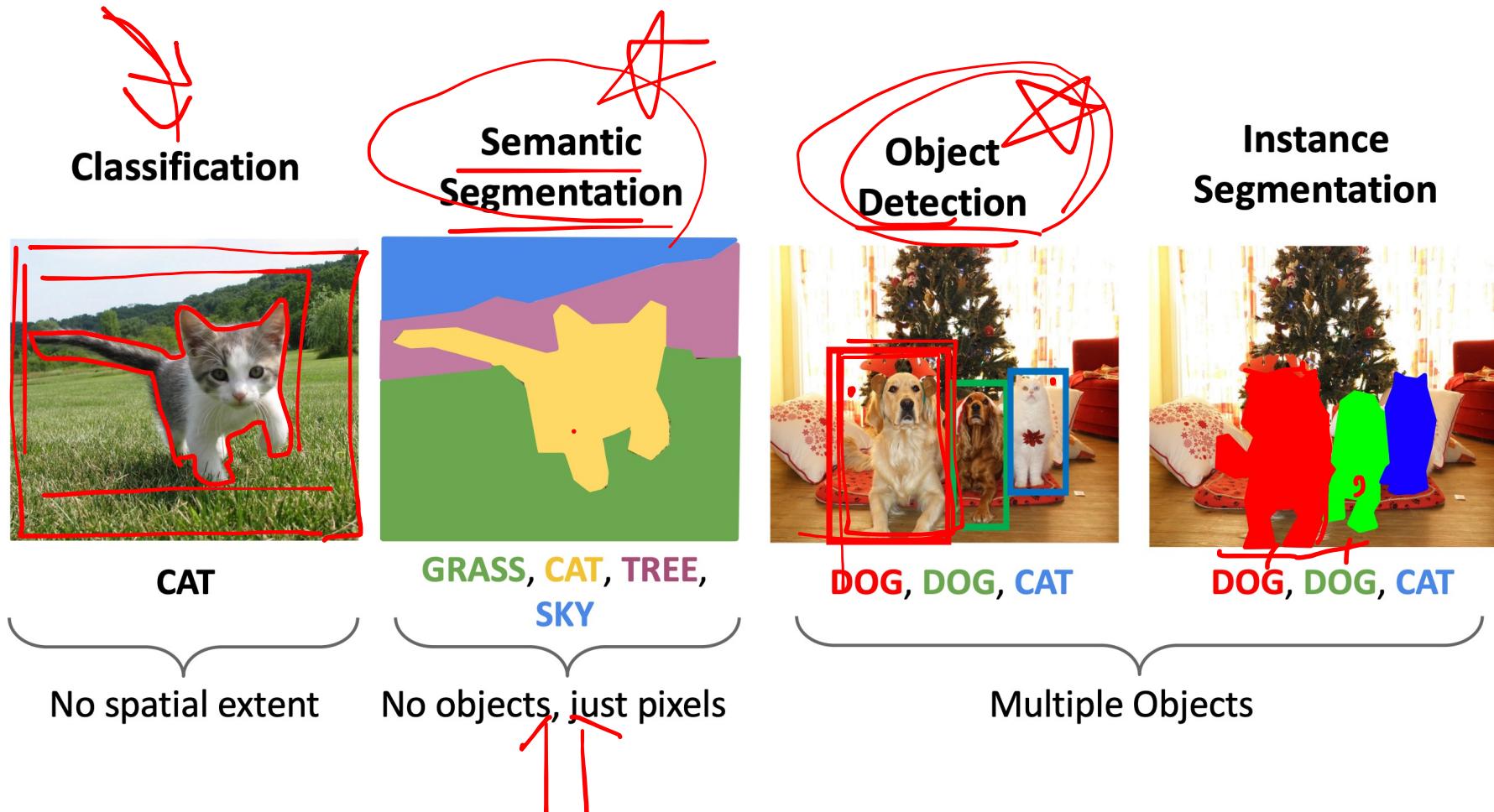


Computer Vision: Image segmentation

Siheng Chen 陈思衡

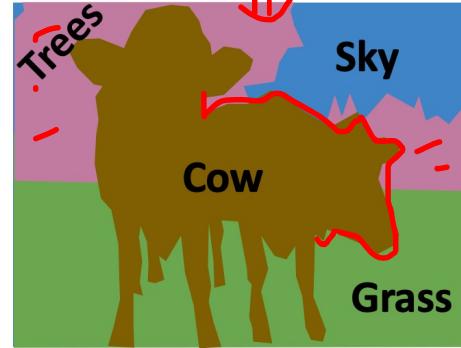
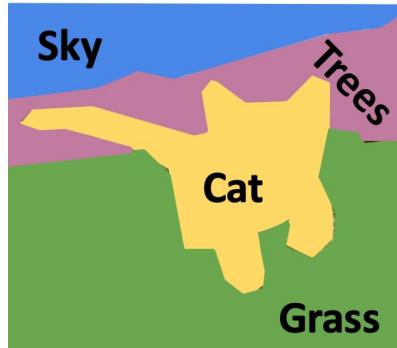
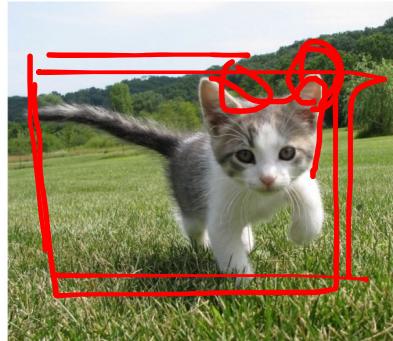
Common computer vision tasks



Semantic segmentation

Image segmentation is the process of partitioning a digital image into multiple image segments, also known as image regions or image objects (sets of pixels)

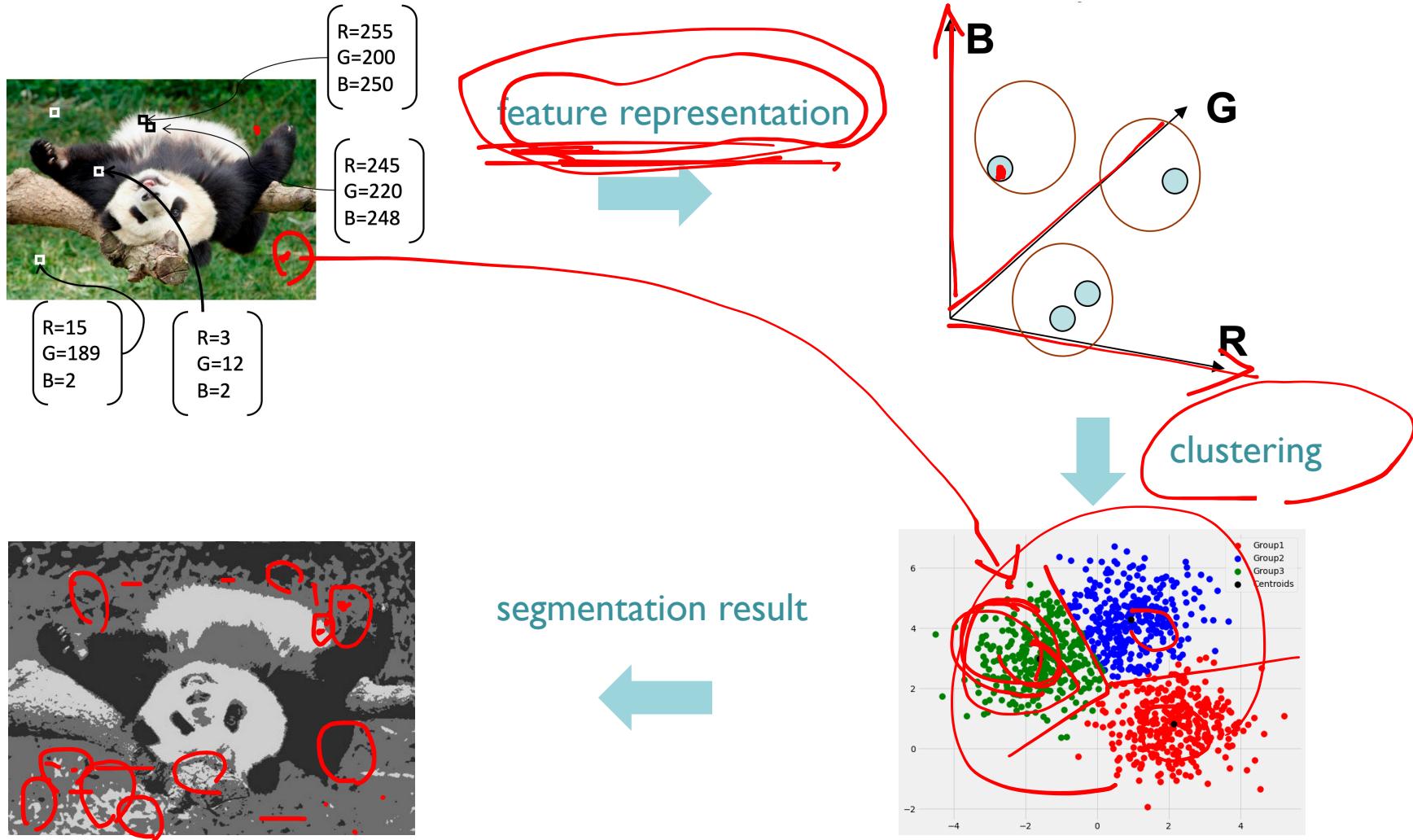
- Label each pixel in the image with a category label
- Semantic: Don't differentiate instances, only care about pixels



Semantic segmentation

- Segmentation as clustering
- Segmentation as energy minimization
- Pixel-wise end2end learning

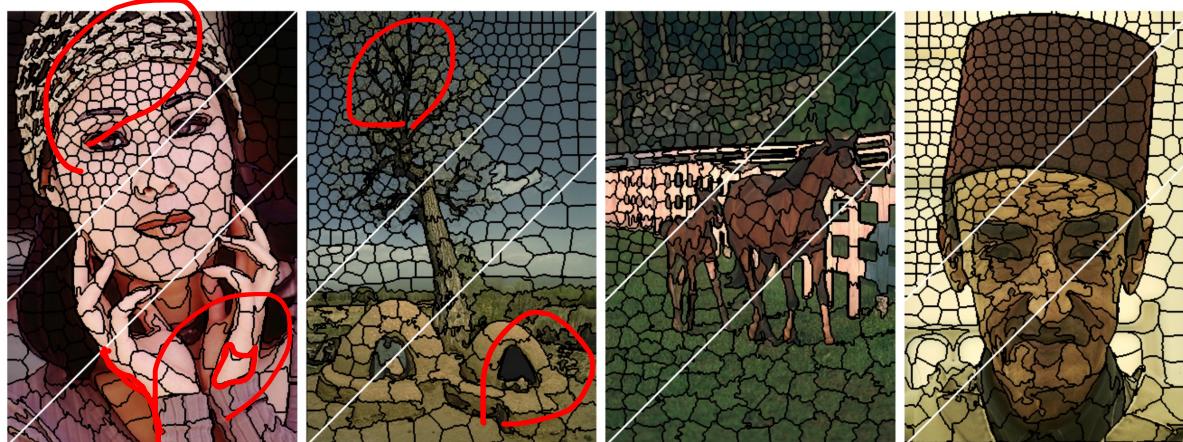
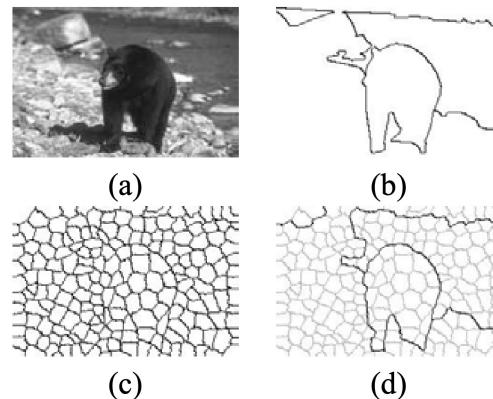
Segmentation as clustering



Segmentation as clustering

Superpixel

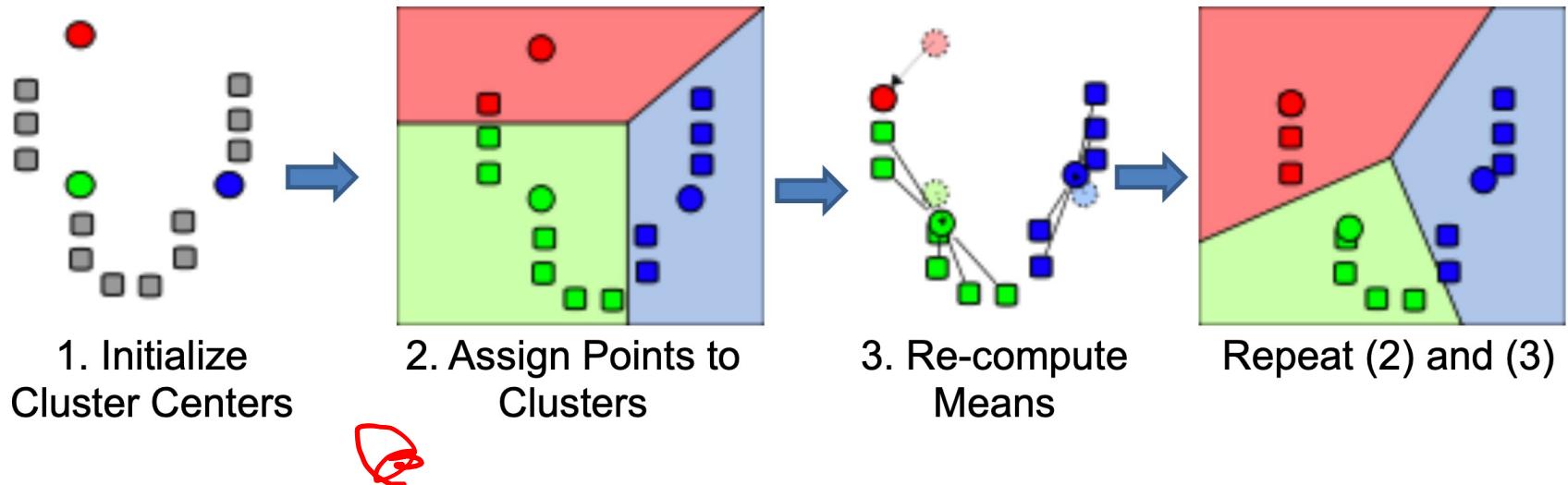
In this section we present a preprocessing stage to group pixels into “superpixels”. The motivations of this preliminary grouping are: (1) pixels are not natural entities; they are merely a consequence of the discrete representation of images; and (2) the number of pixels is high even at moderate resolutions; this makes optimization on the level of pixels intractable. We would like to work with “superpixels” which are local, coherent, and which preserve most of the structure necessary for segmentation at the scale of interest.



