

Computer Vision

CS308 Autumn

Feng Zheng

SUSTech CS Vision Intelligence and Perception



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY



Content

- Brief Review
- Semantic Segmentation
 - Using convolution
 - downsampling and upsampling
 - Image pyramids
 - Skip connections
 - Dilation
 - ROI align

Brief Review



Faster R-CNN

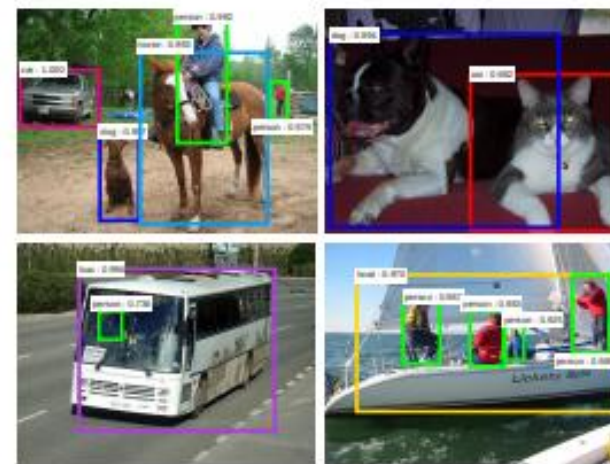
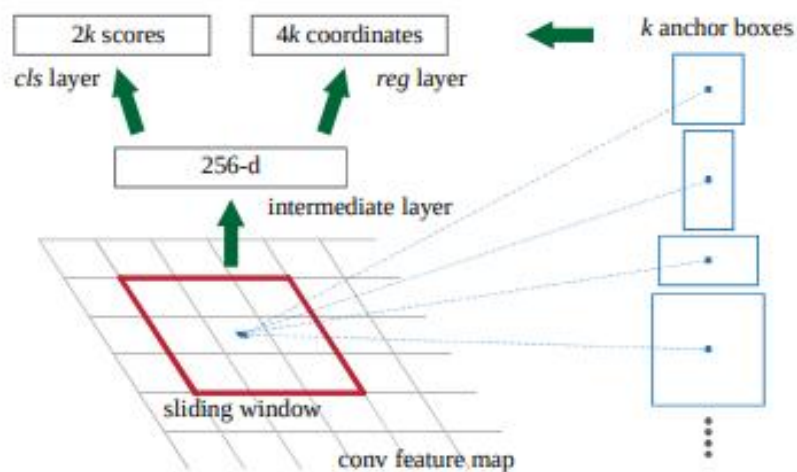
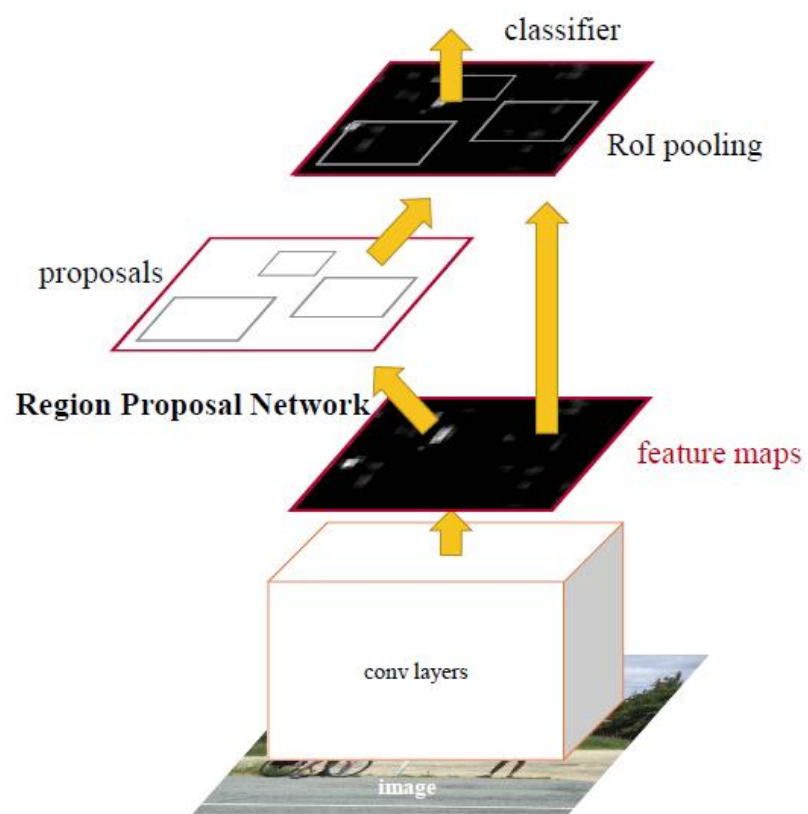
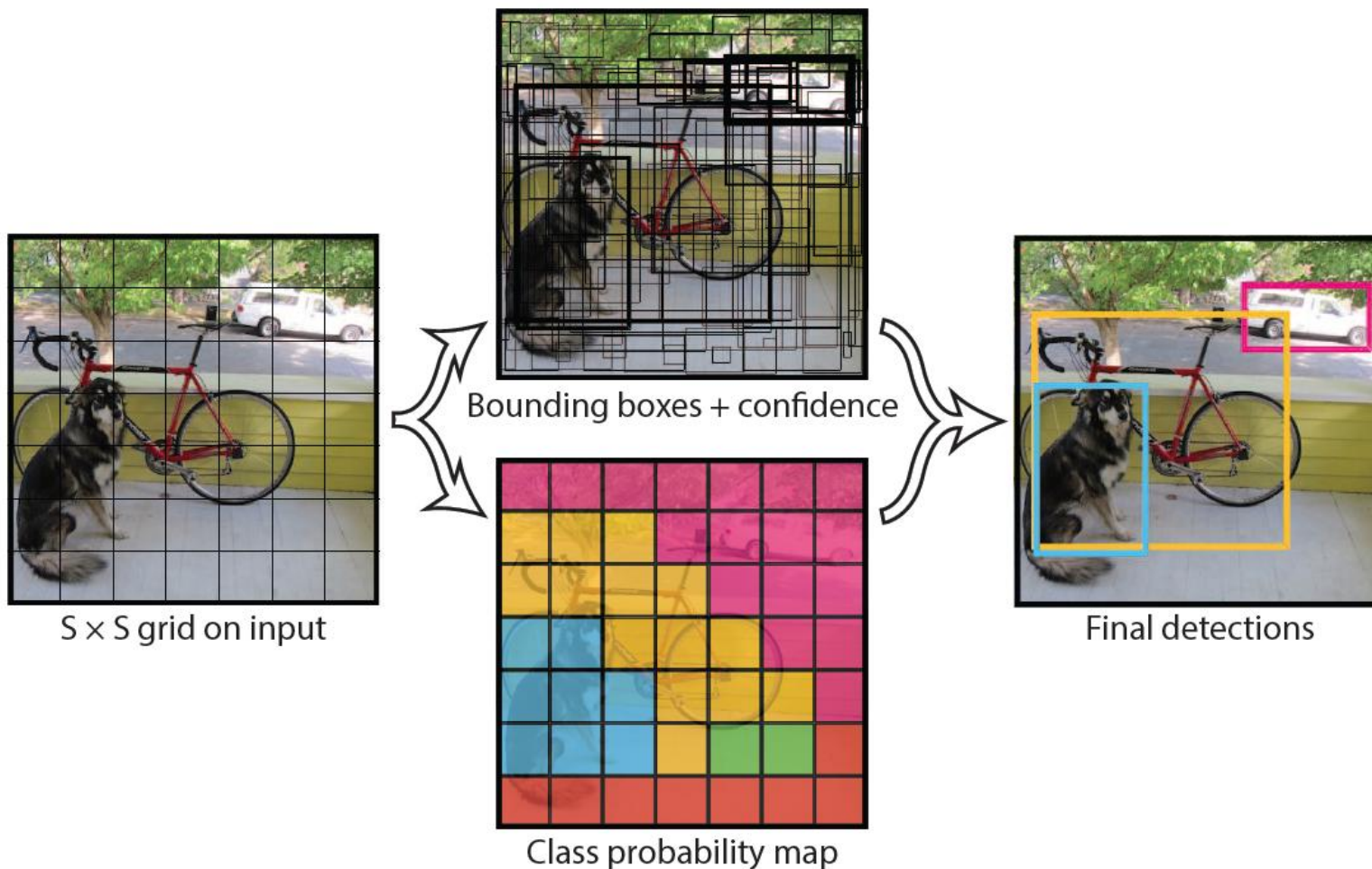


Figure 3: **Left:** Region Proposal Network (RPN). **Right:** Example detections using RPN proposals on PASCAL VOC 2007 test. Our method detects objects in a wide range of scales and aspect ratios.



YOLO Framework



At test time we multiply the conditional class probabilities and the individual box confidence predictions

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

Semantic Segmentation



The Segmentation Task

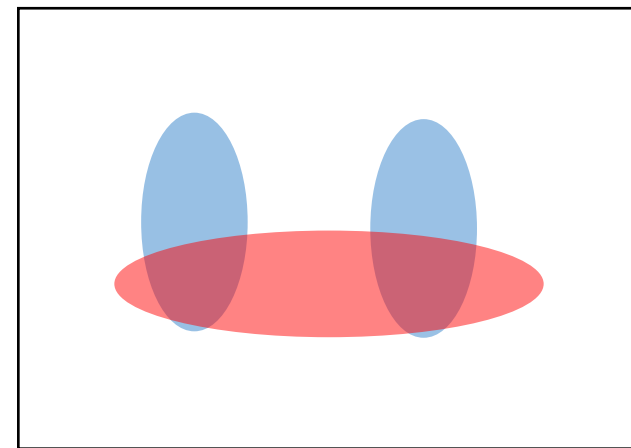


- person
- grass
- trees
- motorbike
- road



Evaluation metric

- Pixel classification!
- Accuracy?
 - Heavily unbalanced
 - Common classes are over-emphasized
- *Intersection over Union*
 - Average across classes and images
- Per-class accuracy
 - Compute accuracy for every class and then average





Things vs Stuff

THINGS

- Person, cat, horse, etc
- Constrained shape
- Individual instances with separate identity
- May need to look at objects



STUFF

- Road, grass, sky etc
- Amorphous, no shape
- No notion of instances
- Can be done at pixel level
- "texture"





Challenges in data collection

- Precise localization is **hard to annotate**
- Annotating every pixel leads to **heavy tails**
 - Common solution: annotate few classes (often things), mark rest as "Other"
 - Weakly-supervised labels: scribble and points
- Common datasets:
 - PASCAL VOC 2012 (~1500 images, 20 categories),
 - COCO (~100k images, 20 categories)
 - The number of classes is limited

