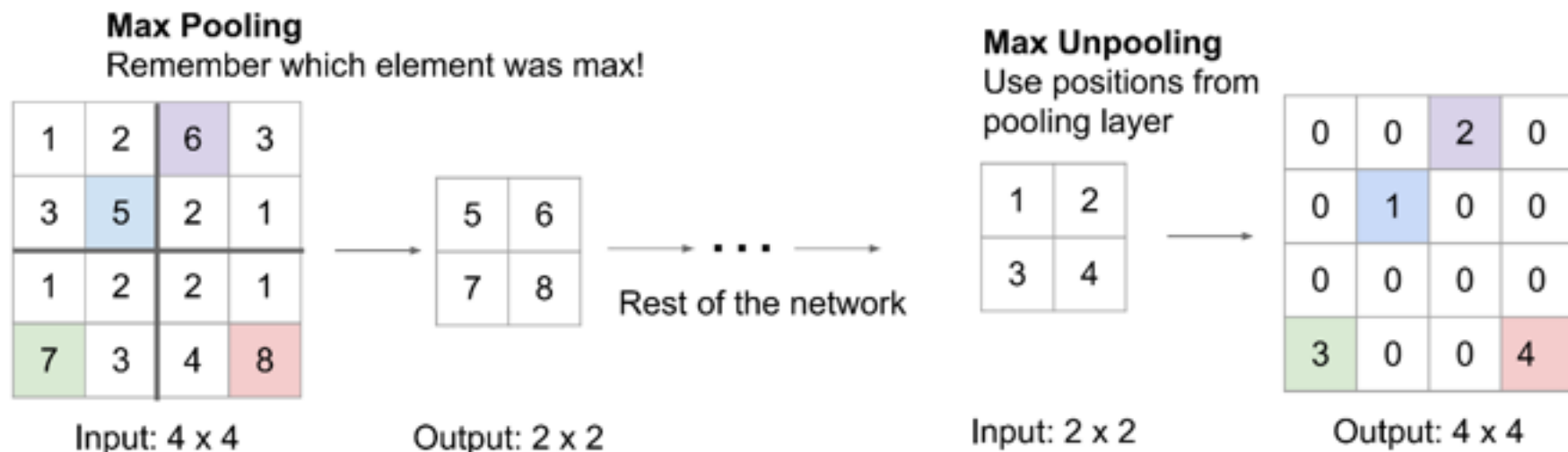# **Unpooling**: a different method of zero insertion

**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

- Unpooling: insert values to upsample where you pooled
- Why does this make sense? The upsampling happens much later in the network…
- And it increases computational overhead and memory to track indices…
- Not very advantageous…

# Back up Slides for Semantic Segmentation



**François Chollet** ✔
@fchollet · · ·

Every single character in Thomas the Tank Engine:

🛢️

(🙂)

8:28 PM · 2/28/23

**41.9K** Views **101** Likes **6** Retweets

**Alexis Taugeron** @ataugeron · 1d · · ·
What about the Troublesome Trucks?
163

**Ben Tseng** @BenjaminTseng · 1d · · ·
That show is the best illustration that
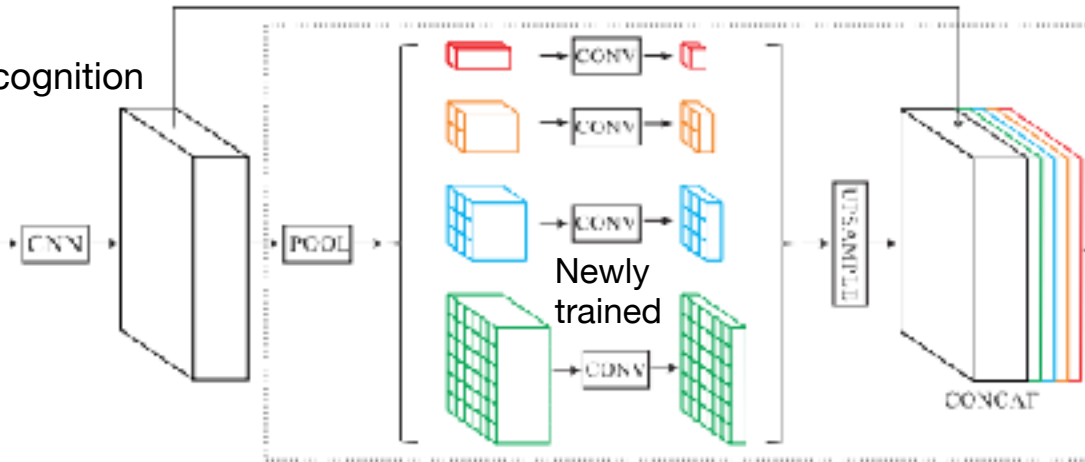sentience in machines won't lead to mass
displacement of human workers
740 ♡ 4

