

# Introduction to Semantic Segmentation

Summer school 2022

Dmitry Sidnev

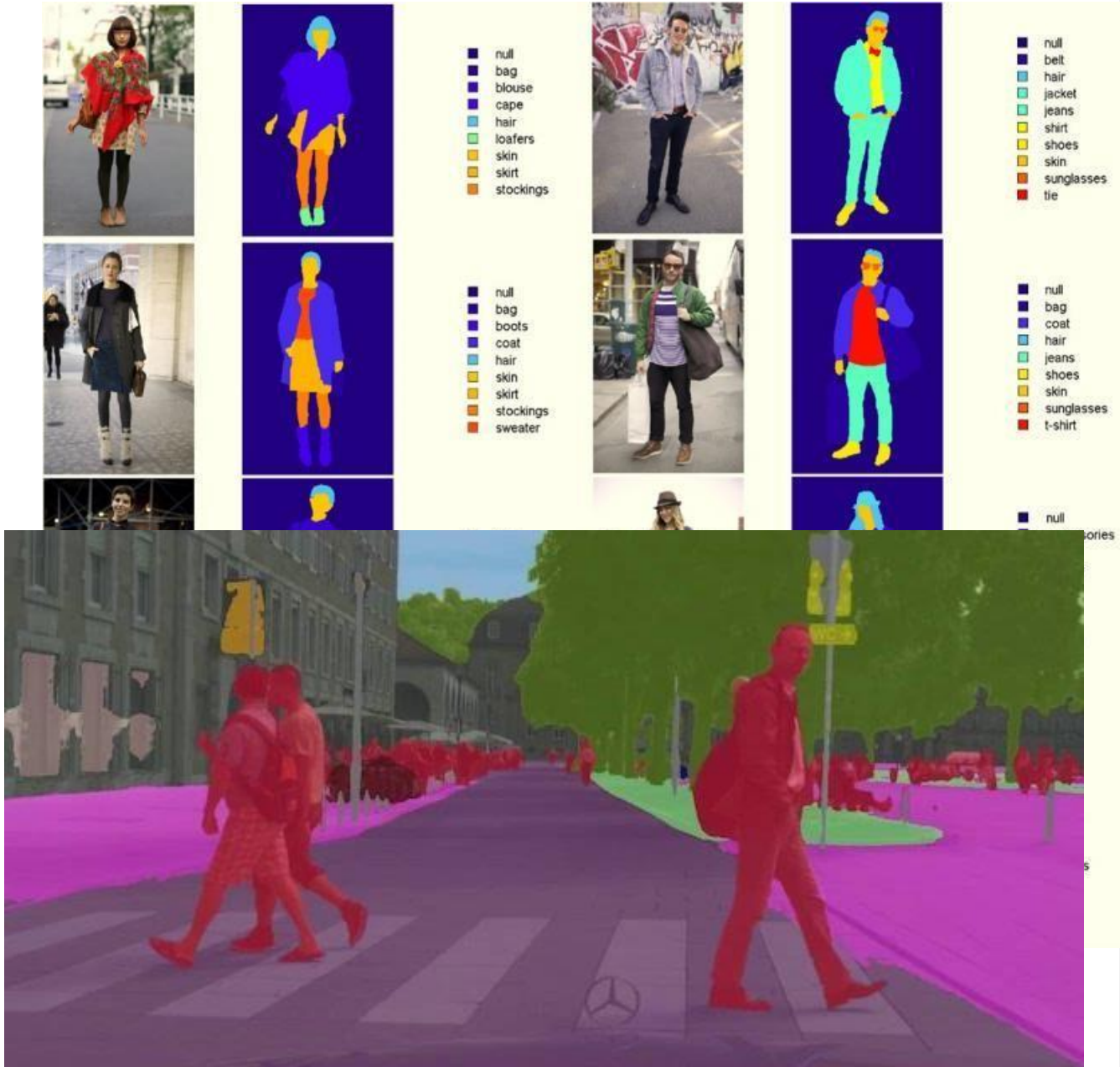
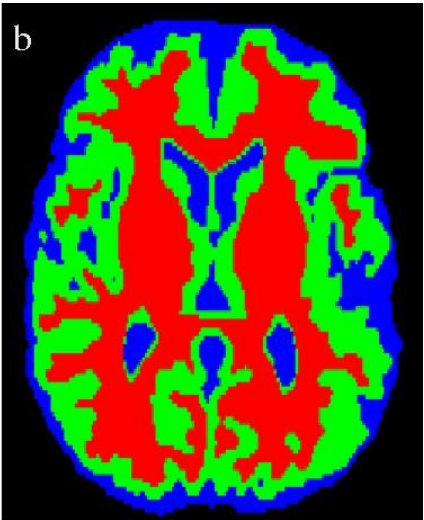
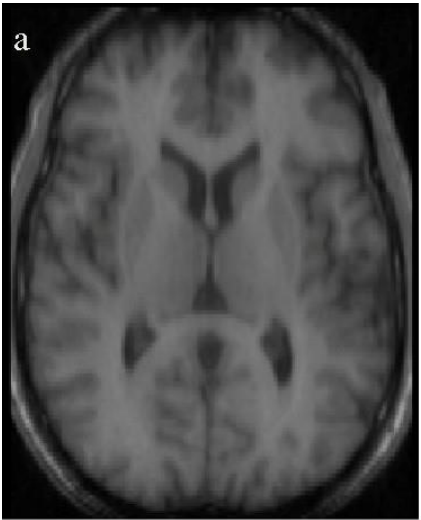
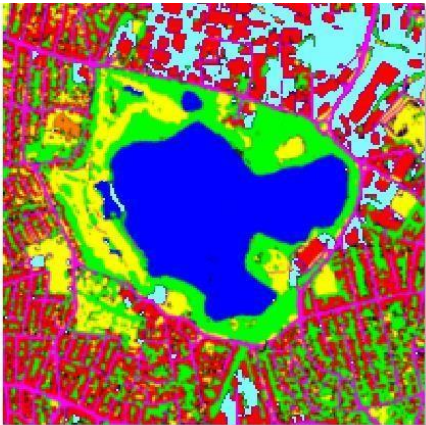
# Agenda

- Problem formulation
- Datasets
- Evaluation metrics
- Architectures
- Loss functions
- Comparison

# Computer vision problems

- Aerospace photos processing
- Medical scan segmentation
- Autonomous driving

# Computer vision problems



# Problem formulation

- Input image:

$$I = \{I_{ij}\}_{0 \leq i < w, 0 \leq j < h}, I_{ij} \in R^c$$

- Set of classes:

$$\mathcal{C} = \{0, 1, \dots, N - 1\}$$

- Mask:

$$M = \{M_{ij}\}_{0 \leq i < w, 0 \leq j < h}, M_{ij} \in \mathcal{C}$$

- Segmentation function:

$$\varphi(R^c) \rightarrow \mathcal{C}$$



# Datasets

Dataset	Train subset	Test subset	Classes
Common objects			
PASCAL VOC 2012 <a href="http://host.robots.ox.ac.uk/pascal/VOC/voc2012">[http://host.robots.ox.ac.uk/pascal/VOC/voc2012]</a>	9 963	1 447	20
ADE20K <a href="http://groups.csail.mit.edu/vision/datasets/ADE20K">[http://groups.csail.mit.edu/vision/datasets/ADE20K]</a>	20 210	2 000	150
MS COCO'15 <a href="http://mscoco.org">[http://mscoco.org]</a>	80 000	40 000	80

# Datasets

Dataset	Train subset	Test subset	Classes
City, streets, cars			
CamVid <a href="http://mi.eng.cam.ac.uk/research/projects/VideoRec/CamVid">[http://mi.eng.cam.ac.uk/research/projects/VideoRec/CamVid]</a>	468	233	11
Cityscapes <a href="https://www.cityscapes-dataset.com">[https://www.cityscapes-dataset.com]</a>	2 975	500	19
KITTI <a href="http://www.cvlibs.net/datasets/kitti">[http://www.cvlibs.net/datasets/kitti]</a>	200	200	4
Interiors			
Sun-RGBD <a href="http://rgbd.cs.princeton.edu">[http://rgbd.cs.princeton.edu]</a>	10 355	2 860	37
NYUDv2 <a href="http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html">[http://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html]</a>	795	645	40



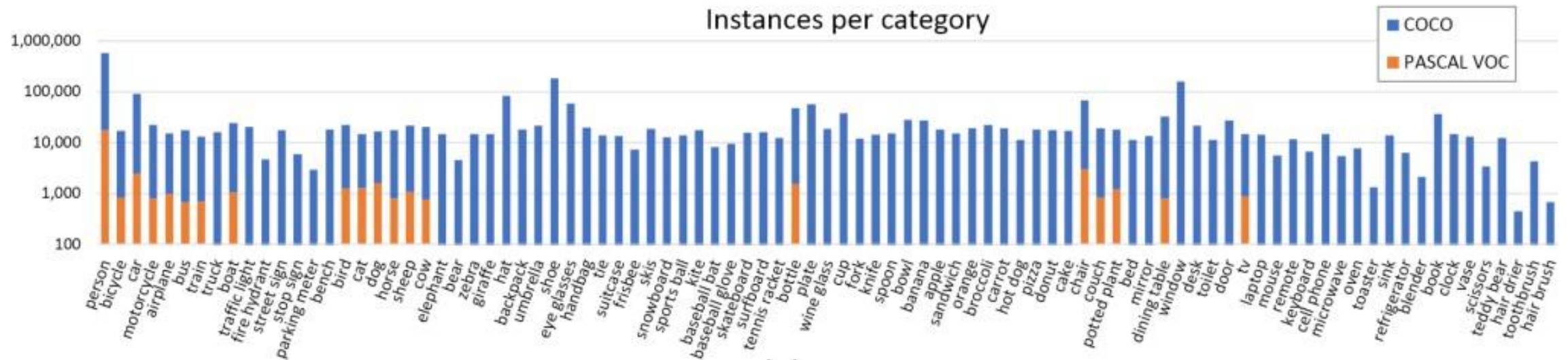
# Datasets: Pascal VOC2012

- Airplane
- Bicycle
- Bird
- Boat
- Bottle
- Bus
- Car
- Cat
- Chair
- Cow
- dining table
- Dog
- Horse
- Motorbike
- Person
- potted plant
- Sheep
- Sofa
- Train
- tv/monitor





# Datasets: MS COCO



Lin T.Y., et al. Microsoft COCO: Common objects in context // Lecture Notes in Computer Science. – Vol. 8693. – 2014. – P. 740-755. [<https://arxiv.org/pdf/1405.0312>].