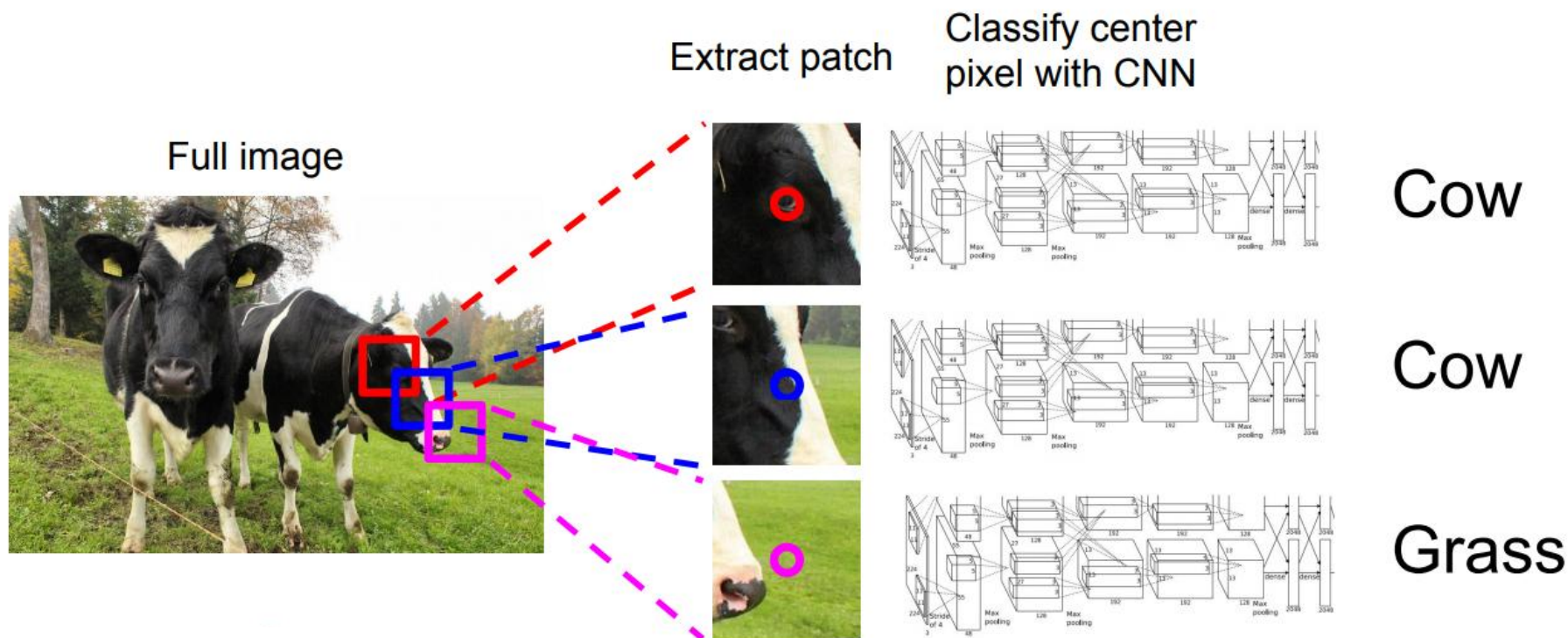


Classical semantic segmentation (Pre-convnet)

- Things
 - Do object detection, then segment out detected objects
- Stuff
 - "Texture classification"
 - Compute histograms of **filter responses**
 - Classify local image patches
- Classical segmentation methods
 - Graph cut
 - Clustering



Using Sliding Window for CNN Features



What's the problem?

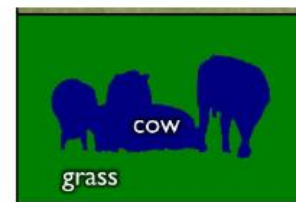
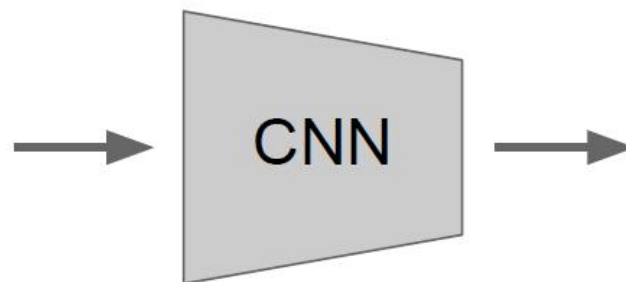
Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



Using Convolution

- We can use pooling:



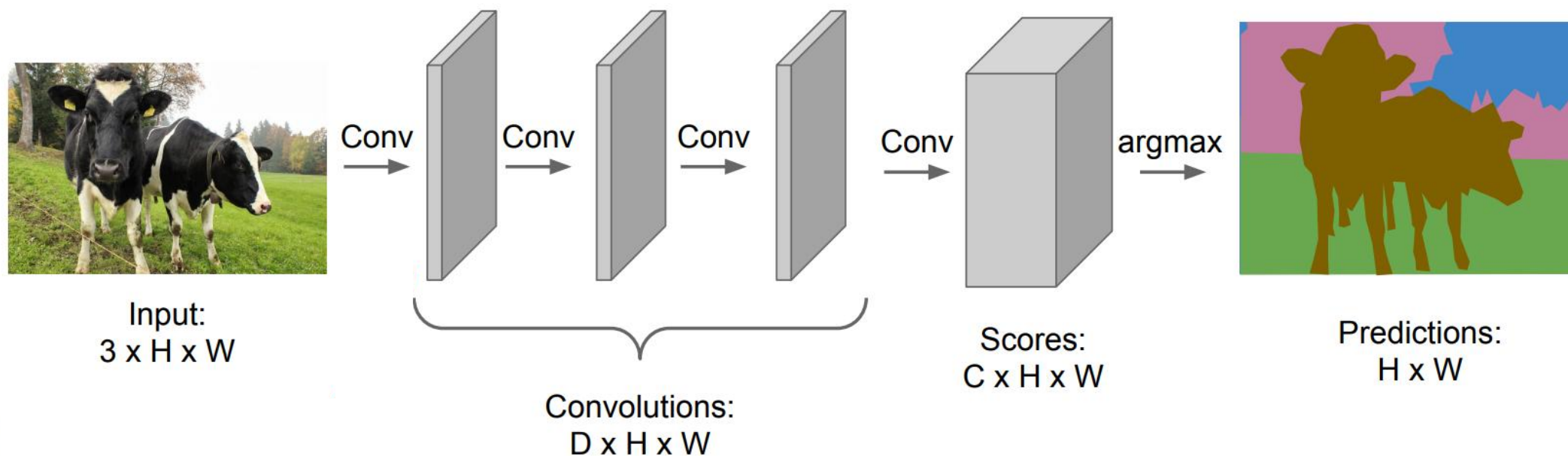
Smaller output
due to pooling

What's the problem?



Using Convolution

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



What's the problem?

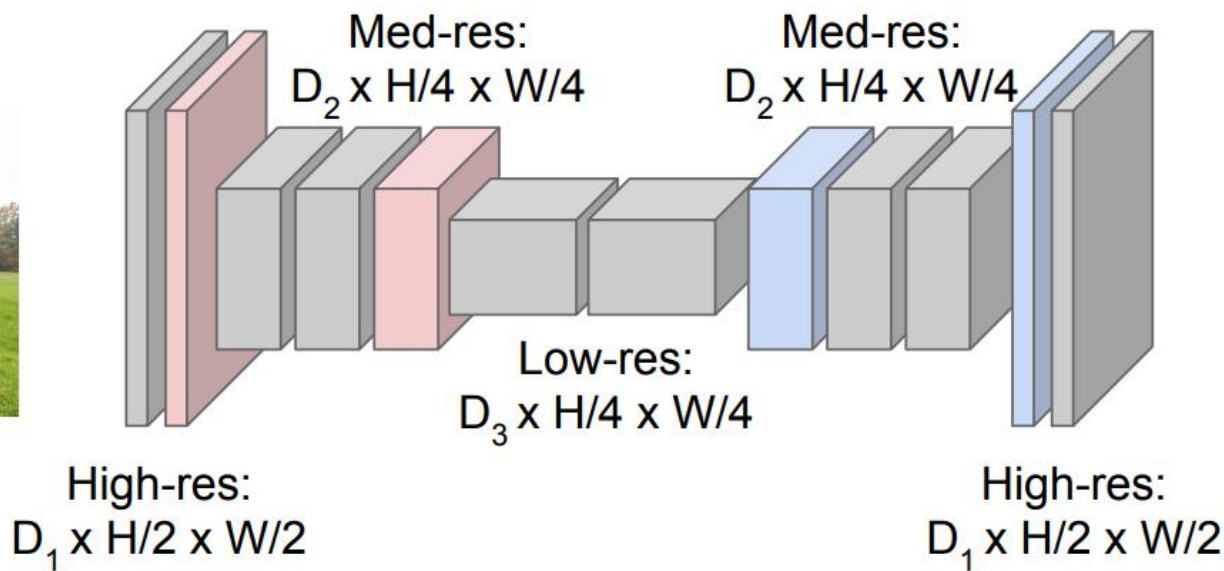


Fully convolutional

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



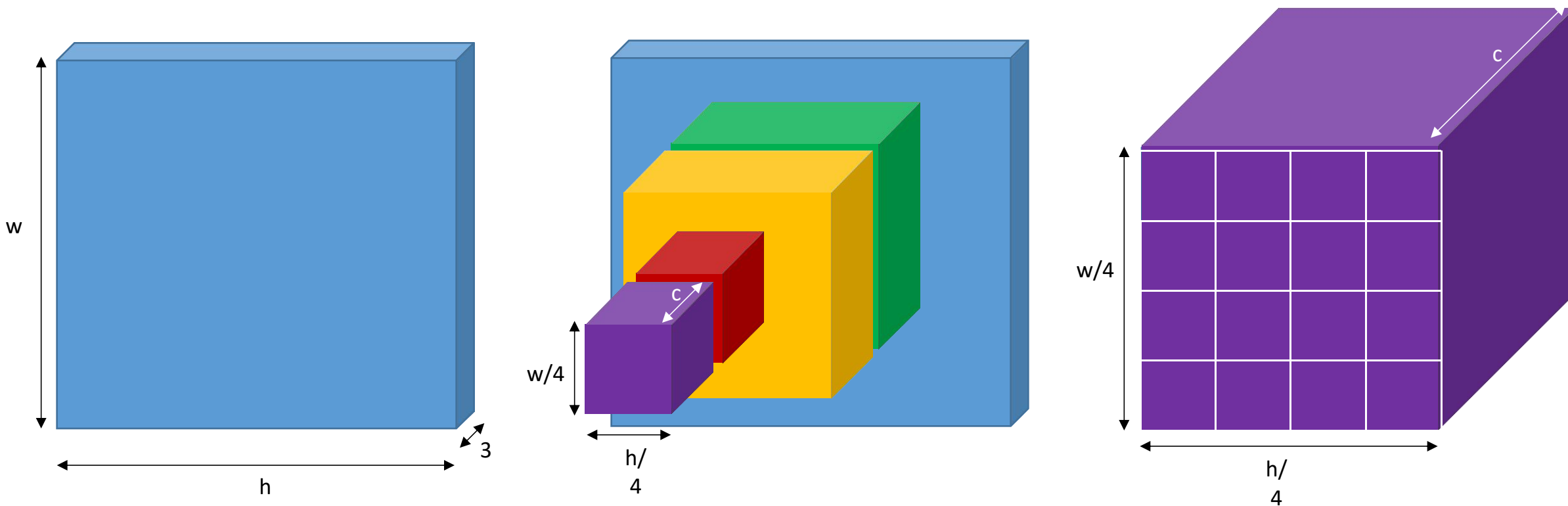
Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

