



СИЛАЙН ОБРАЗОВАНИЕ

Онлайн-образование

Image Segmentation



Ческидова Евгения

Deep Learning engineer
Wolf3d
t.me/fogside

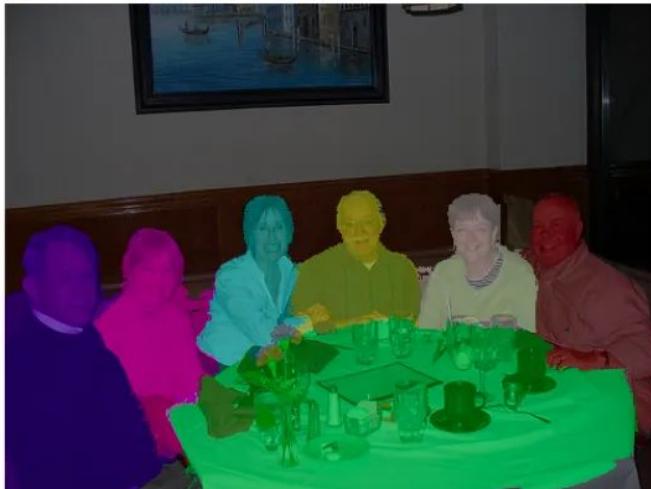
Plan

- Task definition
 - Semantic Segmentation
 - Instance Segmentation
- Applications
- Metrics
- Architectures
- Loss functions
- Frameworks
- Datasets
- Tutorials

Semantic and instance segmentation



Semantic Segmentation



Instance Segmentation

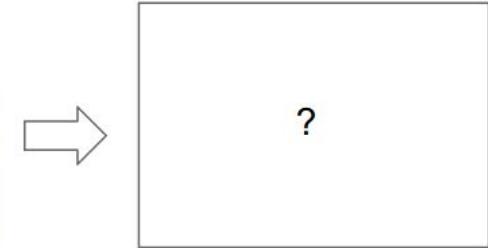
- In **semantic segmentation**, all objects of the same type are marked using one class label while in **instance segmentation** similar objects get their own separate labels.