# Datasets: Citiscapes

- 50 cities
- 5 000 fine annotations
- 20 000 coarse annotations
- 30 classes, 8 groups
- Diversity: daytime, season, weather conditions



The Cityscapes Dataset Homepage [<https://www.cityscapes-dataset.com/examples>].

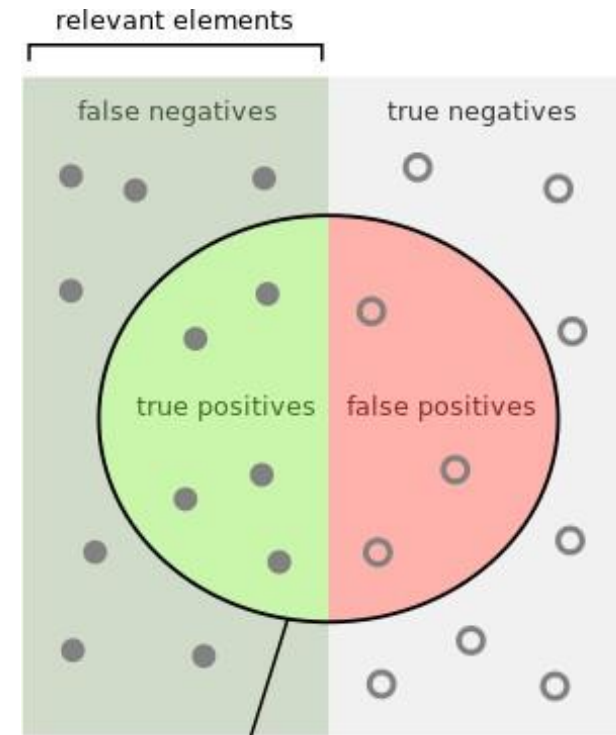
# Evaluation metrics

- Pixel accuracy
- Mean pixel accuracy over classes
- Jaccard index (IoU)
- Dice index

# Pixel accuracy

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

	Prediction							
	True	False						
Ground Truth	<table><tr><td>True</td><td>TP</td><td>FN</td></tr><tr><td>False</td><td>FP</td><td>TN</td></tr></table>	True	TP	FN	False	FP	TN	
True	TP	FN						
False	FP	TN						



How many selected items are relevant?

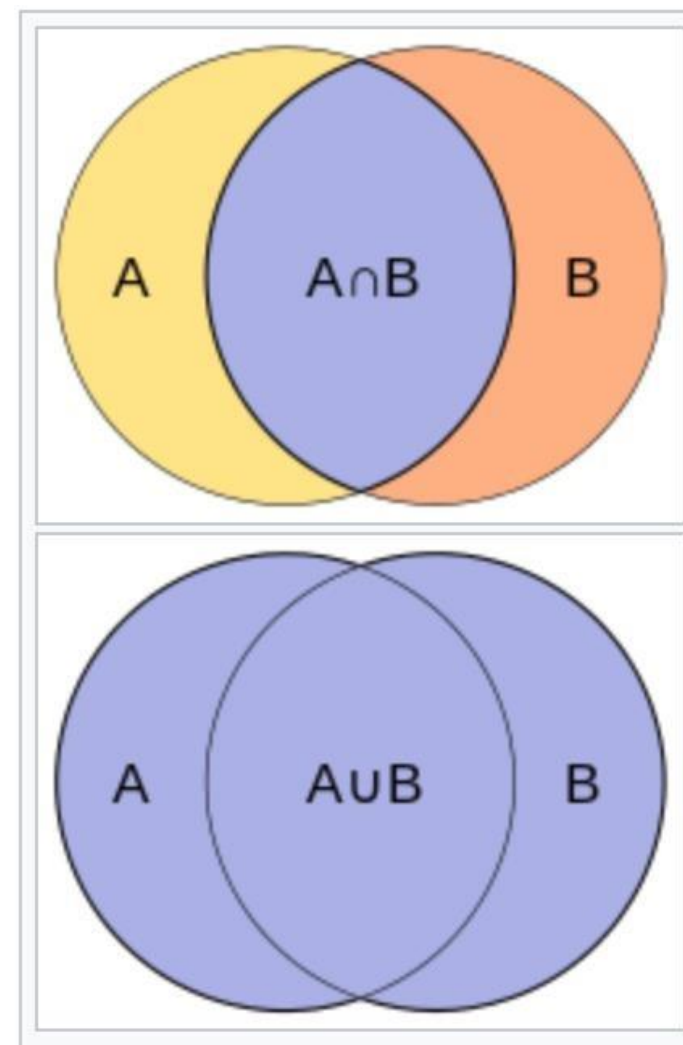
$$Precision = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are selected?

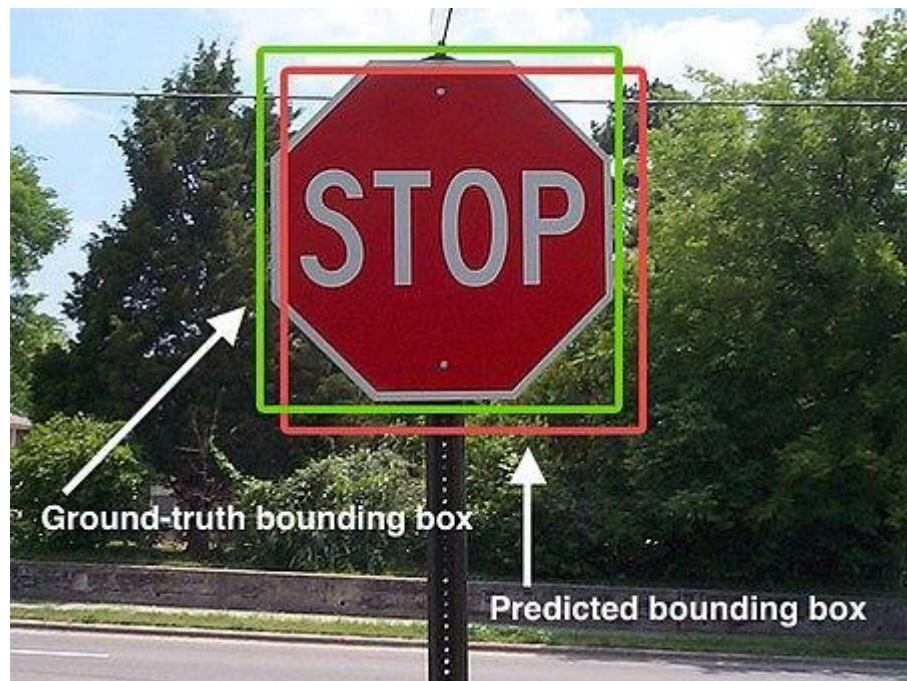
$$Recall = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$


# Jaccard index (IoU) and Dice (F1) index

- $IoU(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{TP}{TP + FN + FP}$
- $F1(A, B) = 2 \frac{|A \cap B|}{|A| + |B|} = \frac{2TP}{2TP + FN + FP}$

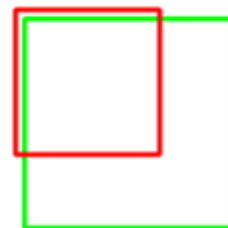


# IoU explanation



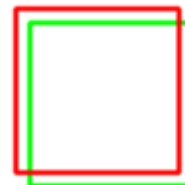
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