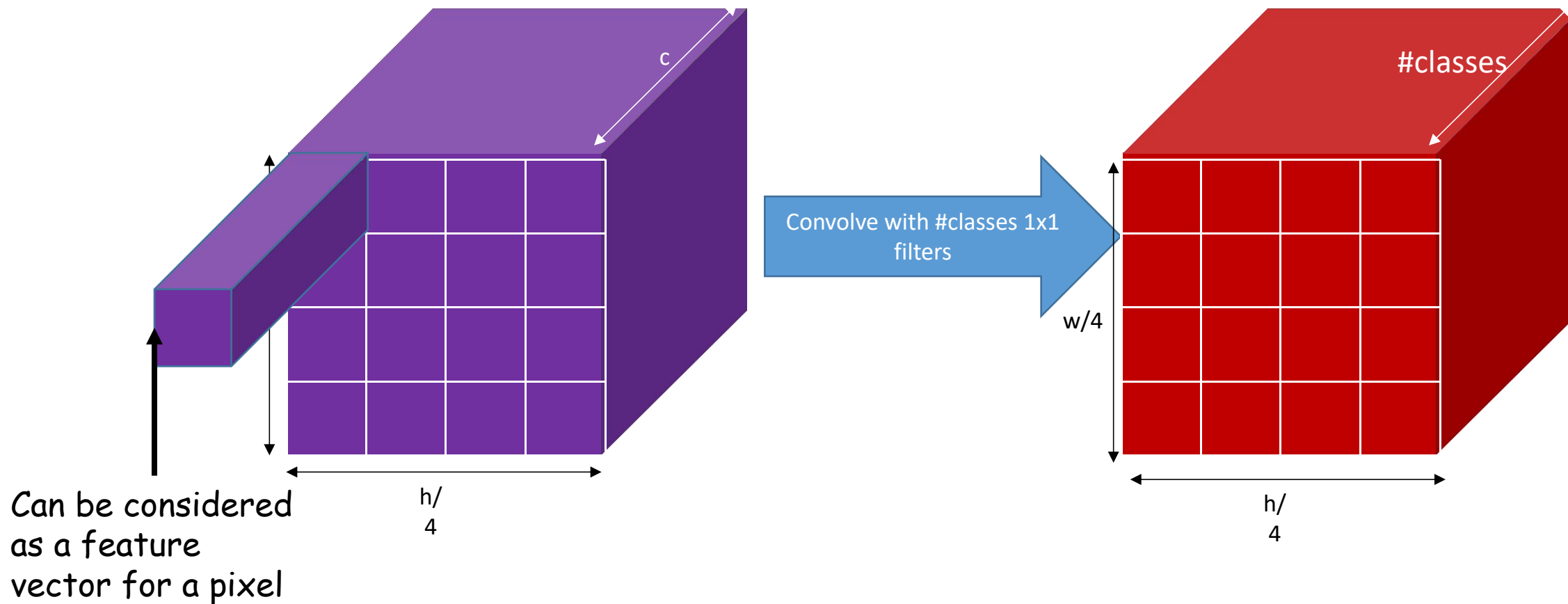


Semantic segmentation using convolutional networks





Semantic segmentation using convolutional networks





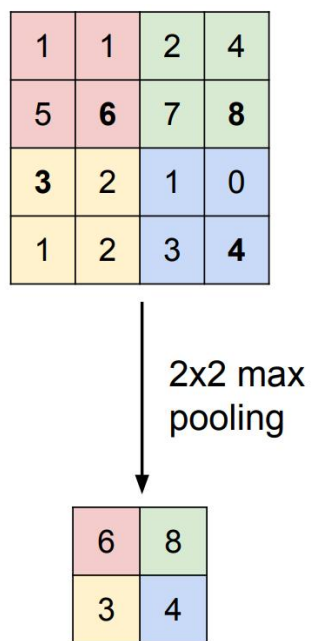
Semantic segmentation using convolutional networks

- Creates a feature map, then **downsample** it and so on...
- The convolution is **not on full size** image.
- All layers are Convolutional, no fully connected layer at the end.
- The **output** size of the last layer will be the **same size as the input** image.
- The **number of channels** is according to the **number of classes** we want to classify (car, road)
- Training is the same (as in previous methods) by **cross-entropy loss for every pixel**.

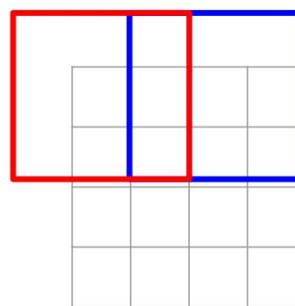


Downsampling

Pooling



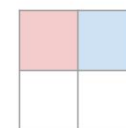
Strided Convolution



Input: 4 x 4



Dot product
between filter
and input



Output: 2 x 2

Filter moves 2 pixels in
the input for every one
pixel in the output

Stride gives ratio between
movement in input and
output

Move 2 pixels



Recall: Strided convolution

Normal convolution: 3x3, pad 1, Stride 1

$$d = \frac{n+2p-f}{s} + 1$$

f – filter size

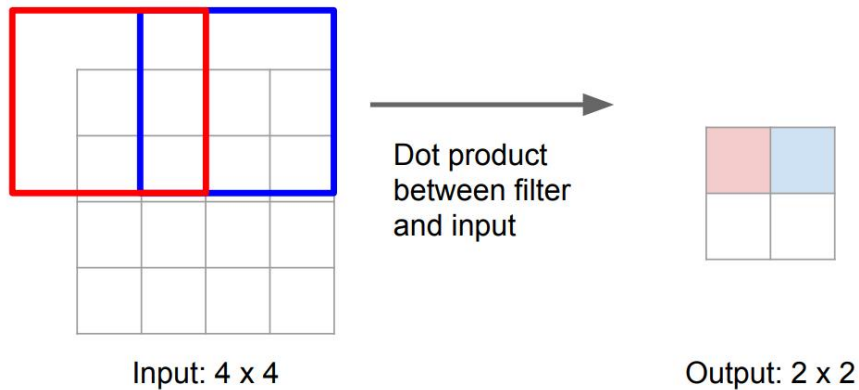
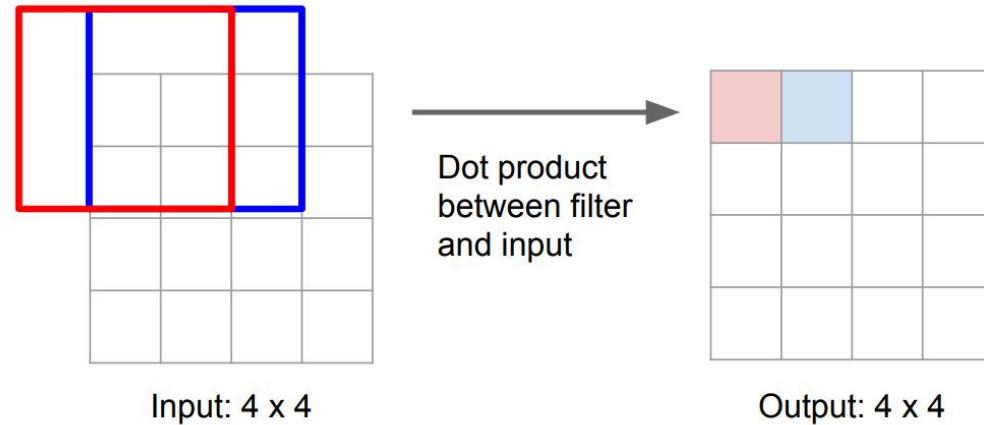
p – padding

n – input image size

s – stride

d – output image size

Normal convolution: 3x3, pad 1, Stride 2



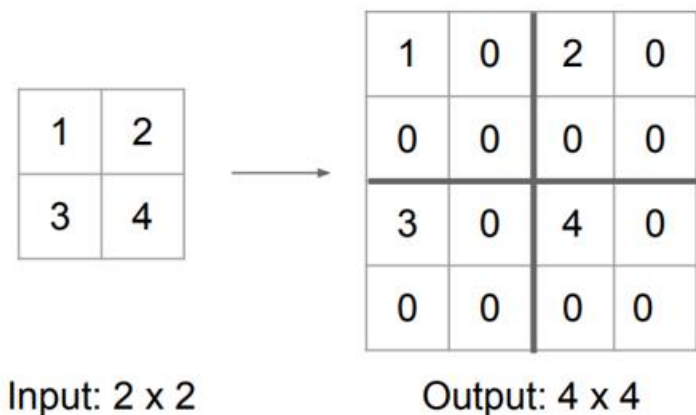
Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

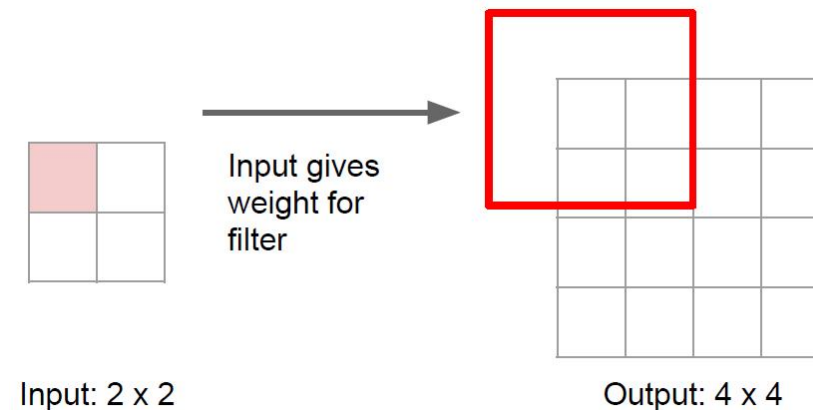


How Does Upsampling Work?

UpPooling



Transpose Convolution





Uppooling Methods

Nearest Neighbor

1	2
3	4

Input: 2 x 2



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

Output: 4 x 4

"Bed of Nails"

1	2
3	4

Input: 2 x 2



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

Output: 4 x 4



Uppooling Methods

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

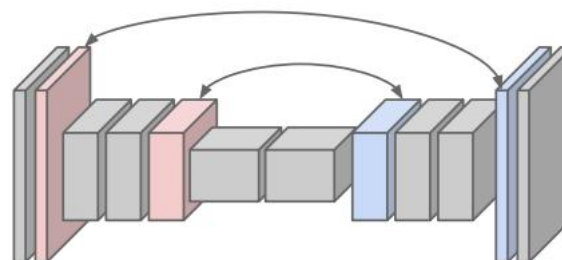
Input: 2 x 2



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

Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers

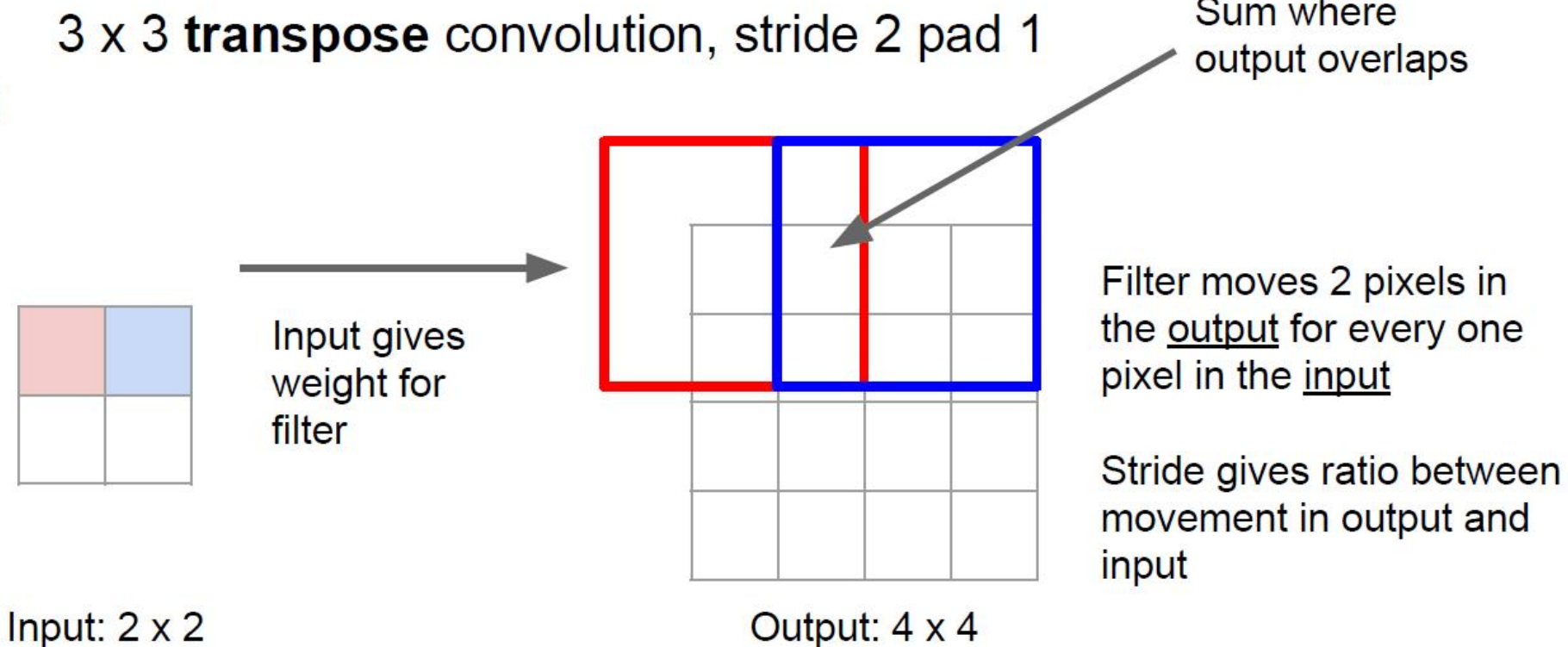




Transpose Convolution

Other names:

