

EE 4146 Data Engineering and Learning Systems

Lecture 13: segmentation with deep learning

Semester A, 2021-2022

Finals

- Open-book and open-notes
- 10 multiple choices and 10 true/false, 10 understanding related questions and 6 calculations.
- No electronic device is allowed except calculator
- Cover all lecture information

Schedules

Week	Date	Topics
1	Sep. 1	Introduction
2	Sep. 8	Data exploration
3	Sep. 15	Feature reduction and selection (HW1 out)
4	Sep. 22	Mid-Autumn Festival
5	Sep. 29	Clustering I: Kmeans based models (HW1 due in this weekend)
6	Oct. 6	Clustering II: Hierarchical/density based/fuzzing clustering
7	Oct. 13	Midterm (no tutorials this week)
8	Oct. 20	Adverse Weather
9	Oct. 27	Linear classifiers
10	Nov. 3	Classification based on decision tree (Tutorial on project) (HW2 out)
11	Nov. 10	Bayes based classifier (Tutorial on codes) (HW2 due in this weekend)
12	Nov. 17	Non-linear Perceptron and Classifier ensemble
13	Nov. 24	Deep learning based models (Quiz)
14		Summary: based on the poll, we will do off-line video recording for the summary and will upload the video before Dec. 1 st .

Course content

■ Course grades

- 2 Homework assignments (10%) (105 submission & 108 submission)
- 1 Project report (10%) ([Deadline Dec. 8](#))
- 1 Midterm (20%) (109 submission)
- 1 Quiz (10%) (Nov. 24)
- Final exam (50%)

■ Remarks:

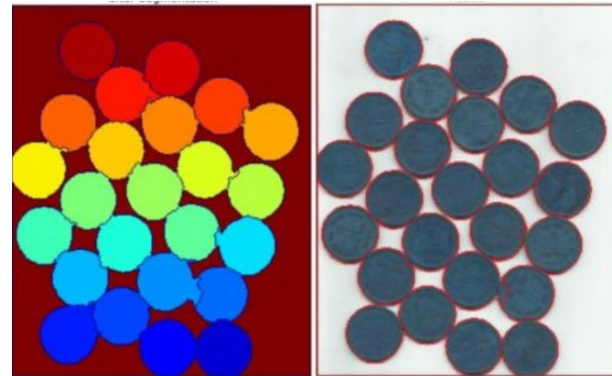
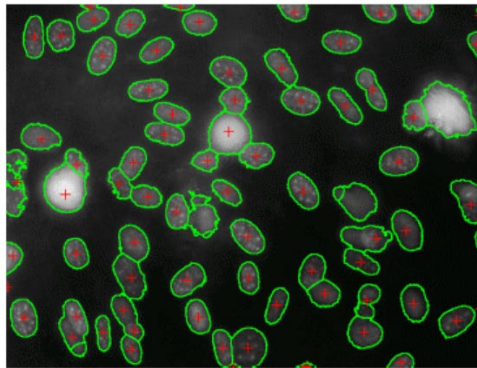
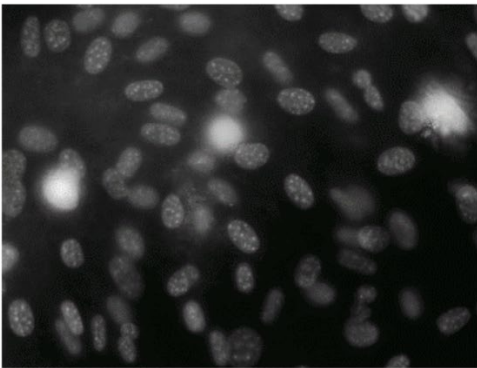
- By university regulation, to pass the course, students are required to achieve at least 30% in course work and 30% in the examination.
- Attend at least two course work, including homework assignments and quiz.

Lecture outline

- Introduction
- Different models
 - Fully Convolutional Network
 - DeconvNet, SegNet
 - U-Net
 - PSPNet
 - DeepLab v1, v2, v3, v3+
 - Transformer
- Loss functions

Image Segmentation

- Segmentation:
 - Split/separate/subdivide an image into regions or objects
 - To facilitate recognition, understand region of interest
- Challenges of Segmentation:
 - The definition of a region/object is problem-dependent
 - One of the most difficult task in image processing
 - Accuracy determines success or failure of application

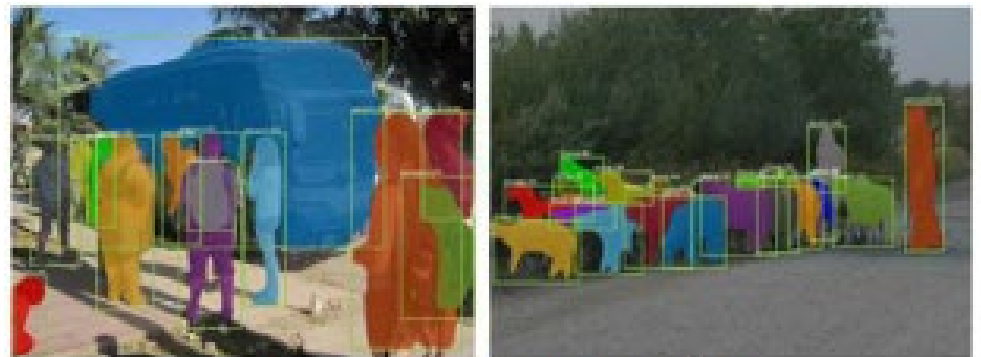


Introduction

- Typically, image segmentation is applied to derive semantics (meaning) known as semantic segmentation.
 - This involves both classification and localization.
 - Localization is at pixel-level unlike object detection where bounding boxes are used.
- Instance segmentation: each instance or occurrence of an object is assigned a different label.



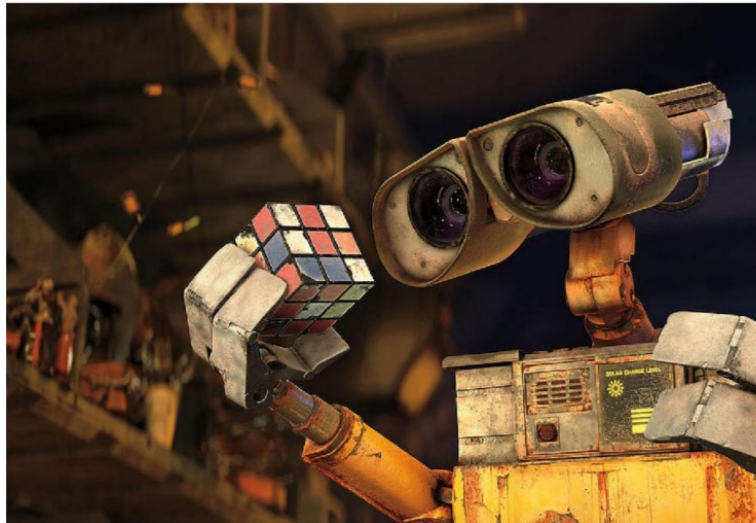
Semantic segmentation



Instance segmentation

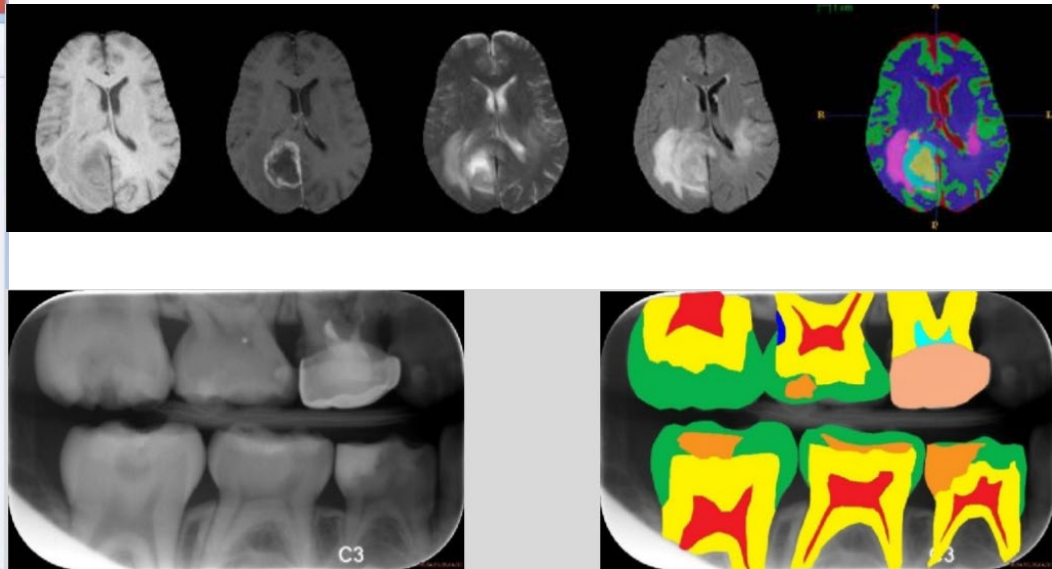
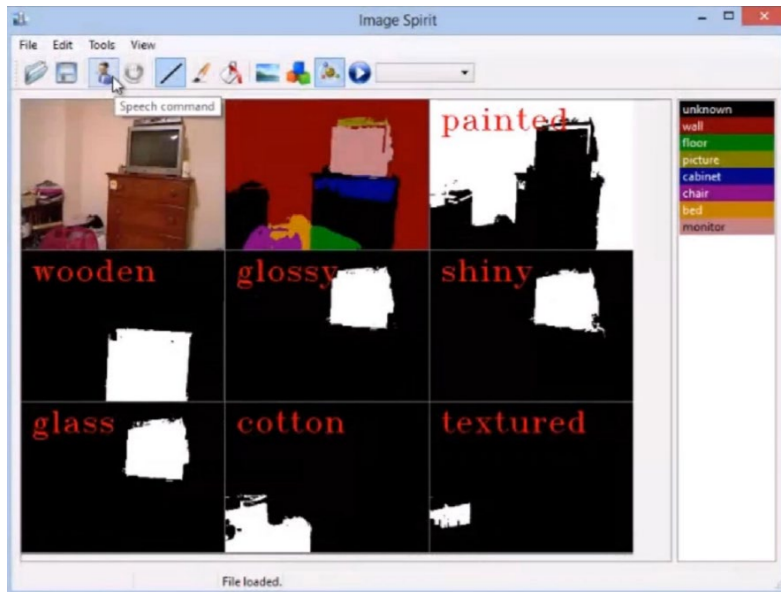
Applications