IoU: 0.4034



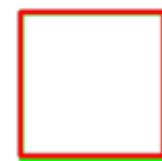
Poor

IoU: 0.7330



Good

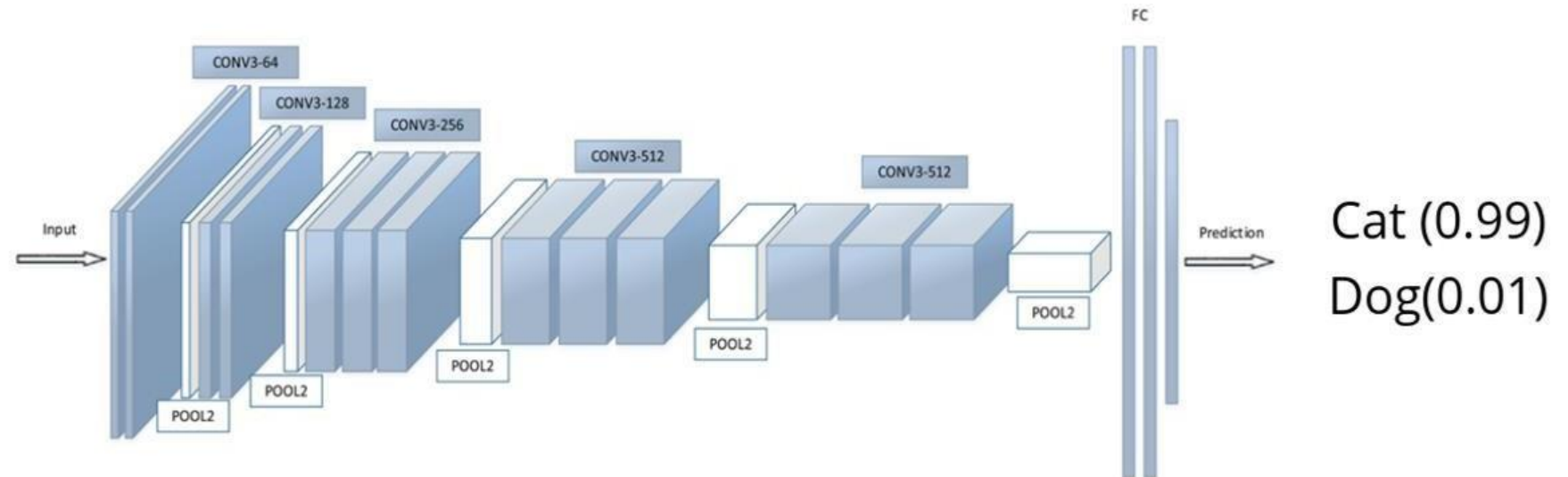
IoU: 0.9264



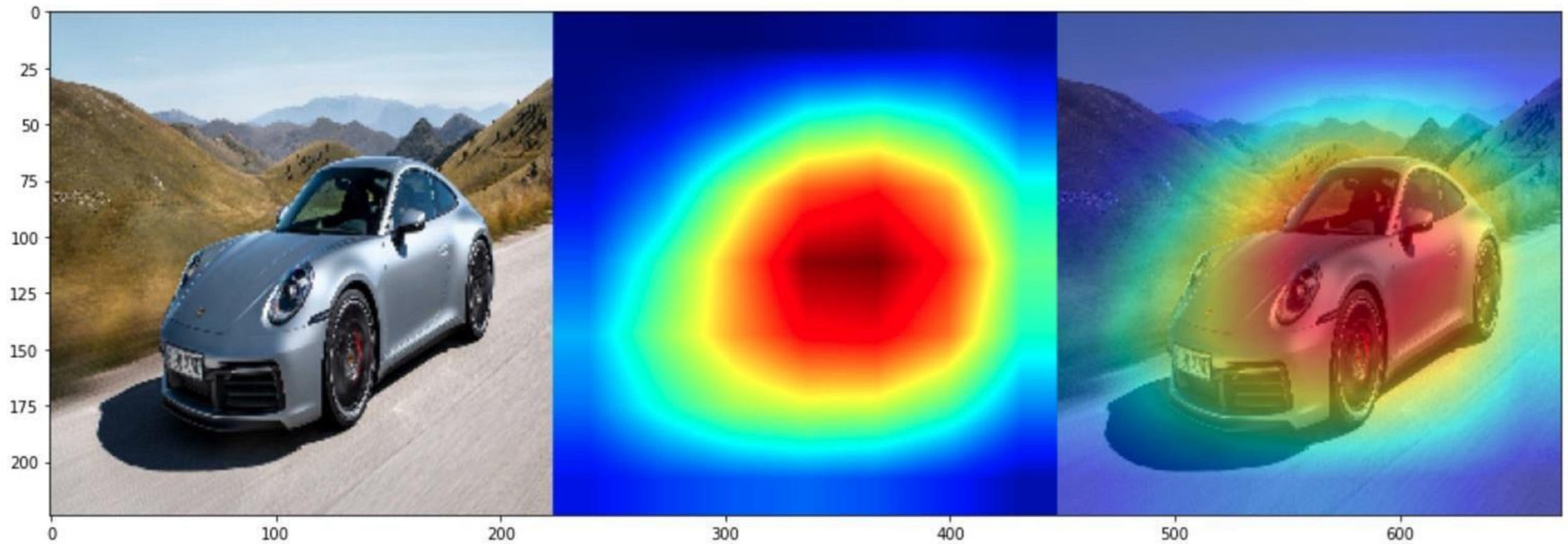
Excellent



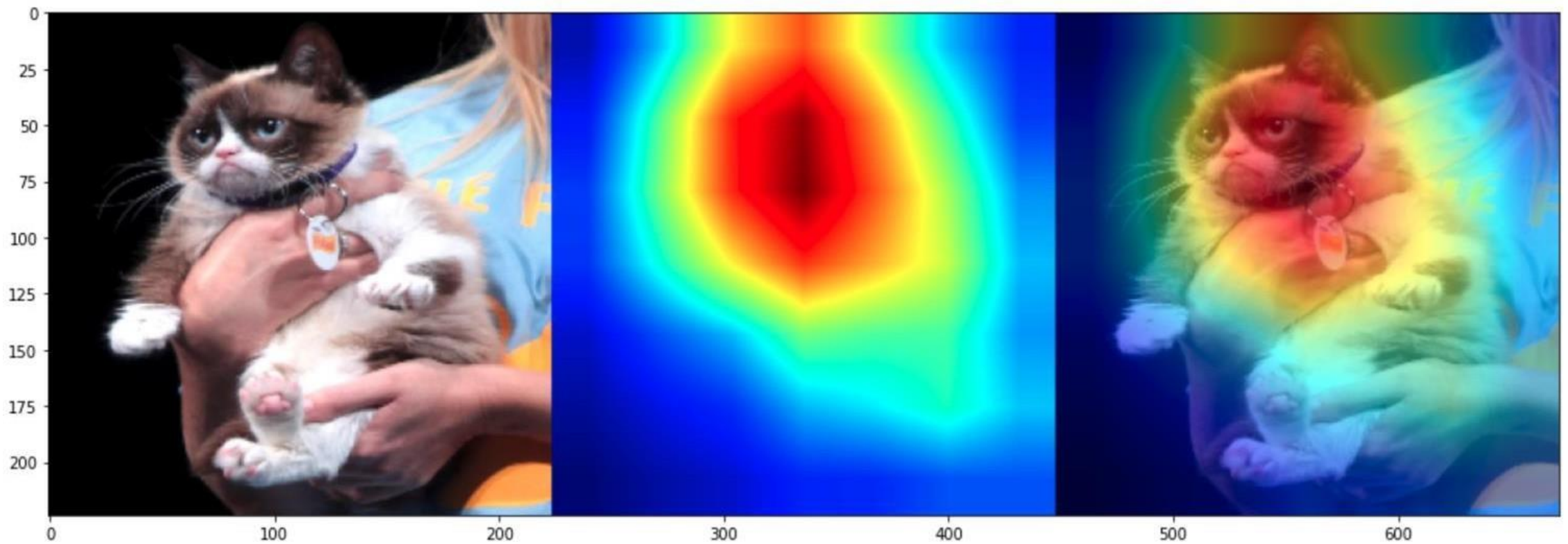