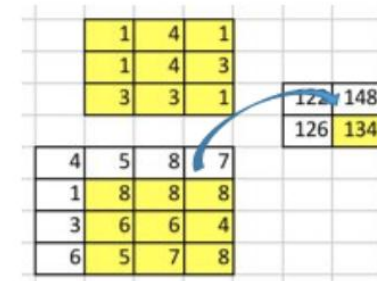
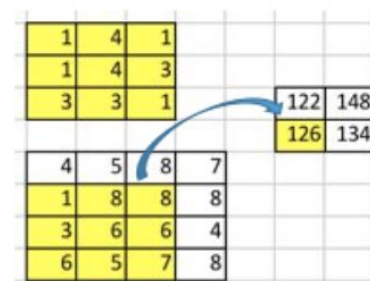
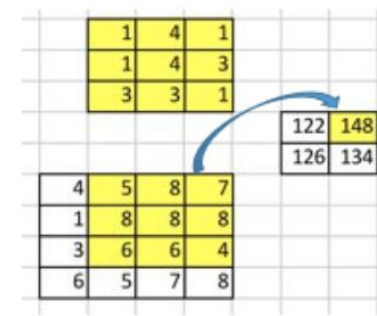
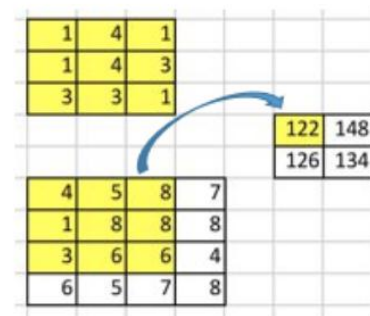
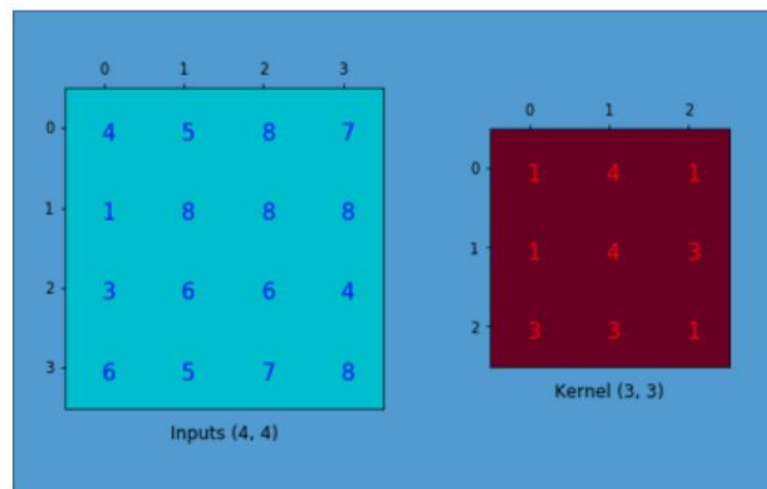
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution



Weighting?