Semantic Segmentation: The Problem



**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.

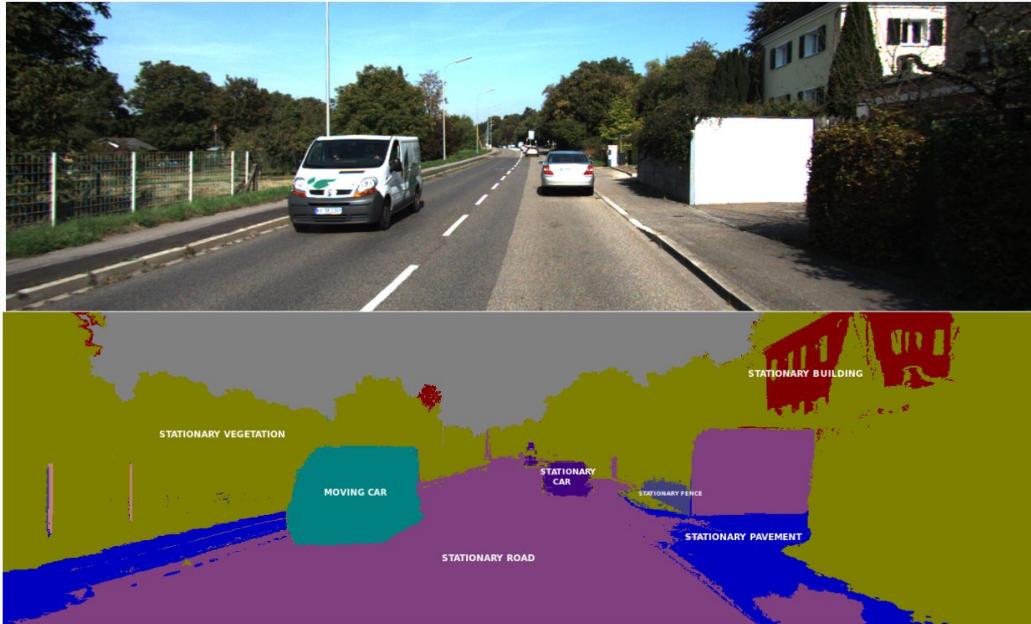




ОНЛАЙН ОБРАЗОВАНИЕ

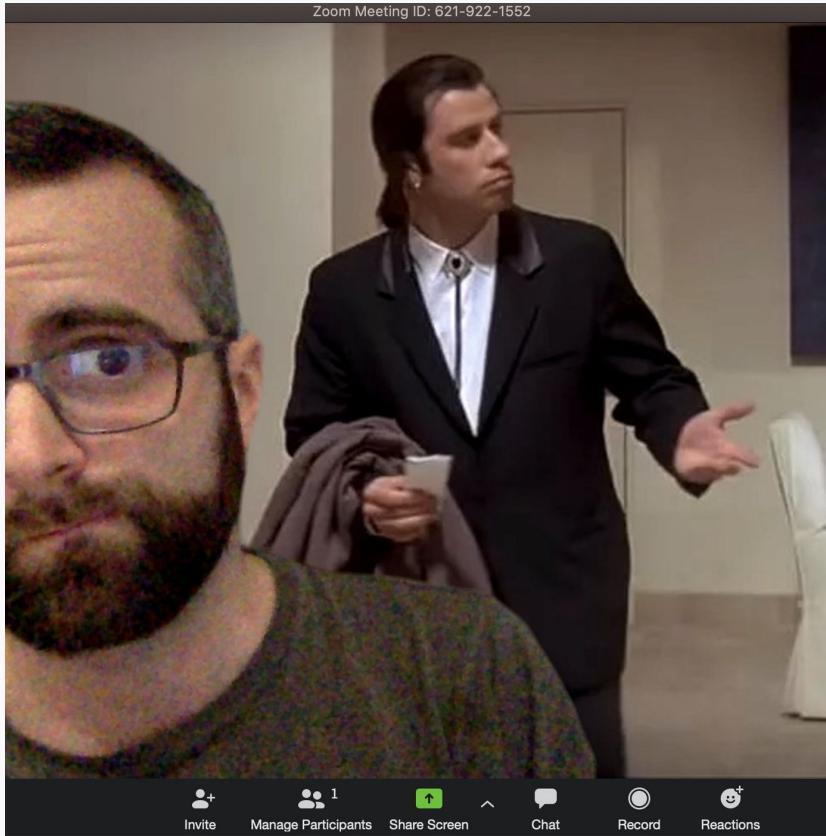
Applications

Applications: Autonomous Driving

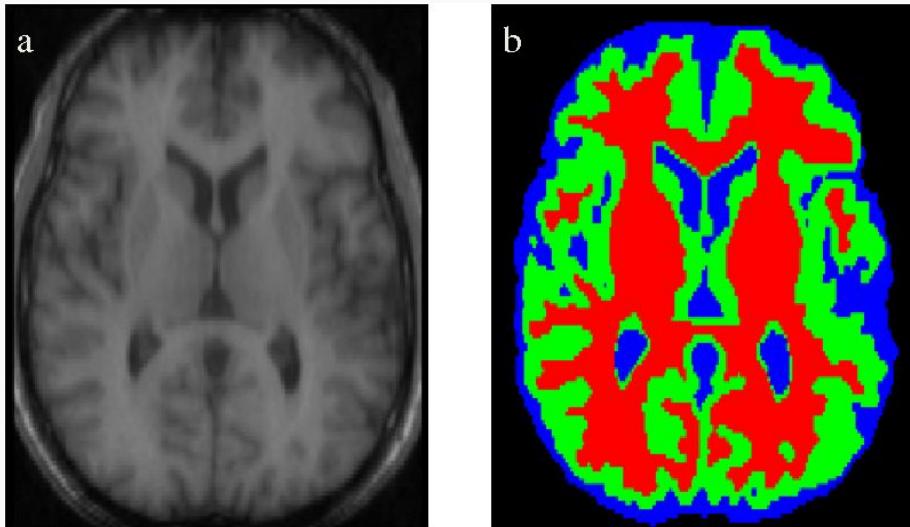


Semantic segmentation for autonomous vehicles

Background replacement



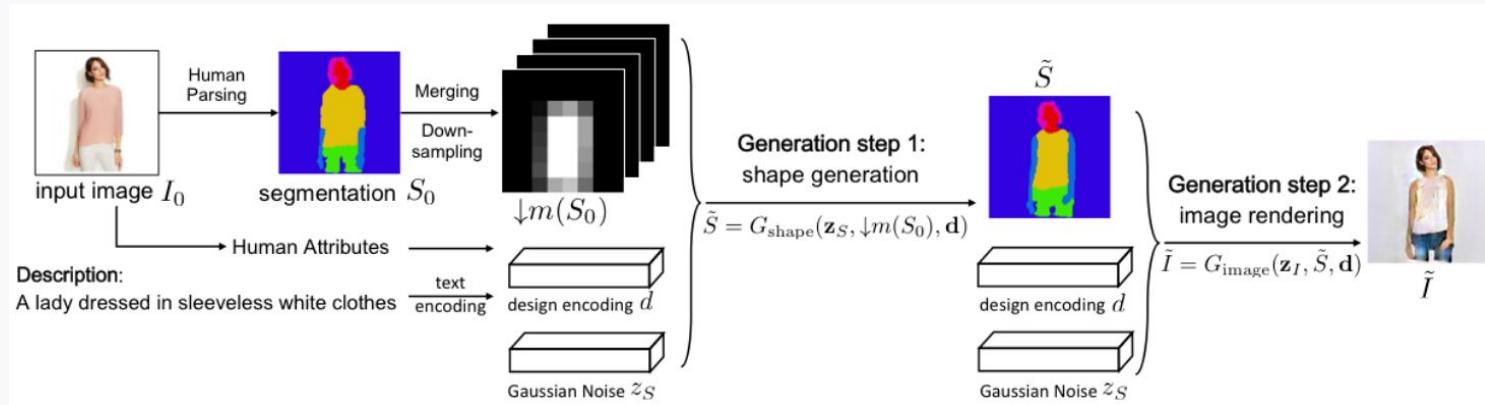
Applications: Medical Image Segmentation



Segmentation of medical scans. ([Source](#))

- 2D segmentation ([U-Net](#))
- 3D segmentation ([V-Net](#))

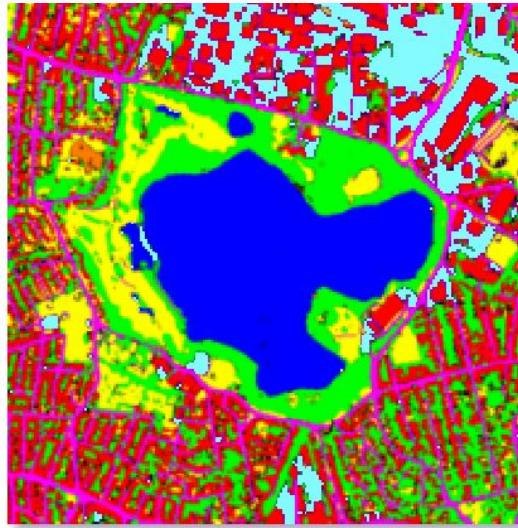
Applications: Fashion Industry



Semantic segmentation used as an intermediate step to redress a human based on text input.