# Architectures: CNN



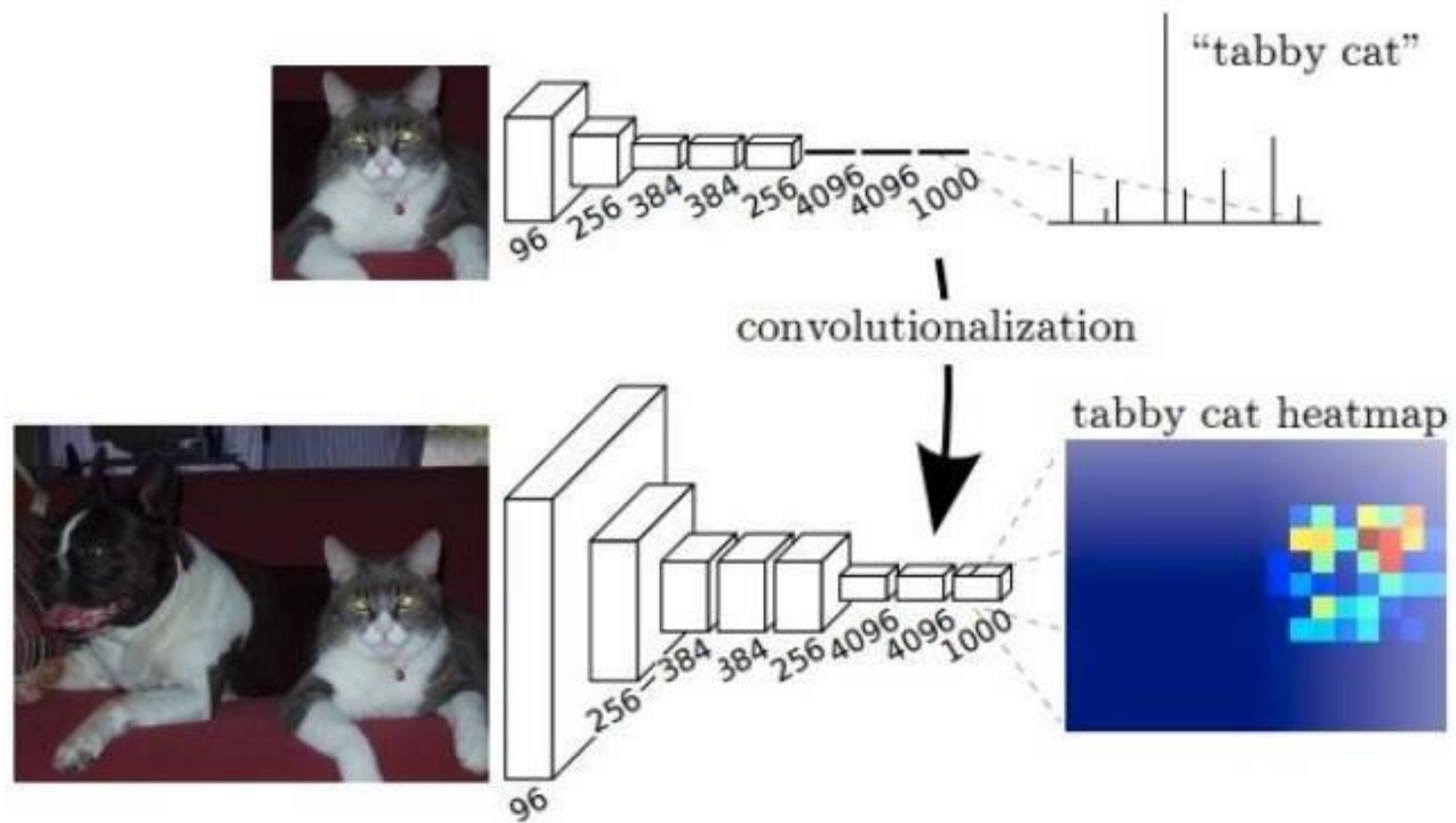
# Architectures: CNN



# Architectures: CNN

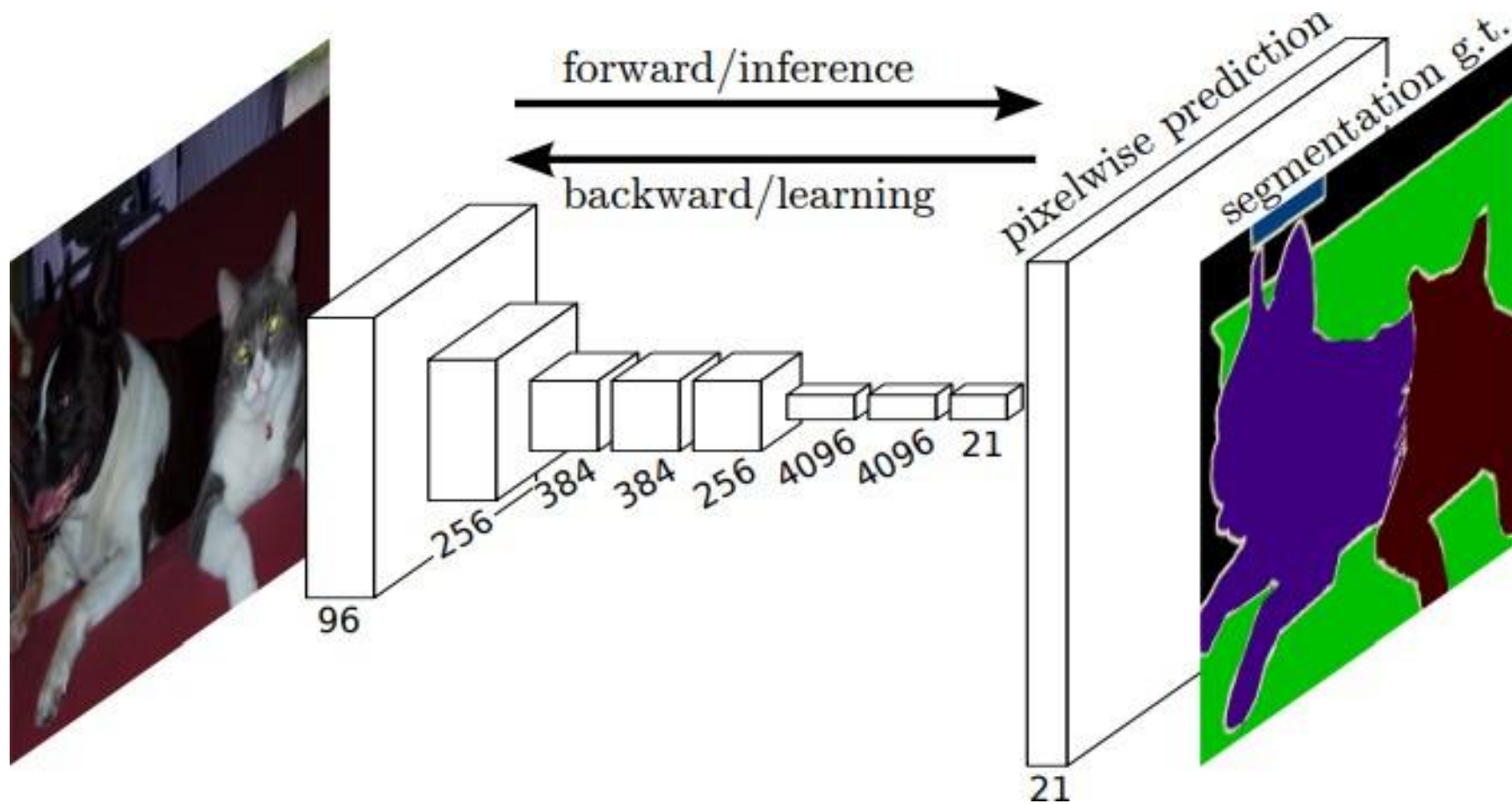


# Architecture: FCN

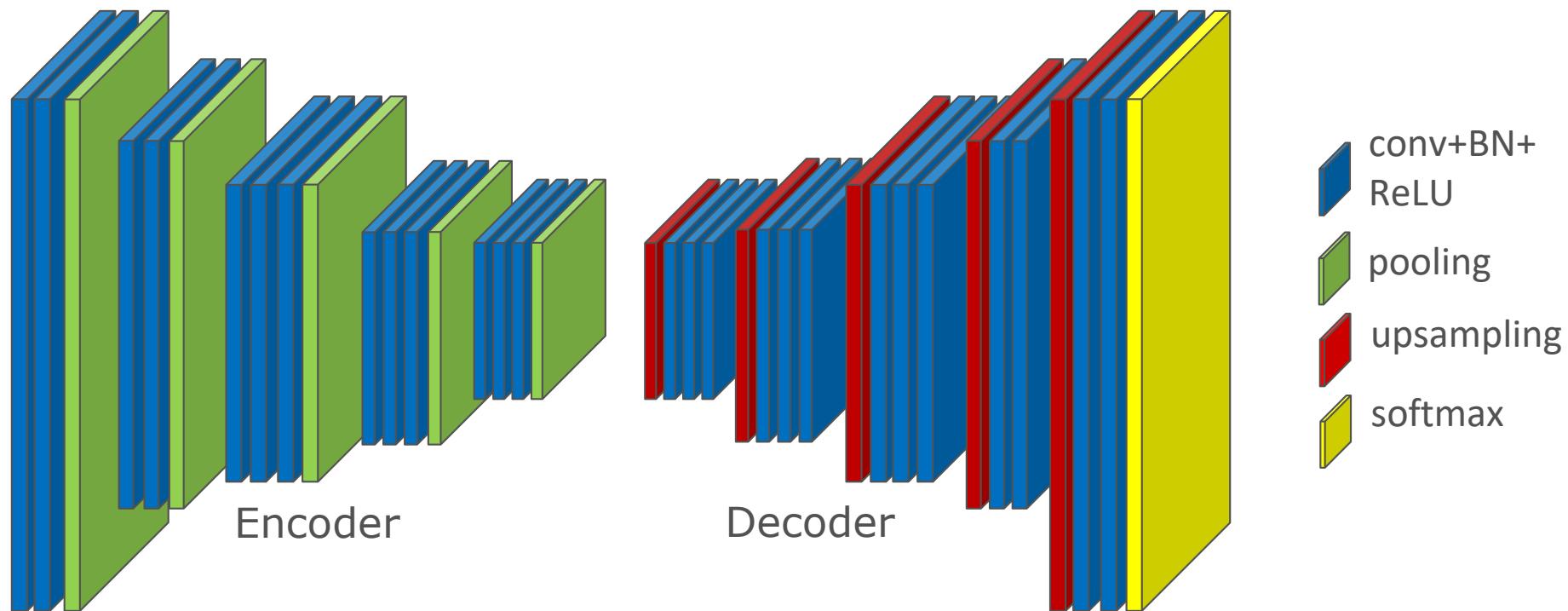




# Architecture: FCN

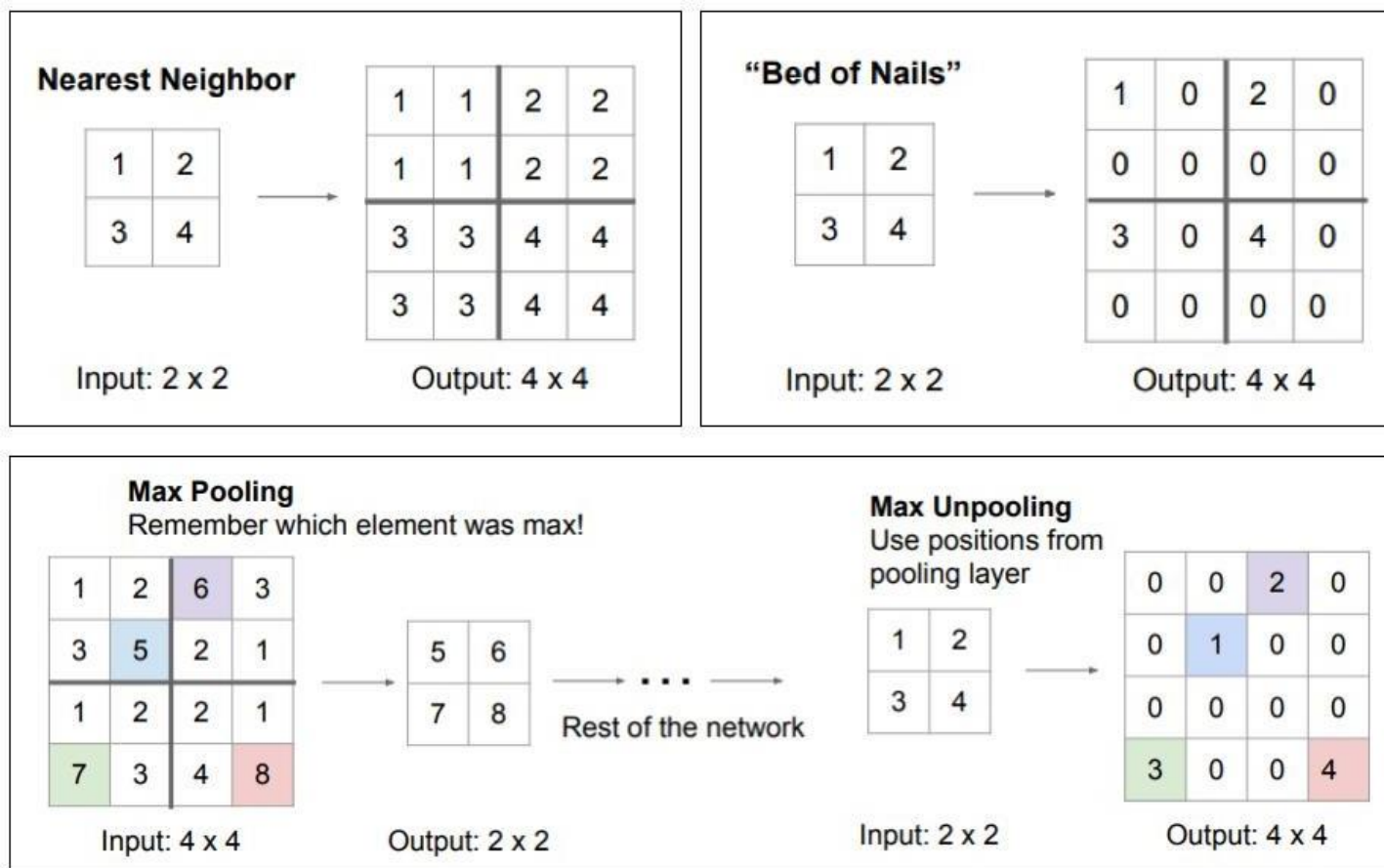