Segmentation as clustering

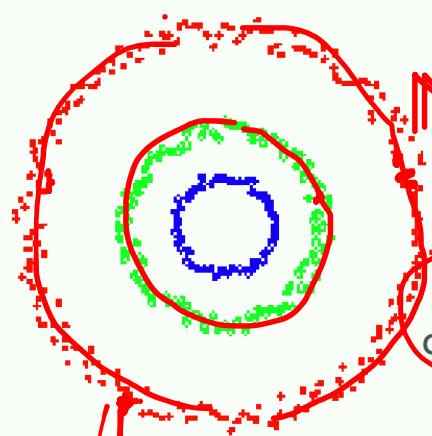
K-means clustering

$$\mathbf{c}^*, \delta^* = \arg \min_{\mathbf{c}, \delta} \sum_{x_i} \sum_{c_j} \delta_{i,j} (x_i - c_j)^2 \quad \delta \in \{0, 1\}^{N \times K}, \delta \mathbf{1} = 1$$

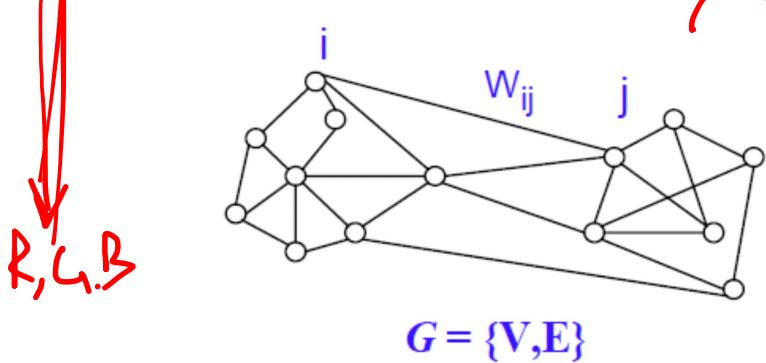
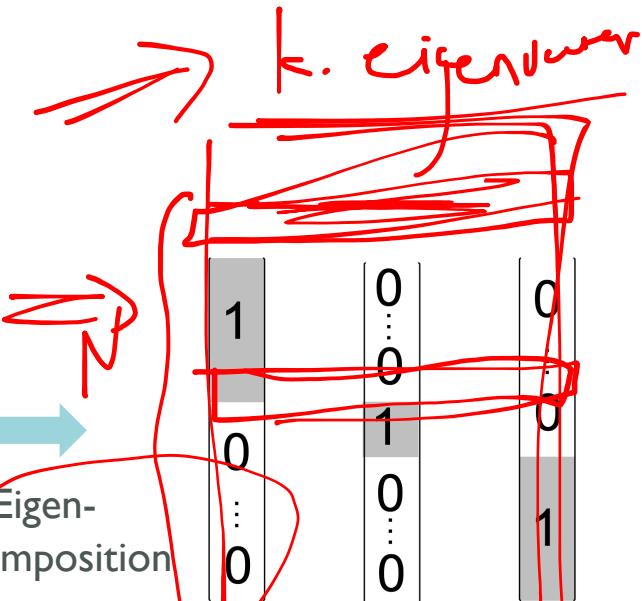


Segmentation as clustering

Spectral clustering



$$L = \begin{bmatrix} L_1 & & & \\ & 0 & & \\ & & L_2 & \\ & & & 0 \\ & & & & L_3 \end{bmatrix}$$



$$\text{minimize } N\text{cut}(A, B) = \frac{\mathbf{f}^T \mathbf{L} \mathbf{f}}{\mathbf{f}^T \mathbf{D} \mathbf{f}}$$

subject to $\mathbf{f}^T \mathbf{D} \mathbf{1} = 0$

$$\mathbf{L}\mathbf{f} = \lambda \mathbf{D}\mathbf{f}$$

Semantic segmentation

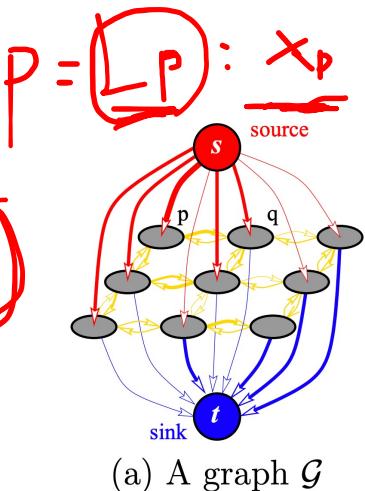
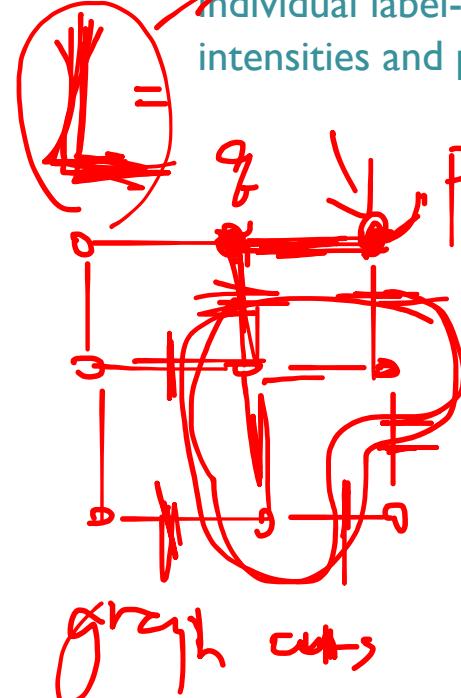
- Segmentation as clustering
- ~~Segmentation as energy minimization~~
- Pixel-wise end2end learning

Segmentation as energy minimization

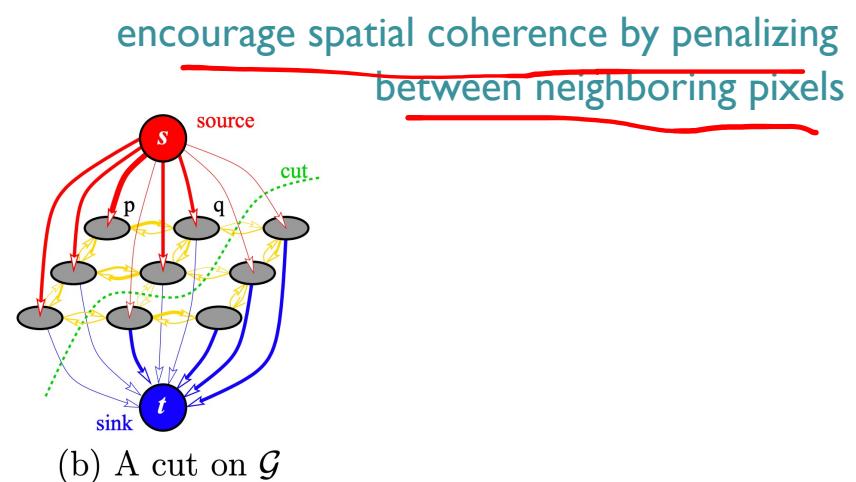
Min-cut/Max-flow algorithms on graphs

$$E(L) = \sum_{p \in \mathcal{P}} D_p(L_p) + \sum_{(p,q) \in \mathcal{N}} V_{p,q}(L_p, L_q)$$

Individual label-preferences of pixels based on observed intensities and pre-specified likelihood function



(a) A graph \mathcal{G}



(b) A cut on \mathcal{G}

graph cuts

Segmentation as energy minimization

Conditional random field (CRF)

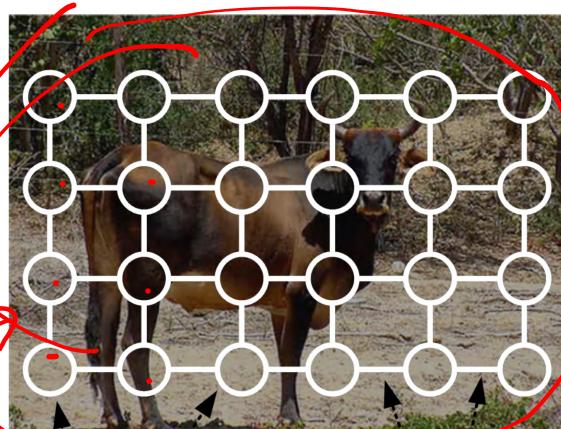
S 7 A

2-11-2014

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left\{ \sum_{k=1}^K \lambda_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}$$

Labels

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \exp \left\{ \sum_{k=1}^K \lambda_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}$$



Unary Term

Pairwise Term

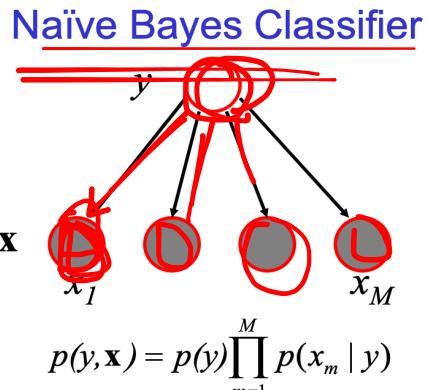


Segmentation Results

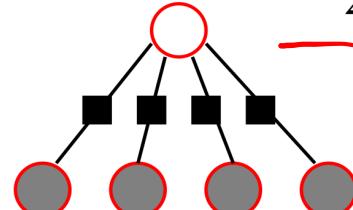
Segmentation as energy minimization

Conditional random field (CRF)

GENERATIVE

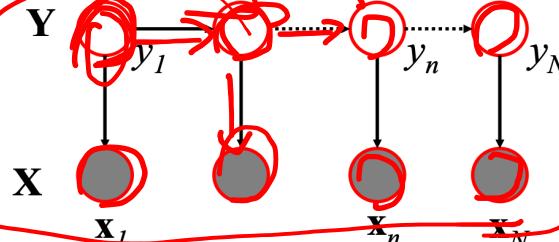


DISCRIMINATIVE



$$p(y | \mathbf{x}) = \frac{\exp \left\{ \sum_{m=1}^M \lambda_m f_m(y, \mathbf{x}) \right\}}{\sum_{y'} \exp \left\{ \sum_{m=1}^M \lambda_m f_m(y', \mathbf{x}) \right\}}$$

Hidden Markov Model



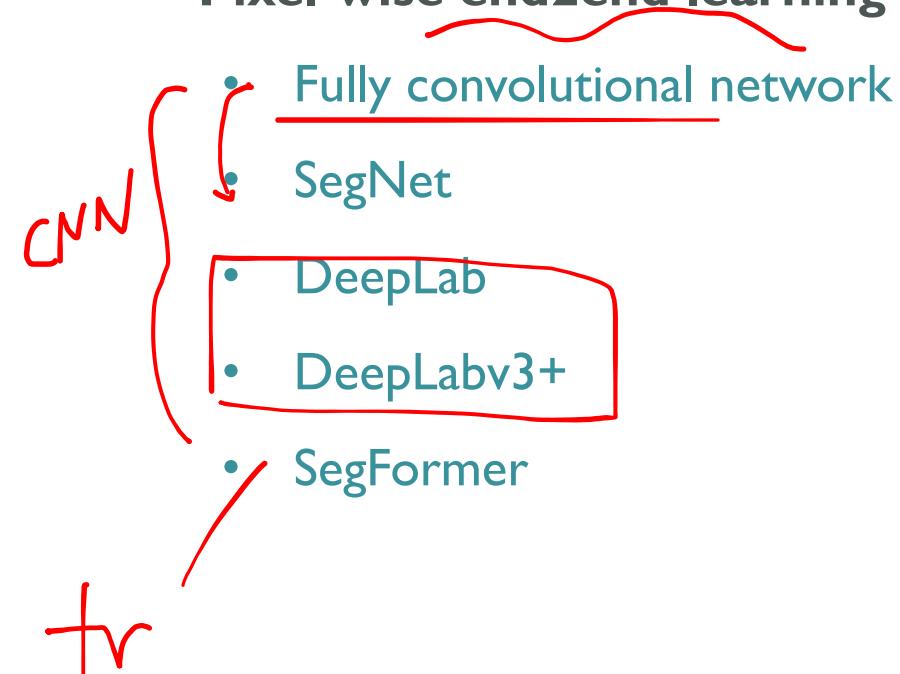
$$p(\mathbf{Y}, \mathbf{X}) = \prod_{n=1}^N p(y_n | y_{n-1}) p(\mathbf{x}_n | y_n)$$

Conditional Random Field

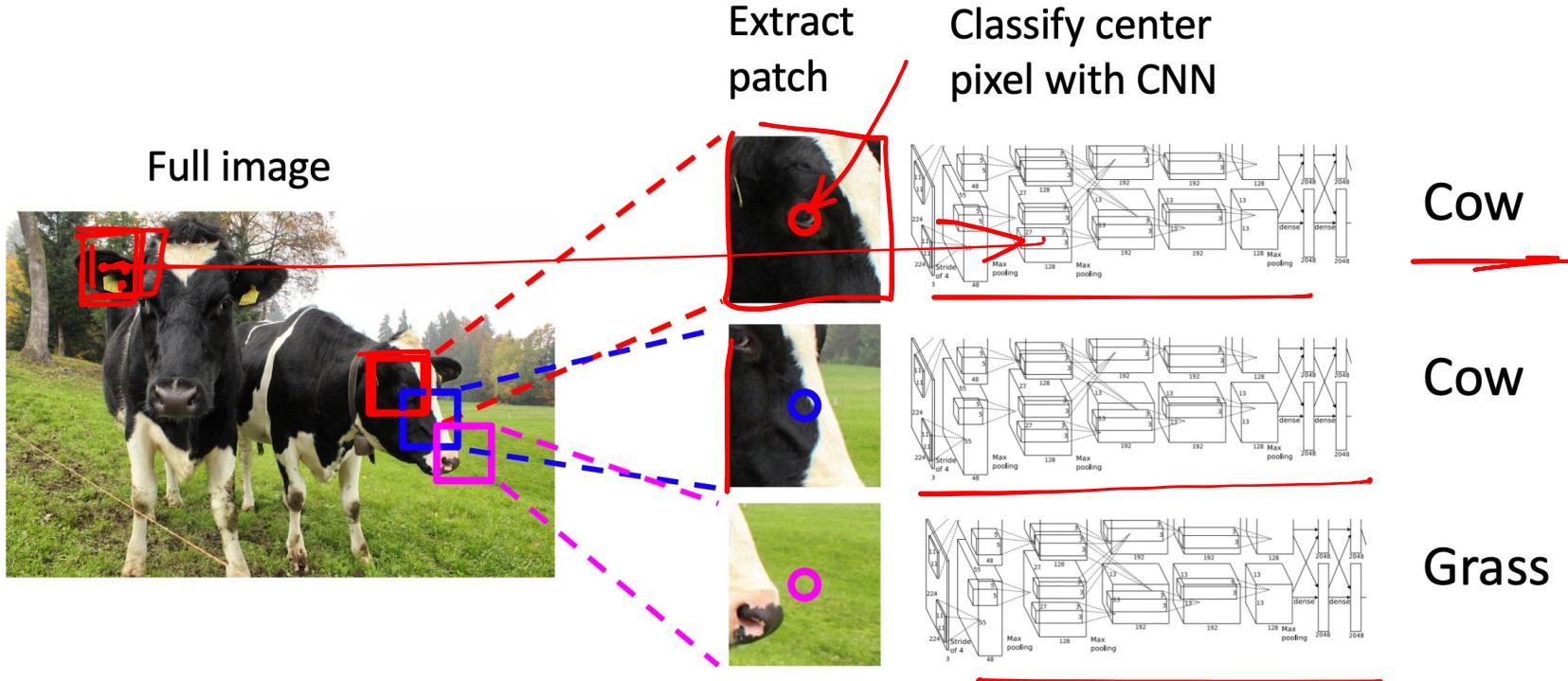
$$p(\mathbf{Y} | \mathbf{X}) = \frac{\exp \left\{ \sum_{m=1}^M \lambda_m f_m(y_n, y_{n-1}, \mathbf{x}_n) \right\}}{\sum_{y'} \exp \left\{ \sum_{m=1}^M \lambda_m f_m(y'_n, y'_{n-1}, \mathbf{x}_n) \right\}}$$