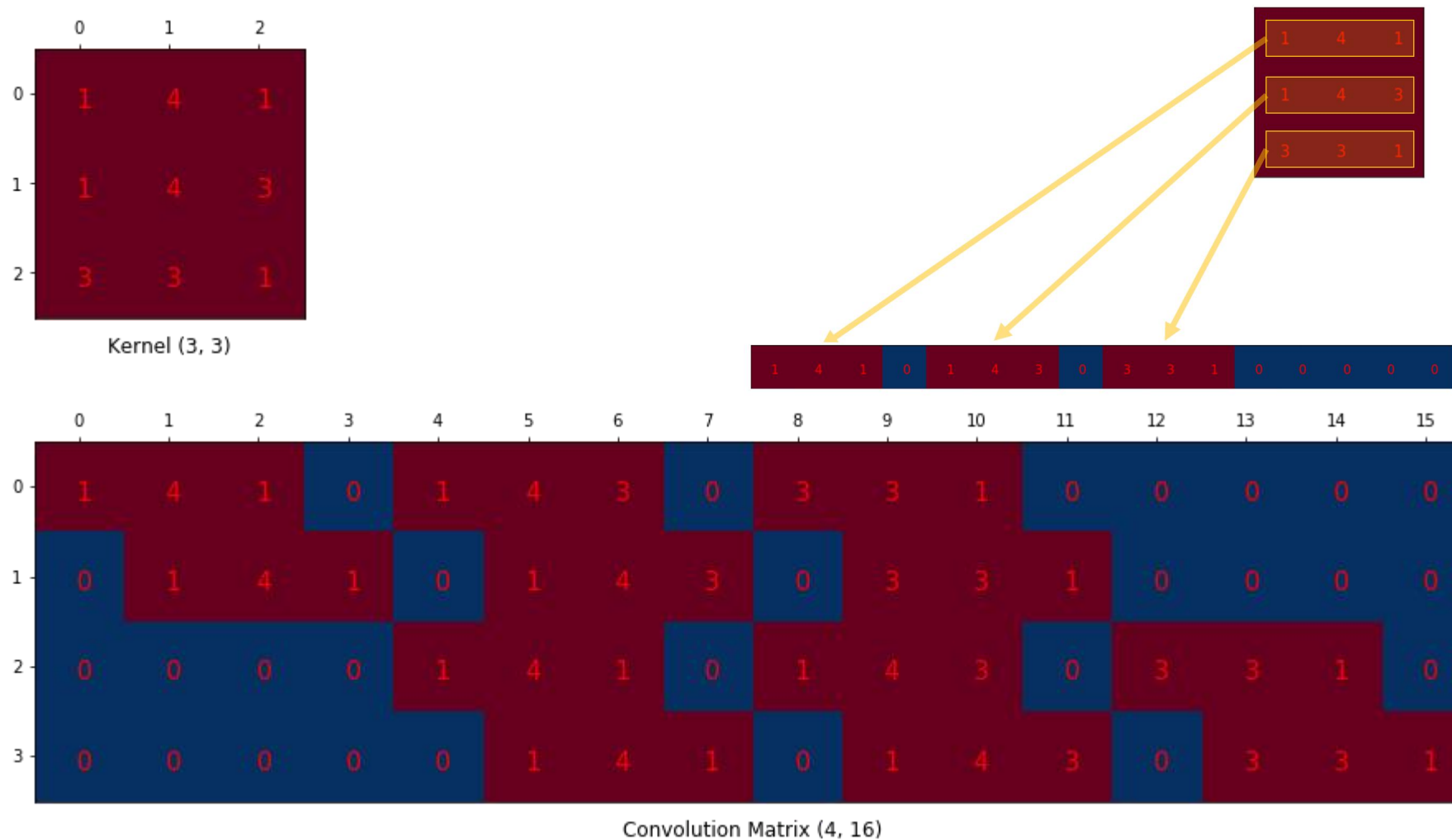


Convolution Matrix



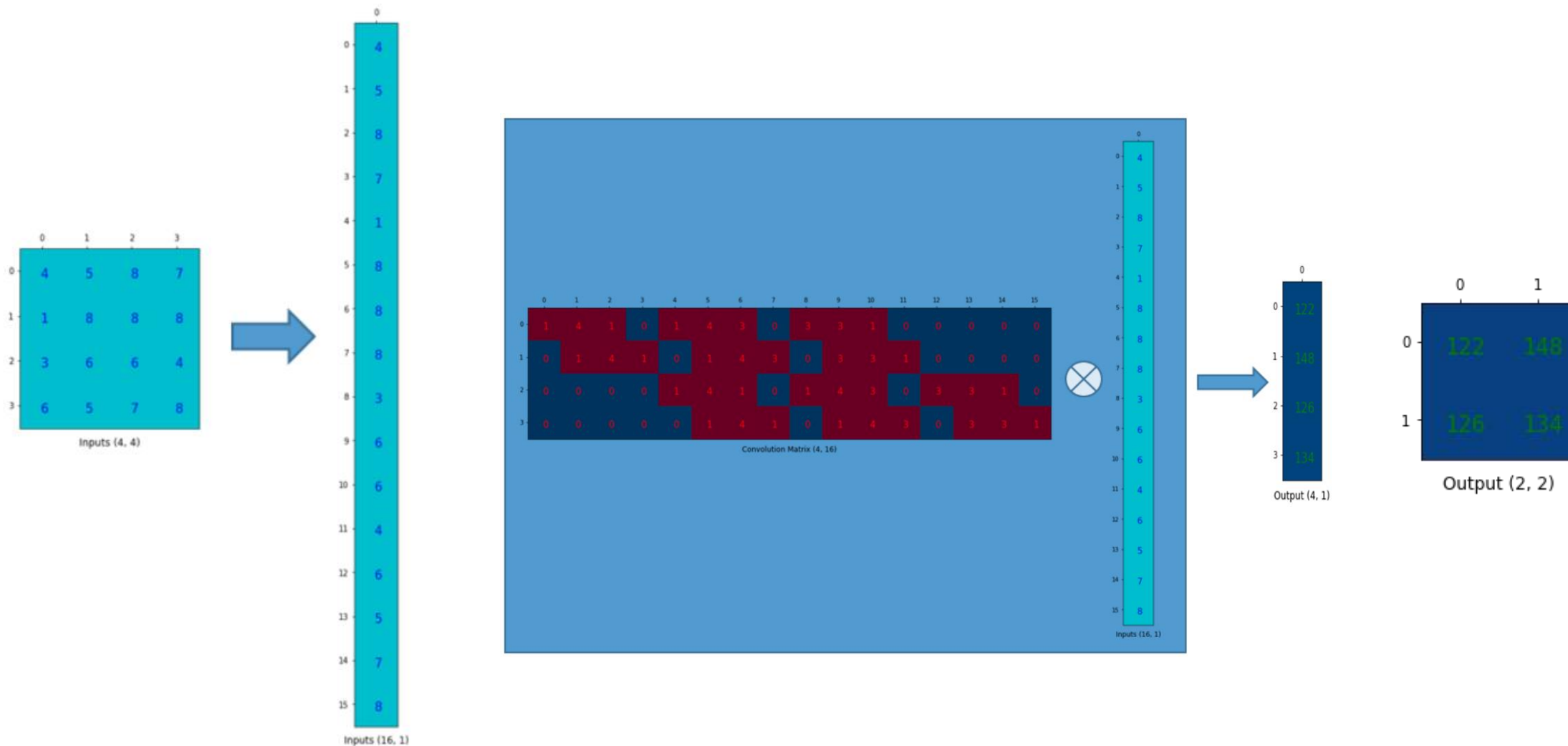


Convolution Matrix



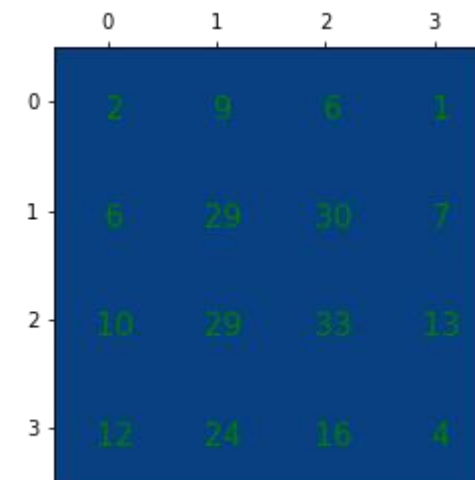
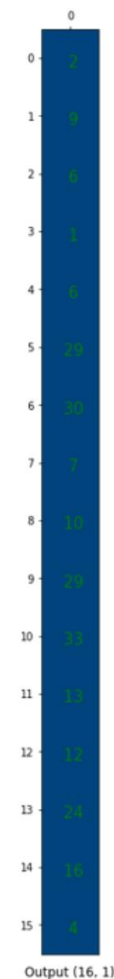
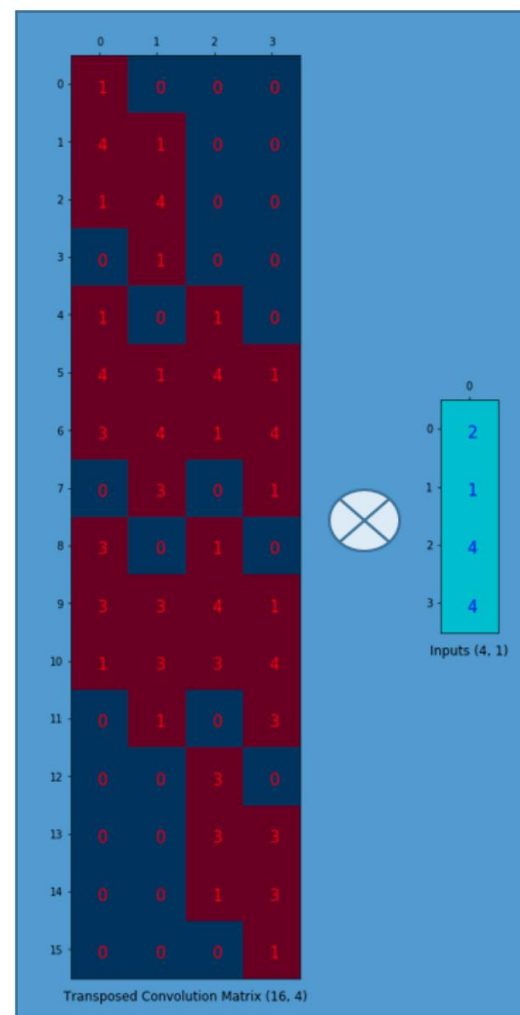
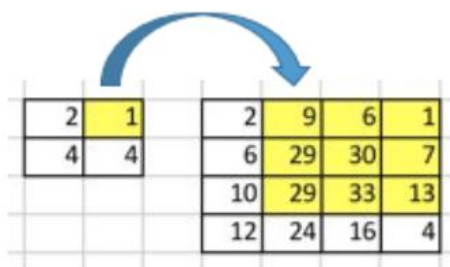
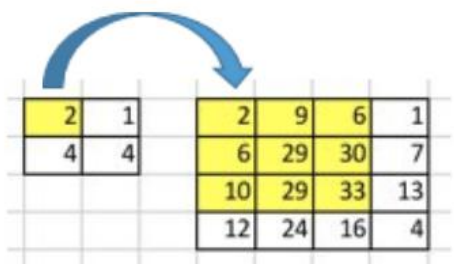


Convolution Matrix





Transposed Convolution Matrix



Output (4, 4)

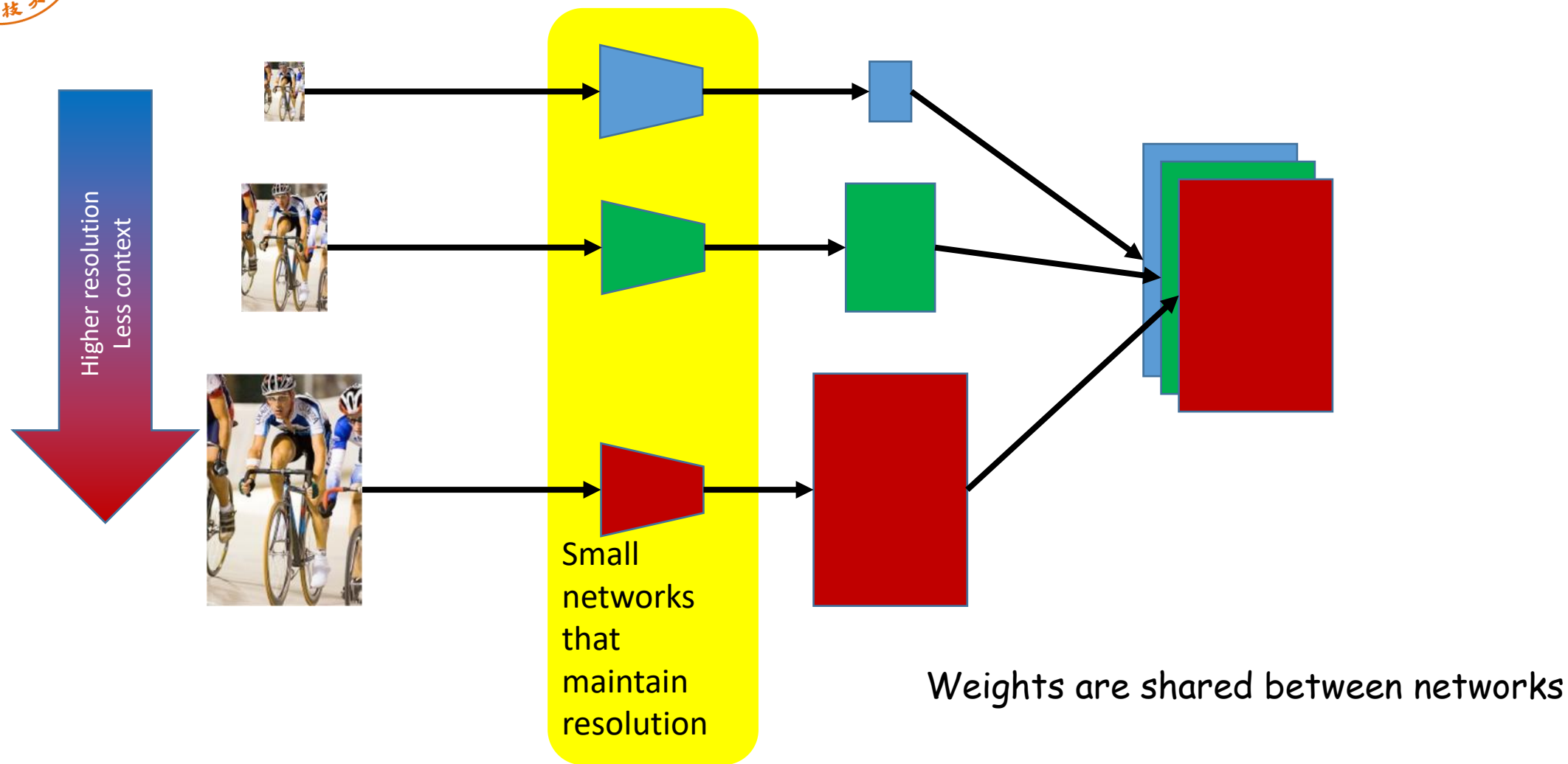


The resolution issue

- Problem: Need **fine** details!
- Shallower network / earlier layers?
 - Deeper networks work better: more **abstract concepts**
 - Shallower network => Not very semantic!
- Remove subsampling?
 - Subsampling allows later layers to capture **larger** and larger patterns
 - Without subsampling => Looks at only a **small** window!