**Pyramid Scene Parsing Network (PSPNet)**



Pre-trained for object recognition
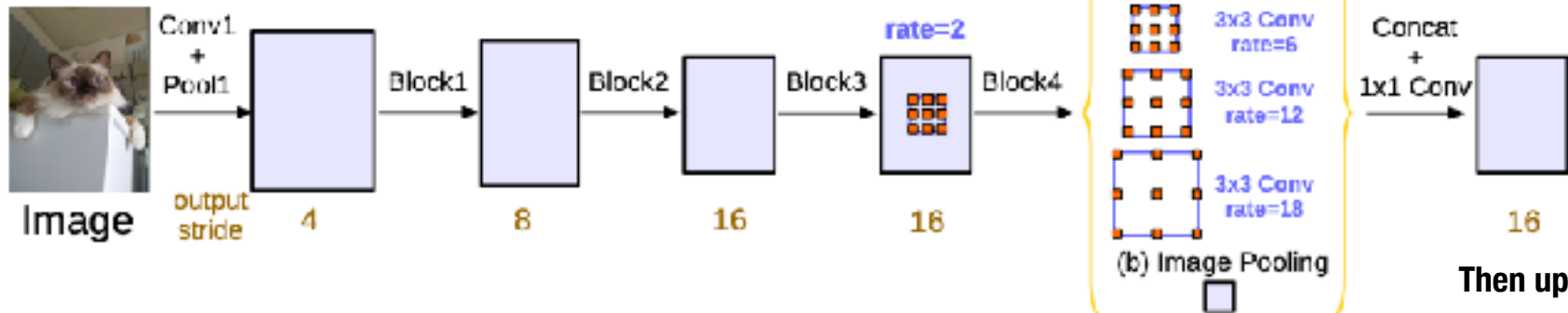
Newly trained

Newly trained

(a) Input Image    (b) Feature Map    (c) Pyramid Pooling Module    (d) Final Prediction

**DeepLabV3: Dilated Convolutions (Atrous Convolutions)**



Then upscaling→

https://towardsdatascience.com/semantic-segmentation-with-deep-learning-a-guide-and-code-e52fc8958823

20
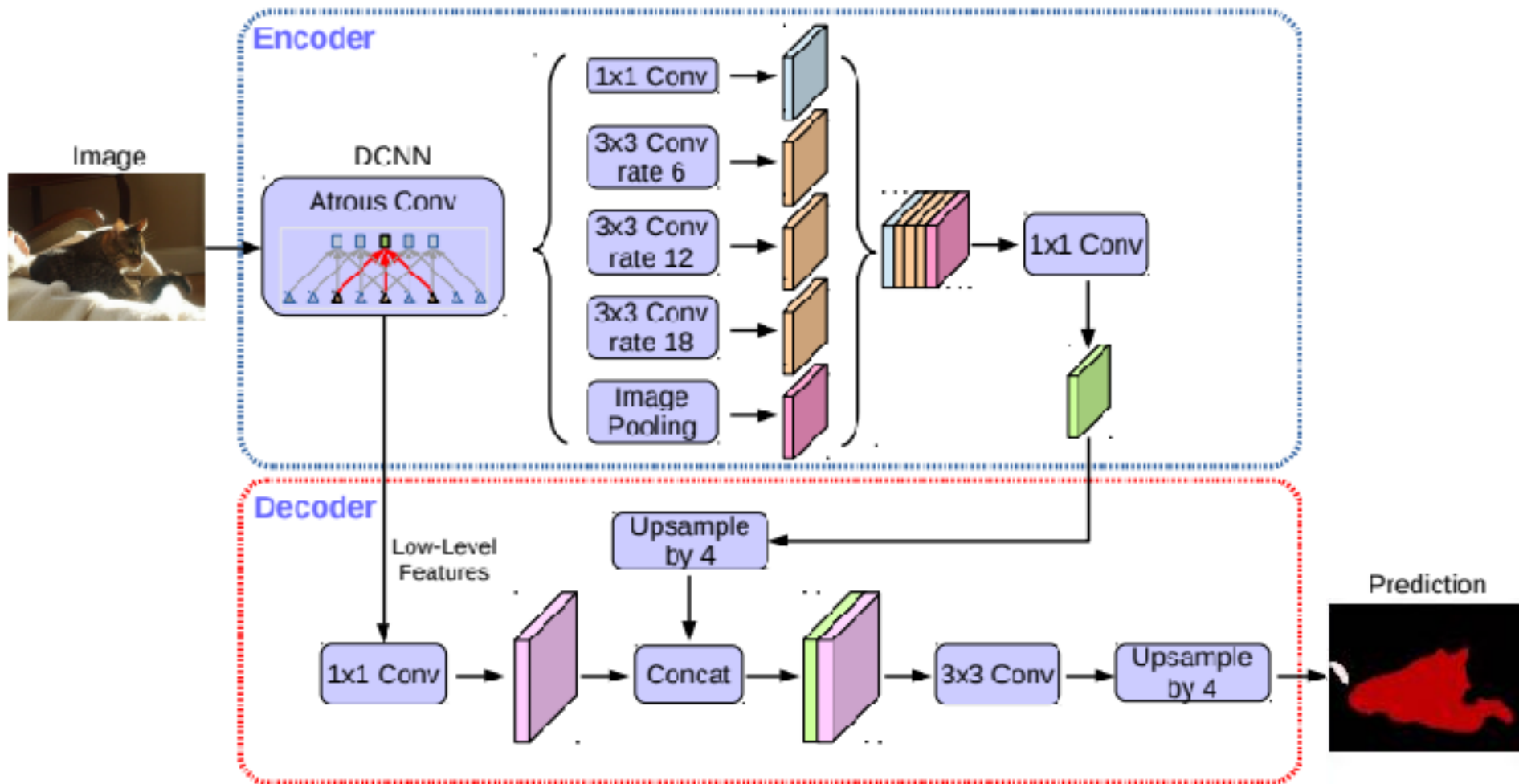
Outputs of convolution are the same size, except for edge effects!
But have advantage of processing at a different scale.