Semantic segmentation

- Segmentation as clustering
- Graph-based segmentation
- Segmentation as energy minimization
- Pixel-wise end2end learning



Semantic segmentation



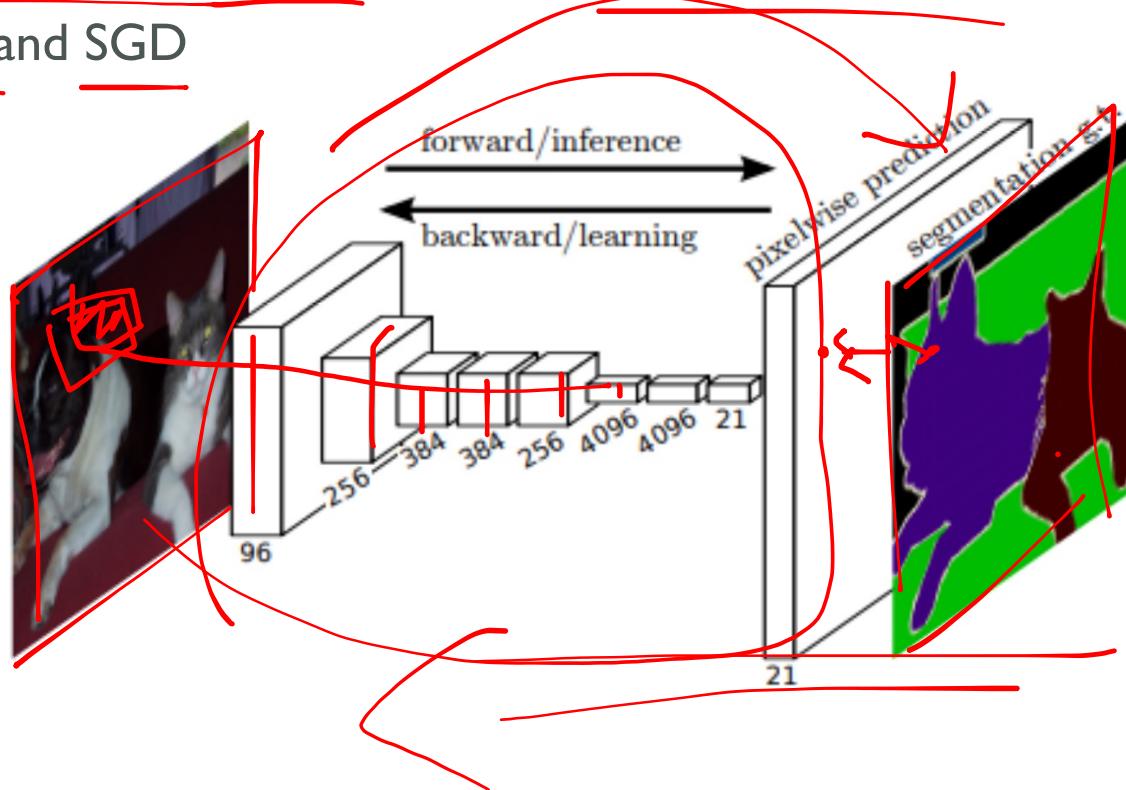
Problem:

Very inefficient! Not reusing shared features between overlapping patches

Fully convolutional network

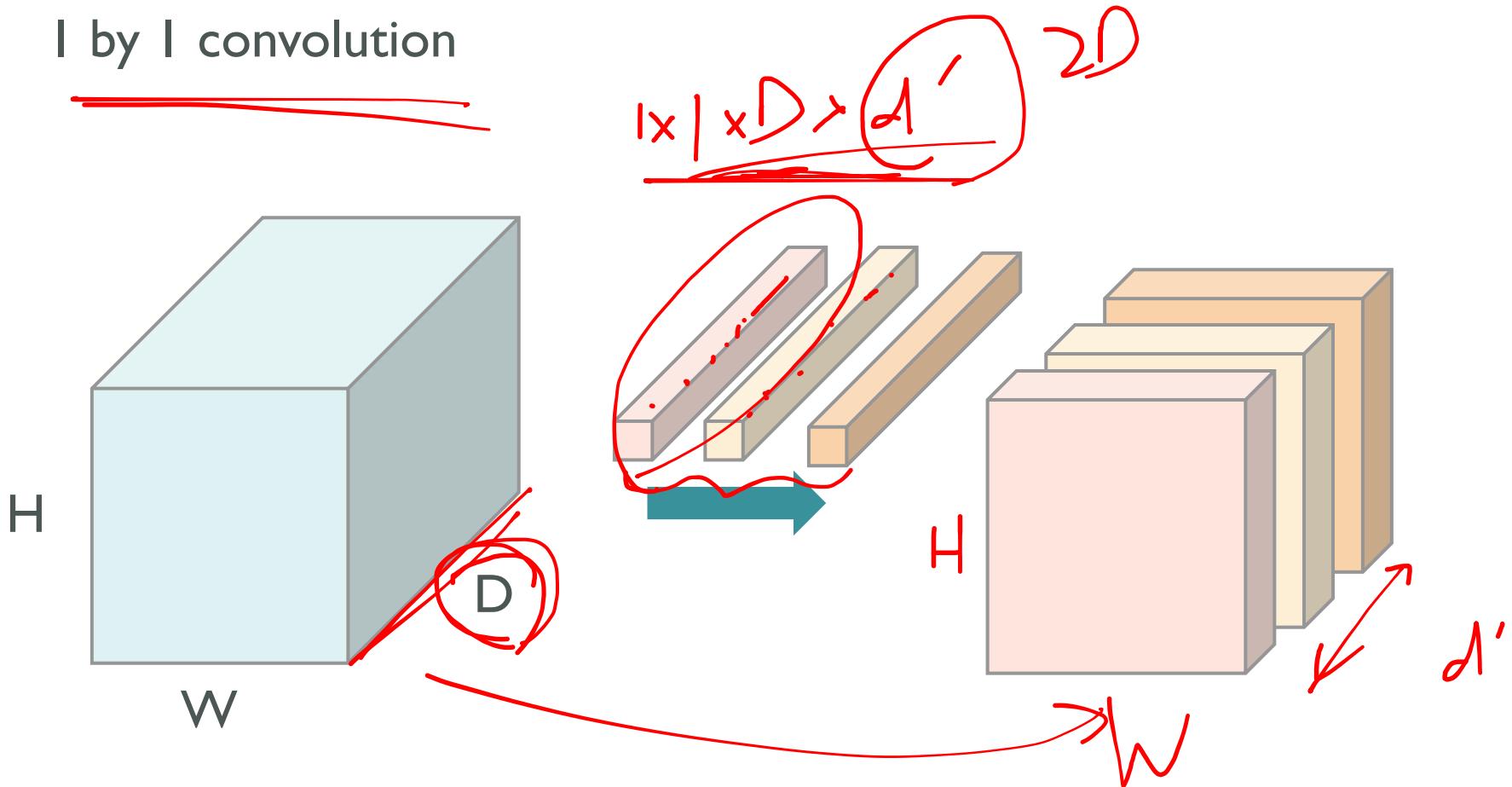
Key ideas

- All convolutional layers.
- Convert the last layer output to the original resolution.
- Softmax-cross entropy loss between pixel-wise predictions and ground truth
- Backprop and SGD



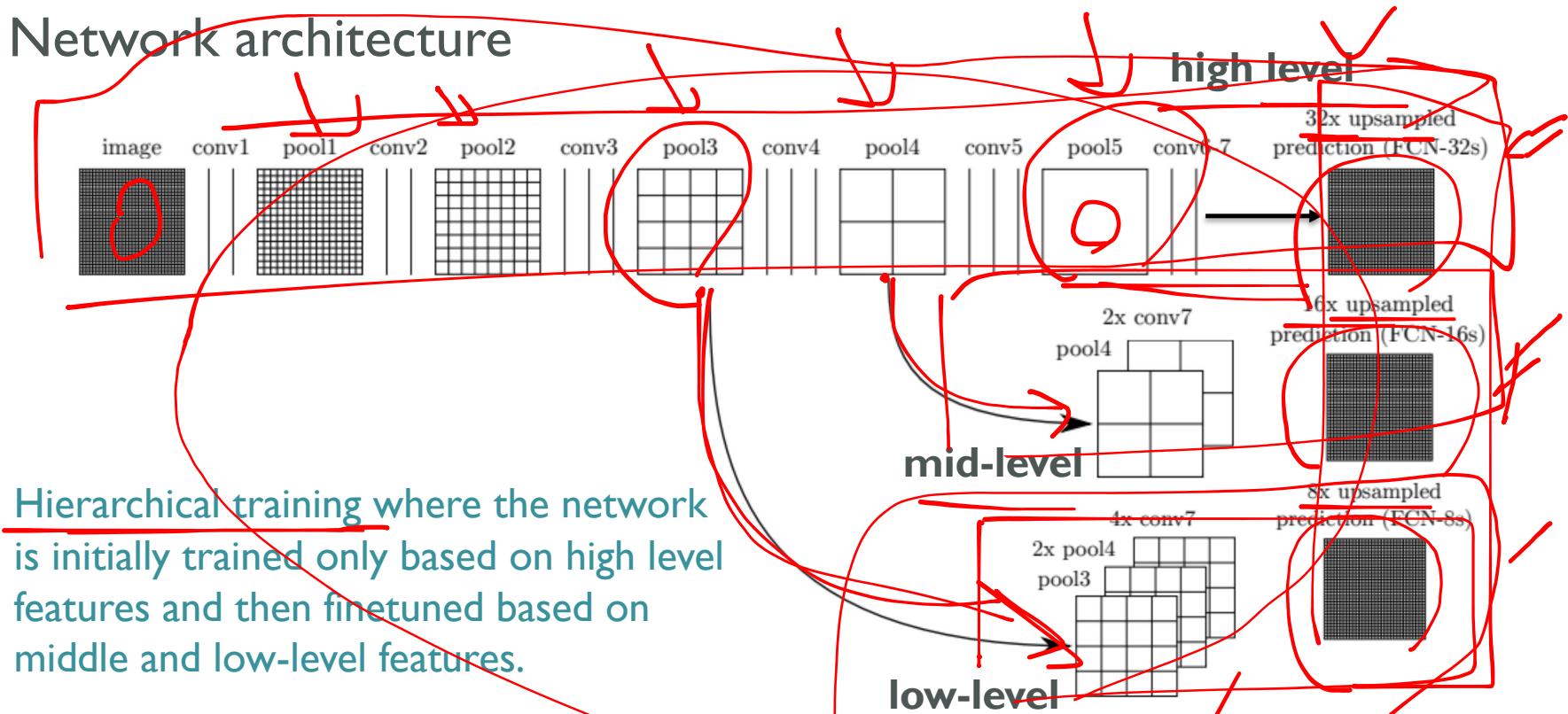
Fully convolutional network

I by I convolution



Fully convolutional network

Network architecture

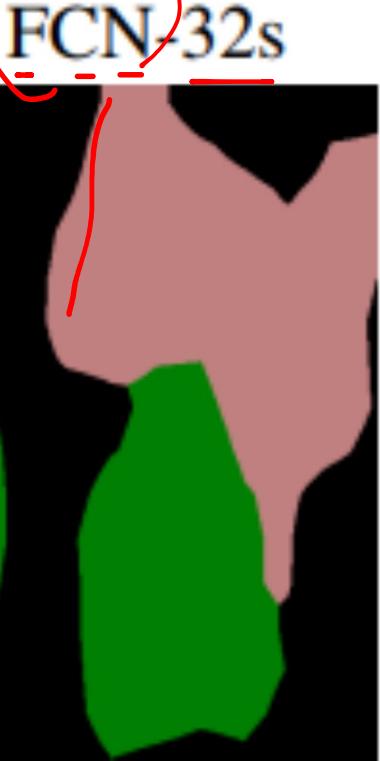


Hierarchical training where the network is initially trained only based on high level features and then finetuned based on middle and low-level features.

This is important because it allows the network to also learn the mid and low-level details of the image, in addition to high level ones.