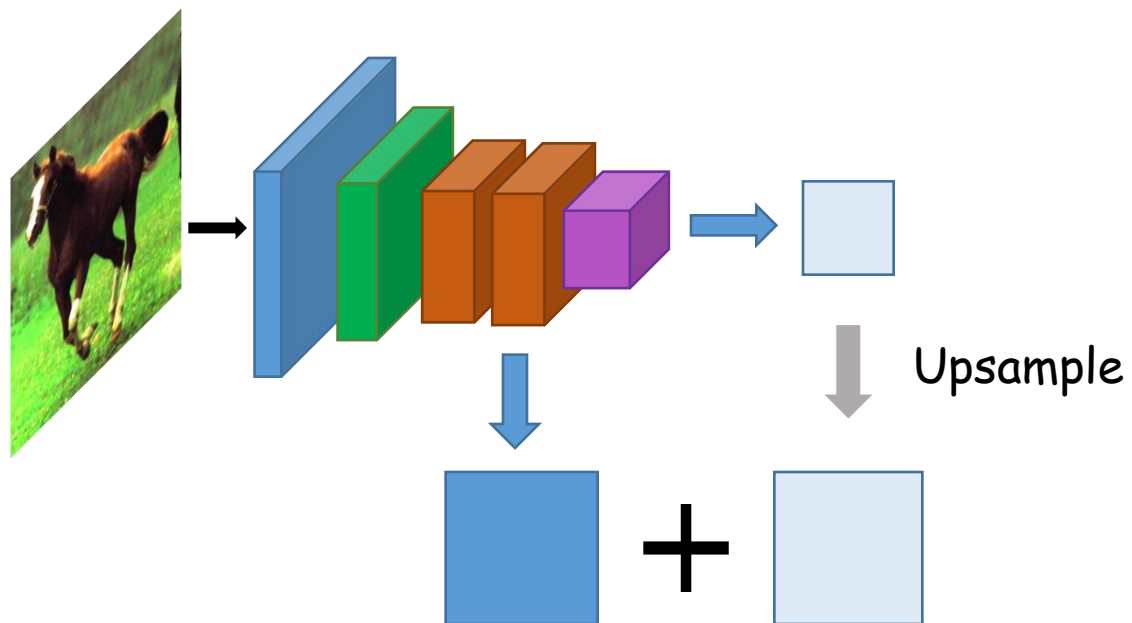


Solution 1: Image pyramids





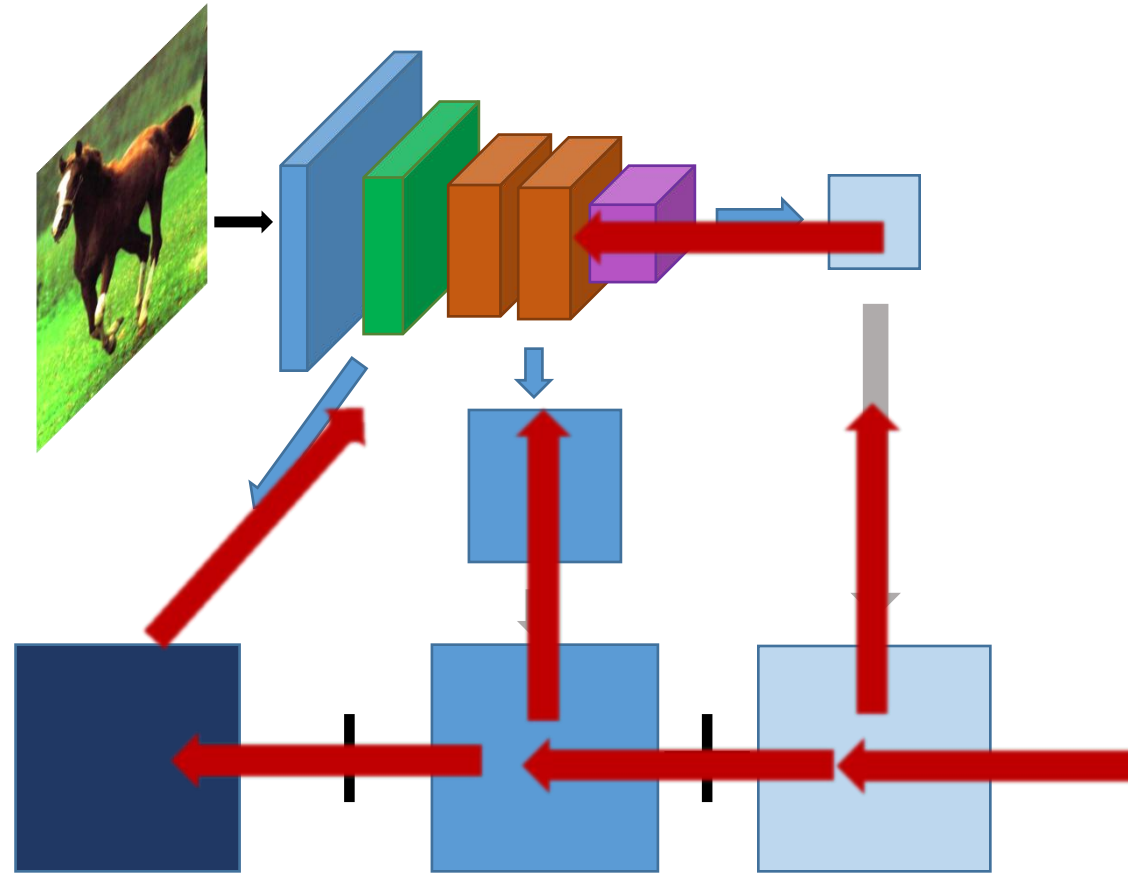
Solution 2: Skip connections



Compute class scores at **multiple layers**, then upsample and add



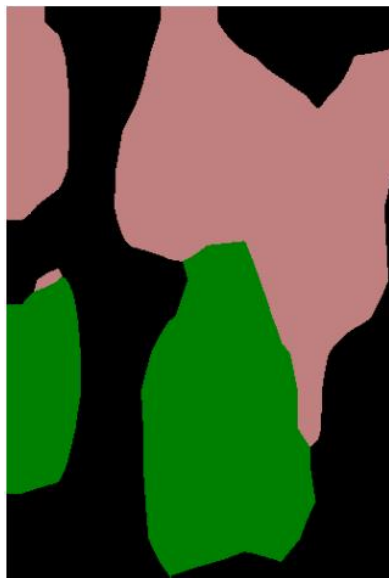
Solution 2: Skip connections



Red arrows indicate
backpropagation



Solution 2: Skip connections



without skip

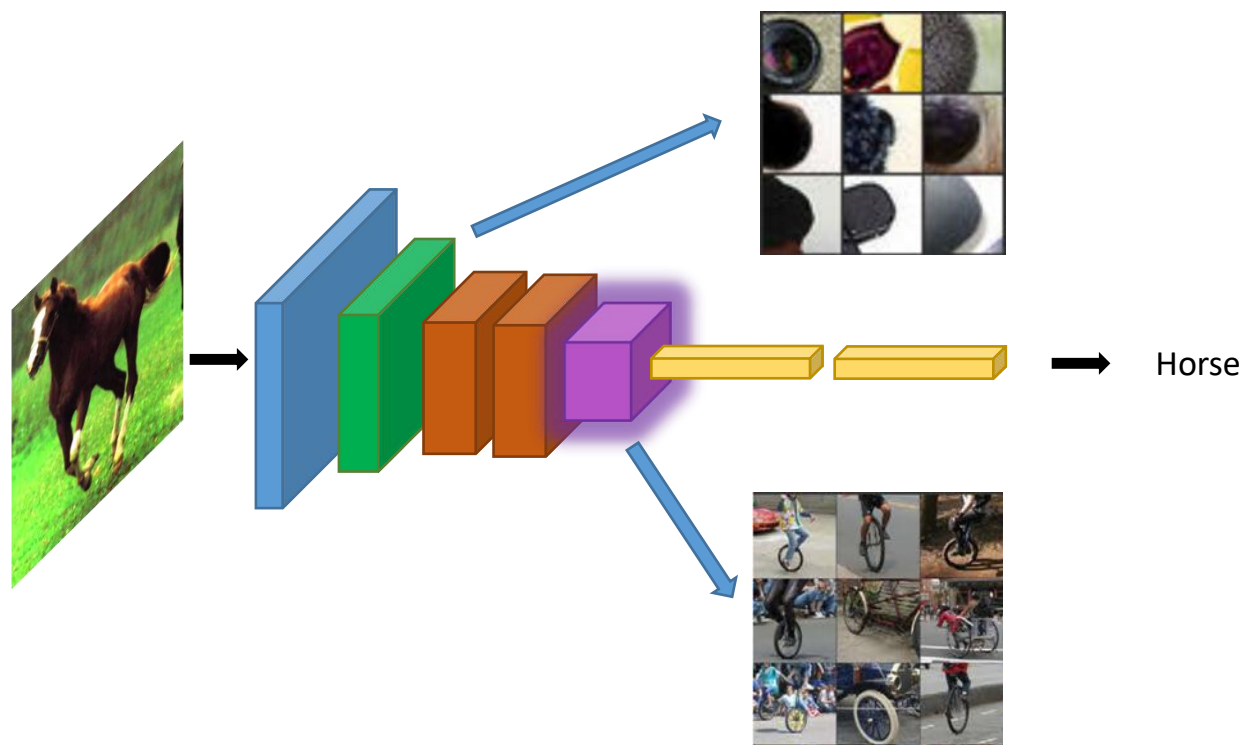


with skip



Solution 2: Skip connections

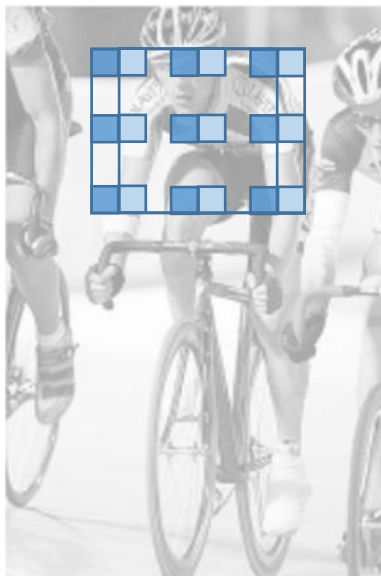
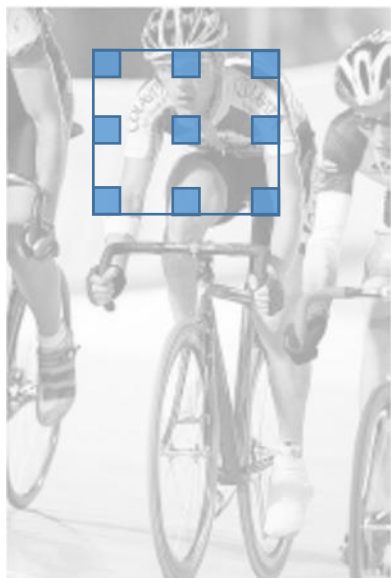
- Problem: early layers not semantic





Solution 3: Dilation

- Need subsampling to allow convolutional layers to **capture large regions with small filters**
 - Can we do this without subsampling?



$$(F * k)(p) = \sum_{s+t=p} F(s) k(t). \quad (F *_l k)(p) = \sum_{s+lt=p} F(s) k(t).$$

*Standard Convolution (Left)
Dilated Convolution (Right)*

*When $l=1$, it is standard convolution.
When $l>1$, it is dilated convolution.*



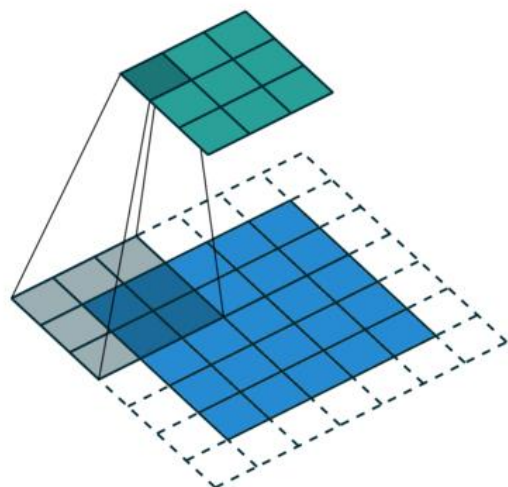
Solution 3: Dilation

- Instead of subsampling by factor of 2: **dilate by factor of 2**
- Dilation can be seen as:
 - Using a much larger filter, but with most entries set to 0
 - Taking a small filter and "exploding"/ "dilating" it
- Not panacea: without subsampling, feature maps are much larger: memory issues
- The idea of Dilated Convolution is come from the **wavelet decomposition**. It is also called "atrous convolution", "algorithme à trous" and "hole algorithm". Thus, any ideas from the past are still useful if we can turn them into the deep learning framework.

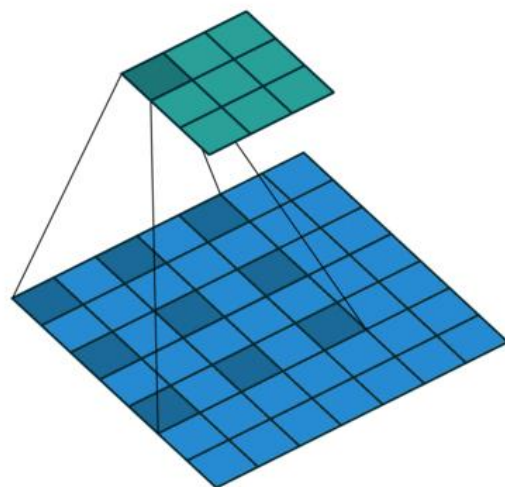


Solution 3: Dilation

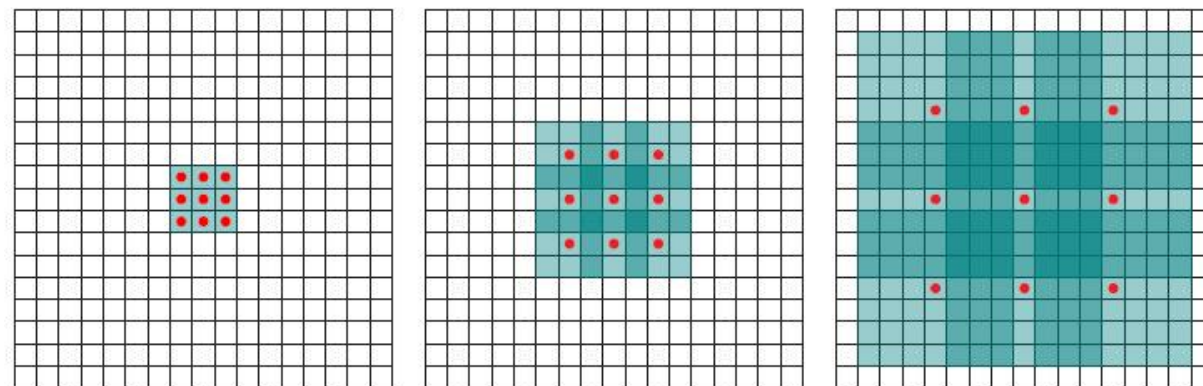
- The figures illustrate an example of **dilated convolution** when $l=2$. We can see that the **receptive field** is **larger** compared with the standard one.



Standard
Convolution ($l=1$)



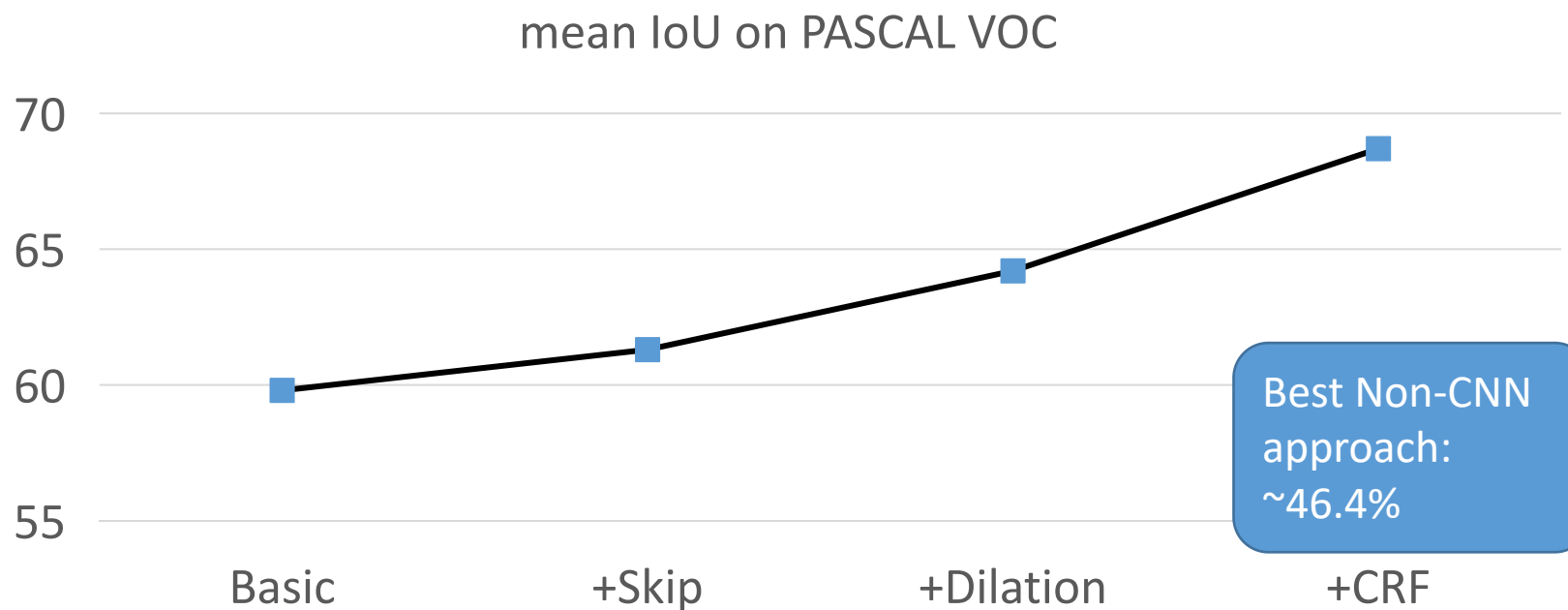
Dilated
Convolution ($l=2$)



$l=1$ (left), $l=2$ (Middle), $l=4$ (Right)



Putting it all together



Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs. Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan Yuille. In *ICLR*, 2015.