- Be Your Own Prada: Fashion Synthesis with Structural Coherence ([2017, Shizhan Zhu et al](#))
- <https://github.com/switchablenorms/DeepFashion2>

Applications: Satellite Image Processing



Applications: Image editing



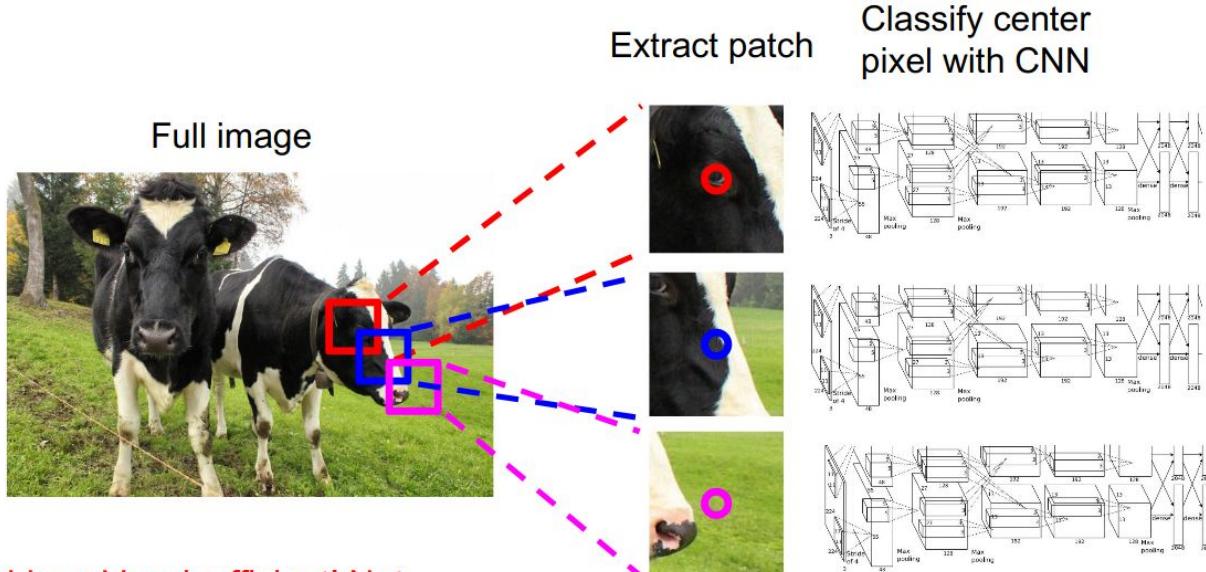
Segmentation used as an intermediate step
for GAN face generation



ОНЛАЙН ОБРАЗОВАНИЕ

Architectures

Semantic Segmentation Idea: Sliding Window

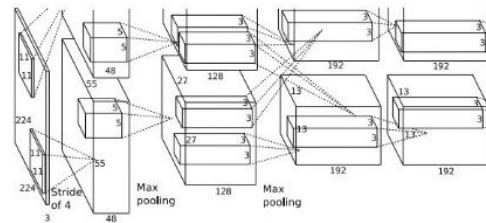


Problem: Very inefficient! Not reusing shared features between overlapping patches

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 Idea: Convolution

Full image

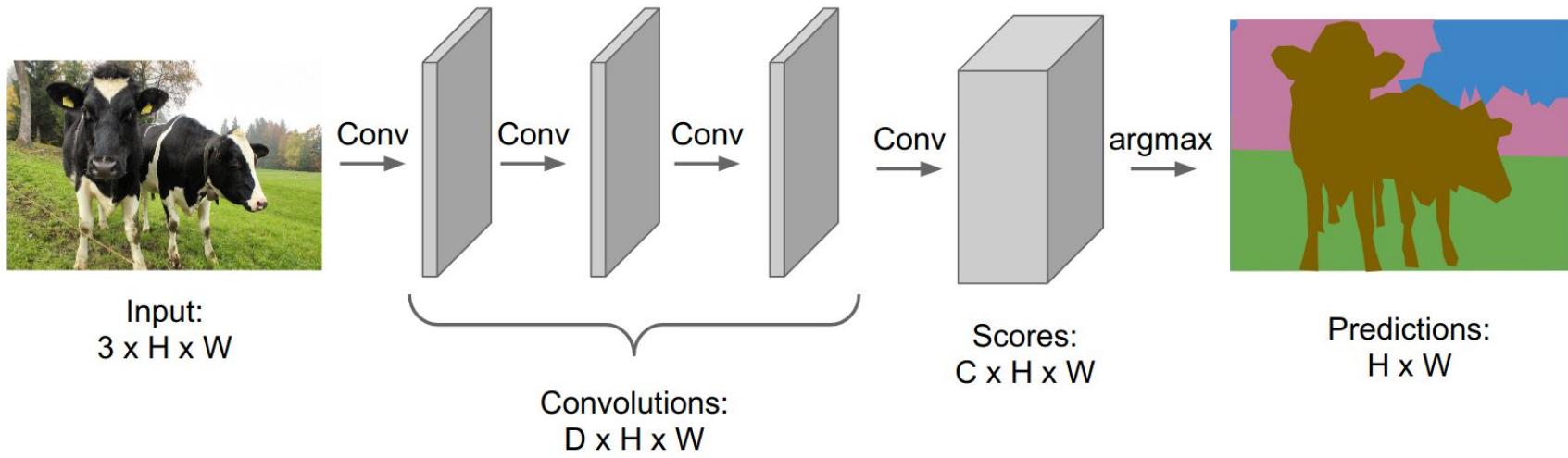


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

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.

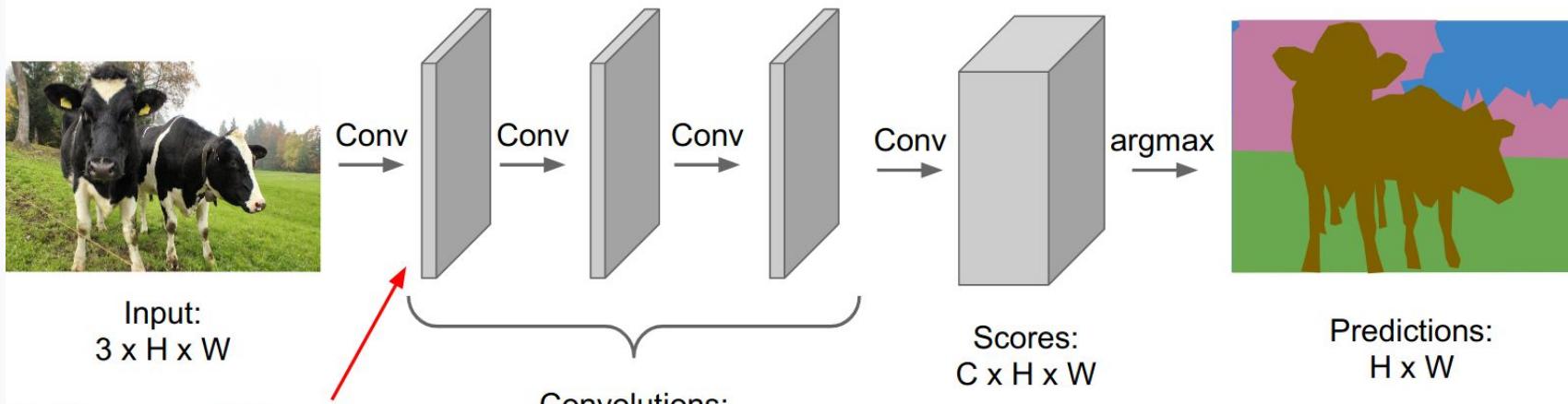
Semantic Segmentation Idea: Fully Convolutional