# Architectures: SegNet

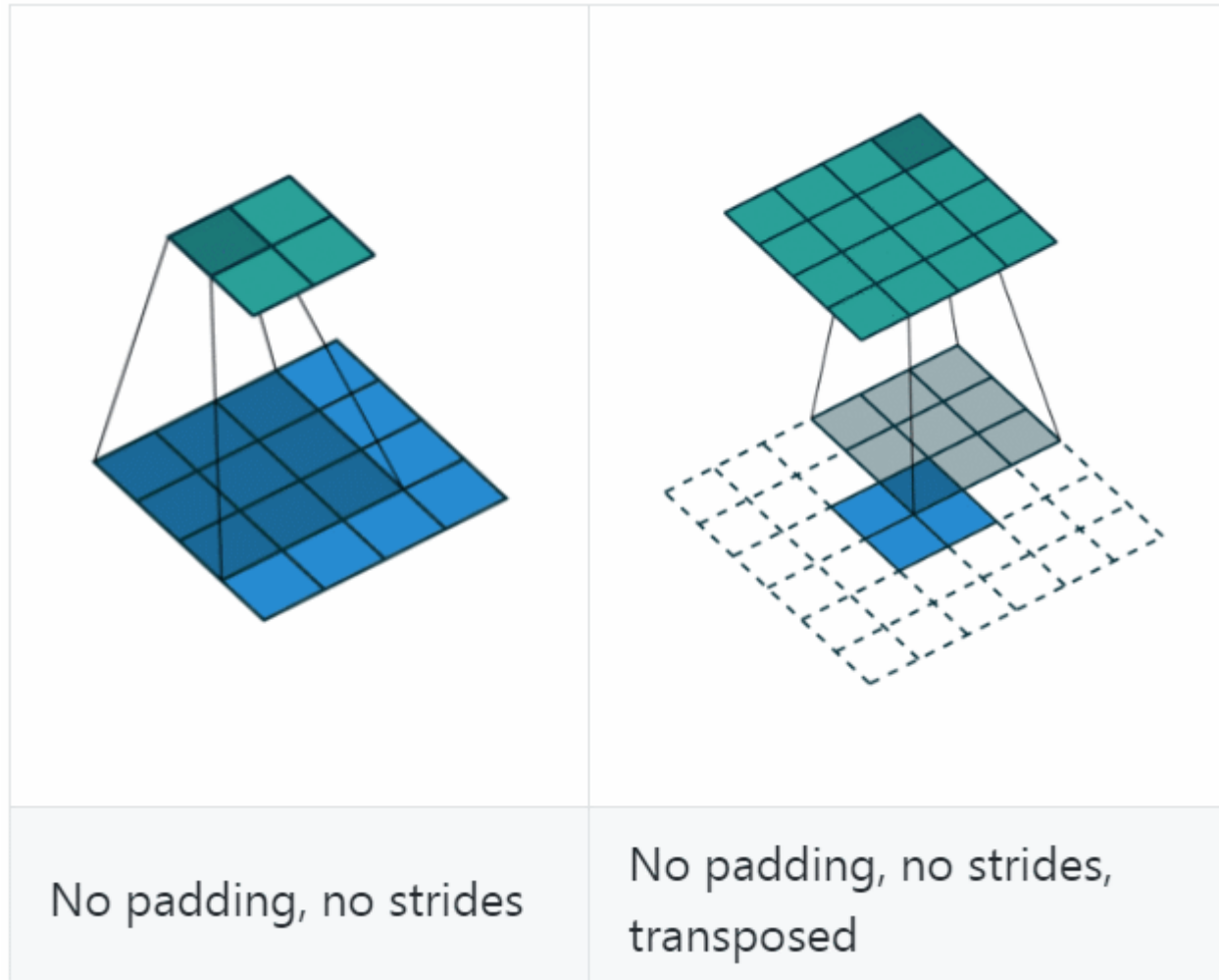




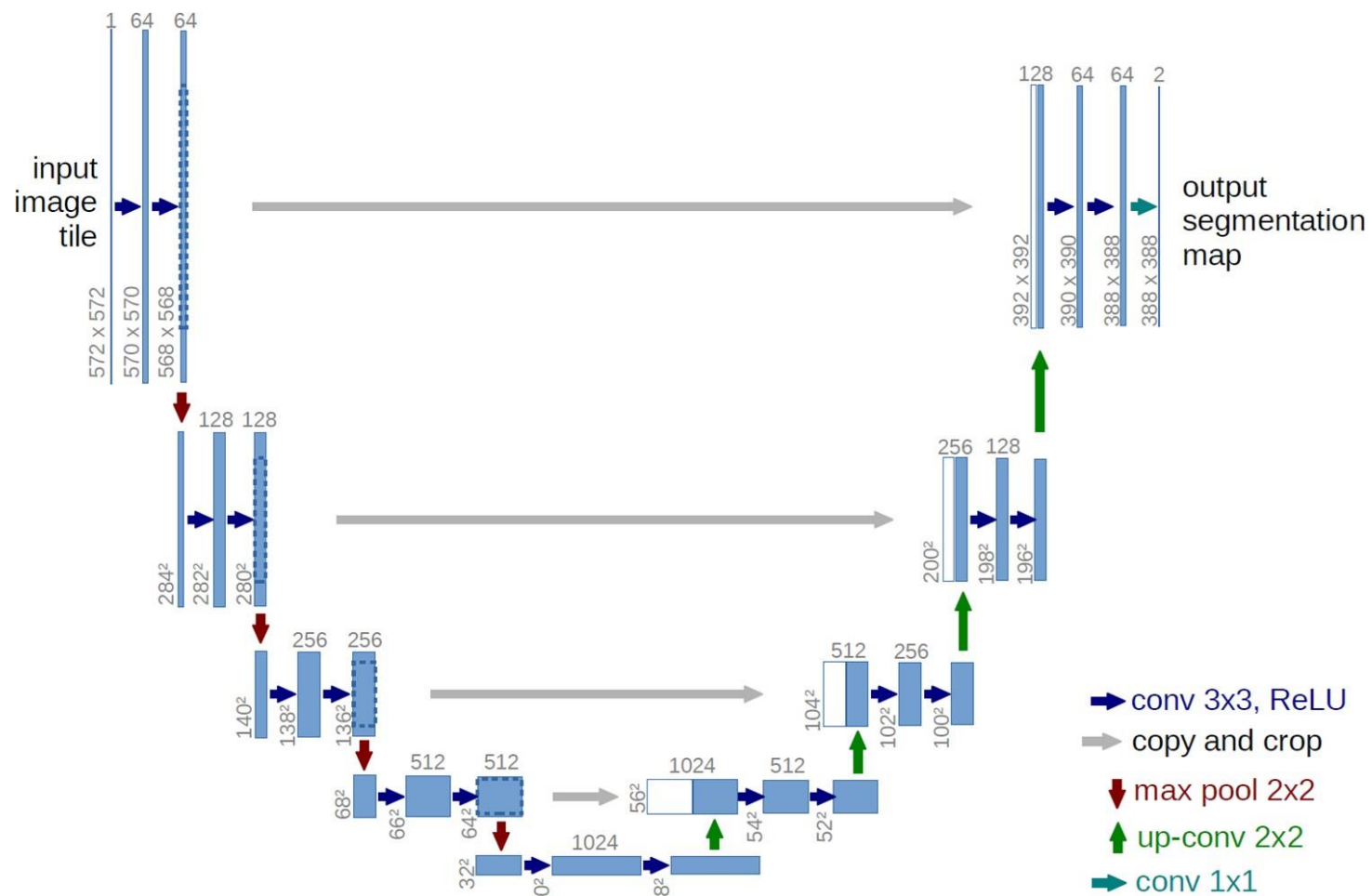
# Architectures: Upsampling



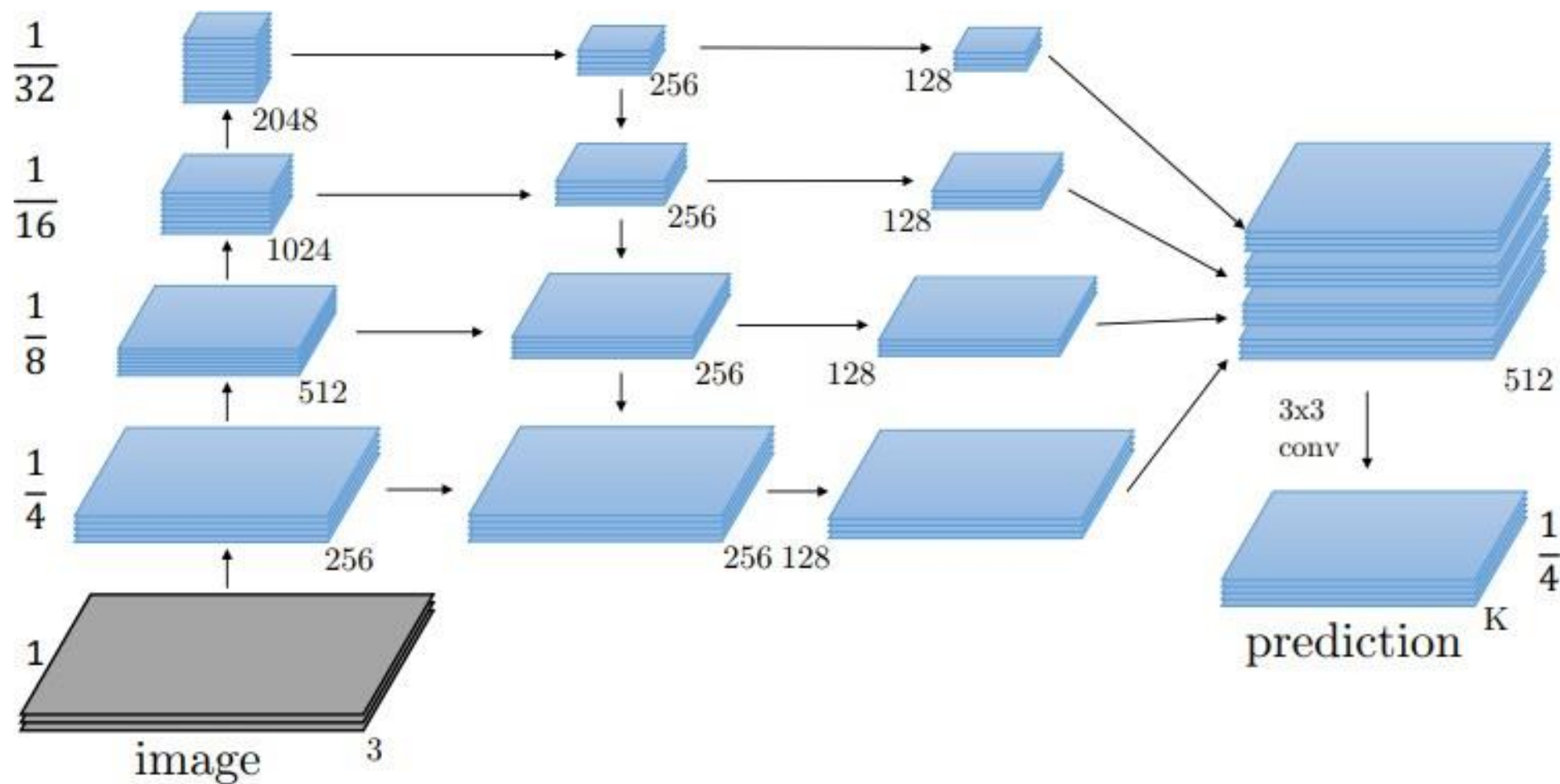
# Architectures: Deconvolution (or transposed convolution)



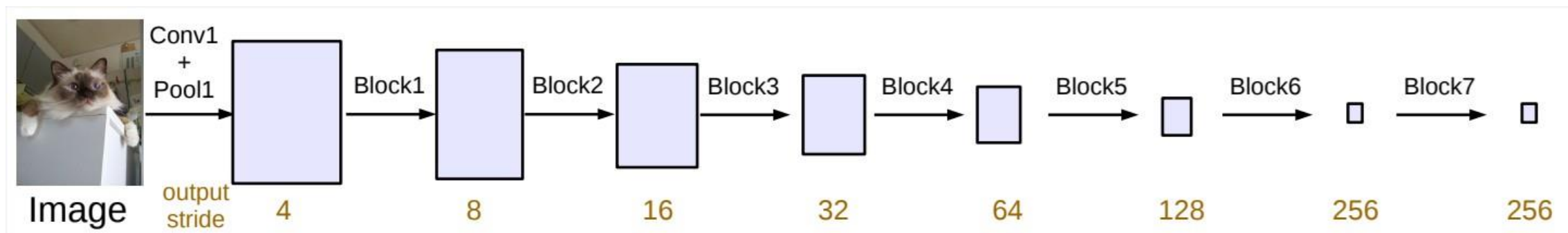