- To let robots segment objects so that they can grasp them
- Road scenes understanding; useful for autonomous navigation of cars and drones



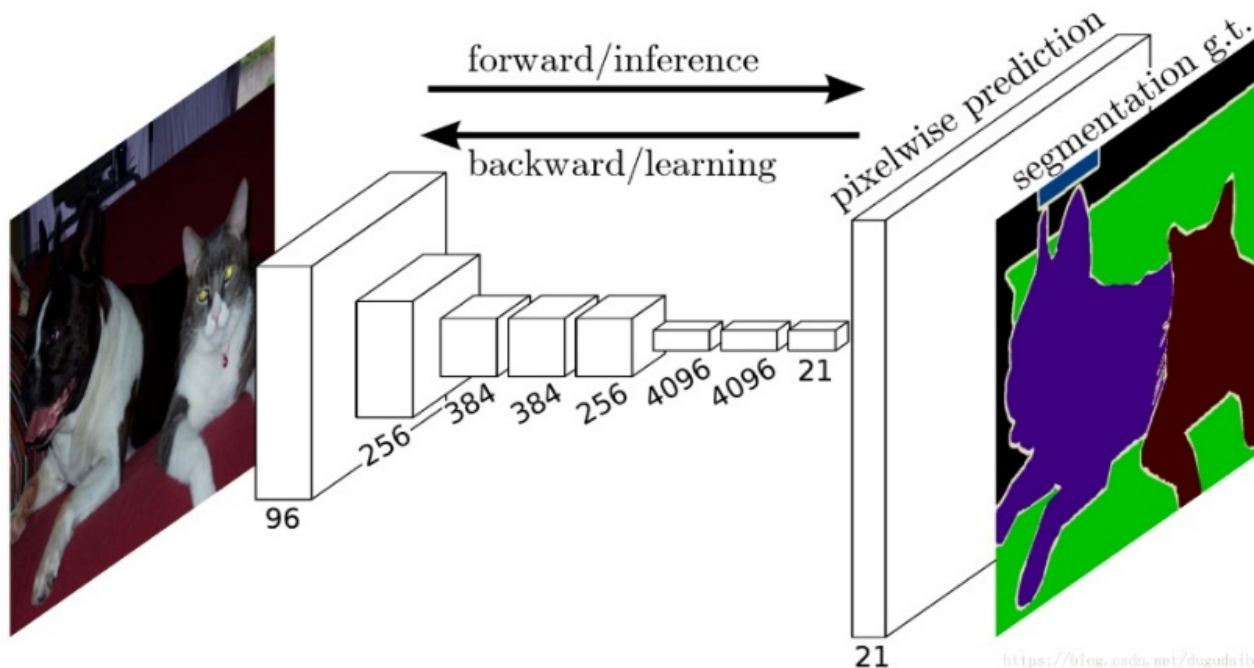
Applications

- Useful tool for editing images
- Medical purposes: e.g. segmenting tumors, dental cavities, ...



Segmentation with deep learning models

- Define CNN architecture
- Define a loss measuring performance (loss function)
 - In the training, loss values will be used to update the model parameters.
- Minimize the loss (optimizer)

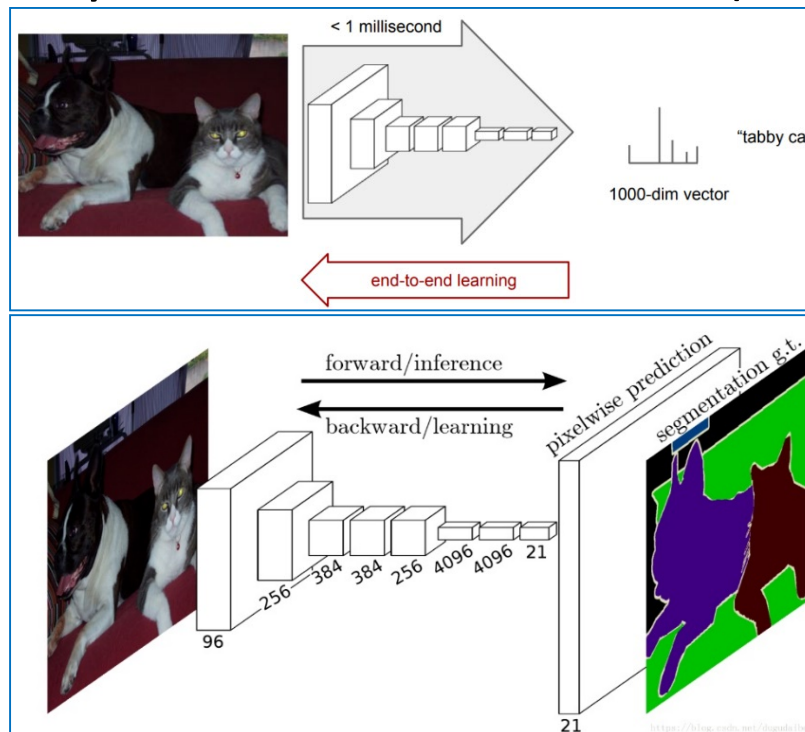


Lecture outline

- Introduction
- Different models
 - Fully Convolutional Network
 - DeconvNet, SegNet
 - U-Net
 - PSPNet
 - DeepLab v1, v2, v3, v3+
- Loss functions

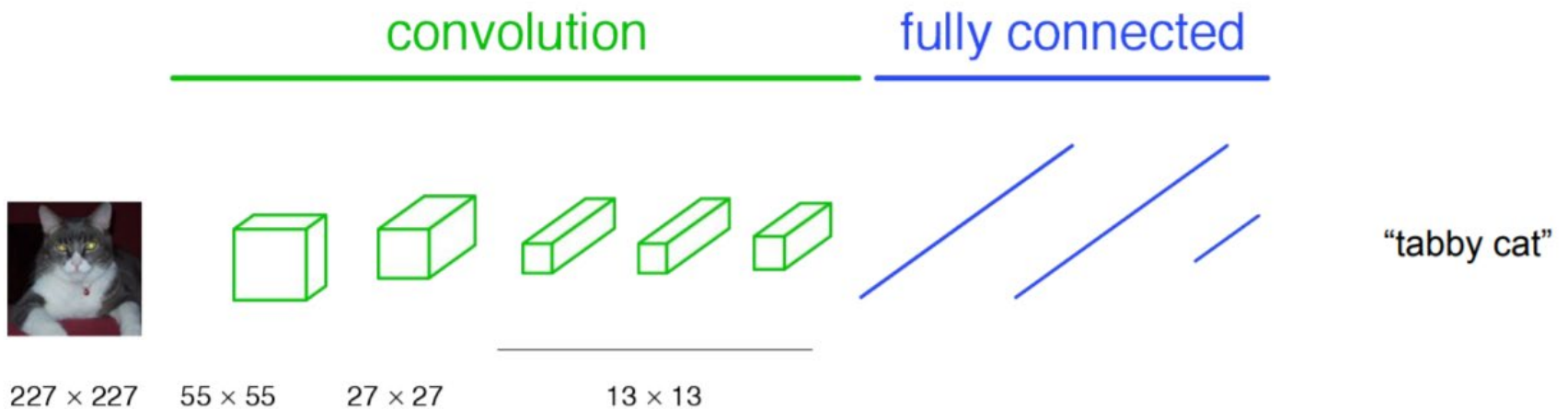
Fully Convolutional Network

- First segmentation architecture proposed using a Convolutional Neural Network (CNN).
- Key idea is to eliminate the use of fully connected layers and replace it with convolution layers.
- Hence called fully convolutional network (FCN)



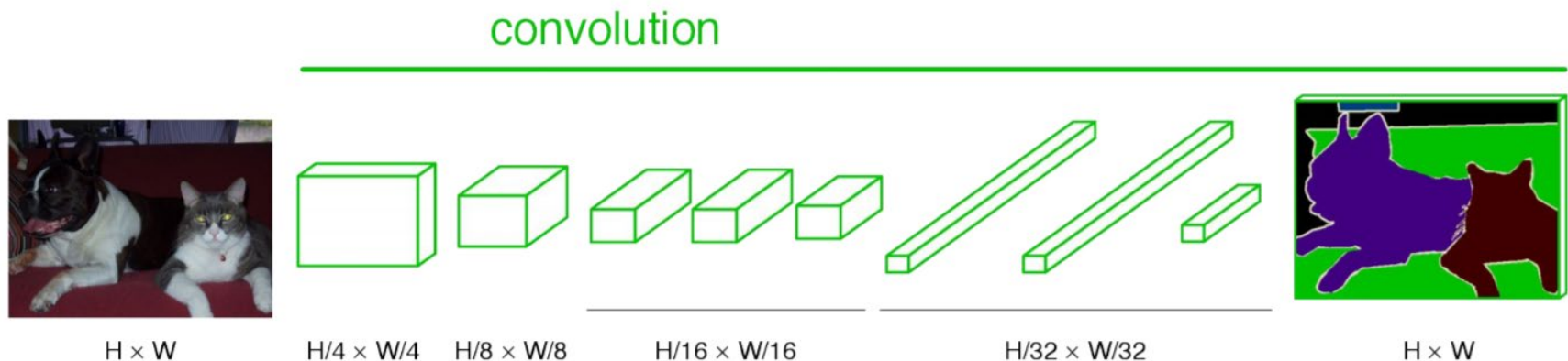
Fully Convolutional Network

- An example of a classification network.
- Convolution layers are typically followed by fully connected layers for classification.