Fully convolutional network

FCN-32s



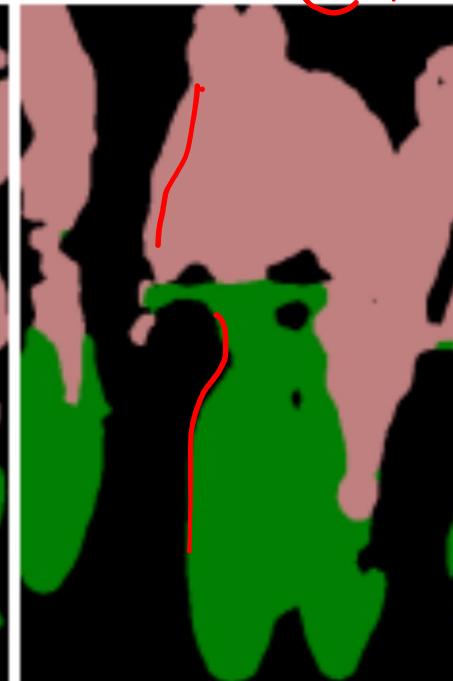
Good

FCN-16s



Better

FCN-8s



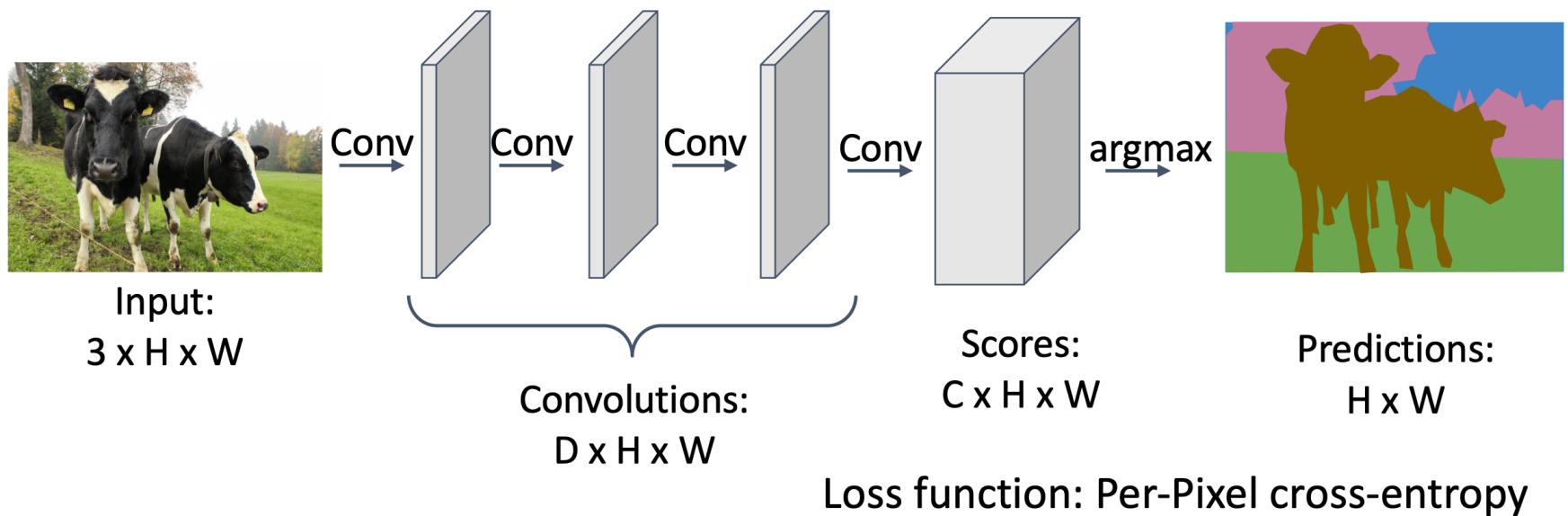
Even better

Ground truth



Fully convolutional network

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

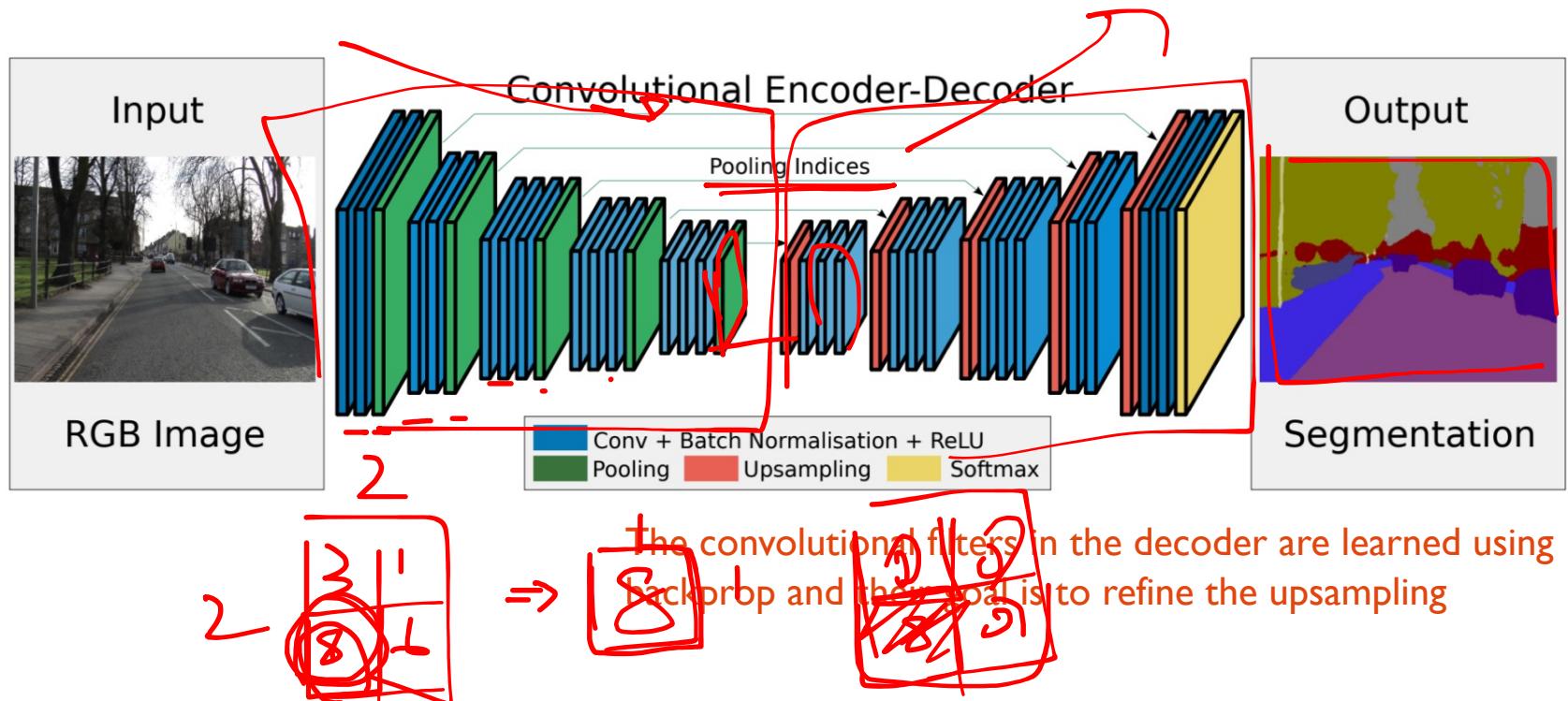


- { **Problem #1:** Effective receptive field size is linear in number of conv layers: With L 3×3 conv layers, receptive field is $1+2L$
- Problem #2:** Convolution on high res images is expensive!

SegNet

Design network as a bunch of convolutional layers, with pooling (downsampling) and upsampling inside the network!

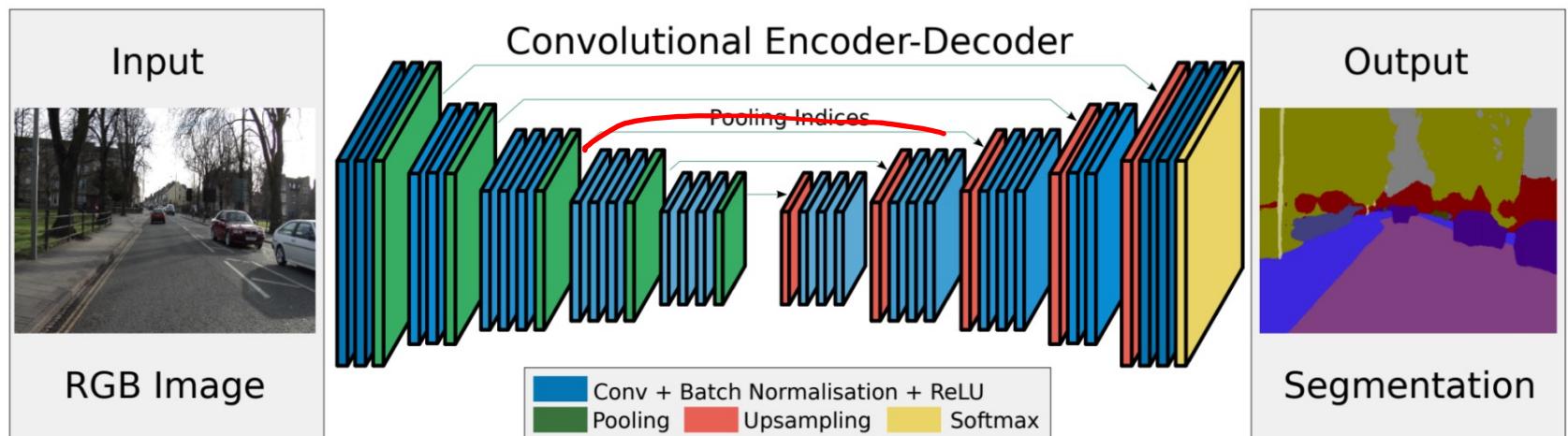
- **Encoder:** Normal convolutional filters + pooling
- **Decoder:** Upsampling + convolutional filters
- **Softmax layer:** The output of the soft-max classifier is a K channel image of probabilities where K is the number of classes



SegNet

Design network as a bunch of convolutional layers, with **pooling** (downsampling) and **upsampling** inside the network!

- **Encoder:** Normal convolutional filters + pooling
- **Decoder:** Upsampling + convolutional filters
- **Softmax layer:** The output of the soft-max classifier is a K channel image of probabilities where K is the number of classes

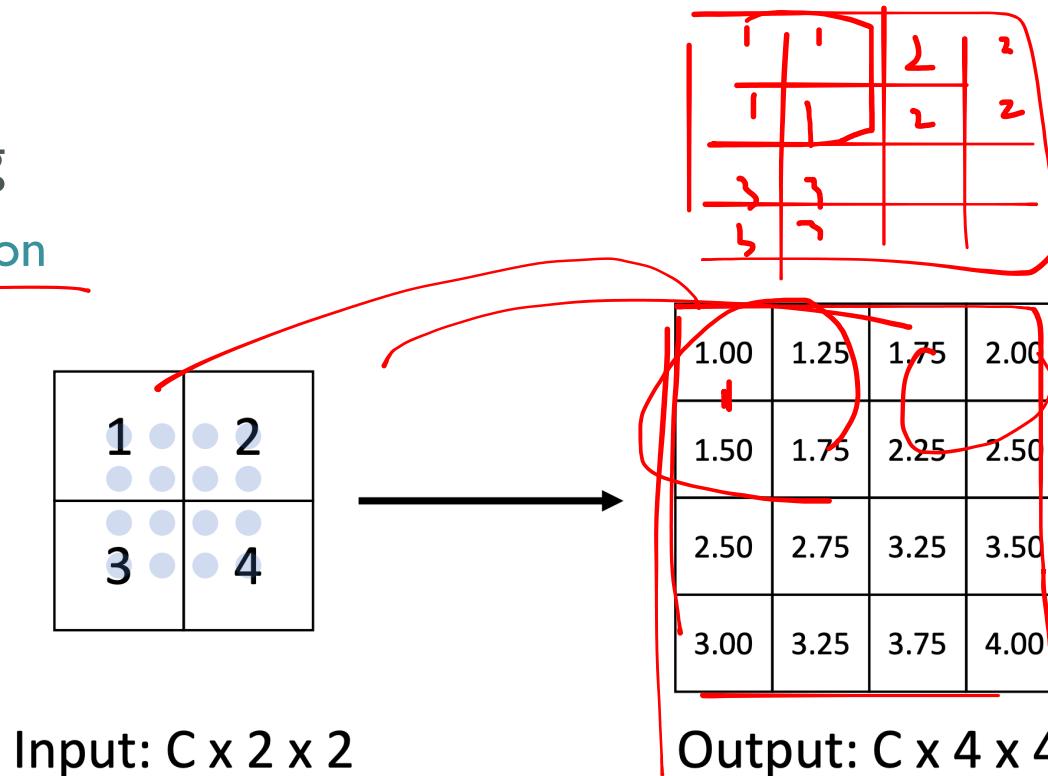


The convolutional filters in the decoder are learned using backprop and their goal is to refine the upsampling

SegNet

Upsampling

- Interpolation



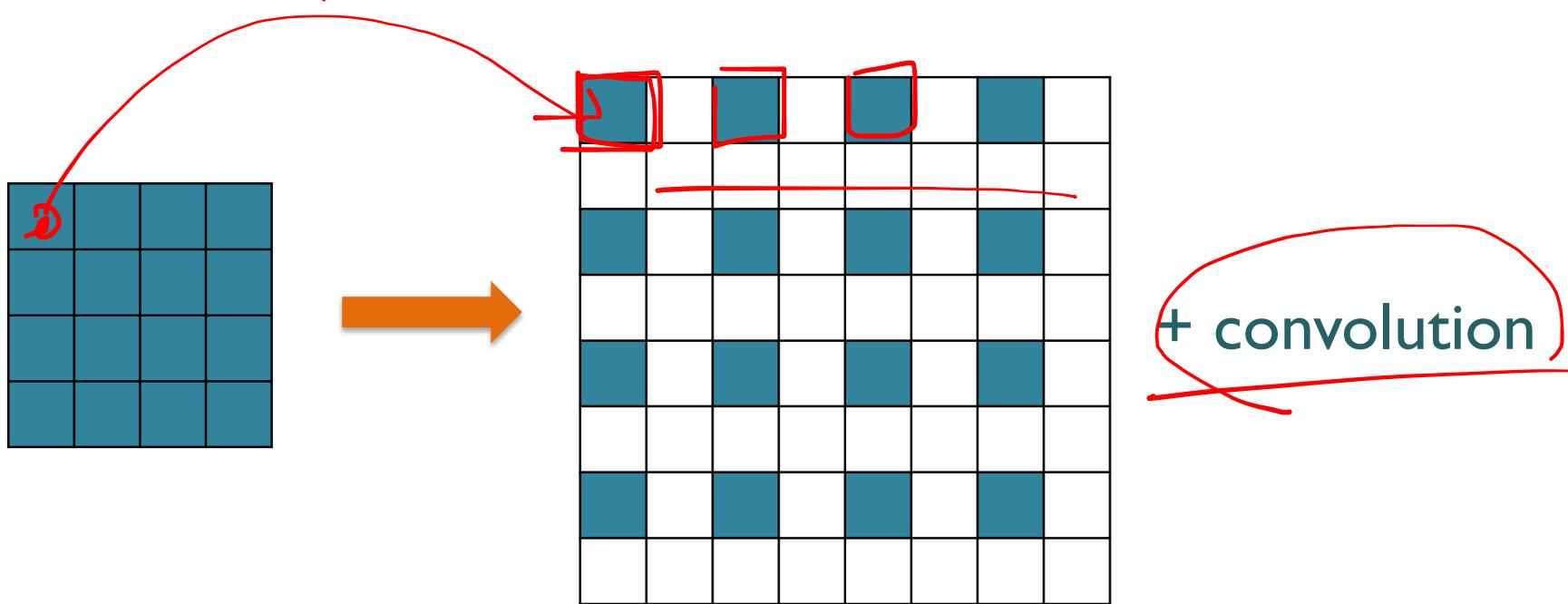
Use two closest neighbors in x and y to construct linear approximations

$$f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$$
$$j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$$