# Architectures: UNet



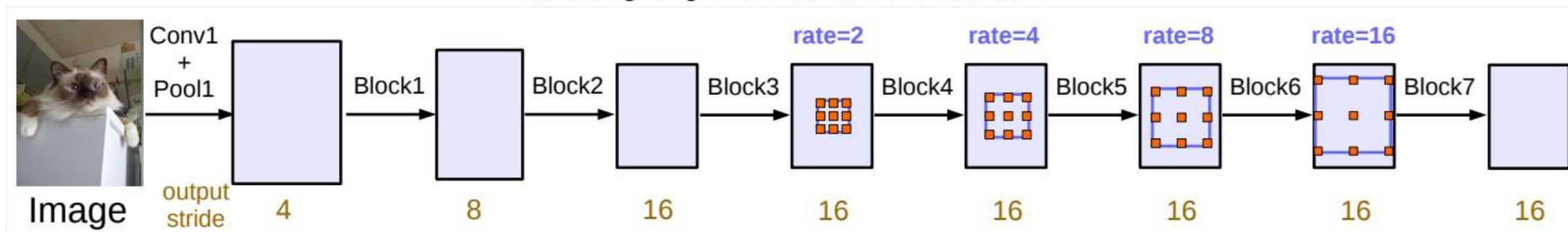
# Architectures: Feature Pyramid Network



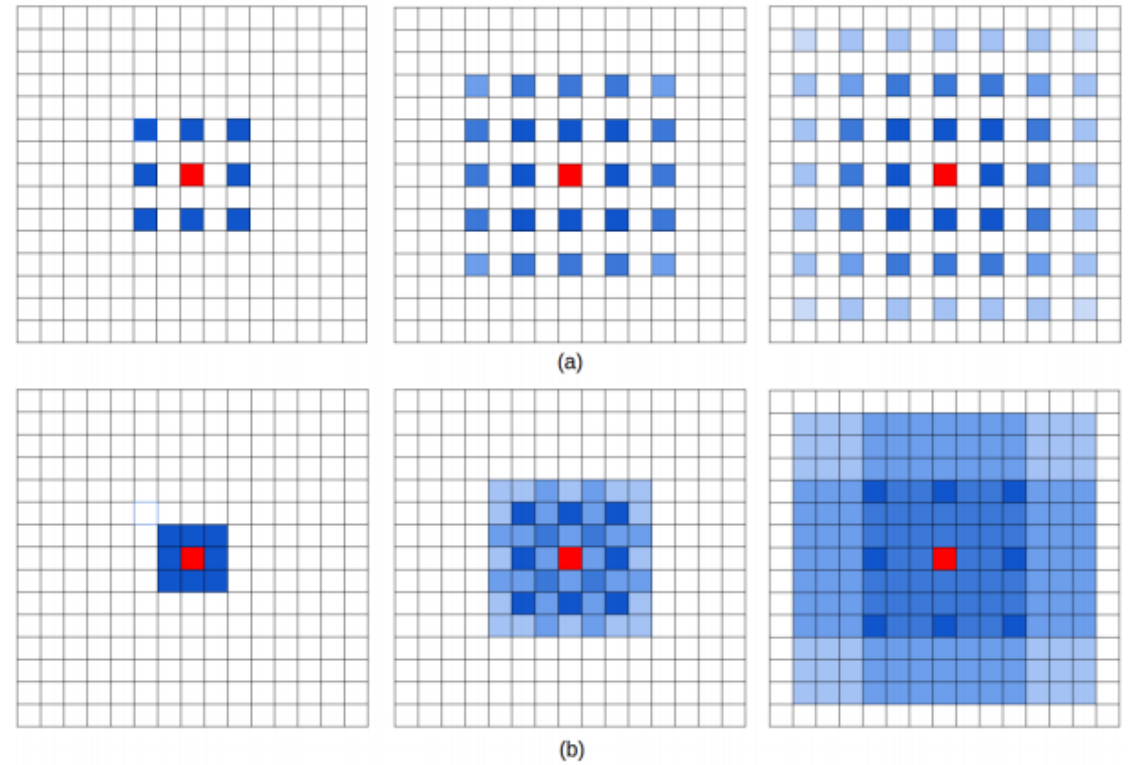
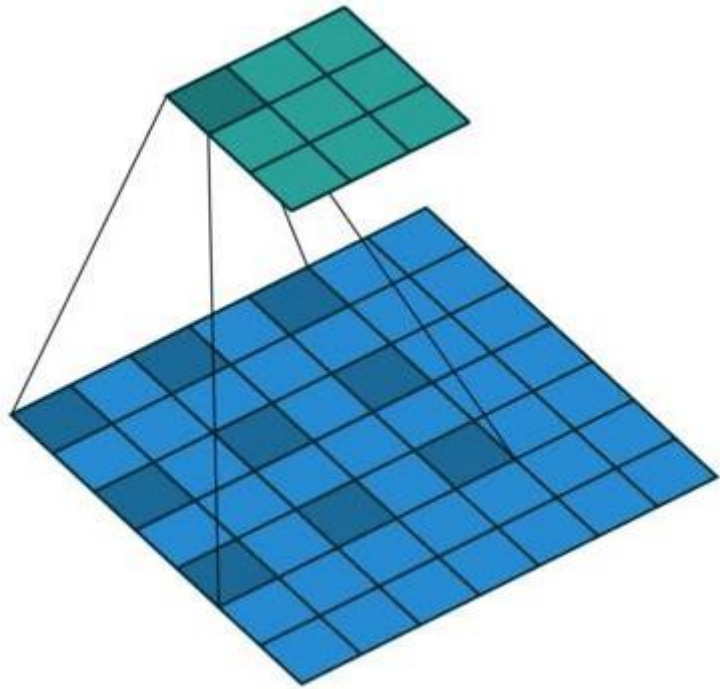
# Architectures: DeepLab v1