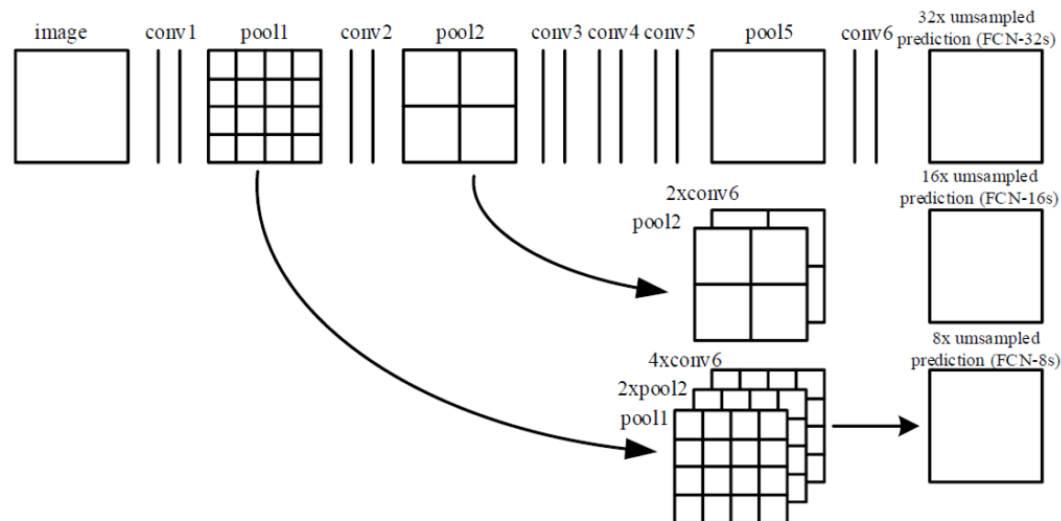
Fully Convolutional Network

- FCN transfers knowledge from VGG16 to perform semantic segmentation.
- The fully connected layers of VGG16 is replaced by convolution layers.
- Upsample the last layer using deconvolution layer to produce output of same size as input.
- This is FCN-32s, as a stride of 32 is used by deconvolution kernel.



Fully Convolutional Network

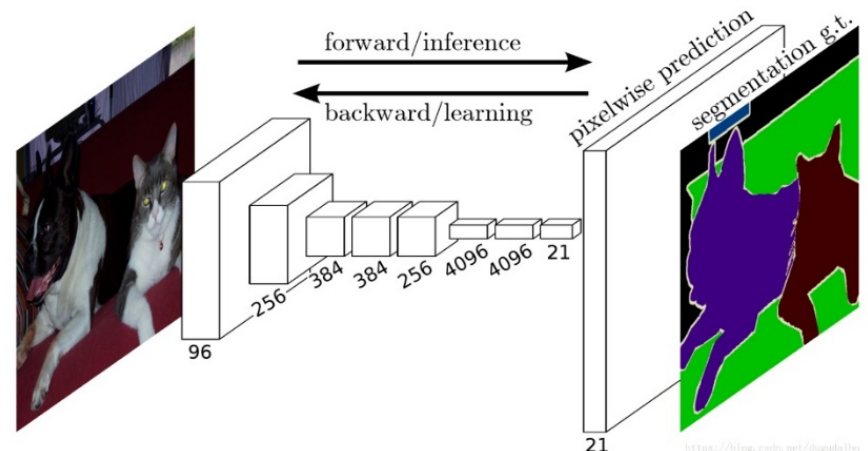
- Other variants do exist such as FCN-8s, FCN-16s
- The deconvolution at the last layer can lose a lot of resolution
- One option is adding “skip” connections from earlier higher-resolution layers. For these variants the last convolution layers are upsampled and fused with earlier pooling layers to produce finer results (via summation).



Fully Convolutional Network

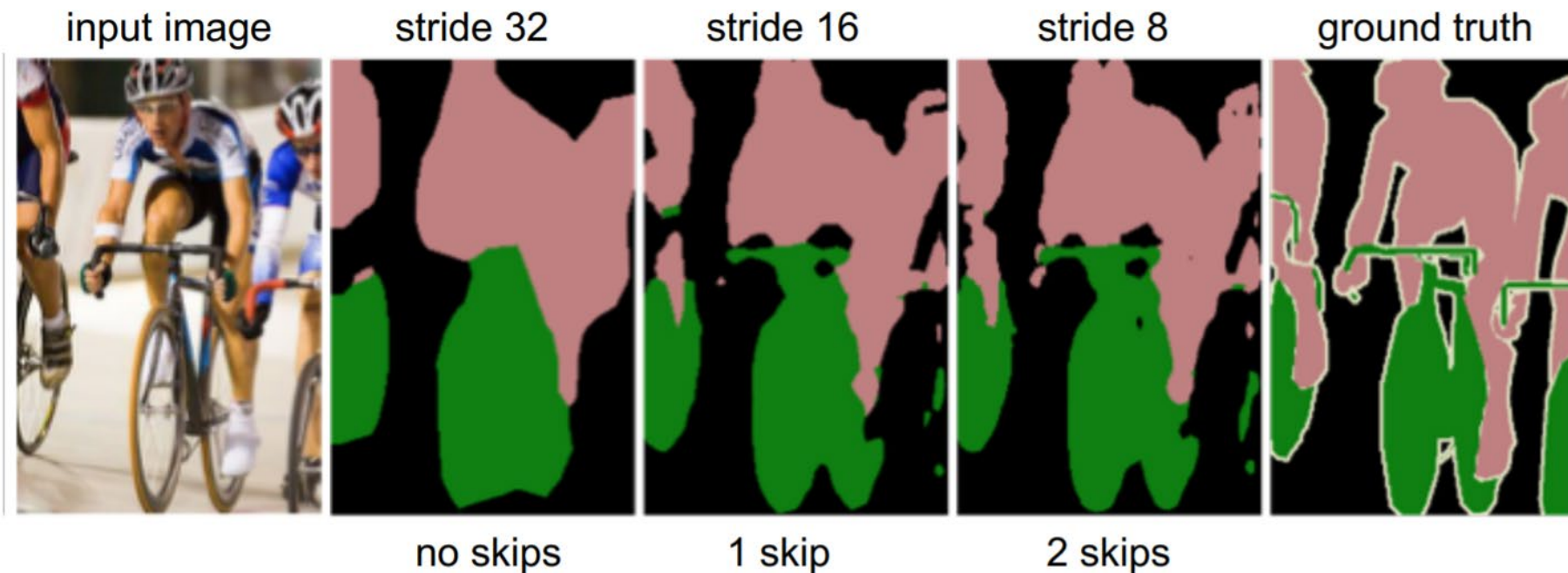
- The output of the last layer is a tensor of depth k , where k is the number of classes.
- **Softmax** is applied such that values are stored between 0 and 1.
- Ground truth mask is stored as one-hot encoding.
- One-hot encoding is a representation of categorical values such that the active class is assigned “1” whereas the rest are “0”.

Label	Car	Person	Building
1	1	0	0
2	0	1	0
3	0	0	1



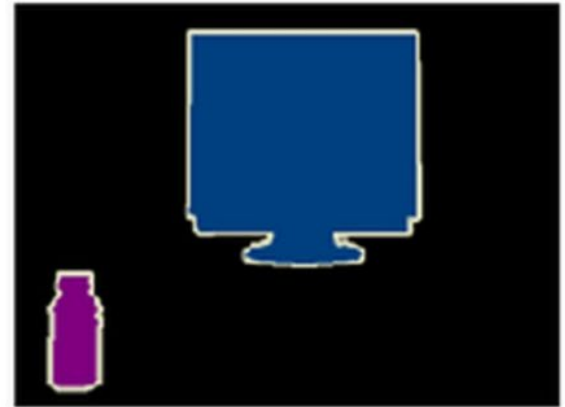
Fully Convolutional Network

- Comparison with different network baselines



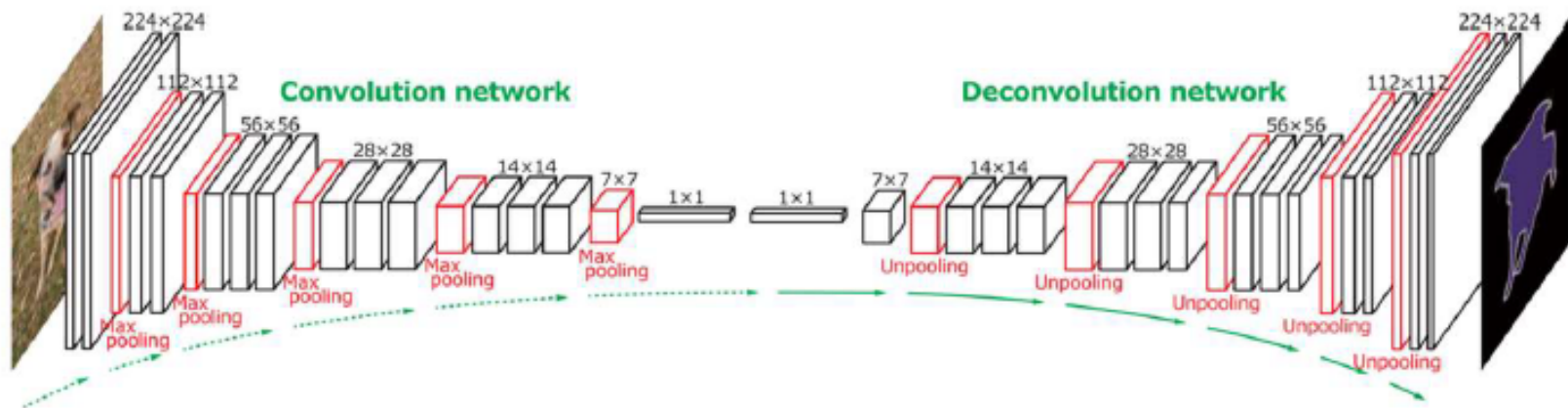
Fully Convolutional Network

- Significantly improved the state of the art in semantic segmentation.
- Poor object delineation: e.g. spatial consistency neglected.



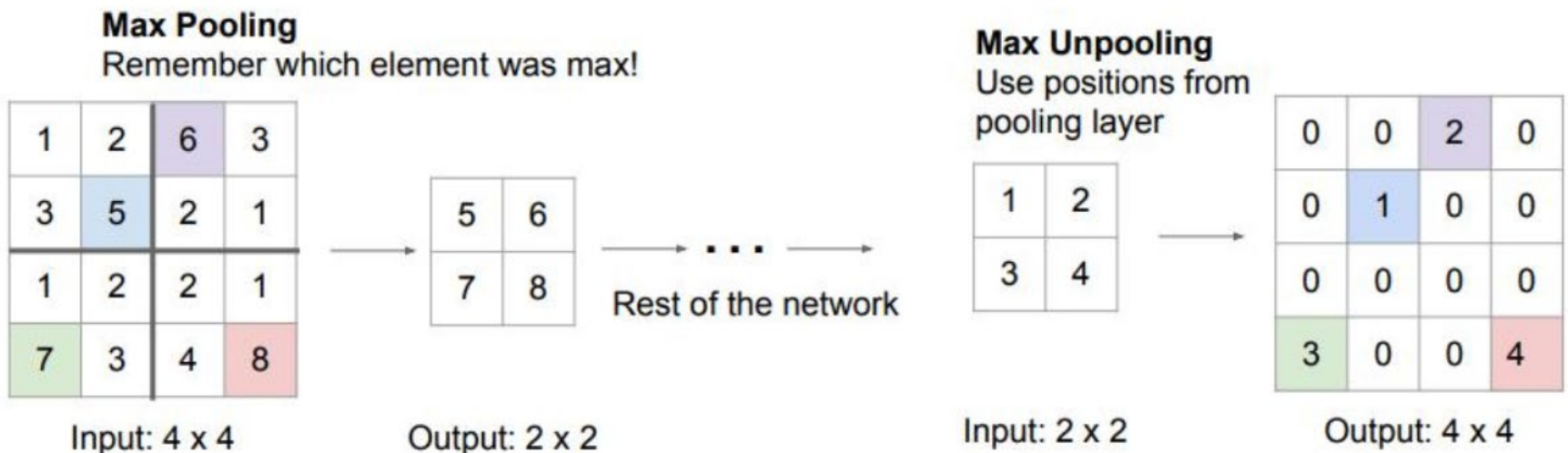
Deconvolutional Network

- Has an **Encoder-Decoder** (Convolution-Deconvolution) structure.
- Decoder side has a stack of unpooling layer and deconvolution layers.
- Compared to FCN, **finer details** are preserved in the output.