SegNet

Upsampling

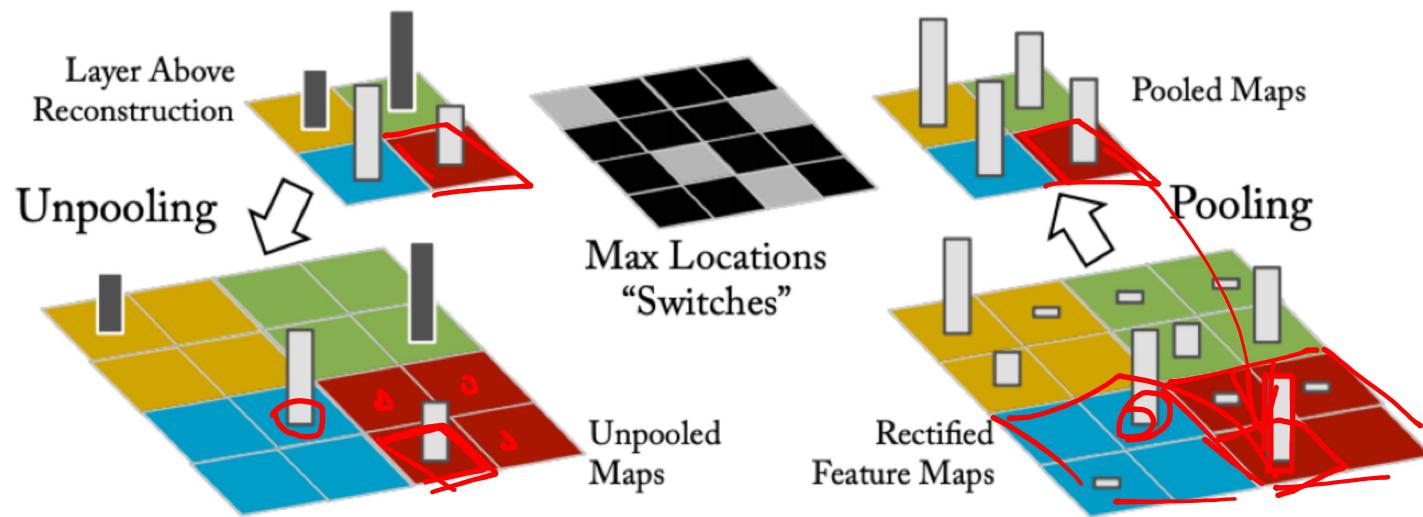
- Interpolation
- Fixed unpooling



SegNet

Upsampling

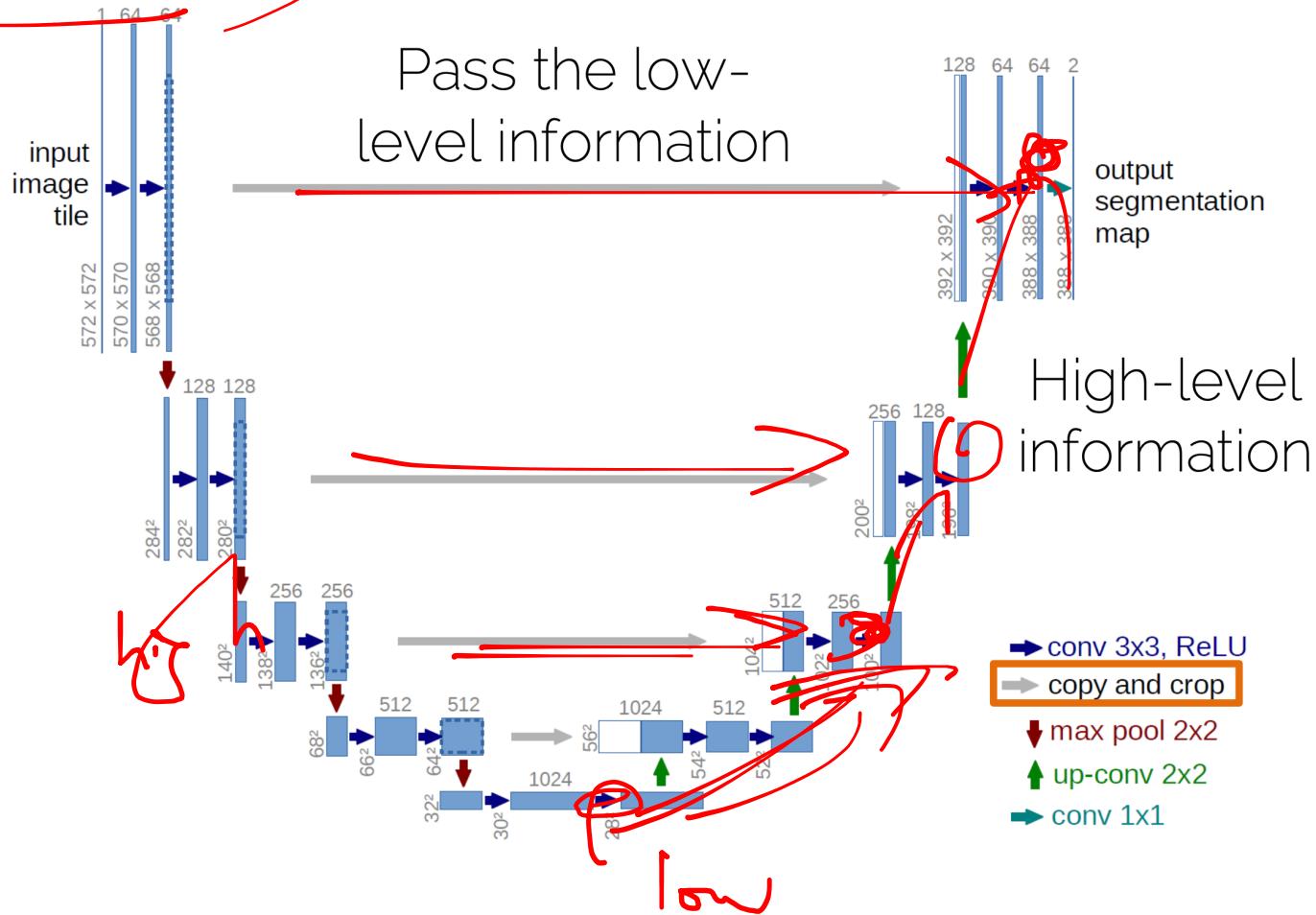
- Interpolation
- Fixed unpooling
- Unpooling in DeconvNet



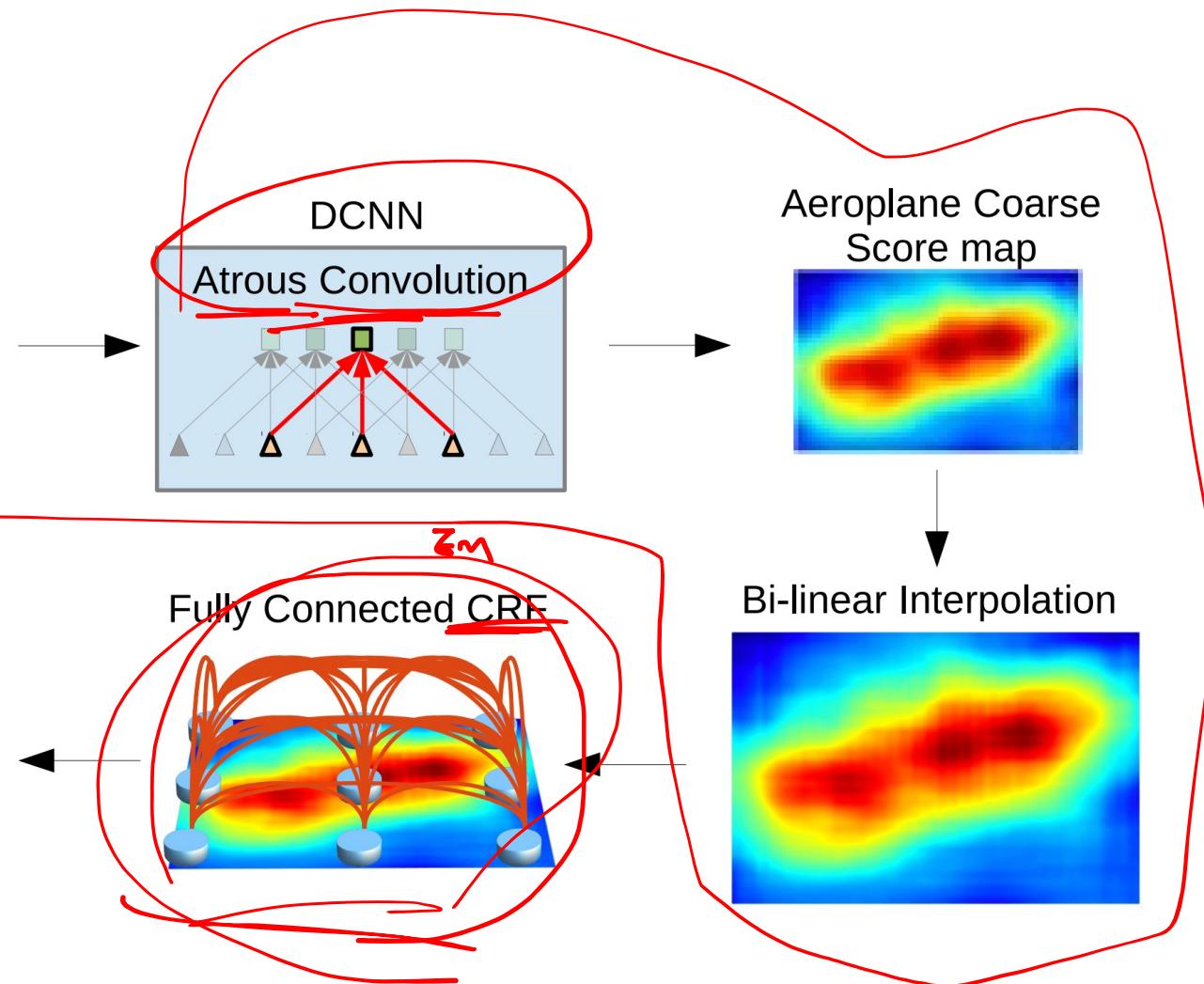
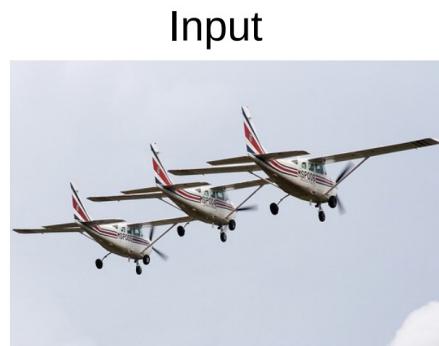
Keep the locations where the max came from

U-Net

Skip Connections



DeepLab



DeepLab

Reduced feature resolution

Solution: Atrous convolutions

Objects exist at multiple scales

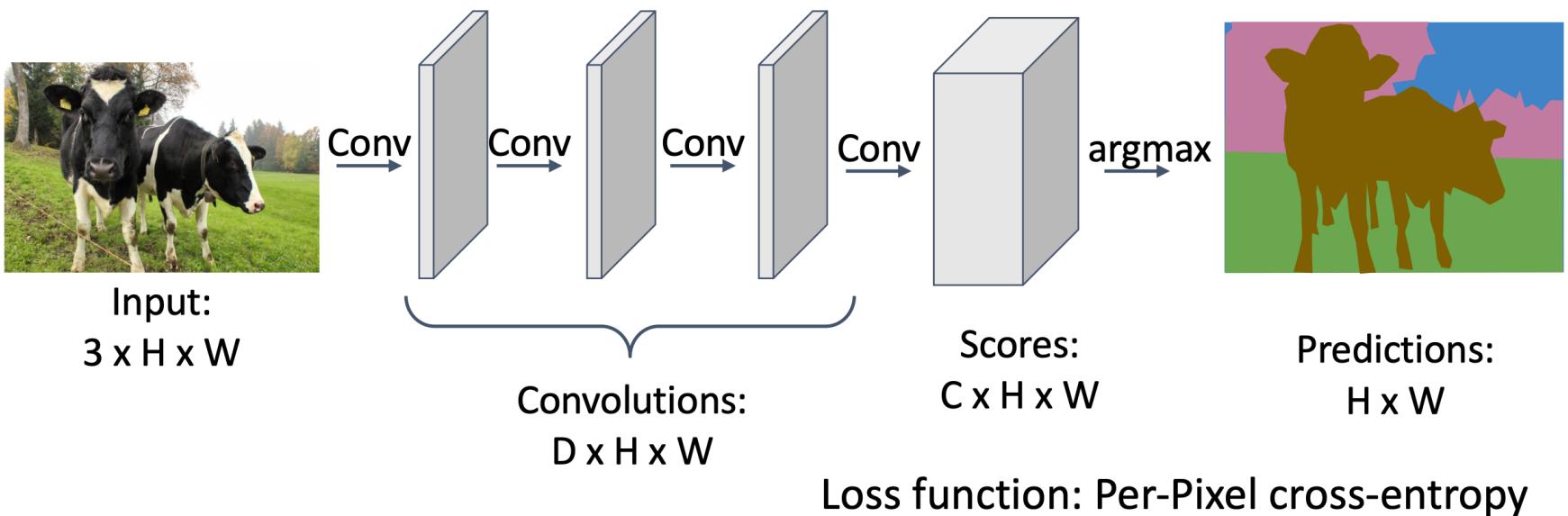
Solution: Pyramid pooling

Poor localization of the edges

Solution: Refinement with Conditional Random Field (CRF)

DeepLab

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Problem #1: Effective receptive field size is linear in number of conv layers: With L 3×3 conv layers, receptive field is $1 + 2L$



Problem #2: Convolution on high res images is expensive!