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



Semantic Segmentation Idea: Fully Convolutional

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



Problem: convolutions at original image resolution will be very expensive ...

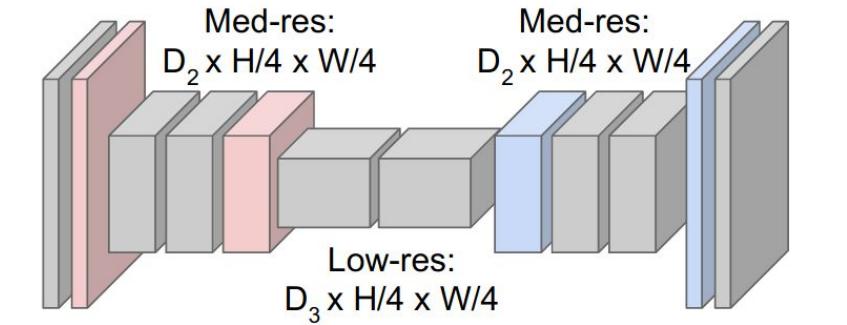
Semantic Segmentation Idea: Fully Convolutional

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:
???



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

In-Network upsampling: “Unpooling”

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

- Keras layer:
https://keras.io/api/layers/reshaping_layers/up_sampling2d/
- source: cs231 2020, [lection 12](#)

In-Network upsampling: “Max Unpooling”

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

5	6
7	8

Rest of the network

Max Unpooling

Use positions from pooling layer

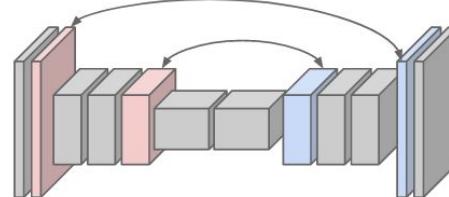
1	2
3	4

Input: 2 x 2

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

Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers



Learnable Upsampling: Transposed Convolutions

Transposed Convolutions

2x2 convolution, stride of 1 and a pad of 0

$$\begin{aligned}4 * 3 &= 12 \\2 * 1 &= 2 \\12 + 2 &= 14\end{aligned}$$

6	14	4
2	17	21
0	1	5

Output

2	4
0	1
3	1
1	5

6	2	
2	10	

	12	4
	4	20

1	0	0
0	0	
0	0	

3	1
1	5

Input

Conv
Kernel



Deconvolution and Checkerboard Artifacts
Use upsampling + convolution instead of deconvolution to avoid it

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

U-net

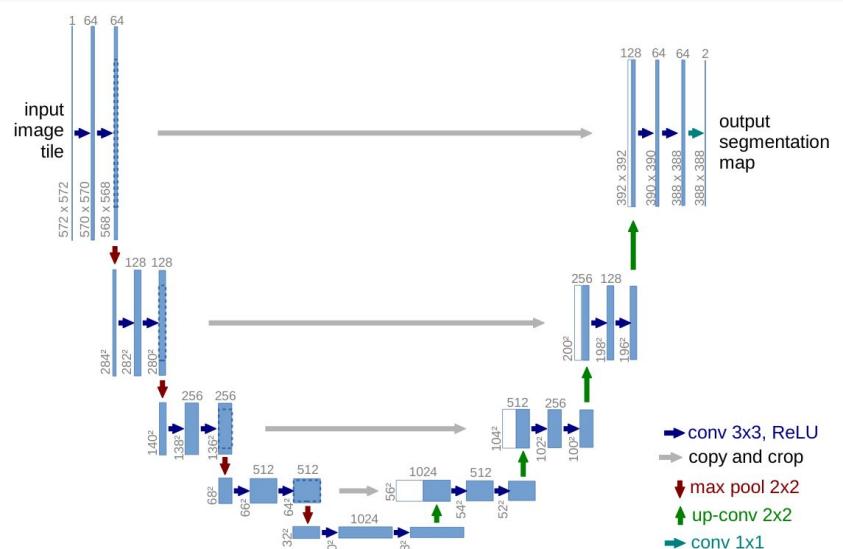


Fig. 1. U-net architecture (example for 32×32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

- Skip connections allows gradients to flow better and provides information from multiple scales of the image size.
- Example:
https://github.com/petrosgk/Kaggle-Carvana-Image-Masking-Challenge/blob/master/model/u_net.py

U-net: pretrained weights usage

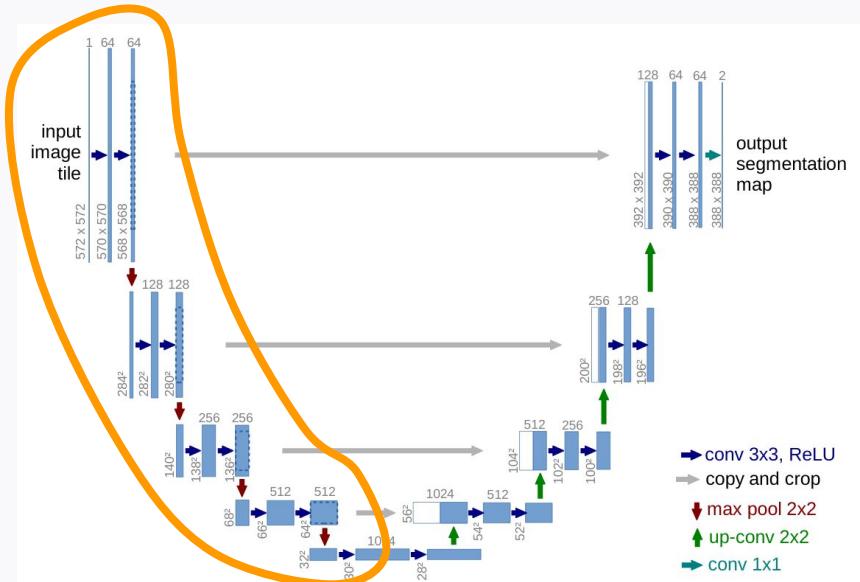


Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

- Weights from pretrained backbone (like ResNet18) could be used to initialize **encoder part of U-net**