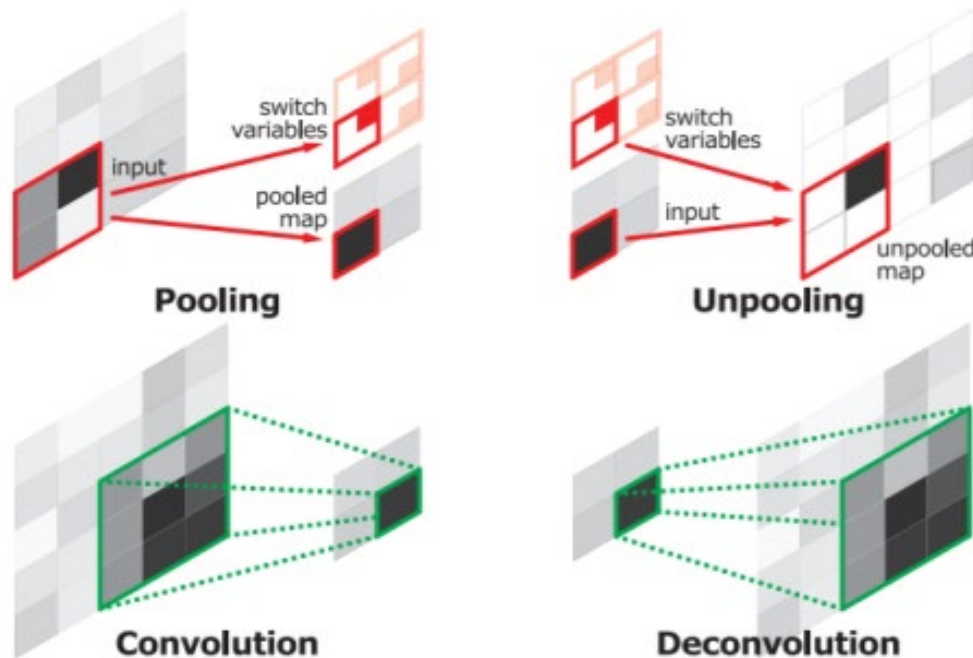
Deconvolutional Network

- Pooling operation (max): selects the max value across a window. Hence, downsamples the input.
- Unpooling operation (max):
 - Stores the indices of max value from the encoder side.
 - On the decoder side, the unpooling layer places back each output back to its pooled location.
 - Due to this operation, the output will be sparse.



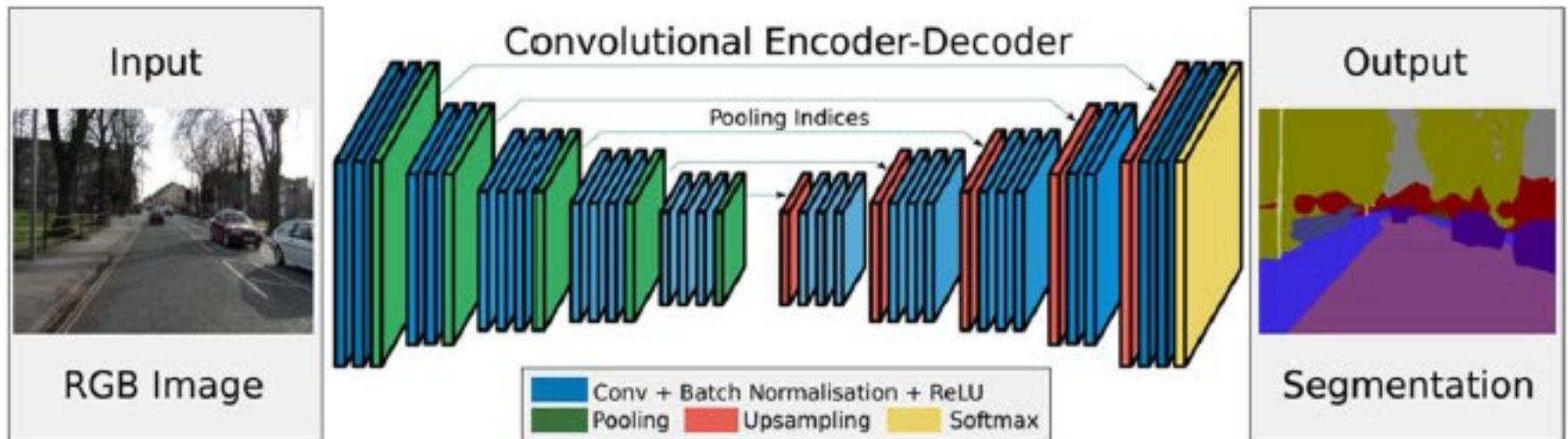
Deconvolutional Network

- The unpooling layer produces sparse output.
- To have dense predictions, deconvolution layers are applied.
 - Convolution - multiple inputs with a filter to produce single output.
 - Deconvolution: single input with a filter to produce multiple outputs.
- Note that deconvolution is called transposed convolution in practice.



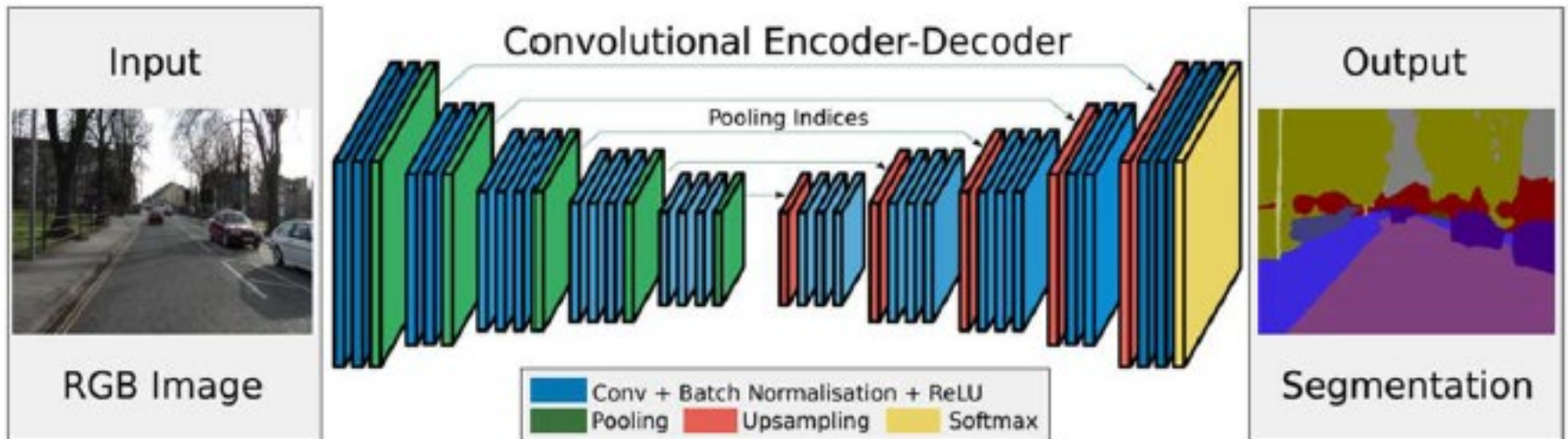
SegNet

- Similar to deconvolutional network but without fully-connected layers.
- Far less parameters as **no fully-connected layers** are used.
- Better performance than deconvolutional network.



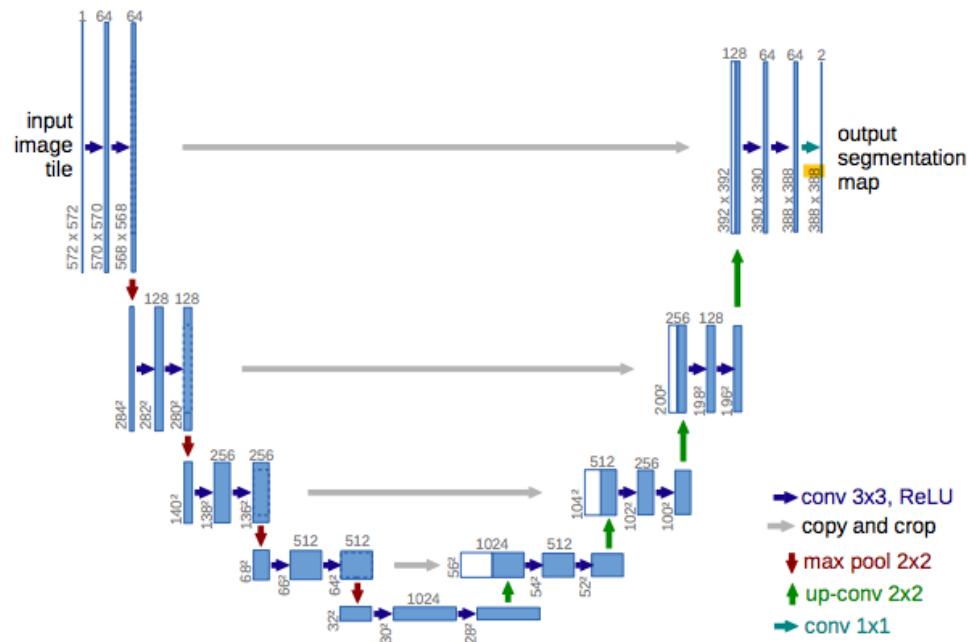
SegNet

- Each encoder: one or more convolutional layers with batch normalization and a ReLU non-linearity, followed by non-overlapping maxpooling
- Use **max-pooling** indices in the decoders to perform upsampling of low resolution feature maps
- Retain high frequency details and also reduce the total number of trainable parameters in the decoders
- Tend to be smooth even without a **Conditional Random Field** based post-processing.(a kind of post-processing model)