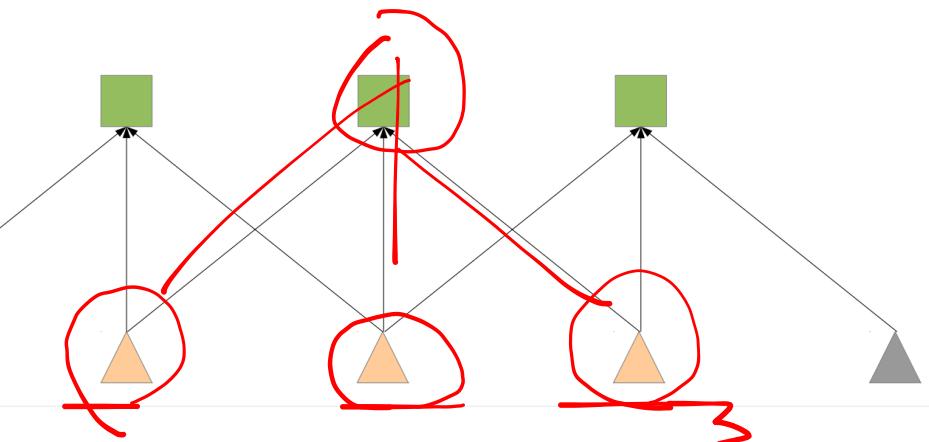
DeepLab

Dilated (atrous) convolutions ID

Output feature

Convolution
kernel = 3
stride = 1
pad = 1

Input feature

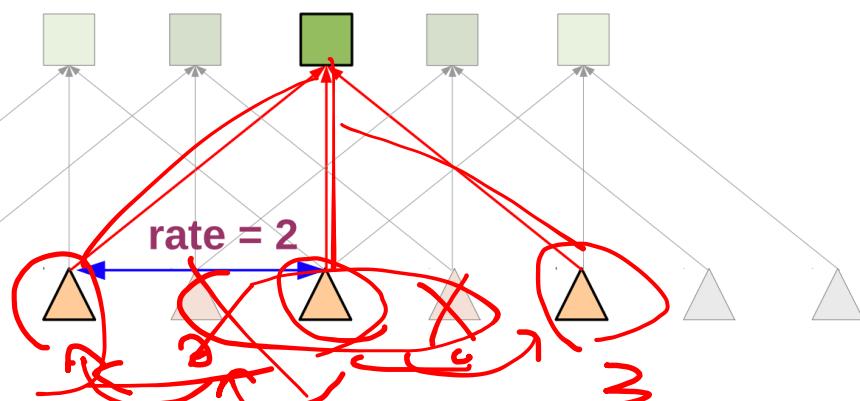


(a) Sparse feature extraction

Convolution

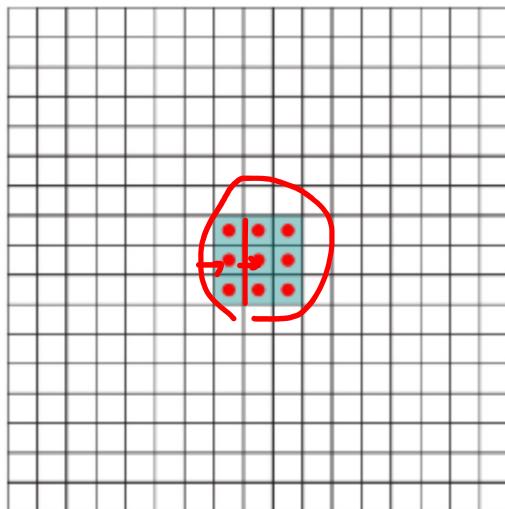
kernel = 3
stride = 1
pad = 2

rate = 2
~~insert 1 zero~~

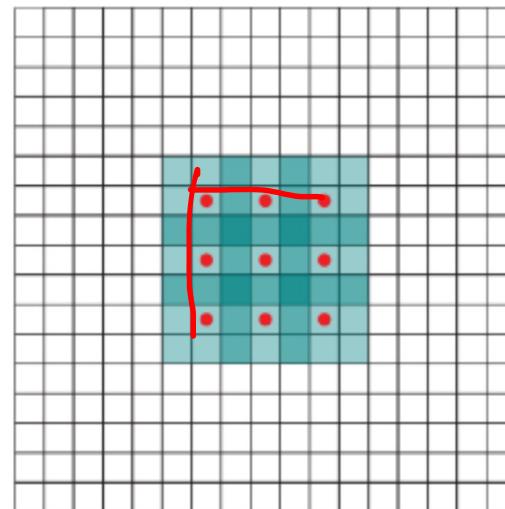


DeepLab

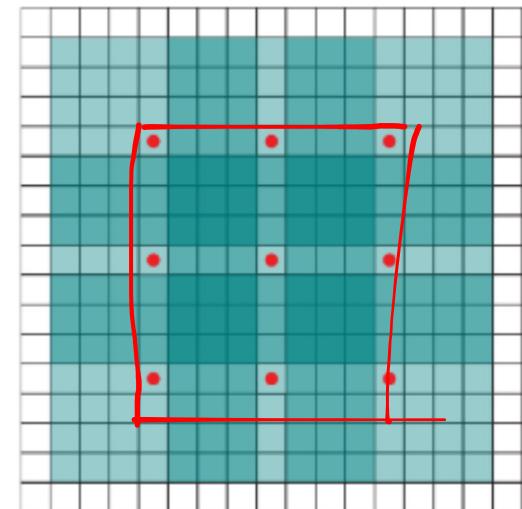
Dilated (atrous) convolutions 2D



(a)



(b)



(c)

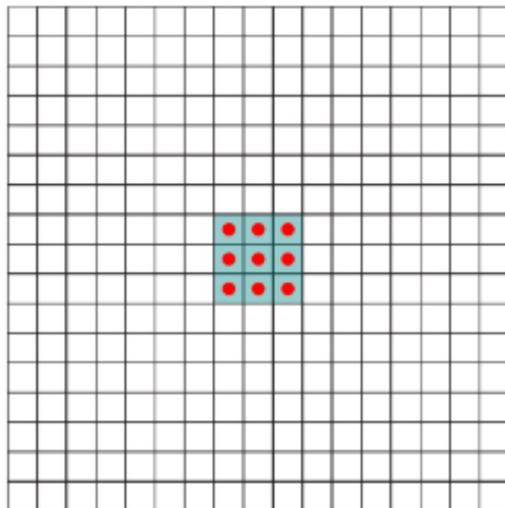
(a) the dilation parameter is 1, and each element produced by this filter has reception field of 3×3 .

(b) the dilation parameter is 2, and each element produced by it has reception field of 7×7 .

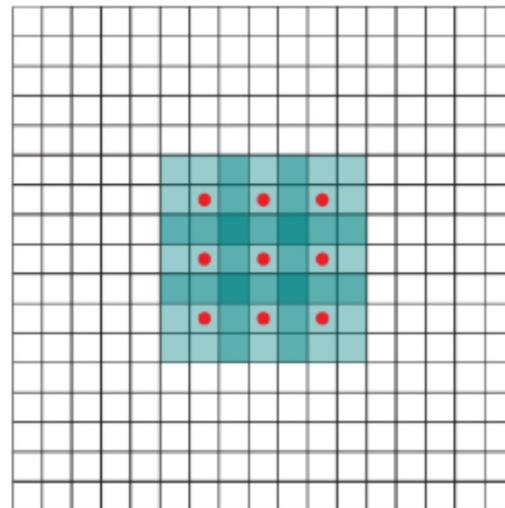
(c) the dilation parameter is 4, and each element produced by it has reception field of 15×15

DeepLab

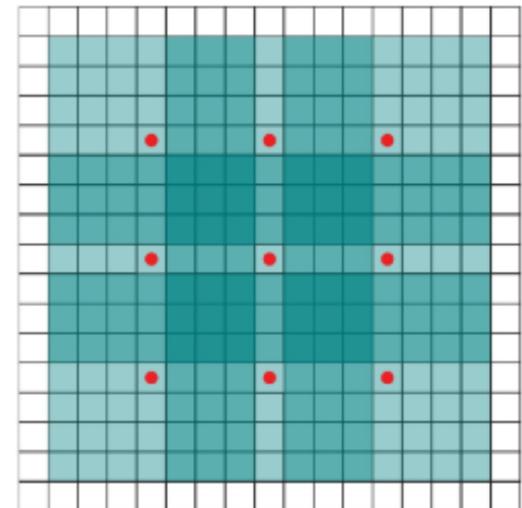
Dilated (atrous) convolutions 2D



(a)



(b)

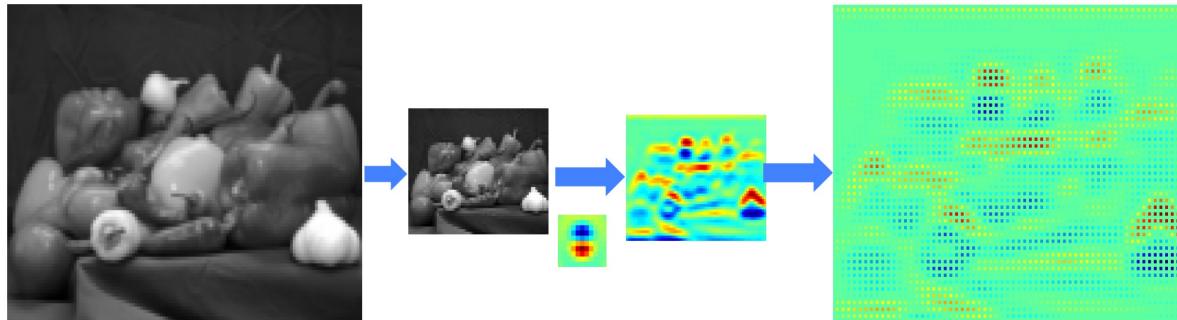


(c)

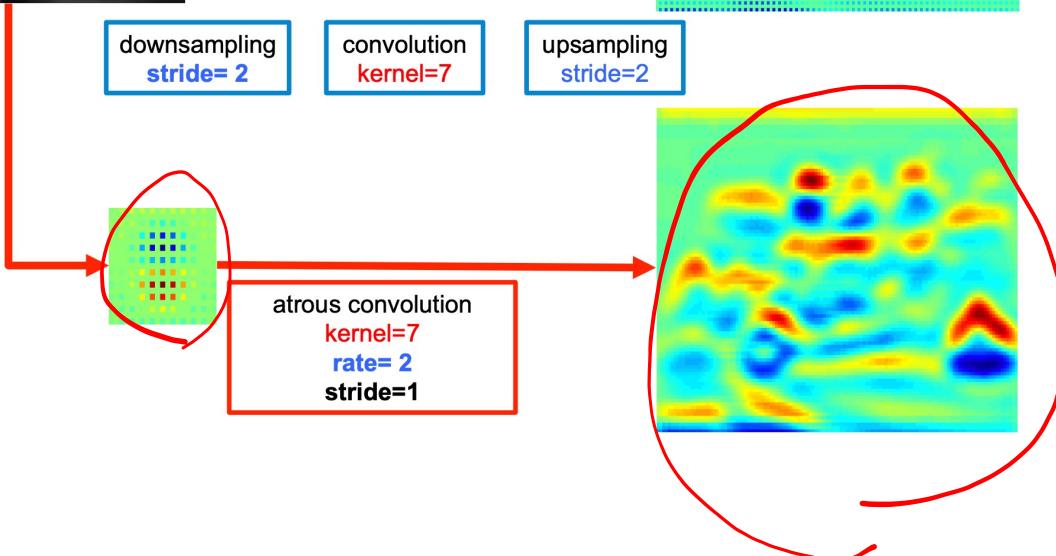
Each layer has the same number of parameters, but the receptive field grows exponentially while the number of parameters grows linearly.

DeepLab

Dilated (atrous) convolutions



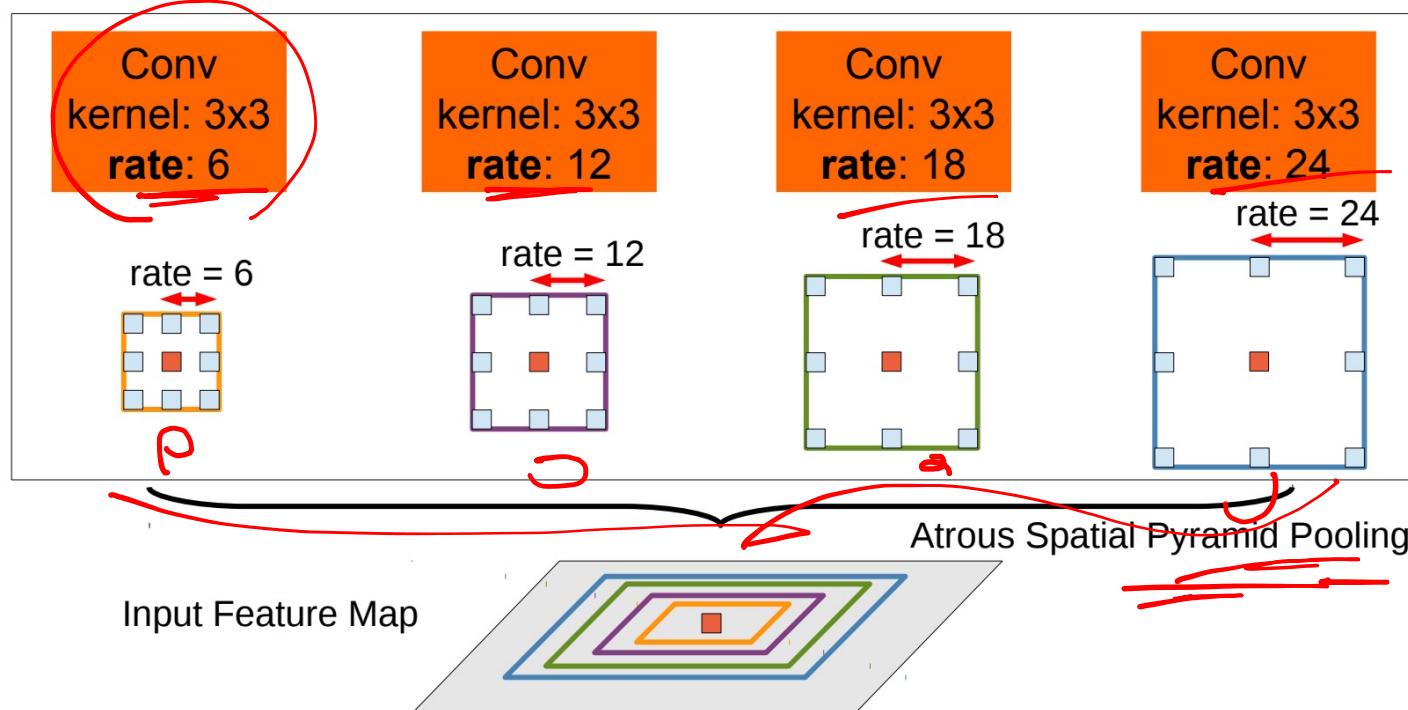
Sparse feature extraction
with standard convolution
on a low resolution input
feature map



Dense feature extraction
with atrous convolution with
rate r=2, applied on a high
resolution input feature map.

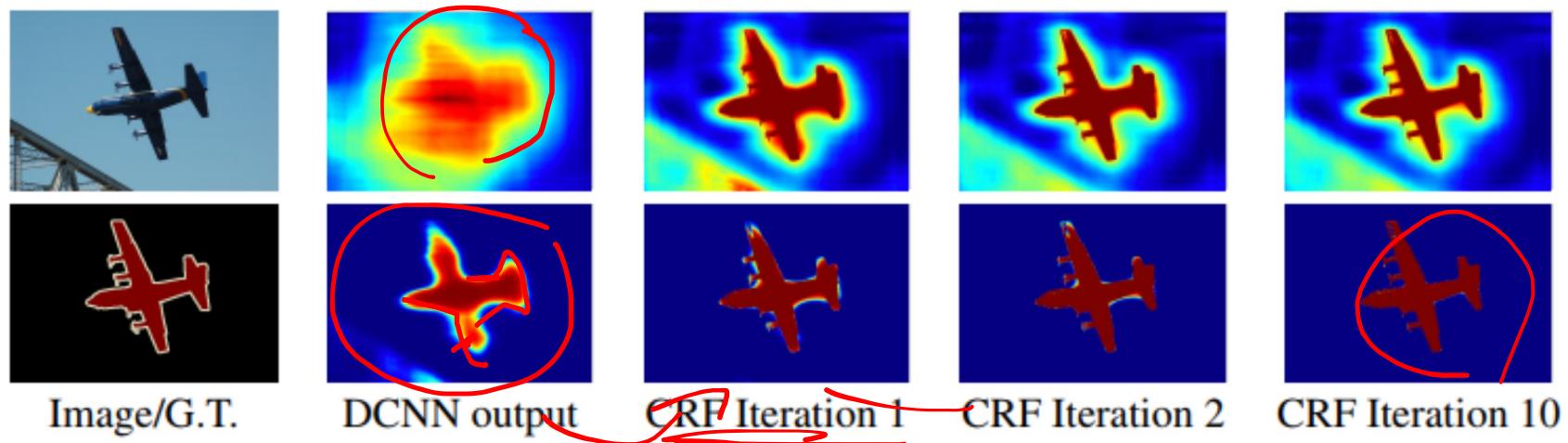
DeepLab

Atrous Spatial Pyramid Pooling



DeepLab

Conditional Random Field (CRF)



Score map (input before softmax function) and belief map (output of softmax function) for Aeroplane. The image shows the score (1st row) and belief (2nd row) maps after each mean field iteration. The output of last DCNN layer is used as input to the mean field inference.