V-net

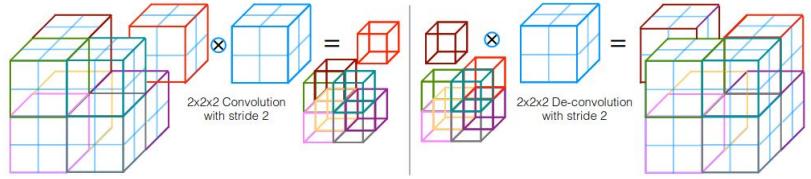
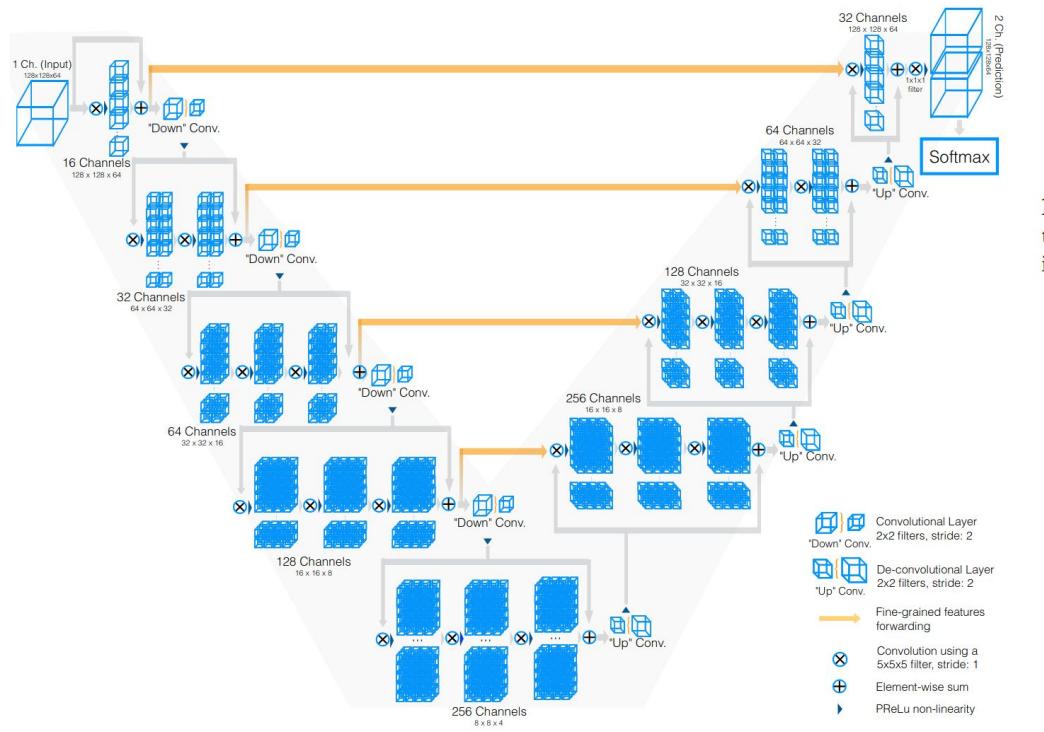


Fig. 3. Convolutions with appropriate stride can be used to reduce the size of the data. Conversely, de-convolutions increase the data size by projecting each input voxel to a bigger region through the kernel.

● 3D-convolutions

Tiramisu Model

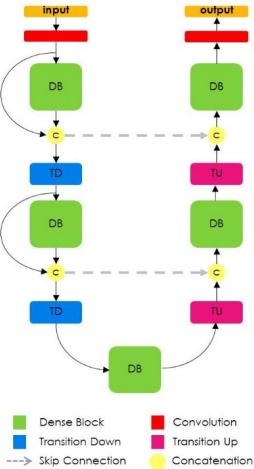


Figure 1. Diagram of our architecture for semantic segmentation. Our architecture is built from dense blocks. The diagram is composed of a downsampling path with 2 Transitions Down (TD) and an upsampling path with 2 Transitions Up (TU). A circle represents concatenation and arrows represent connectivity patterns in the network. Gray horizontal arrows represent skip connections, the feature maps from the downsampling path are concatenated with the corresponding feature maps in the upsampling path. Note that the connectivity pattern in the upsampling and the downsampling paths are different. In the downsampling path, the input to a dense block is concatenated with its output, leading to a linear growth of the number of feature maps, whereas in the upsampling path, it is not.