Other additions

Method	mean IoU (%)
VGG16 + Skip + Dilation	65.8
ResNet101	68.7
ResNet101 + Pyramid	71.3
ResNet101 + Pyramid + COCO	74.9

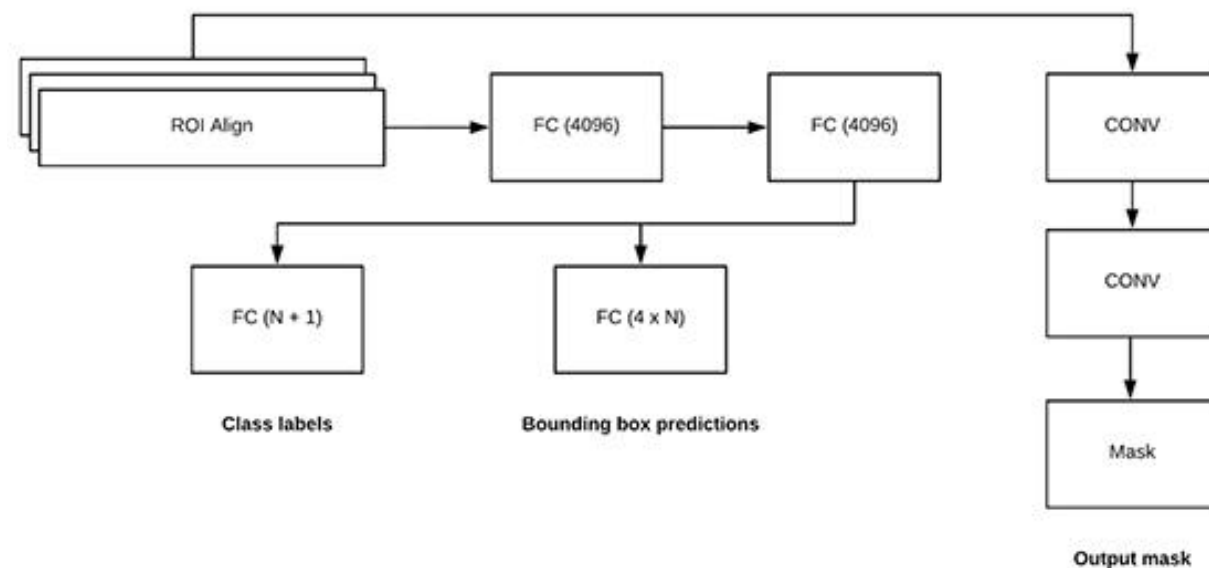
DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan Yuille. Arxiv 2016.



Mask R-CNN

- Faster R-CNN
 - Region Proposal Network (RPN) \leftarrow Selective Search algorithm
- Mask R-CNN
 - Replacing the **ROI Pooling** module with a more accurate **ROI Align module**
 - Inserting an additional branch out of the ROI Align module

Mask R-CNN



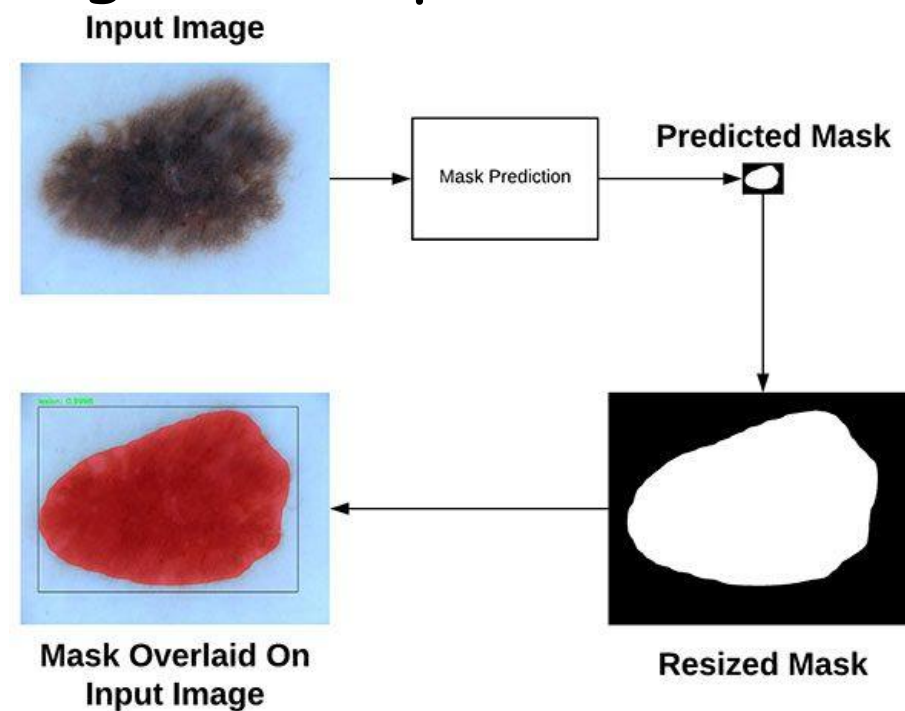
Mask R-CNN: Kaiming He, Georgia Gkioxari, Piotr Dollar and Ross Girshick
<https://arxiv.org/pdf/1703.06870.pdf>



Mask R-CNN

- Each of the 300 selected ROIs go through three parallel branches of the network:
 - Label prediction
 - Bounding box prediction
 - Mask prediction
- A multi-task loss
 - Average binary cross-entropy loss
 - Per pixel: K classes

$$L = L_{cls} + L_{box} + L_{mask}$$





Mask R-CNN

• ROI Align

- Removes this quantization which is causes this misalignment
- For each bin, you regularly sample 4 locations and do bilinear interpolation
- Result are not sensitive to exact sampling location or the number of samples

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.1	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.1	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

0.8	0.6
0.9	0.6

0.88	0.6
0.9	0.6

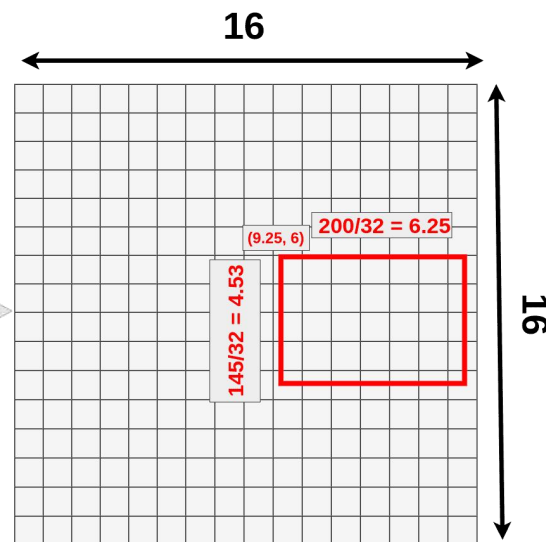
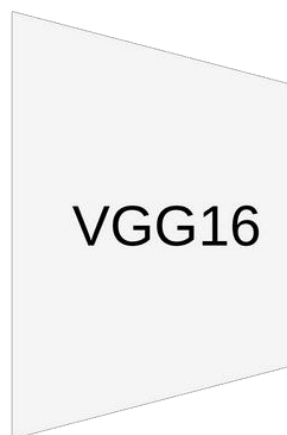
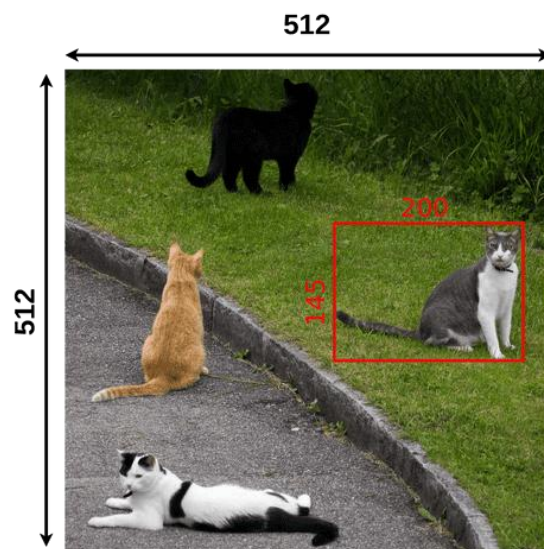
	AP	AP ₅₀	AP ₇₅	AP ^{bb}	AP ^{bb} ₅₀	AP ^{bb} ₇₅
<i>RoIPool</i>	23.6	46.5	21.6	28.2	52.7	26.9
<i>RoIAlign</i>	30.9	51.8	32.1	34.0	55.3	36.4
	+7.3	+ 5.3	+10.5	+5.8	+2.6	+9.5



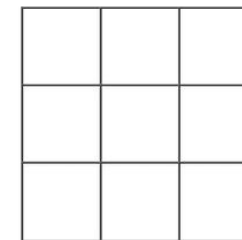
ROI Align

- An image input of size 512x512x3
- A 16x16x512 feature map
- The proposed RoIs (145x200 box)
- 3*3 RoI Pooling

- (9.25, 6) - top left corner
- 6.25 - width
- 4.53 - height



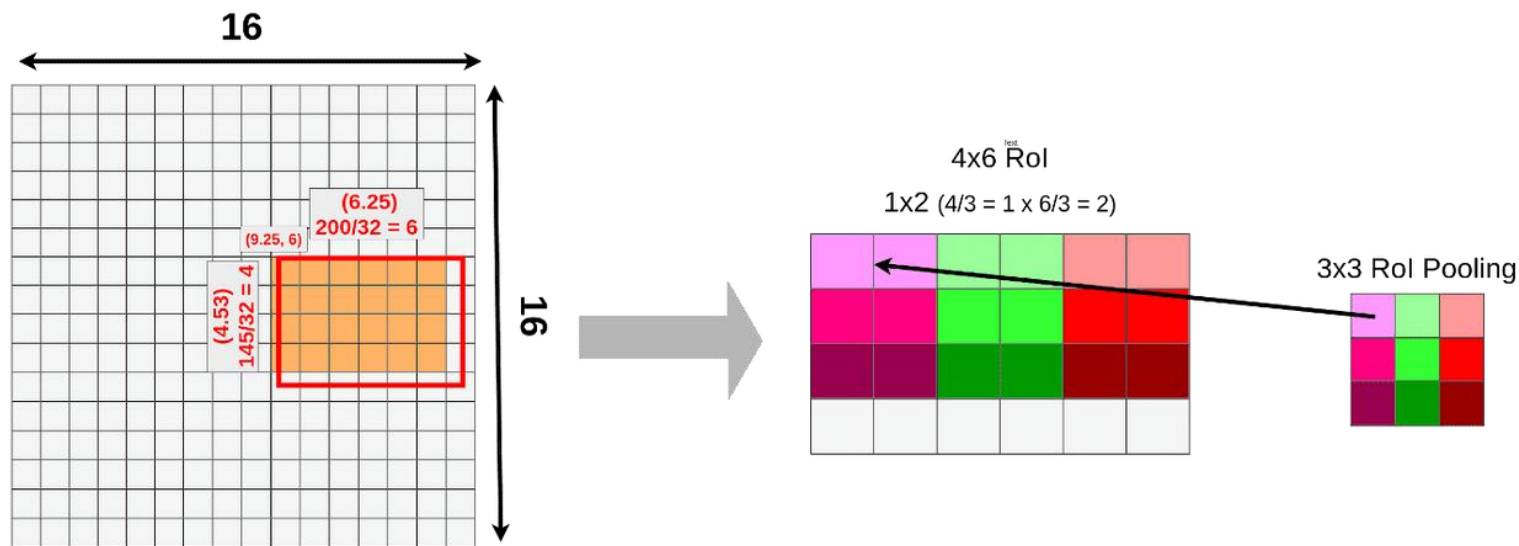
3x3 RoI Pooling





ROI Align

- The main difference between RoI Pooling and RoI Align is quantization (applying quantization twice)
 - First time in the mapping process
 - Second time during the pooling process