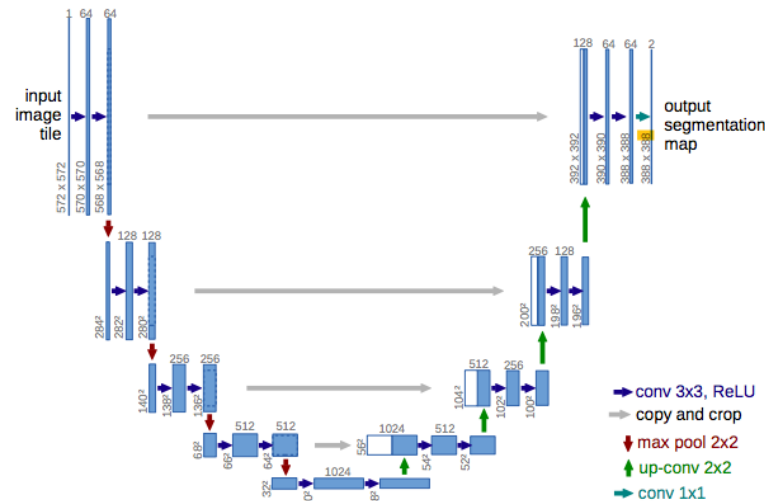
U-Net

- Introduced for ISBI challenge of segmentation neuronal structures.
- Has “U” structure as the name implies.
- At every three layers, outputs are cropped.
- On the decoder side, the network concatenates the outputs from encoder and then upsamples it.
- There is also a “W” net for segmentation!



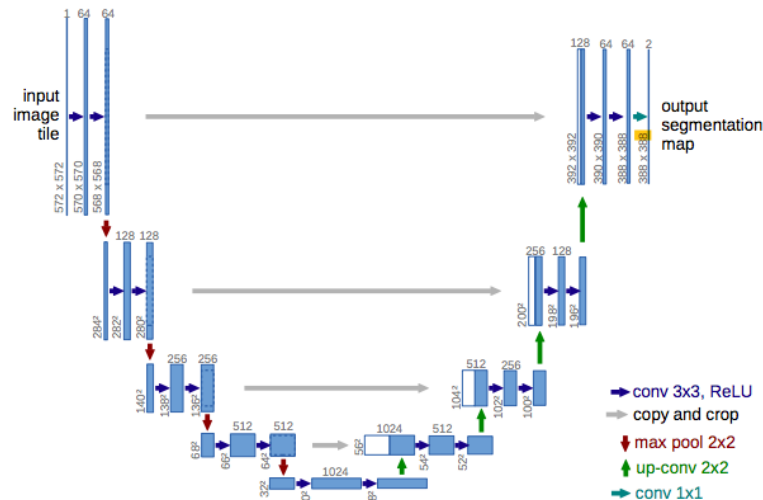
U-Net

- Contracting/downsampling path with 4 blocks
 - 3x3 Convolution Layer + activation function (with batch normalization)
 - 3x3 Convolution Layer + activation function (with batch normalization)
 - 2x2 Max Pooling
 - Purpose: capture the context of the input image in order to be able to do segmentation. This coarse contextual information will then be transferred to the upsampling path by means of skip connections.



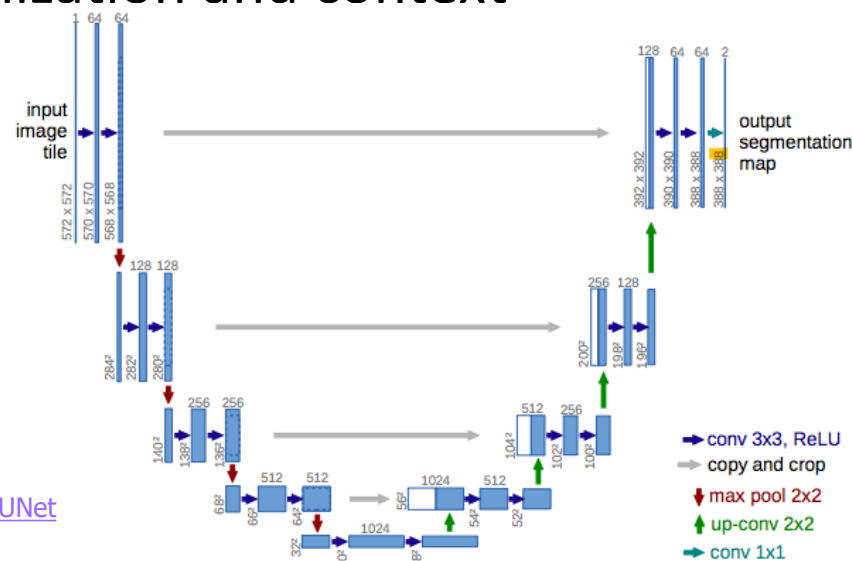
U-Net

- Expanding/upsampling path with 4 blocks.
 - Deconvolution layer with stride 2
 - Concatenation with the corresponding cropped feature map from the contracting path
 - 3x3 Convolution layer + activation function (with batch normalization)
 - 3x3 Convolution layer + activation function (with batch normalization)
- Purpose: enable precise localization combined with contextual information from the contracting path.



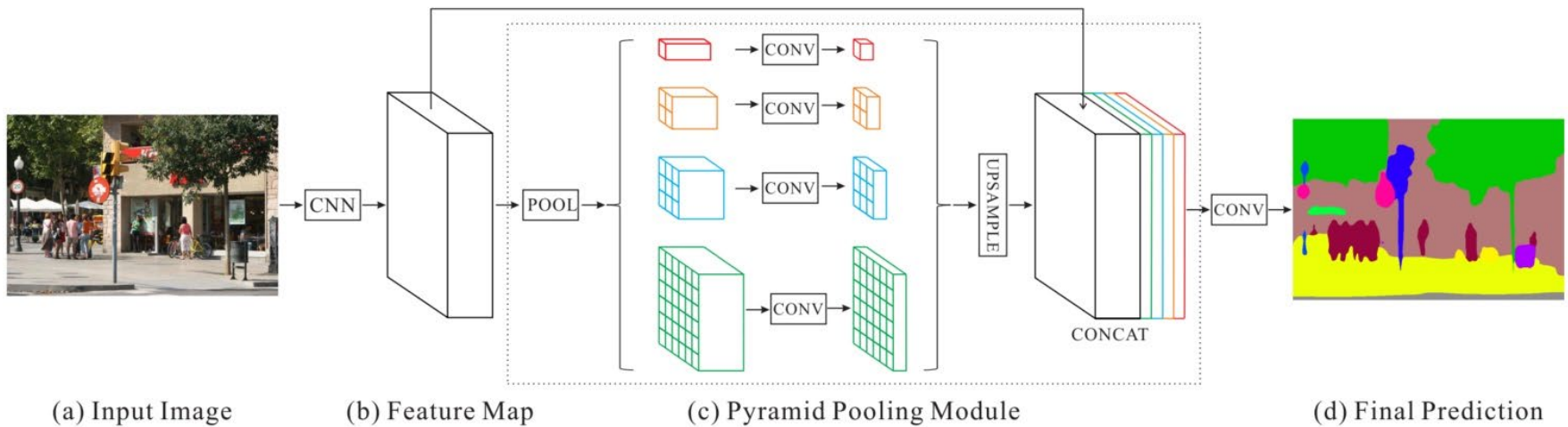
Unet

- Bottleneck
 - This part of the network is between the contracting and expanding paths. The bottleneck is built from simply 2 convolutional layers (with batch normalization), with dropout.
- The U-Net combines the location information from the downsampling path with the contextual information in the upsampling path to finally obtain a general information combining localization and context



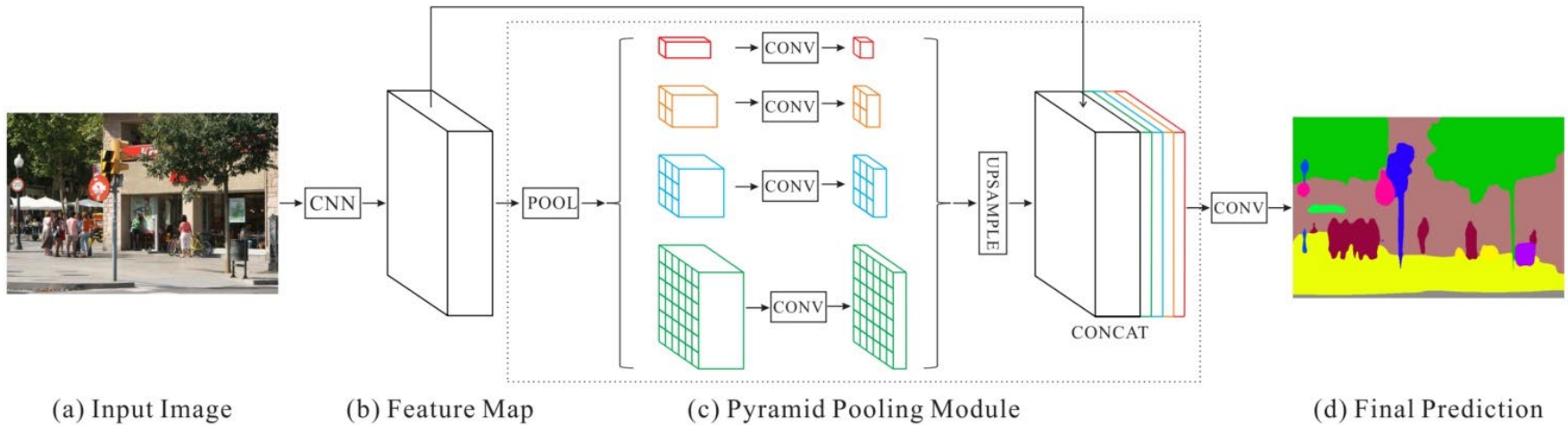
Pyramid Scene Parsing Network (PSPNet)

- Motivation: Errors occur because of contextual relations between of objects.
- Eg: car on the road (not in the sky!)
- To take context into account, Pyramid pooling module is proposed.