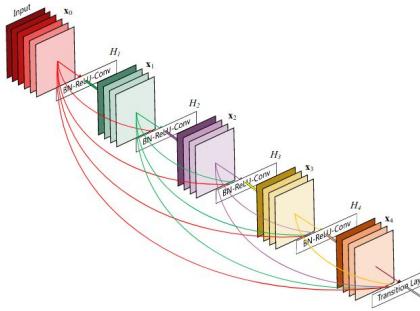


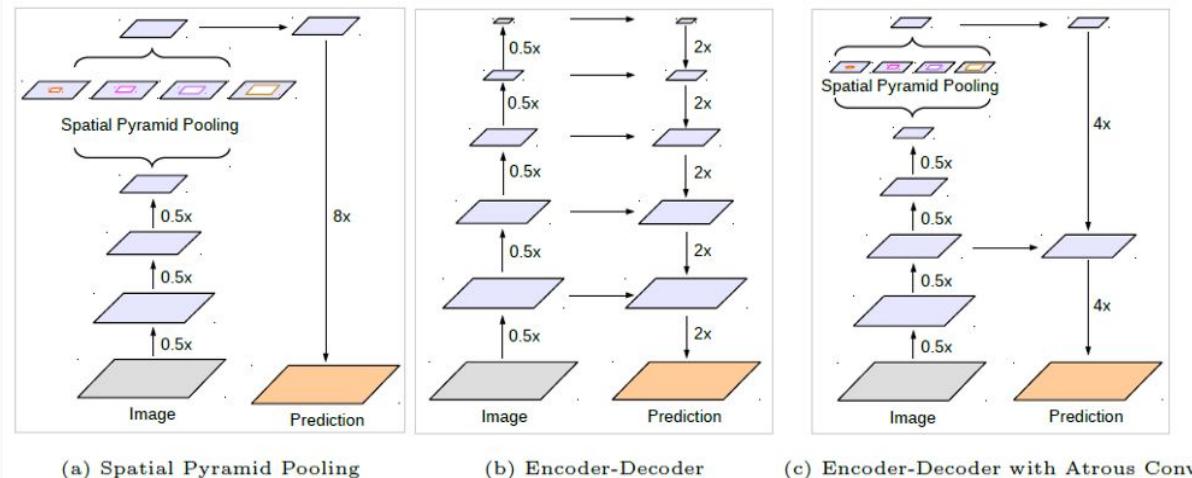
Figure 1: A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.

Dense block ([source](#))

- Uses dense blocks
- The resultant network is extremely parameter efficient and can better access features from older layers
- Due to the nature of the concatenation operations in several ML frameworks, it is not very memory efficient

DeepLab by Google

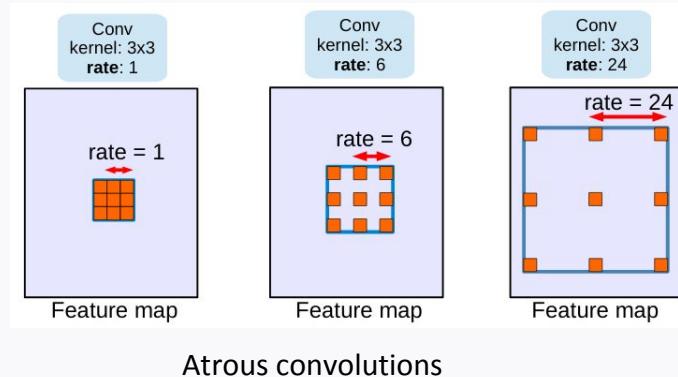
- [DeepLabv1](#) (2015 ICLR), [DeepLabv2](#) (2018 TPAMI), and [DeepLabv3](#) (arXiv)



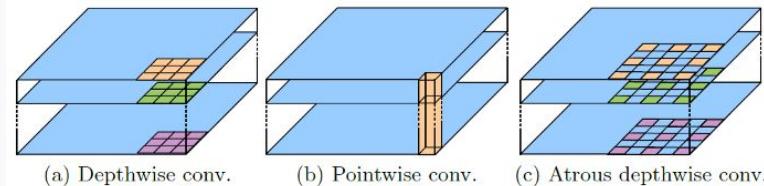
- (a): With Atrous Spatial Pyramid Pooling (ASPP), able to encode multi-scale contextual information.
(b): With Encoder-Decoder Architecture, the location/spatial information is recovered.
(c): DeepLabv3+ makes use of (a) and (b) and Atrous Separable Convolution

Deeplab repo: <https://github.com/tensorflow/models/tree/master/research/deeplab>

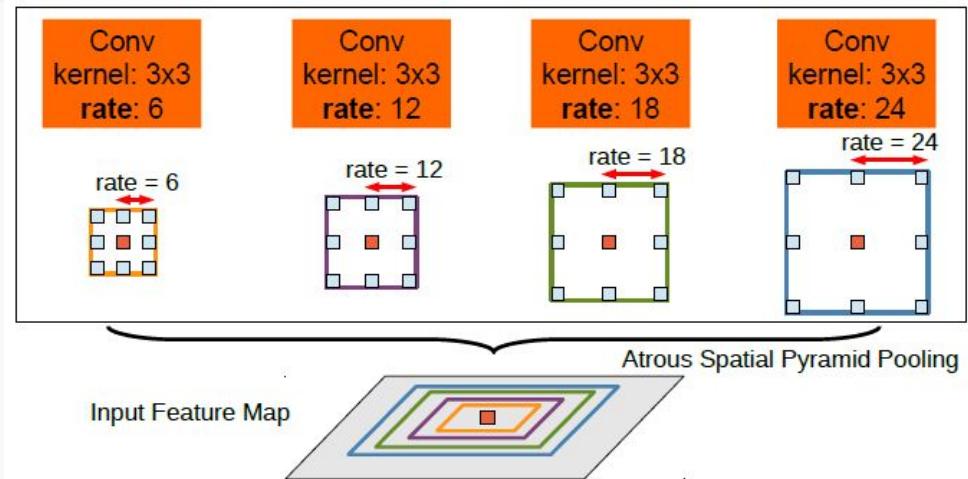
DeepLabV3+: Atrous (dilated) Convolutions