DeepLab

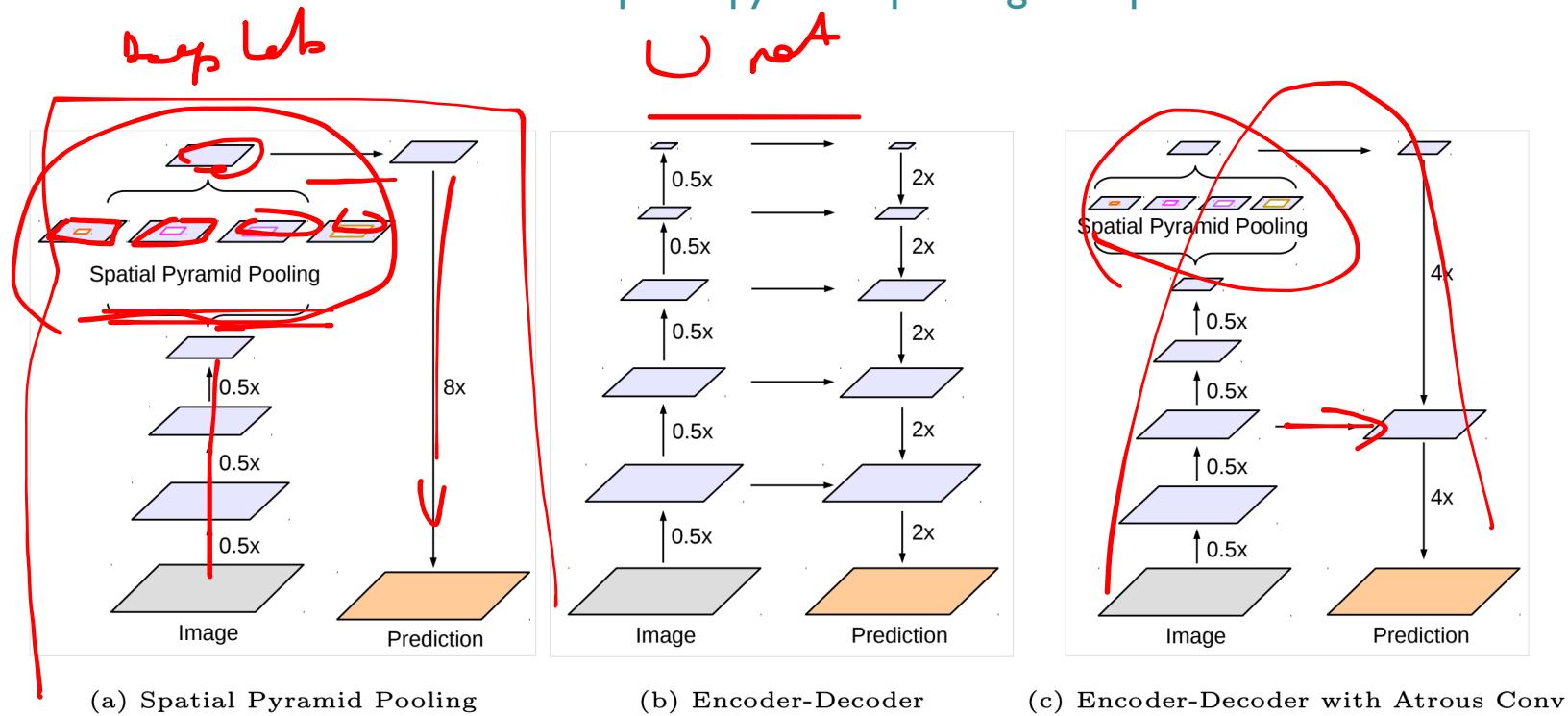


Fig. 6: PASCAL VOC 2012 *val* results. Input image and our DeepLab results before/after CRF.

DeepLabv3

Multiresolution

atrous convolutions + spatial pyramid pooling + skip connection

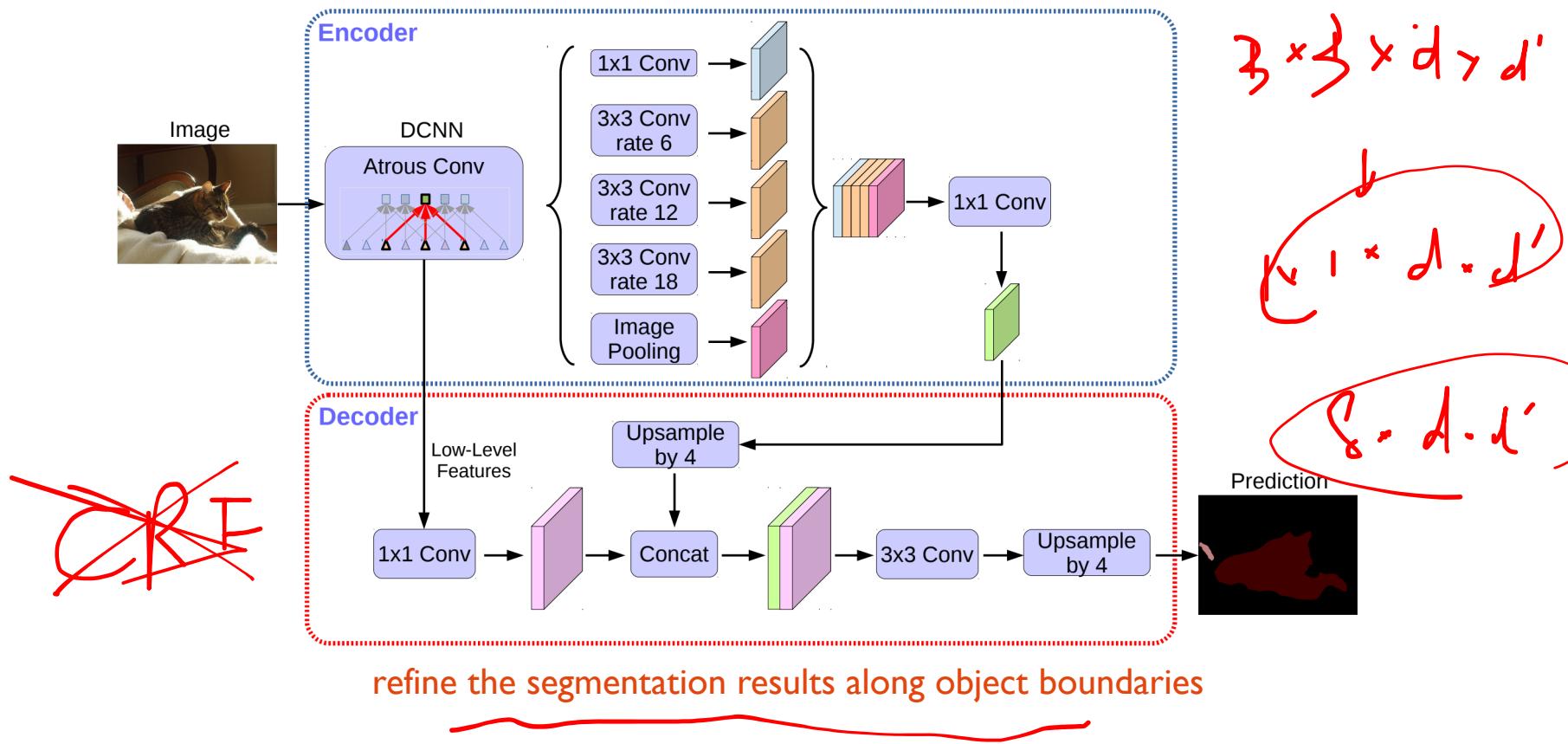


Combine atrous convolutions and spatial pyramid pooling with an encoder-decoder module.

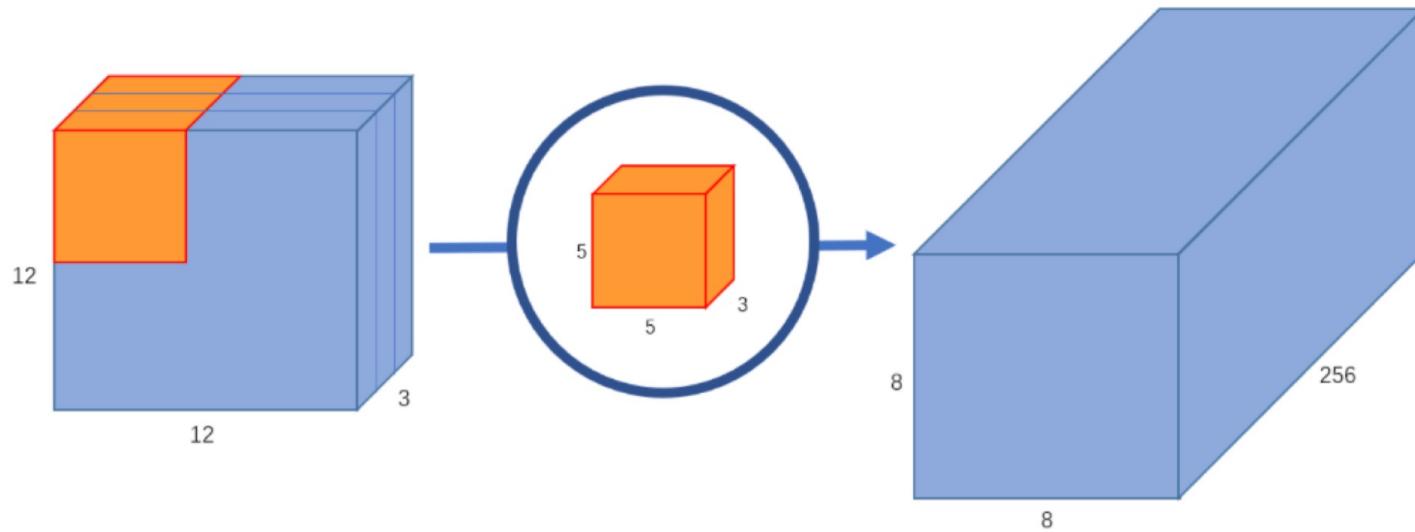
DeepLabv3

Network architecture

Encode multi-scale contextual information by applying atrous convolution at multiple scales

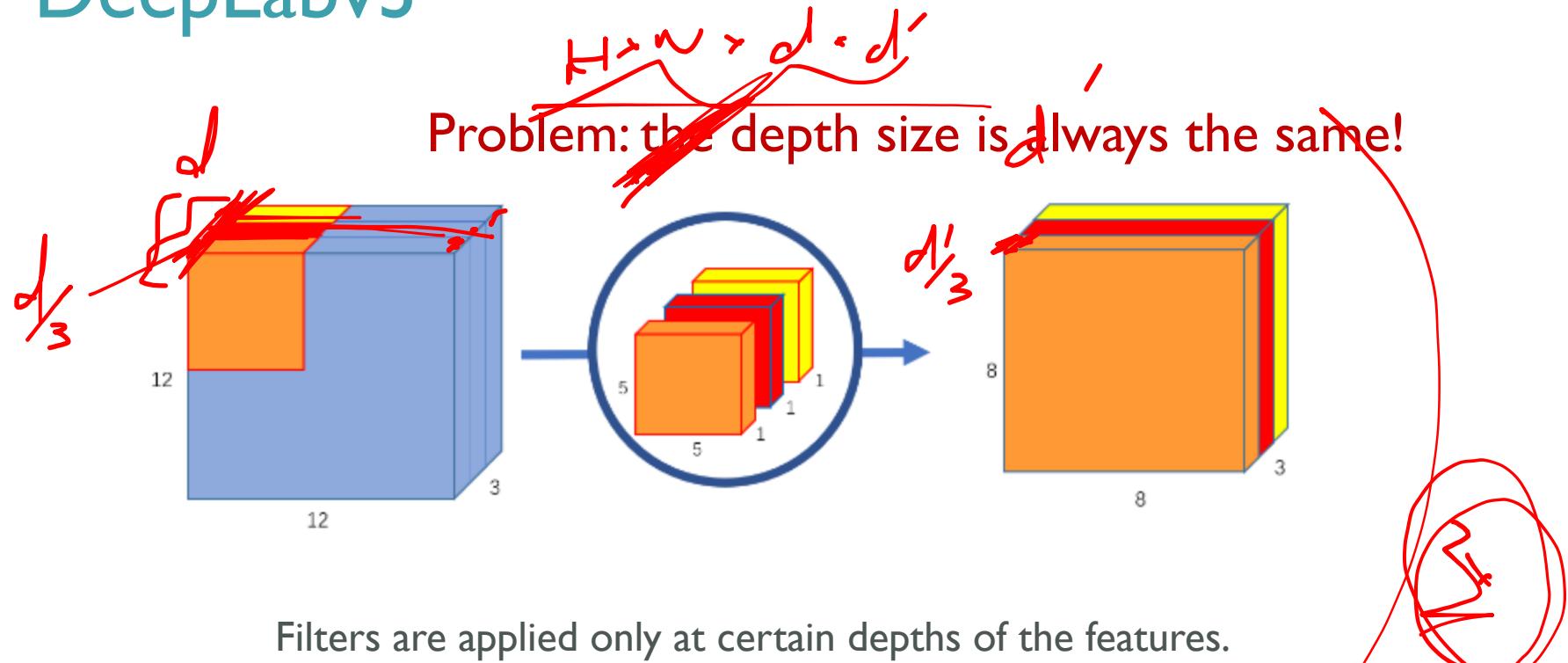


DeepLabv3



Normal convolutions act on all channels.