Quantization when mapping and pooling



ROI Align

- RoI Align is not using quantization for data pooling



1. Divide mapped RoI (6.25×4.53) by 3



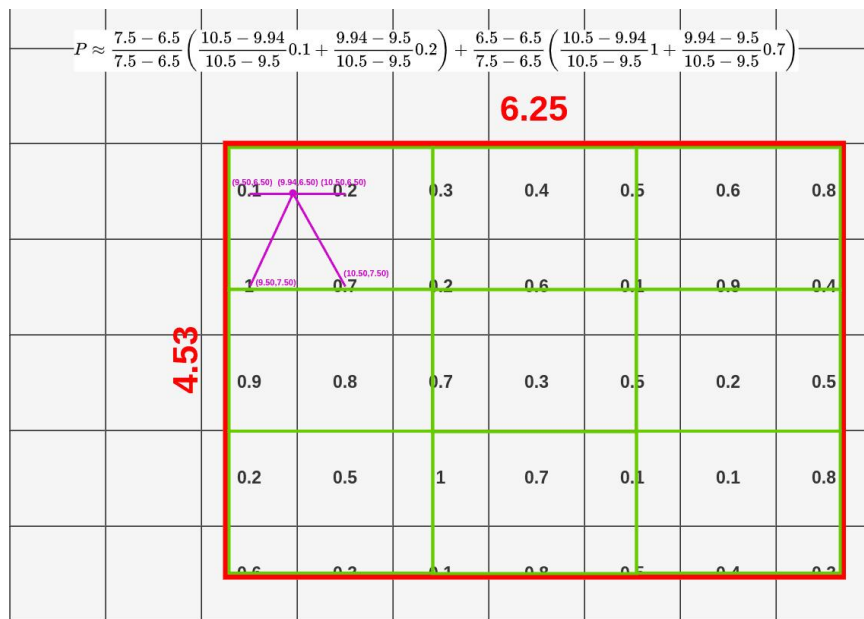
2. Create four sampling points inside that box (divide box by 3)



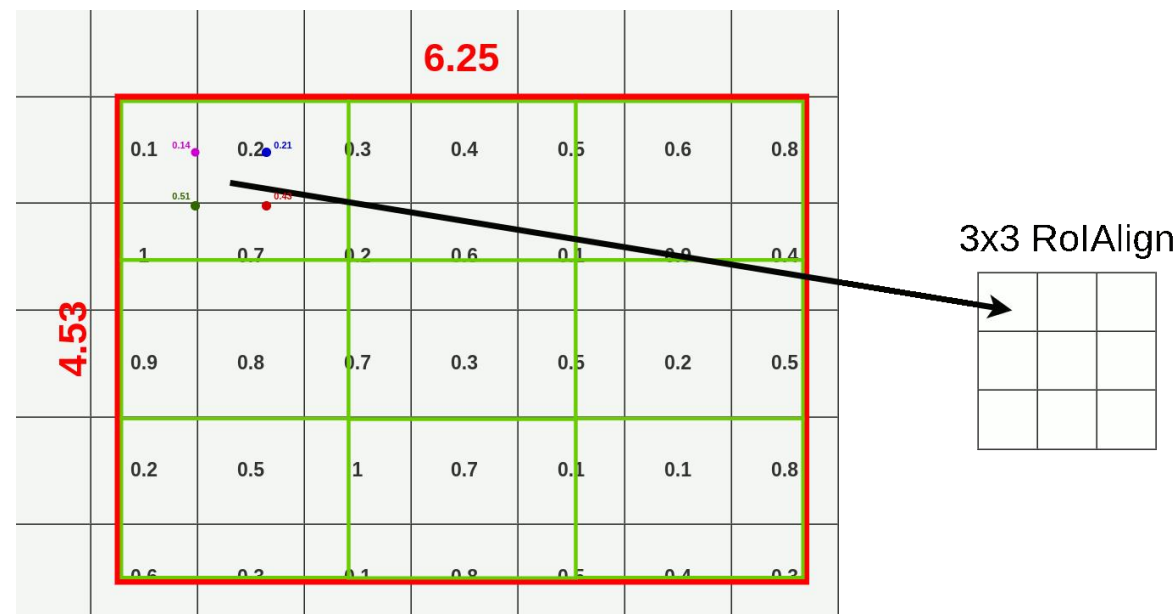
ROI Align

- RoI Align is not using quantization for data pooling

$$P \approx \frac{y_2 - y}{y_2 - y_1} \left(\frac{x_2 - x}{x_2 - x_1} Q_{11} + \frac{x - x_1}{x_2 - x_1} Q_{21} \right) + \frac{y - y_1}{y_2 - y_1} \left(\frac{x_2 - x}{x_2 - x_1} Q_{12} + \frac{x - x_1}{x_2 - x_1} Q_{22} \right)$$



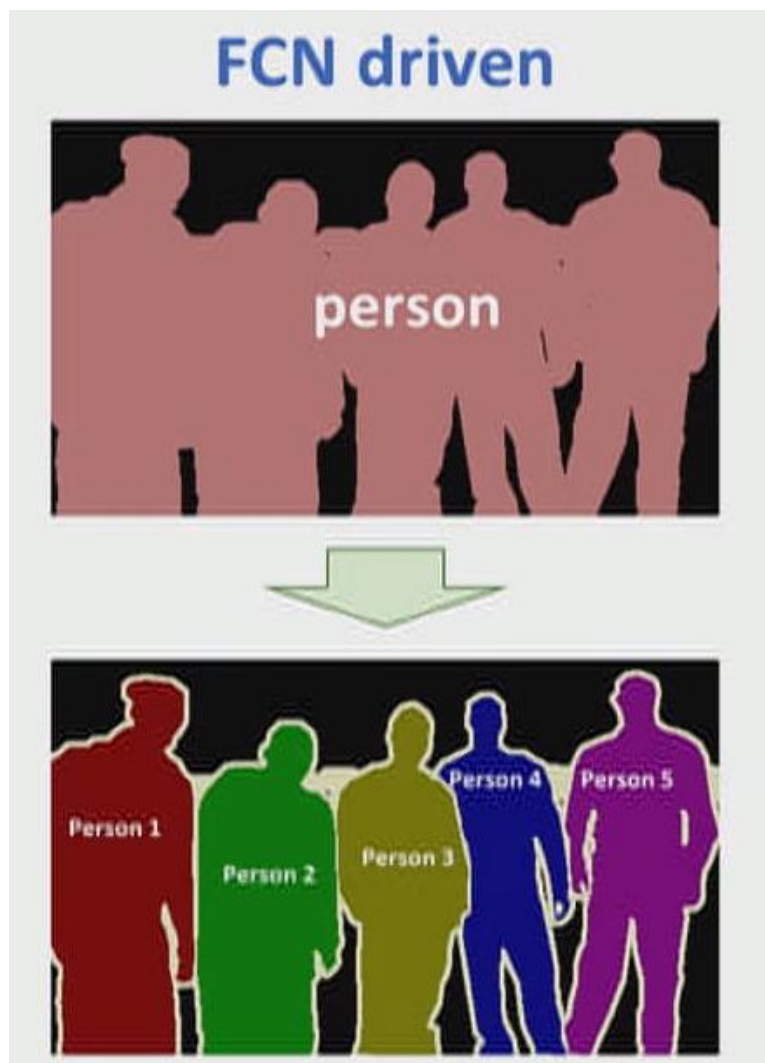
3. Bilinear Interpolation equation



4. RoIAlign pooling process



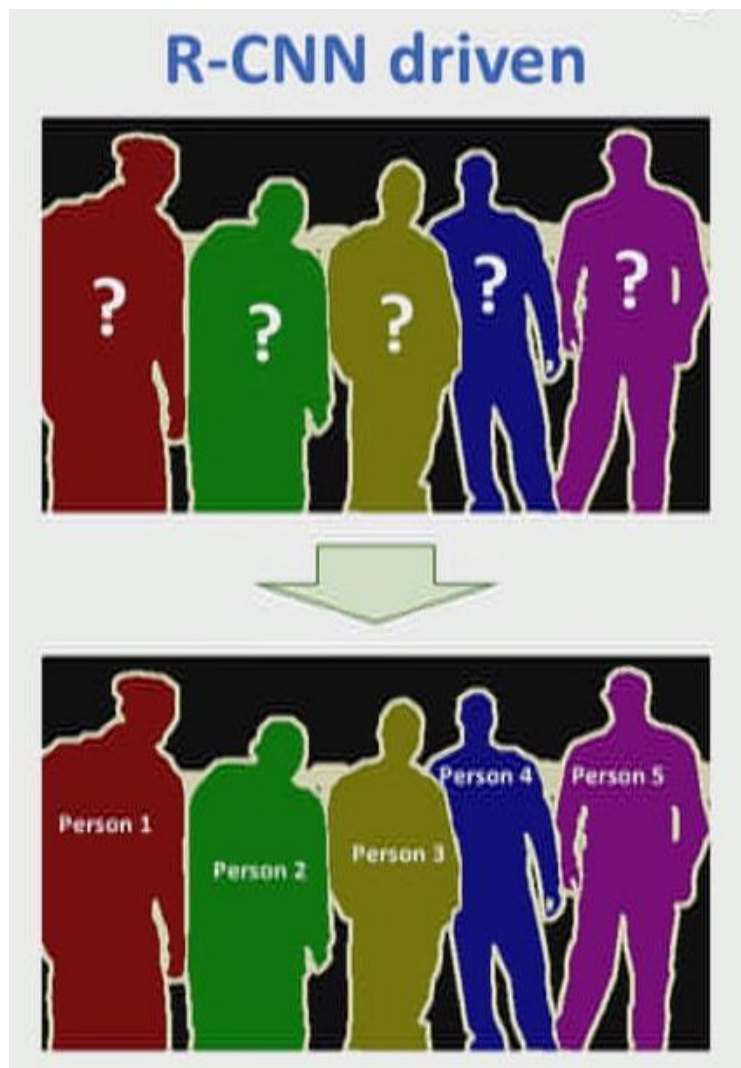
Instance Segmentation Methods



FCN driven methods:
start from a segmentation result,
then learn how divide the results into
individual instances.



Instance Segmentation Methods



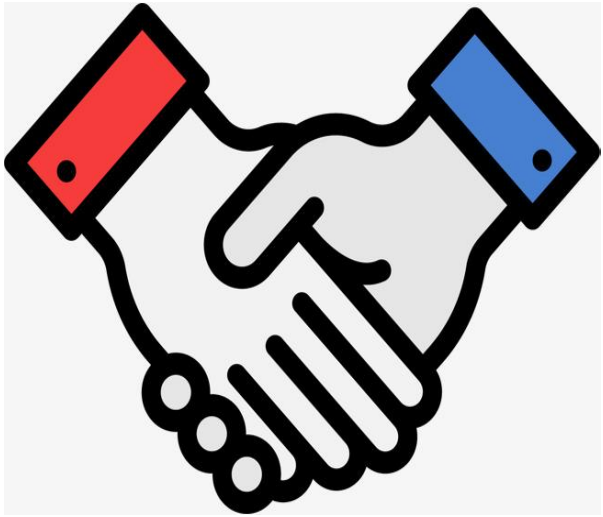
RCNN driven methods:
start from segmentation level
proposals, then train a classifier to
classify these proposals into
semantic categories.

Conclusions



Conclusion

- Semantic Segmentation
 - Using convolution
 - downsampling and upsampling
 - Image pyramids
 - Skip connections
 - Dilation
 - ROI align



Thanks



zhengf@sustc.edu.cn