Atrous convolutions



Depthwise Separable Convolution Using
Atrous Convolution



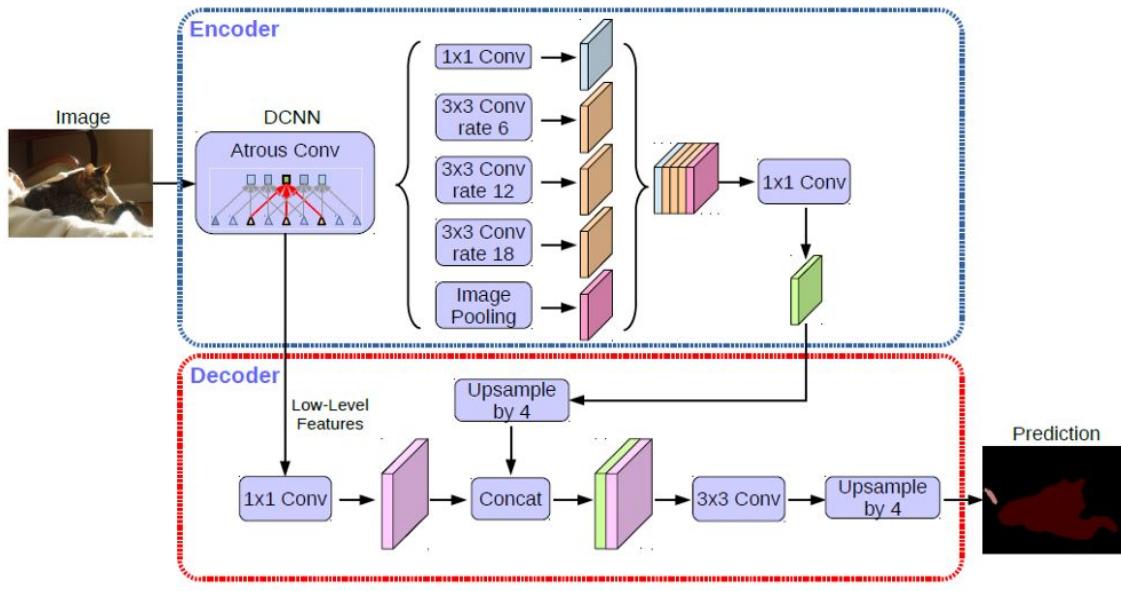
Atrous Spatial Pyramid Pooling (ASPP)

(rate - 1) -- это расстояние между 2-мя пикселями, на которые применяется фильтр.

Иллюстрация не точная по соотношению площадей и указанного rate.

Более подробно смотреть в [этой заметке](#).

DeepLabV3+: architecture



Encoder-Decoder architecture

- DeepLabv3 as Encoder
- Atrous Separable Convolution instead of Atrous Convolution
- Modified Aligned Xception as Backbone

More information in [the article](#)

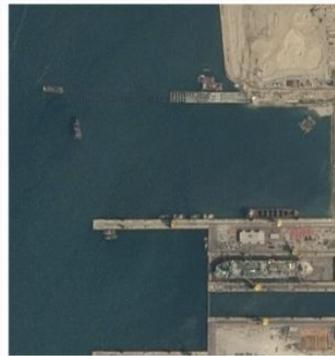


ОНЛАЙН ОБРАЗОВАНИЕ

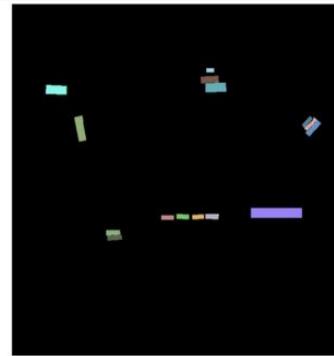
Metrics

Pixel Accuracy

Percent of pixels in your image that are classified correctly



Input image



ground truth mask



your prediction
that gives 95%
accuracy

Class imbalance issue

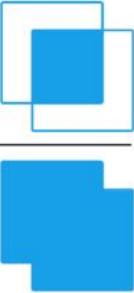
Intersection-Over-Union (IoU, Jaccard Index)

Jaccard Similarity Index is the most intuitive ratio between the intersection and union

- ranges from 0–1 (0–100%)
- mean IoU for multi-class segmentation

How to calculate?

Ships

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


$$\text{Area of Overlap} = 0$$

$$\text{Area of Union} = (5+0)-0 = 5$$

$$\text{Area of Overlap}/\text{Area of Union} = 0\%$$

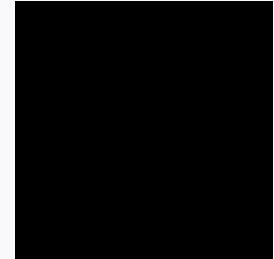
Background

$$\text{Area of Overlap} = 95,$$

$$\text{Area of Union} = (95+100)-95 = 100$$

$$\text{Area of Overlap}/\text{Area of Union} = 95\%$$

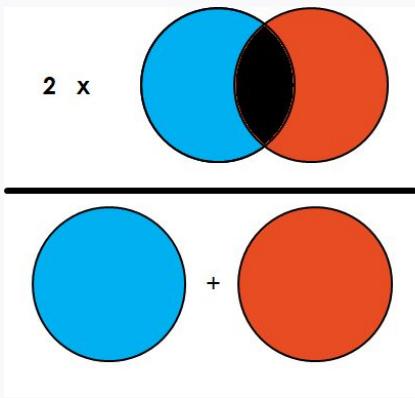
$$\text{Mean IoU} = (\text{Ships} + \text{Background})/2 = \\ (0\%+95\%)/2 = \textcolor{red}{47.5\%}$$



$$IoU = \frac{TP}{(TP + FP + FN)}$$

Dice Similarity Coefficient

$2 * \text{the Area of Overlap} / (\text{total number of pixels in both images})$



How to calculate?

Let total Number of Pixels for both images combined = 200

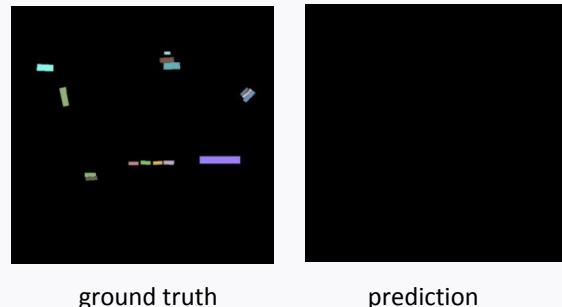
Ships: Area of Overlap = 0

$$(2 * \text{Area of Overlap}) / (\text{total pixels combined}) = 0 / (5+0) = 0$$

Background: Area of Overlap = 95

$$(2 * \text{Area of Overlap}) / (\text{total pixels combined}) = 95 * 2 / (100+95) \approx 0.97$$

$$\text{Dice} = (\text{Ships} + \text{Background}) / 2 = \\ (0\% + 97\%) / 2 = \mathbf{48.5\%}$$