Pyramid Scene Parsing Network (PSPNet)

- Take output from last layer of a CNN,
 - Pool it at different levels (**global average pooling, $2*2, 3*3, 6*6$**).
 - Convolution followed by each pooling
 - Upsample all the outputs and concat the results.
 - Convolution to produce the final outputs.



DeepLab

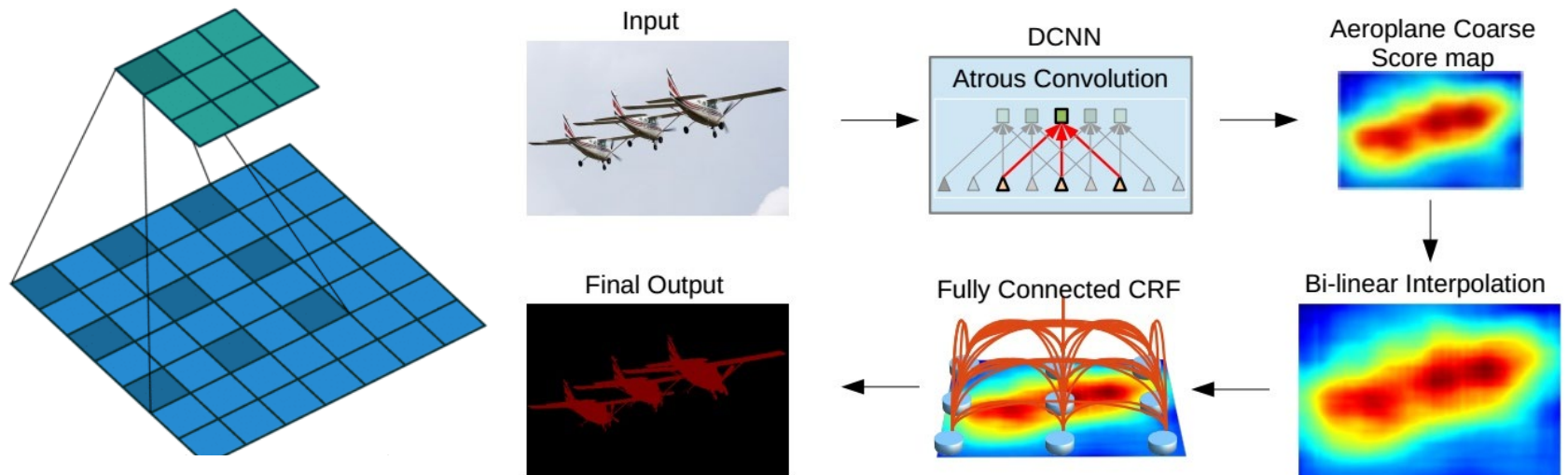
- DeepLab is a state-of-art deep learning model for semantic image segmentation, where the goal is to assign semantic labels (e.g. person, dog, cat) to every pixel in the input image.
- DeepLabv1, DeepLabv2, DeepLabv3, DeepLabv3+,

<https://www.tensorflow.org/lite/models/segmentation/overview>

DeepLab v1

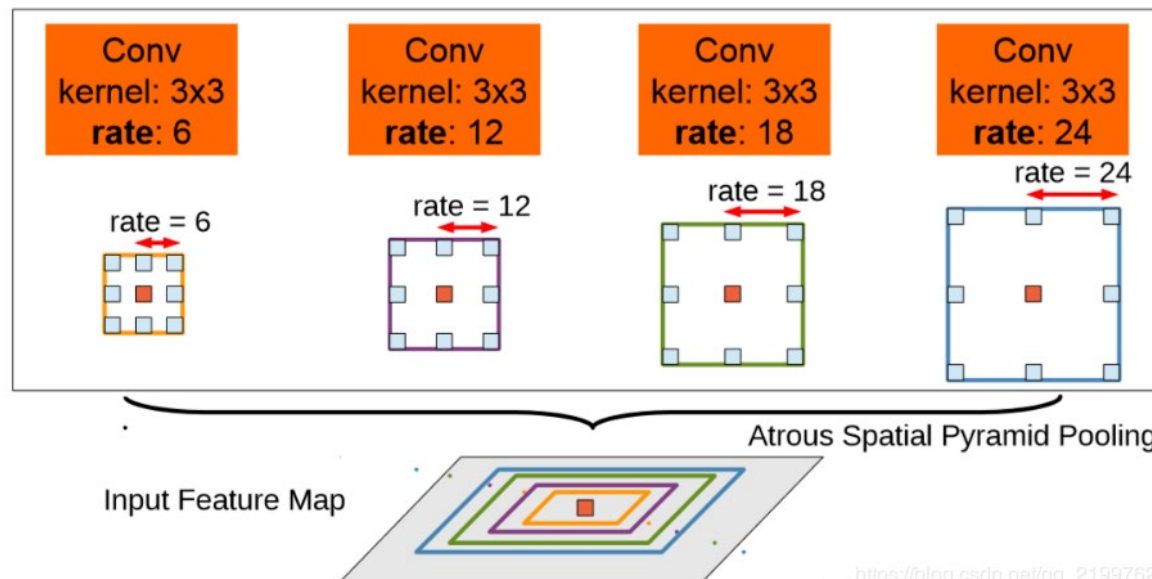
■ DeepLab v1

- Use pretrained **VGG-16** to generate coarse map.
- Atrous convolution enlarges the field of view of filters to incorporate larger context without increasing the number of parameters or the amount of computation.
- Bilinear interpolation to upsample result.
- Apply **Conditional Random Field (CRF)** to refine the result.



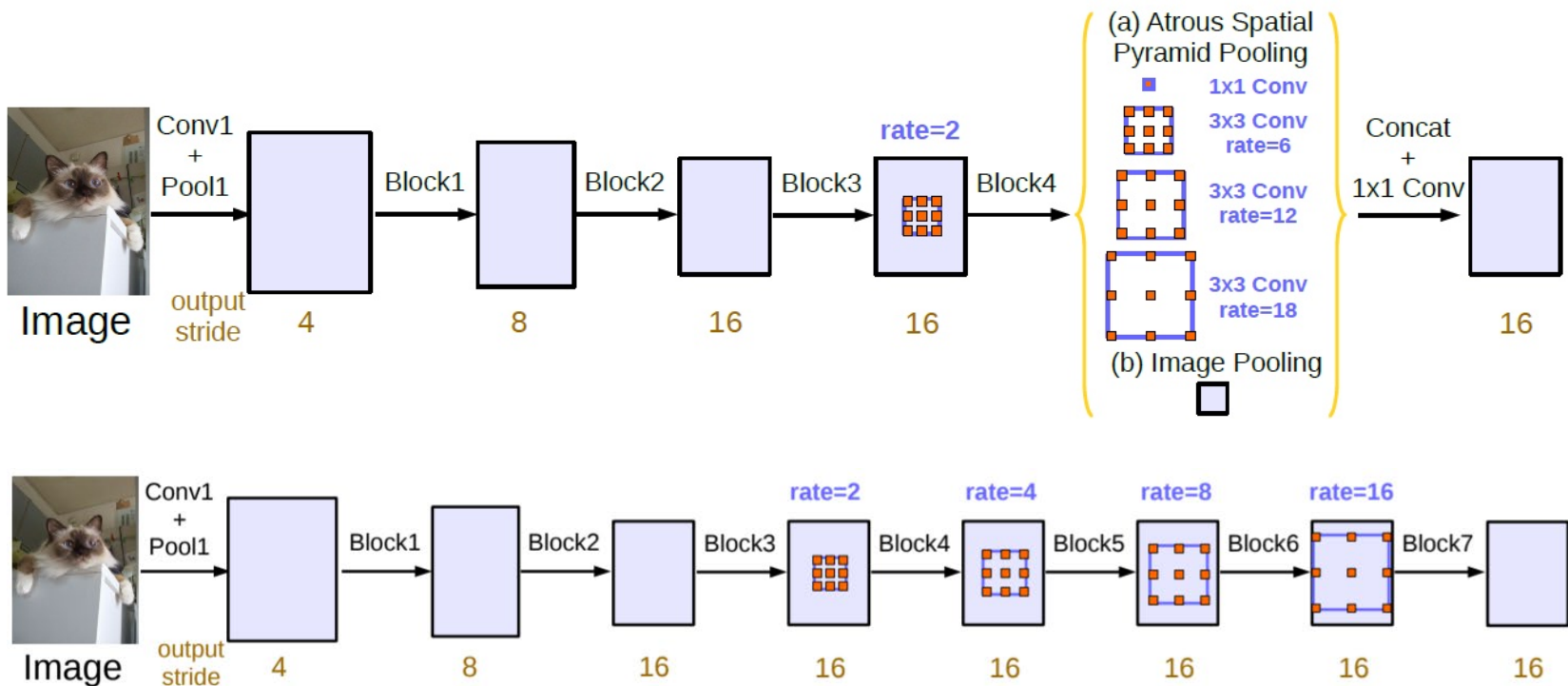
DeepLab v2

- DeepLab v2
- Motivation: existence of object at multiple scales.
 - Use [ResNet](#) instead of VGG.
 - Introduce [Atrous Spatial Pyramid Pooling](#) (ASPP): the outputs at the last convolution layer is dilated at different rates in parallel branches and fused together. ASPP helps to account for different object scales which can improve the accuracy.