(a) Going deeper without atrous convolution.

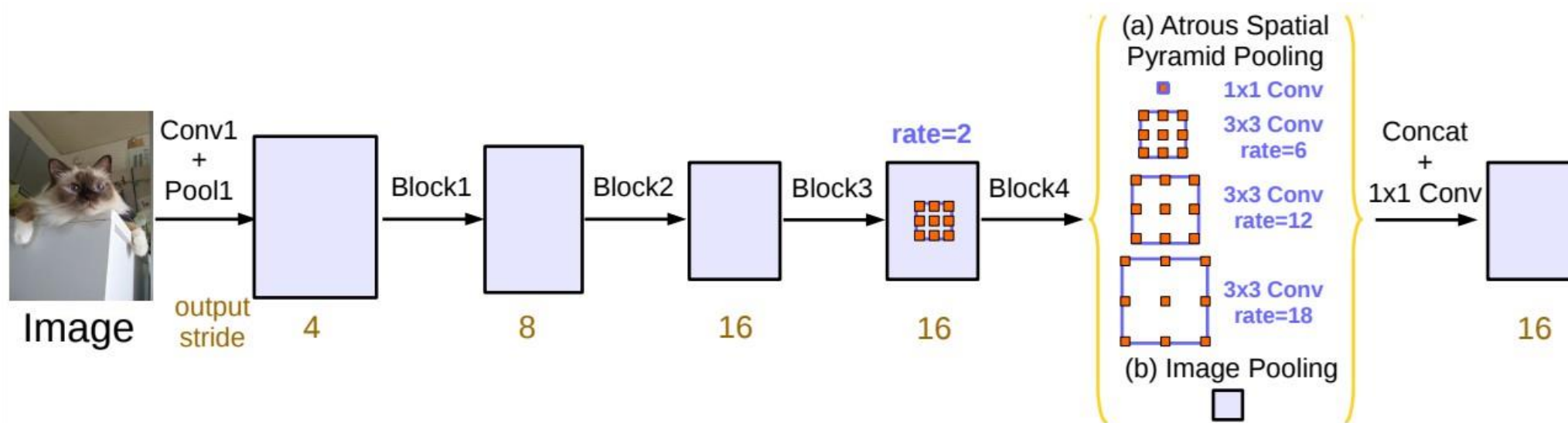


# Architectures: Atrous convolutions

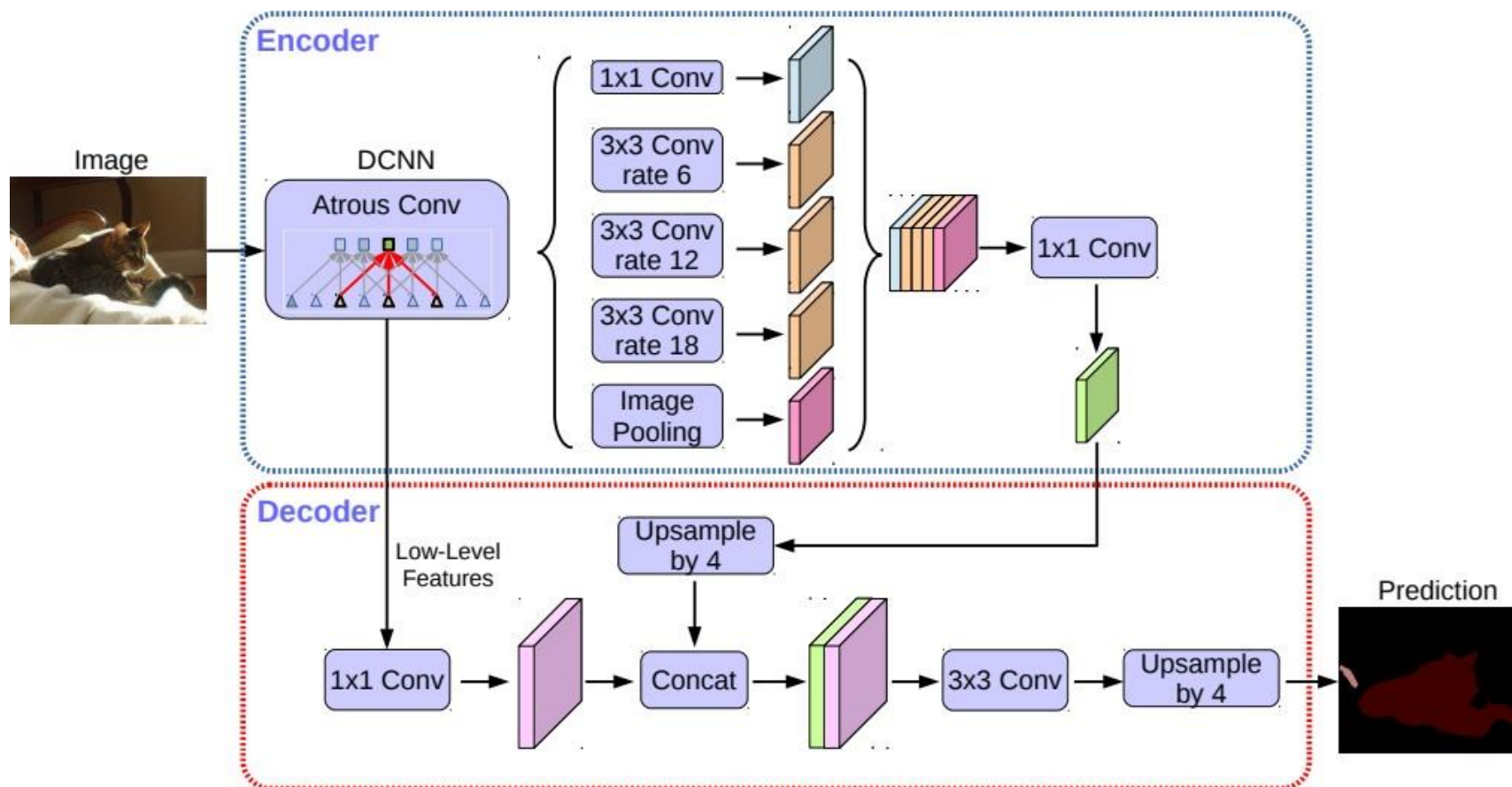




# Architectures: DeepLab v2



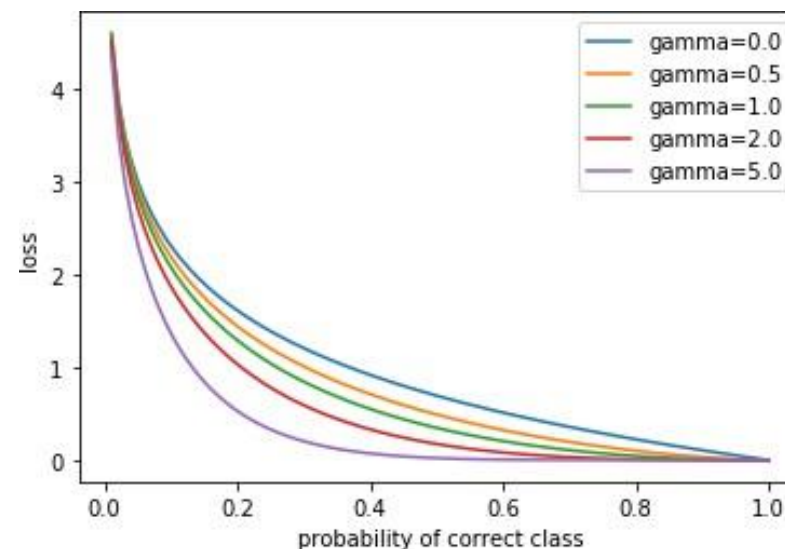
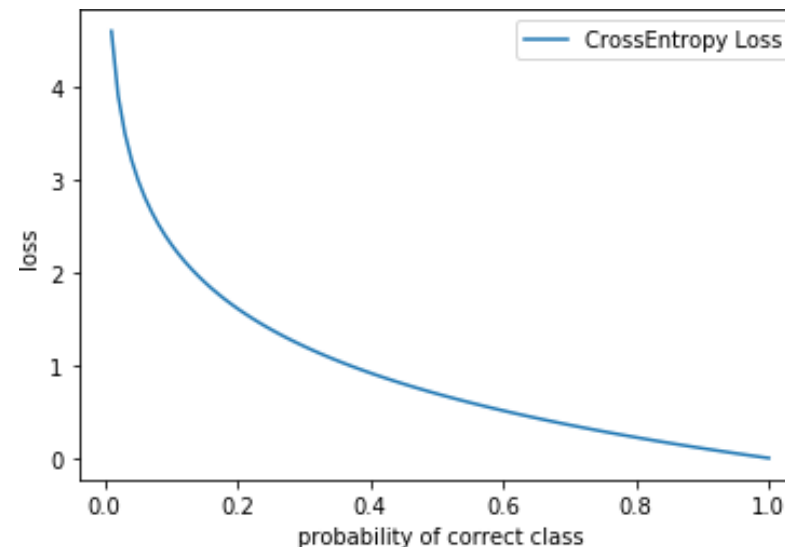
# Architectures: DeepLab v3+