https://towardsdatascience.com/review-dilated-convolution-semantic-segmentation-9d5a5bd768f5

https://github.com/tensorflow/models/tree/master/research/deeplab

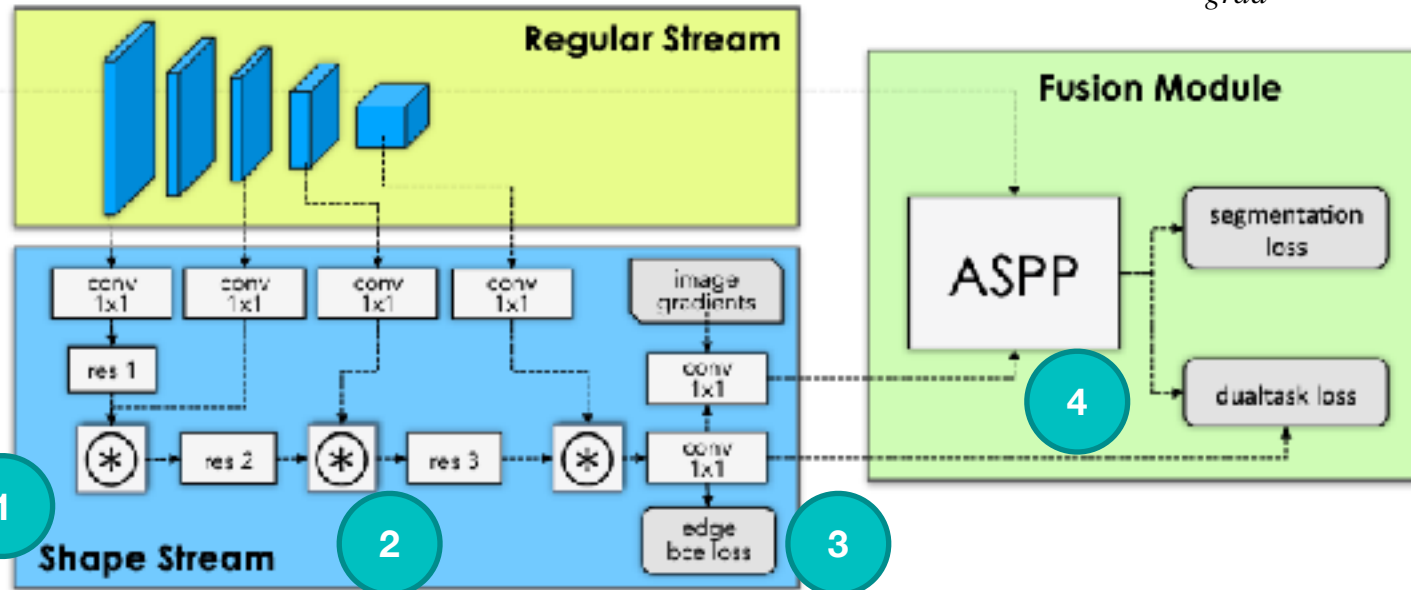https://towardsdatascience.com/semantic-segmentation-with-deep-learning-a-guide-and-code-e52fc8958823

# Gated-SCNN (Gate Shape CNN)

**1** Shape stream employs Traditional Image Processing for edge detection (**image gradients**)

**2** Uses activations to "gate" the image gradient. $\sigma(A) \odot I_{grad}$



Figure 2: **GSCNN architecture.** Our architecture constitutes of two main streams. The regular stream and the shape stream. The regular stream can be any backbone architecture. The shape stream focuses on shape processing through a set of residual blocks, Gated Convolutional Layers (GCL) and supervision. A fusion module later combines information from the two streams in a multi-scale fashion using an Atrous Spatial Pyramid Pooling module (ASPP). High quality boundaries on the segmentation masks are ensured through a Dual Task Regularizer.