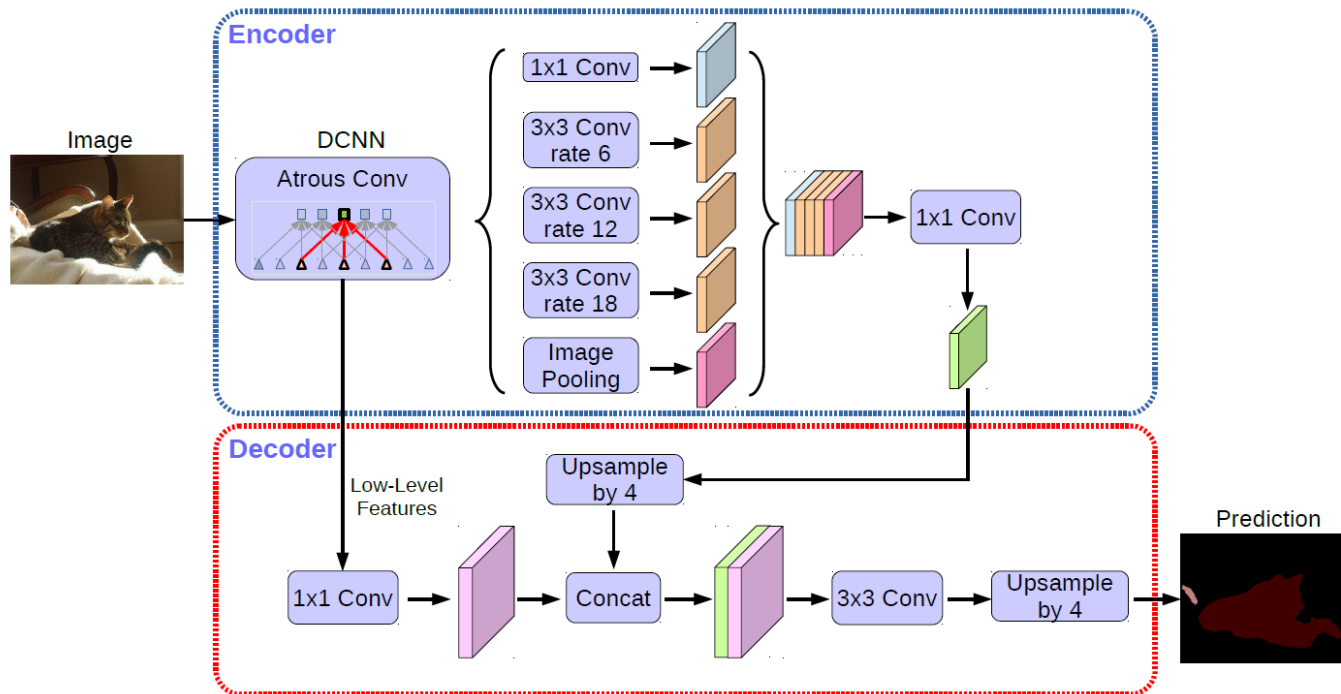
DeepLab v3

- Motivation: capture sharper object boundaries
- Can be used for different framework
- Design modules which employ atrous convolution in **cascade** or in **parallel** to capture multi-scale context by adopting multiple atrous rates
- Different Atrous Spatial Pyramid Pooling: image pooling and 1*1 conv



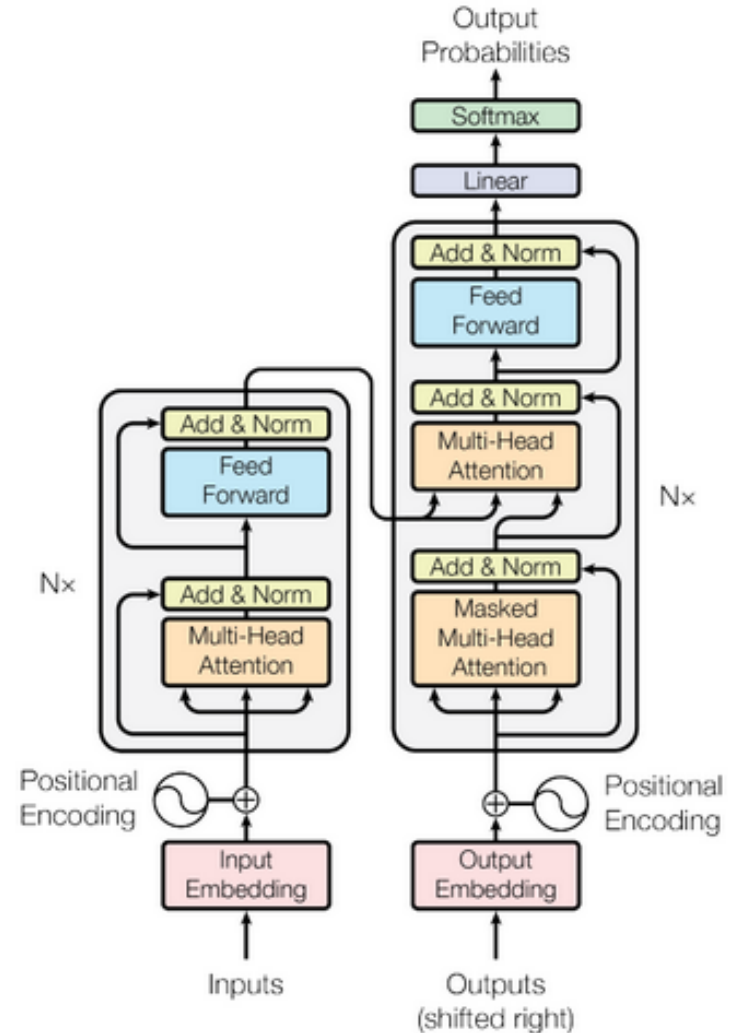
DeepLab v3+

- Add a simple yet effective decoder module to further refine the segmentation results especially along object boundaries.
- DeepLab v3 + Decoder = DeepLab v3+



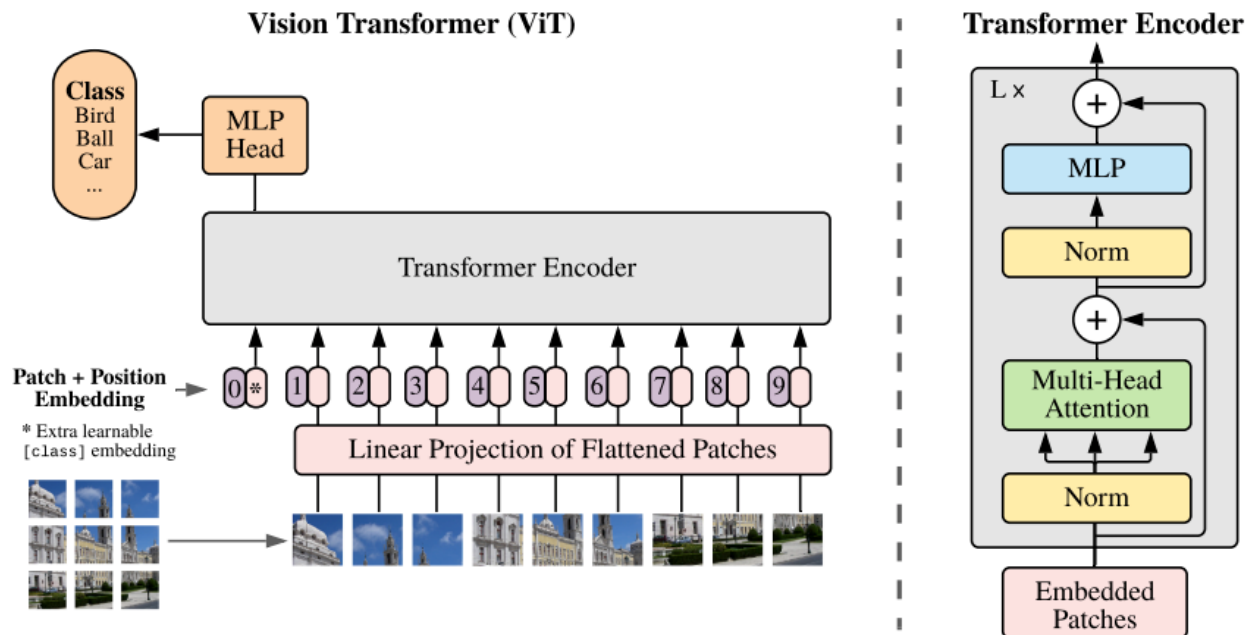
Original transformer

- It is used primarily in the field of natural language processing (NLP)
- A transformer is a deep learning model that adopts the mechanism of self-attention, differentially weighting the significance of each part of the input data.
- Transformers are designed to handle sequential input data
- Transformers do not necessarily process the data in order.
- Rather, the attention mechanism provides context for any position in the input sequence.



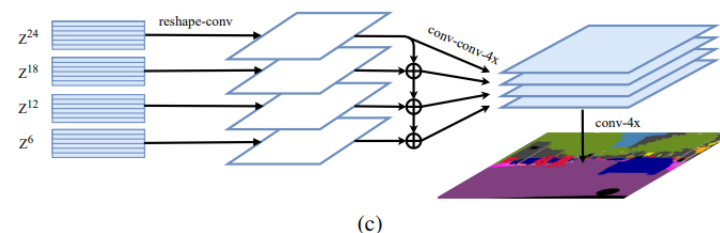
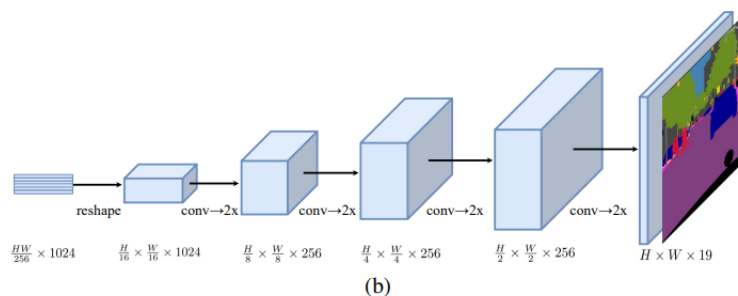
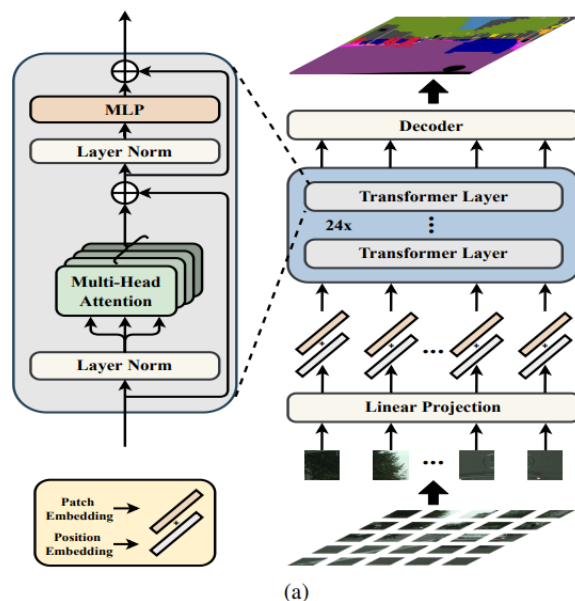
Original ViT

- Vision Transformer (ViT), a vision model based as closely as possible on the Transformer architecture originally designed for text-based tasks.
 - ViT represents an input image as a sequence of image patches, similar to the sequence of word embeddings used when applying Transformers to text, and directly predicts class labels for the image.
 - ViT demonstrates excellent performance when trained on sufficient data, outperforming a comparable state-of-the-art CNN with four times fewer computational resources.