# Loss functions: Cross entropy

$$L_{CE}(p, y) = - \sum_{c=1}^M y_c \log(p_c)$$

$$L_{CE}(p, y) = - \sum_{c=1}^M y_c (1 - p_c) \log p_c$$

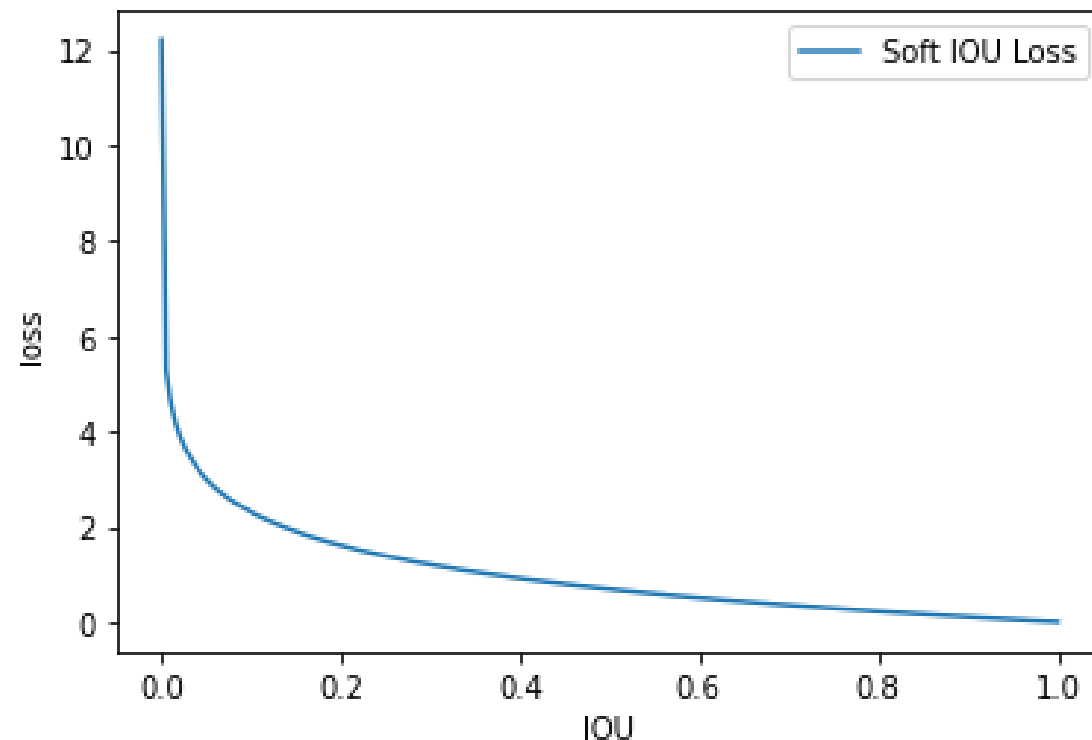


# Loss functions: IoU

$$IoU(A, B) = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$IoU(p, y) = \frac{\sum_{i=1}^N p_i y_i}{\sum_{i=1}^N p_i + \sum_{i=1}^N y_i - \sum_{i=1}^N p_i y_i}$$

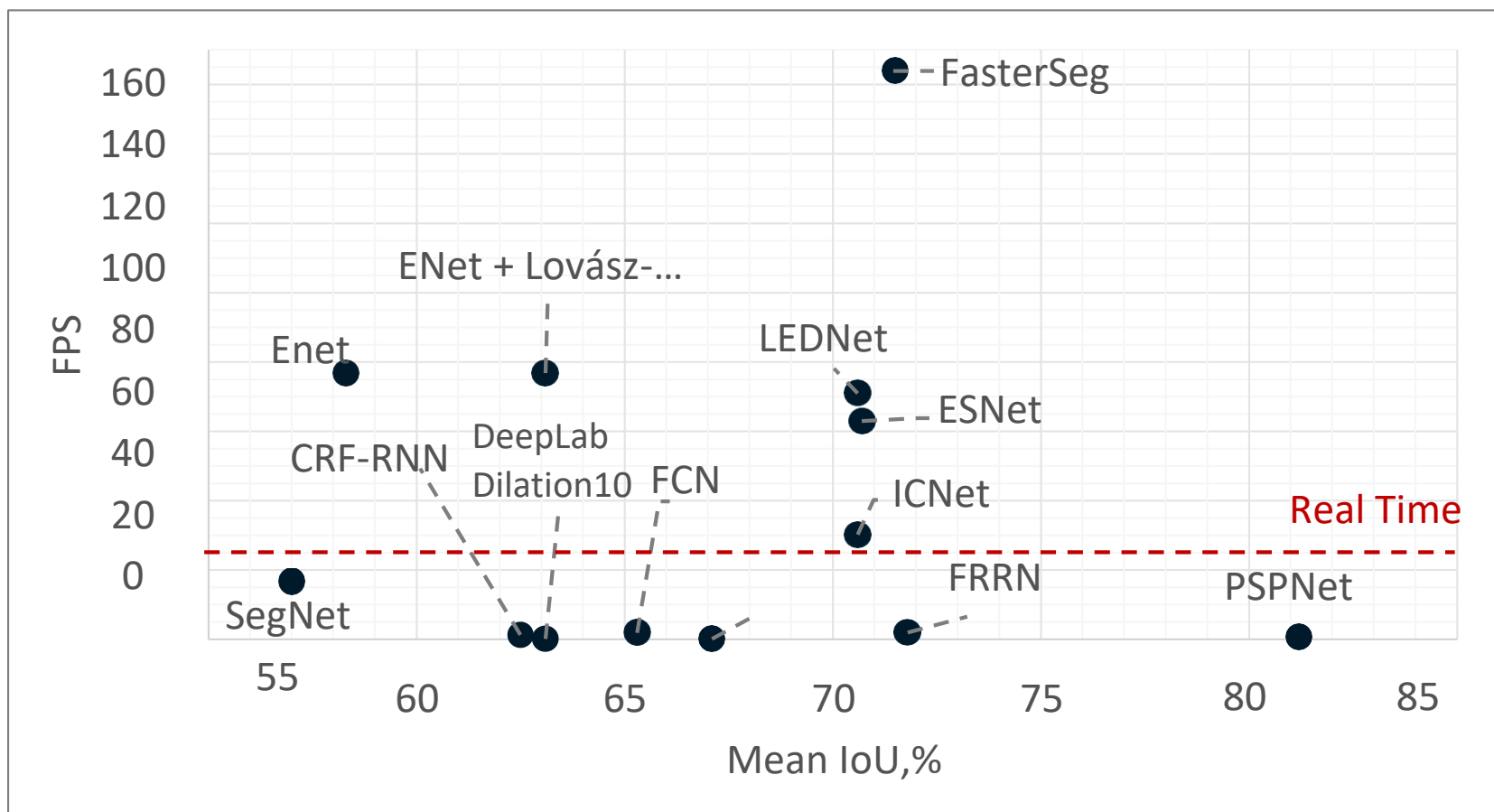
$$L_{IoU} = -\log \left( \frac{\sum_{i=1}^N p_i y_i}{\sum_{i=1}^N p_i + \sum_{i=1}^N y_i - \sum_{i=1}^N p_i y_i} \right)$$



# Comparison

Model	Year	Mean IoU, %	FPS	Latency, ms
DeepLab	2014	63.1	0.25	4000
SegNet	2015	57.0	16.7	60
CRF-RNN	2015	62.5	1.4	700
Dilation10	2015	67.1	0.25	4000
ENet	2016	58.3	76.9	13
FCN	2016	65.3	2	500
FRRN	2016	71.8	2.1	469
ICNet	2017	70.6	30.3	33
PSPNet	2017	81.2	0.78	1288
ENet + Lovász-Softmax	2018	63.1	76.9	13
LEDNet	2019	70.6	71	14
ESNet	2019	70.7	63	16
FasterSeg	2019	71.5	163.9	6.1

# Comparison





# Useful links

- UNet: <https://arxiv.org/abs/1505.04597>
- DeepLab: <https://arxiv.org/abs/1606.00915>
- DeepLabV3: <https://arxiv.org/abs/1706.05587>
- DeepLabV3+: <https://arxiv.org/abs/1802.02611>
- SegNet: <https://arxiv.org/abs/1511.00561>
- FCN: <https://arxiv.org/abs/1411.4038>
- Grad-CAM: <https://arxiv.org/abs/1610.02391>
  
- <https://github.com/mrgloom/awesome-semantic-segmentation>
- Kaggle: <https://www.kaggle.com/>
- ODS (@bes): <https://ods.ai/> <https://opendatascience.slack.com>
- Deep Learning Book: <https://www.deeplearningbook.org/>

AI is coming...