**3** Also uses Labeled Boundaries in BCE Edge Loss Function

**4** Merges segmentation with edges for finer masks. Concatenate + atrous convolution

https://heartbeat.fritz.ai/a-2019-guide-to-semantic-segmentation-ca8242f5a7fc

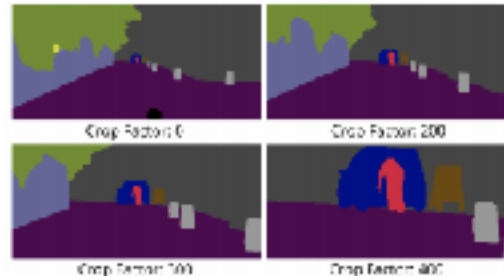Figure 3: Illustration of the crops used for the distance-based evaluation.

Figure 4: Predictions at diff. crop factors.

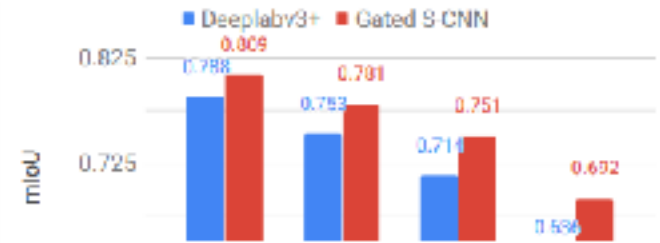Figure 5: **Distance-based evaluation**: Comparison of mIoU at different crop factors.

| Method | road | s.walk | build. | wall | fence | pole | t-light | t-sign | veg | terrain | sky | person | rider | car | truck | bus | train | motor | bike | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LRR [18] | 97.7 | 79.9 | 90.7 | 44.4 | 48.6 | 58.6 | 68.2 | 72.0 | 92.5 | 69.3 | 94.7 | 81.6 | 60.0 | 94.0 | 43.6 | 56.8 | 47.2 | 54.8 | 69.7 | 69.7 |
| DeepLabV2 [9] | 97.9 | 81.3 | 90.3 | 48.8 | 47.4 | 49.6 | 57.9 | 67.3 | 91.9 | 69.4 | 94.2 | 79.8 | 59.8 | 93.7 | 56.5 | 67.5 | 57.5 | 57.7 | 68.8 | 70.4 |
| Piecewise [32] | 98.0 | 82.6 | 90.6 | 44.0 | 50.7 | 51.1 | 65.0 | 71.7 | 92.0 | 72.0 | 94.1 | 81.5 | 61.1 | 94.3 | 61.1 | 65.1 | 53.8 | 61.6 | 70.6 | 71.6 |
| PSP-Net [58] | 98.2 | 85.8 | 92.8 | 57.5 | 65.9 | 62.6 | 71.8 | 80.7 | 92.4 | 64.5 | 94.8 | 82.1 | 61.5 | 95.1 | 78.6 | 88.3 | 77.9 | 68.1 | 78.0 | 78.8 |
| DeepLabV3+ [11] | 98.2 | 84.9 | 92.7 | 57.3 | 62.1 | 65.2 | 68.6 | 78.9 | 92.7 | 63.5 | 95.3 | 82.3 | 62.8 | 95.4 | 85.3 | 89.1 | 80.9 | 64.6 | 77.3 | 78.8 |
| Ours (GSCNN) | 98.3 | 86.3 | 93.3 | 55.8 | 64.0 | 70.8 | 75.9 | 83.1 | 93.0 | 65.1 | 95.2 | 85.3 | 67.9 | 96.0 | 80.8 | 91.2 | 83.3 | 69.6 | 80.4 | 80.8 |

Table 1: Comparison in terms of IoU vs state-of-the-art baselines on the Cityscapes val set.

**mIoU == mean Intersection over Union** $= \dfrac{\text{Area of Overlap}}{\text{Area of Union}}$

Lecture Notes for

# Neural Networks and Machine Learning

FCN Learning

**Next Time:**
Fully Convolutional Objects
**Reading:** None