DeepLabv3

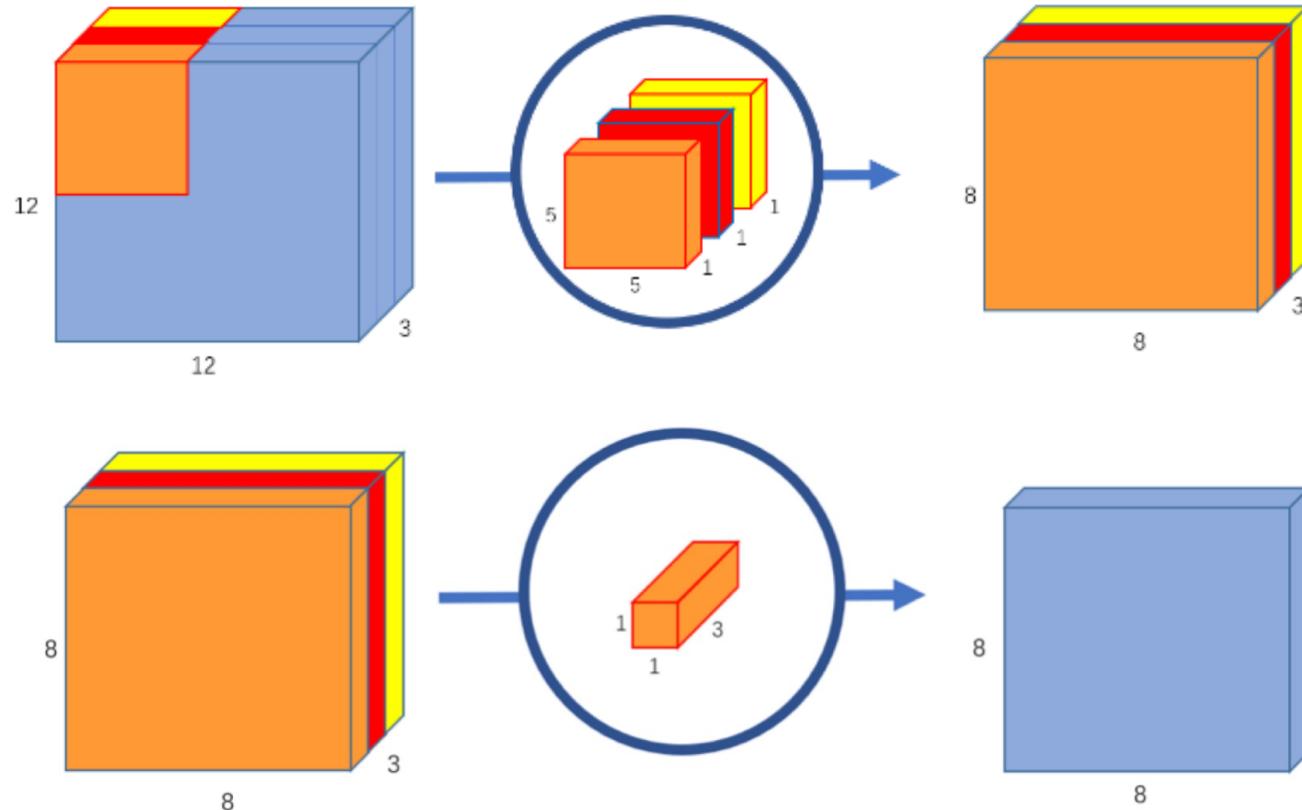


Filters are applied only at certain depths of the features.
Normal convolutions have groups set to 1, the convolutions
used in this image have groups set to 3.

$$H \times W \times \frac{d}{3} < \frac{d'}{3} \times 3$$

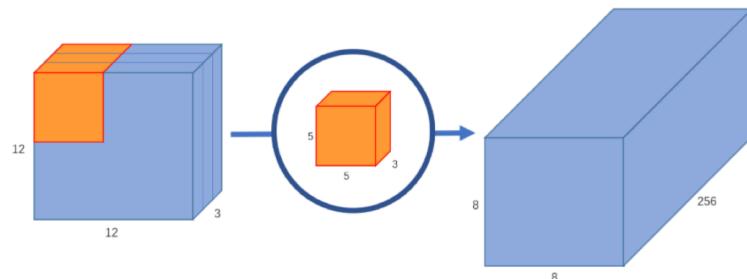
DeepLabv3

Problem: the depth size is always the same!



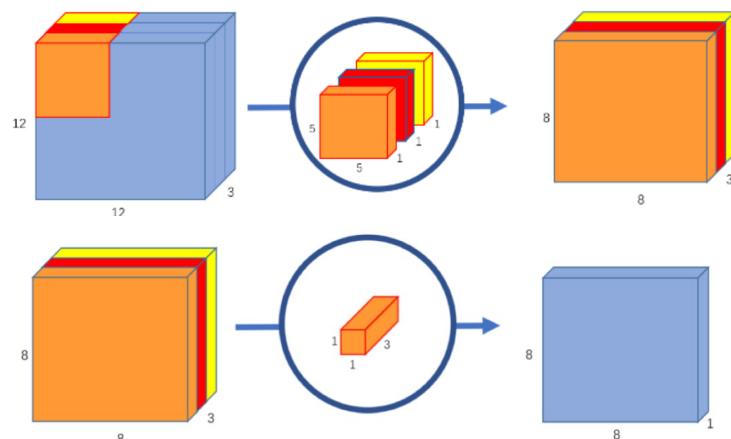
Solution: 1×1 convs!

DeepLabv3



Original convolution
256 kernels of size 5x5x3

Multiplications:
 $256 \times 5 \times 5 \times 3 \times (8 \times 8 \text{ locations}) = 1,228,800$



Depth-wise convolution
3 kernels of size 5x5x1

Multiplications:
 $5 \times 5 \times 3 \times (8 \times 8 \text{ locations}) = 4800$

1x1 convolution
256 kernels of size 1x1x3

Multiplications:
 $256 \times 1 \times 1 \times 3 \times (8 \times 8 \text{ locations}) = 49,152$

DeepLabv3



Semantic segmentation

SegFormer

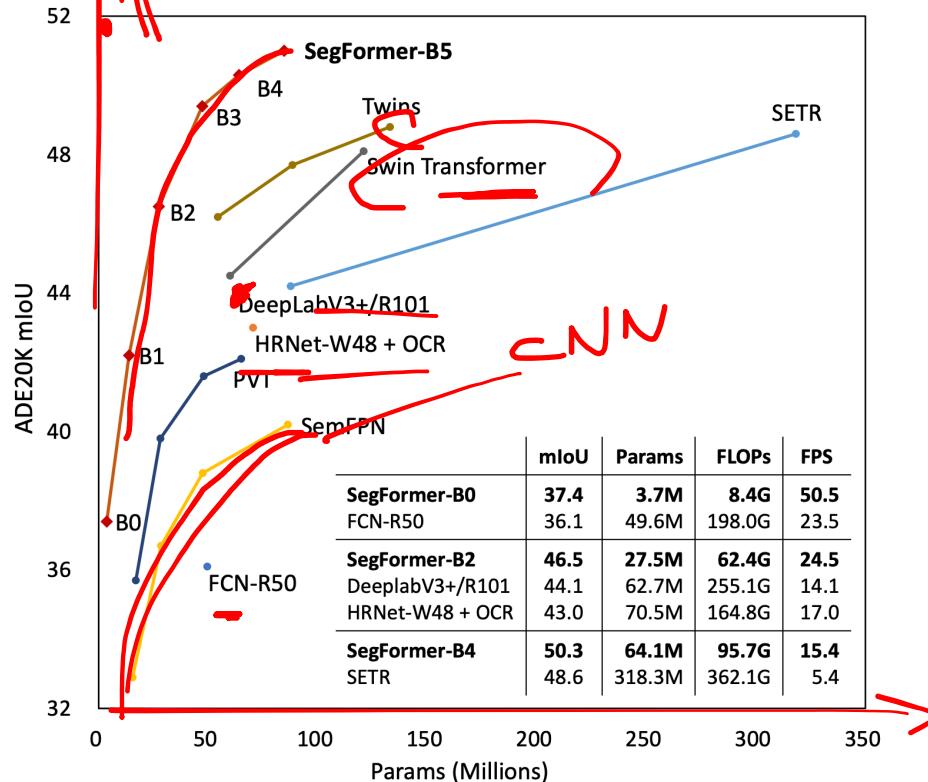
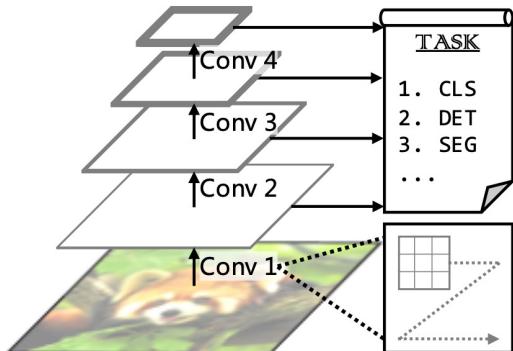


Figure 1: Performance vs. model efficiency on ADE20K. All results are reported with single model and single-scale inference. SegFormer achieves a new state-of-the-art 51.0% mIoU while being significantly more efficient than previous methods.

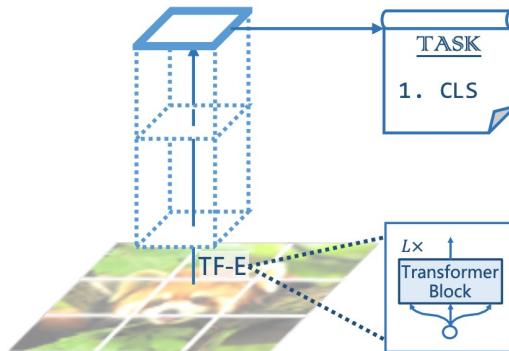
Semantic segmentation

SegFormer



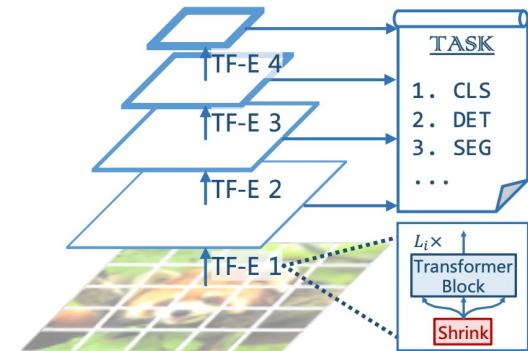
(a) CNNs: VGG [54], ResNet [22], etc.

Many CNN backbones use a pyramid structure for dense prediction tasks such as object detection, instance and semantic segmentation



(b) Vision Transformer [13]

a "columnar" structure specifically designed for image classification (CLS)

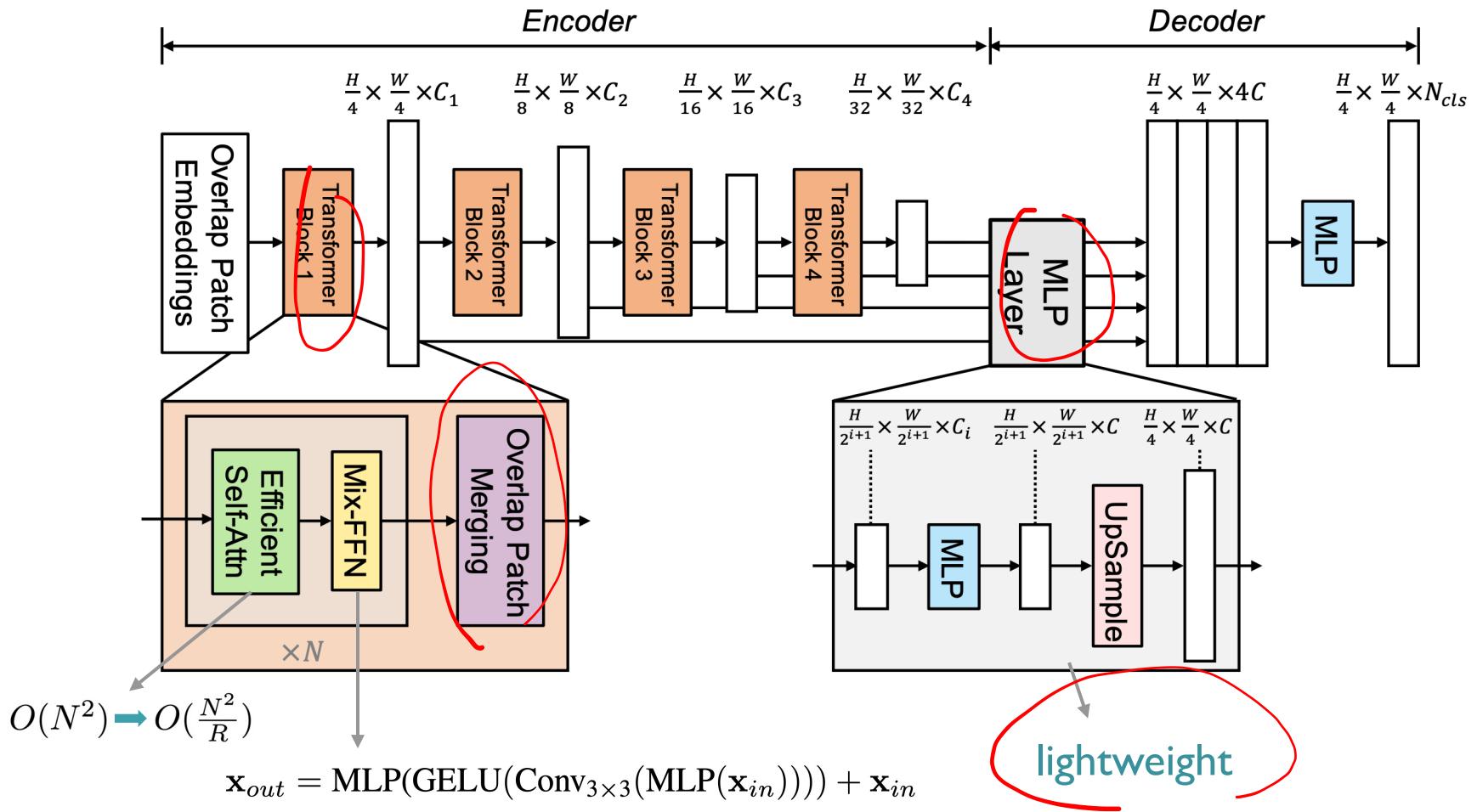


(c) ~~Pyramid~~ Vision Transformer (ours)

by incorporating the pyramid structure from CNNs, the Pyramid Vision Transformer (PVT) can be used as a versatile backbone for many computer vision tasks, broadening the scope and impact of ViT

Semantic segmentation

SegFormer



Semantic segmentation

SegFormer

- For semantic segmentation, maintaining large receptive field to include context information has been a central issue
- decoder design benefits from the non-local attention in Transformers and leads to a larger receptive field without being complex

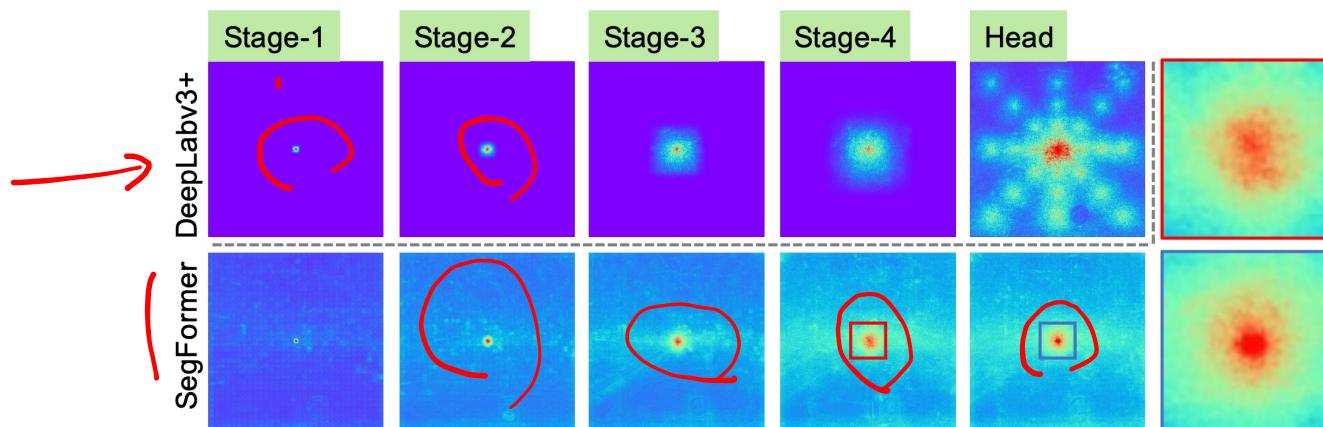


Figure 3: **Effective Receptive Field (ERF) on Cityscapes** (average over 100 images). Top row: Deeplabv3+. Bottom row: SegFormer. ERFs of the four stages and the decoder heads of both architectures are visualized. Best viewed with zoom in.

Next lecture: Video understanding



Thank you very much!

sihengc@sjtu.edu.cn