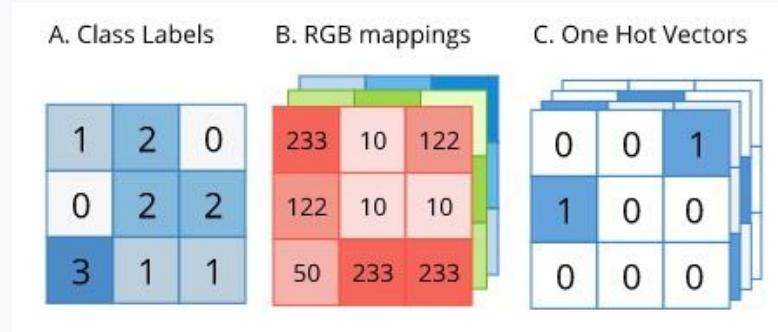
$$\text{Dice} = \frac{2 \times TP}{(TP + FP) + (TP + FN)}$$



ОНЛАЙН ОБРАЗОВАНИЕ

Loss functions

Pixel-wise Softmax with Cross Entropy



- Labels in a one-hot form
- Softmax must be applied pixel-wise on the predicted output before applying cross entropy, as each pixel can belong to any one of our target classes

Focal Loss

Good choice in cases
with extreme class
imbalance

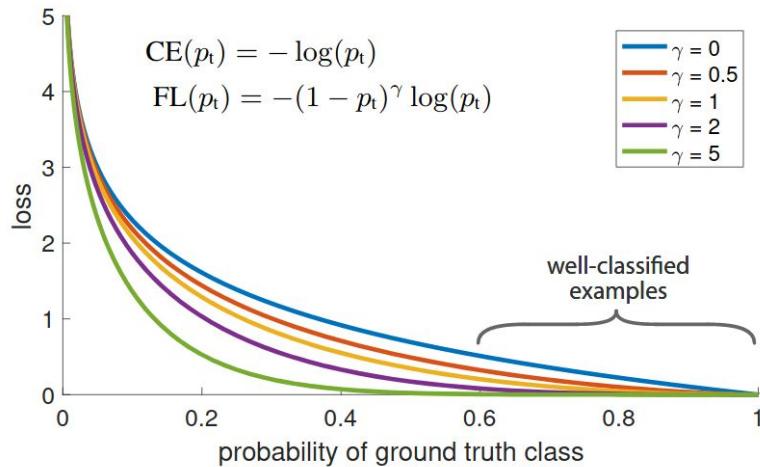


Figure 1. We propose a novel loss we term the *Focal Loss* that adds a factor $(1 - p_t)^\gamma$ to the standard cross entropy criterion. Setting $\gamma > 0$ reduces the relative loss for well-classified examples ($p_t > .5$), putting more focus on hard, misclassified examples. As our experiments will demonstrate, the proposed focal loss enables training highly accurate dense object detectors in the presence of vast numbers of easy background examples.

Dice Loss

Good choice in cases
with extreme class
imbalance

$$D = \frac{2 \sum_i^N p_i g_i}{\sum_i^N p_i^2 + \sum_i^N g_i^2}$$

Dice coefficient between **p** (predicted)
and **g** (ground truth) *binary* volumes
ranging between 0 and 1

$$\frac{\partial D}{\partial p_j} = 2 \left[\frac{g_j \left(\sum_i^N p_i^2 + \sum_i^N g_i^2 \right) - 2p_j \left(\sum_i^N p_i g_i \right)}{\left(\sum_i^N p_i^2 + \sum_i^N g_i^2 \right)^2} \right]$$

Derivative of dice coefficient computed with
respect to the j-th voxel of the prediction

Image Segmentation Datasets

- **Common Objects in COnText—Coco Dataset**

COCO is a large-scale object detection, segmentation, and captioning dataset. The [dataset](#) contains 91 classes. It has 250,000 people with key points. Its download size is 37.57 GiB. It contains 80 object categories. It is available under the [Apache 2.0 License](#) and can be downloaded from [here](#).

- **PASCAL Visual Object Classes (PASCAL VOC)**

PASCAL has 9963 images with 20 different classes. The training/validation set is a 2GB tar file. The dataset can be downloaded from the [official website](#).

- **The Cityscapes Dataset**

This dataset contains images of city scenes. It can be used to evaluate the performance of vision algorithms in urban scenarios. The dataset can be downloaded from [here](#).

- **The Cambridge-driving Labeled Video Database—CamVid**

This is a motion-based segmentation and recognition dataset. It contains 32 semantic classes. This [link](#) contains further explanations and download links to the dataset.

- [The Berkeley Segmentation Dataset and Benchmark](#)
- [ADE20K](#)
- [DeepFashion2](#)

Frameworks

- For instance segmentation [detectron2](#)
- For semantic segmentation:
 - [DeepLab](#)
 - [Catalyst.Segmentation](#) (Unet)
 - [SMP](#)

Tutorials

- [Detectron2 -- instance segmentation with pretrained mask-RCNN and custom balloon dataset](#)
- [DeepLab on ADE20k](#)
- [DeepLab train on your own dataset](#)
 - Background segmentation example -- collecting data, making labels, preparing tfdata and training
 - But it's possible to solve background segmentation just with inference of pretrained model like [this tutorial](#) shows.
- [Advanced segmentation tutorial with Catalyst](#) [*Highly recommended*]



Pretrained models from SMP

- PyTorch (you can use them in Catalyst as well):
 - https://github.com/qubvel/segmentation_models.pytorch
- Keras/Tensorflow
 - https://github.com/qubvel/segmentation_models



Tools for annotations

- [LabelMe](#): One of the most known tools. The UI was a bit too slow, though, especially when zooming in on large images.
- [RectLabel](#): Simple and easy to work with. Mac only.
- [LabelBox](#): Pretty good for larger labeling projects and has options for different types of labeling tasks.
- [VGG Image Annotator \(VIA\)](#): Fast, light, and really well designed. This is the one I ended up using.
- [COCO UI](#): The tool used to annotate the COCO dataset.