SEgmentation TRansformer (SETR)

- With the global context modeled in every layer of the transformer, this encoder can be combined with a simple decoder to provide a powerful segmentation model
 - Decompose an image into a grid of fixed-sized patches, forming a sequence of patches
 - With a linear embedding layer applied to the flattened pixel vectors of every patch to obtain a sequence of feature embedding vectors as the input to a transformer
 - Feed learned features to a decoder to recover the original image resolution



(b) progressive upsampling (resulting in a variant called SETRPUP); and (c) multi-level feature aggregation

Swin Transformer

- Challenges in adapting Transformer from language to vision arise from differences between the two domains, such as large variations in the scale of visual entities and the high resolution of pixels in images compared to words in text.
 - The shifted windowing scheme brings greater efficiency by limiting self-attention computation to non-overlapping local windows while also allowing for cross-window connection.
 - This hierarchical architecture has the flexibility to model at various scales and has linear computational complexity with respect to image size.

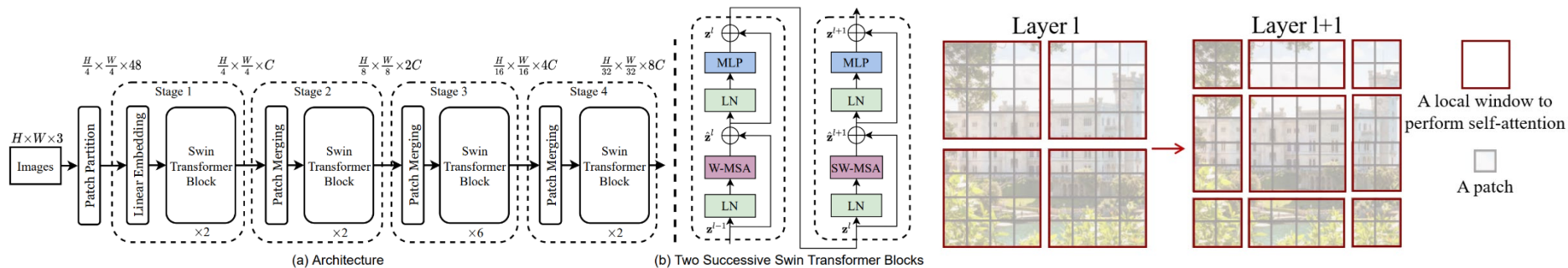


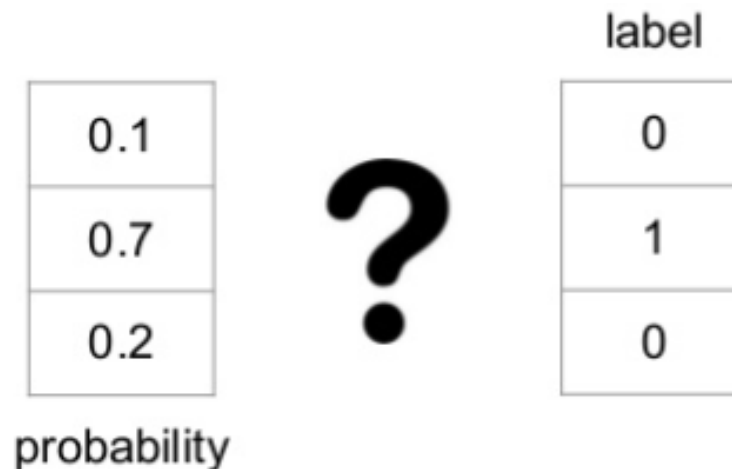
Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

Lecture outline

- Introduction
- Different models
 - Fully Convolutional Network
 - DeconvNet, SegNet
 - U-Net
 - PSPNet
 - DeepLab v1, v2, v3, v3+
- Loss functions

Loss functions for segmentation

- In a supervised deep learning, the loss functions measure the quality of a particular set of parameters based on how well the output of the network agrees with the ground truth labels in the training data.
- Loss functions are used to **guide the training process** in order to find a set of parameters that reduce the value of the loss function.



Loss functions

- Definition
- Let $\mathbf{P}(Y = 0) = p$ and $\mathbf{P}(Y = 1) = 1 - p$
- The predictions are given by the logistic/sigmoid function

$$\mathbf{P}(\hat{Y} = 0) = \frac{1}{1+e^{-x}} = \hat{p} \text{ and } \mathbf{P}(\hat{Y} = 1) = 1 - \frac{1}{1+e^{-x}} = 1 - \hat{p}$$

- **Cross entropy (CE)** can be defined as follows

$$\text{CE}(p, \hat{p}) = -(p \log(\hat{p}) + (1 - p) \log(1 - \hat{p}))$$

- This loss examines each pixel individually, comparing the class predictions (depth-wise pixel vector) to our one-hot encoded target vector.

Weighted cross entropy

- Weighted cross entropy is a variant of CE where all positive examples get weighted by some coefficients.
- It is used in the case of class imbalance.
- For example, when you have an image with 10% black pixels and 90% white pixels, regular CE won't work very well.
- It is defined as follows

$$\text{WCE}(p, \hat{p}) = -(\beta p \log(\hat{p}) + (1 - p) \log(1 - \hat{p}))$$

Balanced cross entropy

- Balanced cross entropy (BCE) is similar to WCE. The only difference is that we weight also the negative examples.
- BCE can be defined as follows:
- $$\text{BCE}(p, \hat{p}) = -(\beta p \log(\hat{p}) + (1 - \beta)(1 - p) \log(1 - \hat{p}))$$

Focal loss

- Easily classified negatives comprise the majority of the loss and dominate the gradient. While α balances the importance of positive/negative examples, it does not differentiate between **easy/hard** examples.
- Focal loss (FL) tries to down-weight the contribution of easy examples, so that the CNN focuses more on hard examples.
- FL can be defined as follows:

$$\text{FL}(p, \hat{p}) = -(\alpha(1 - \hat{p})^\gamma p \log(\hat{p}) + (1 - \alpha)\hat{p}^\gamma(1 - p) \log(1 - \hat{p}))$$

Dice Loss / F1 score

- The Dice coefficient is similar to the Jaccard Index (Intersection over Union, IoU):

$$\text{DC} = \frac{2TP}{2TP + FP + FN} = \frac{2|X \cap Y|}{|X| + |Y|}$$

$$\text{IoU} = \frac{TP}{TP + FP + FN} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

- The dice coefficient can also be defined as a loss function:

$$\text{DL}(p, \hat{p}) = 1 - \frac{2p\hat{p} + 1}{p + \hat{p} + 1}$$

- where $p \in \{0, 1\}$ and $0 \leq \hat{p} \leq 1$.

Dice

■ Example

$$|A \cap B| = \begin{bmatrix} 0.01 & 0.03 & 0.02 & 0.02 \\ 0.05 & 0.12 & 0.09 & 0.07 \\ 0.89 & 0.85 & 0.88 & 0.91 \\ 0.99 & 0.97 & 0.95 & 0.97 \end{bmatrix} * \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \xrightarrow{\text{element-wise multiply}} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0.89 & 0.85 & 0.88 & 0.91 \\ 0.99 & 0.97 & 0.95 & 0.97 \end{bmatrix} \xrightarrow{\text{sum}} 7.41$$

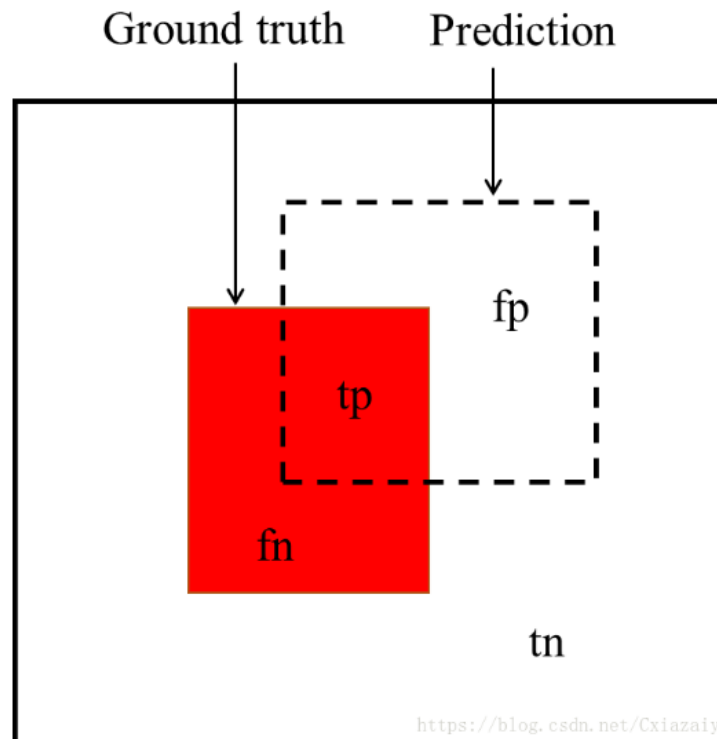
prediction target

$$|A| = \begin{bmatrix} 0.01 & 0.03 & 0.02 & 0.02 \\ 0.05 & 0.12 & 0.09 & 0.07 \\ 0.89 & 0.85 & 0.88 & 0.91 \\ 0.99 & 0.97 & 0.95 & 0.97 \end{bmatrix}^{2 \text{ (optional)}} \xrightarrow{\text{sum}} 7.82$$

$$|B| = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}^{2 \text{ (optional)}} \xrightarrow{\text{sum}} 8$$

Tversky loss

- Tversky index (TI) is a generalization of Dice's coefficient. TI adds a weight to FP (false positives) and FN (false negatives).
- To give FNs higher weights than FPs in training our network for highly imbalanced data, where detecting small lesions is crucial.



Tversky loss

- Tversky index (TI) is a generalization of Dice's coefficient. TI adds a weight to FP (false positives) and FN (false negatives).

$$\text{TI}(p, \hat{p}) = \frac{p\hat{p}}{p\hat{p} + \beta(1-p)\hat{p} + (1-\beta)p(1-\hat{p})}$$

- Let $\beta = \frac{1}{2}$
- which is just the regular Dice coefficient.

$$\text{TI}(p, \hat{p}) = \frac{2p\hat{p}}{2p\hat{p} + (1-p)\hat{p} + p(1-\hat{p})} = \frac{2p\hat{p}}{\hat{p} + p}$$