

2D Object Detection and Segmentation

Jiayuan Gu
2019.10.29

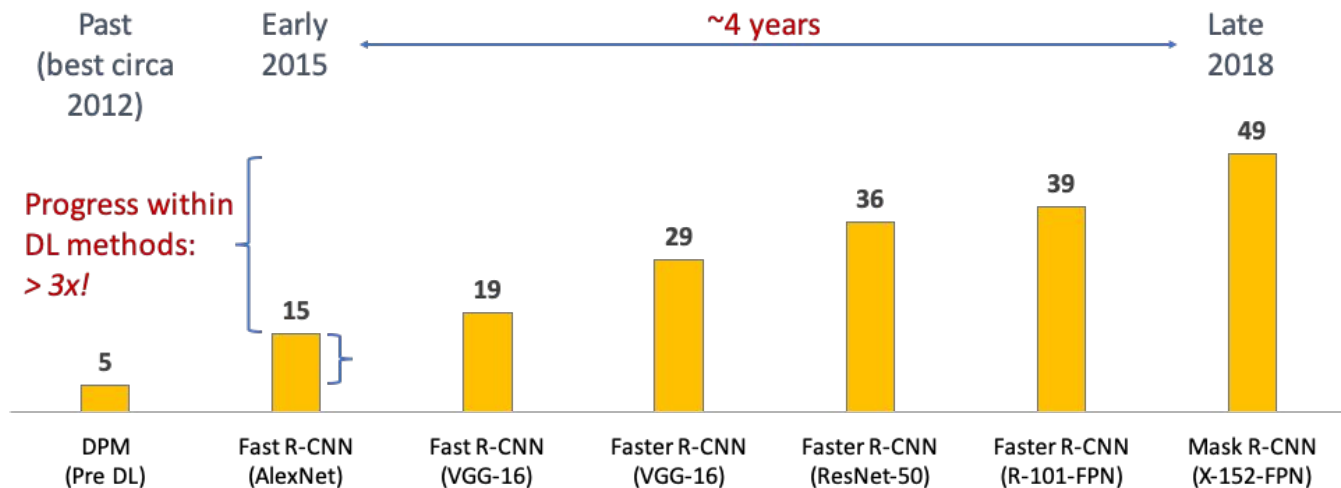
Outline

- What is object detection and segmentation?
- Object detection: R-CNN
- Segmentation: U-Net



Why do we focus on R-CNN?

COCO Object Detection Average Precision (%)



Resources

Tutorials of CVPR

- <http://deeplearning.csail.mit.edu>
- <https://sites.google.com/view/cvpr2018-recognition-tutorial>
- <http://feichtenhofer.github.io/cvpr2019-recognition-tutorial>

Github Repo

<https://github.com/facebookresearch/detectron2>

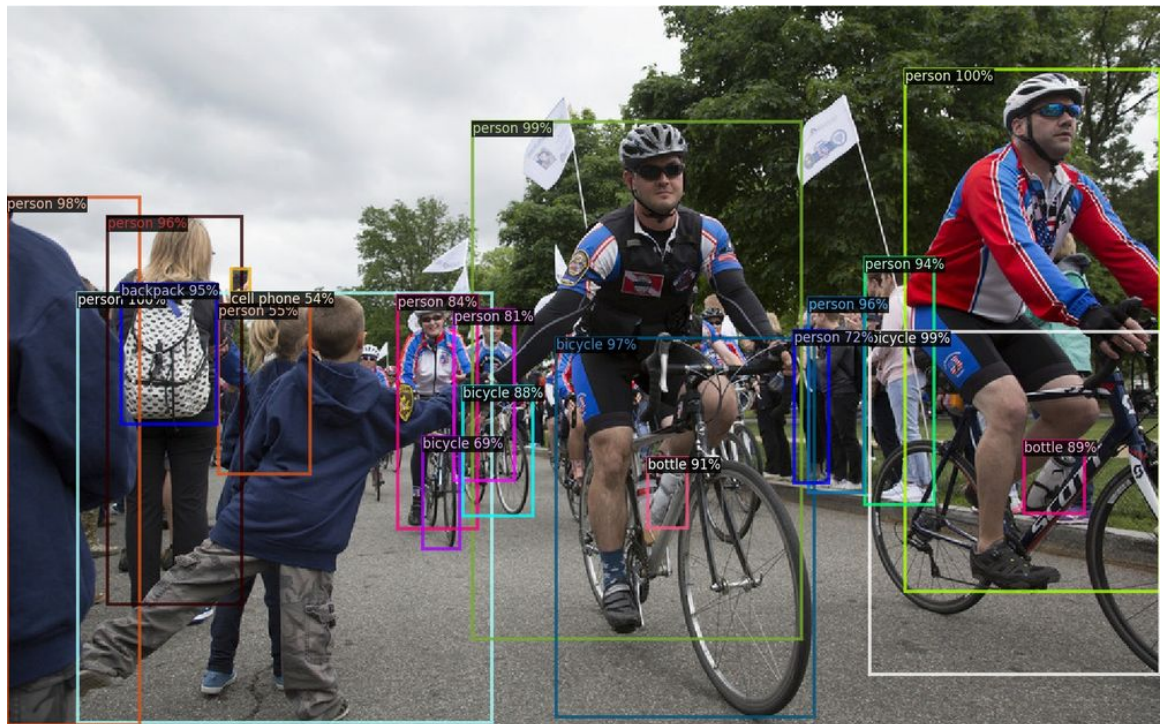


The background is a solid dark blue color. In the top right corner, there is a decorative pattern of overlapping triangles in various shades of blue and white, creating a geometric, stepped effect.

Background

How do we represent objects

- Bounding box



How do we represent objects

- Bounding box
- Instance mask



How do we represent objects

- Bounding box
- Instance mask
- Keypoint

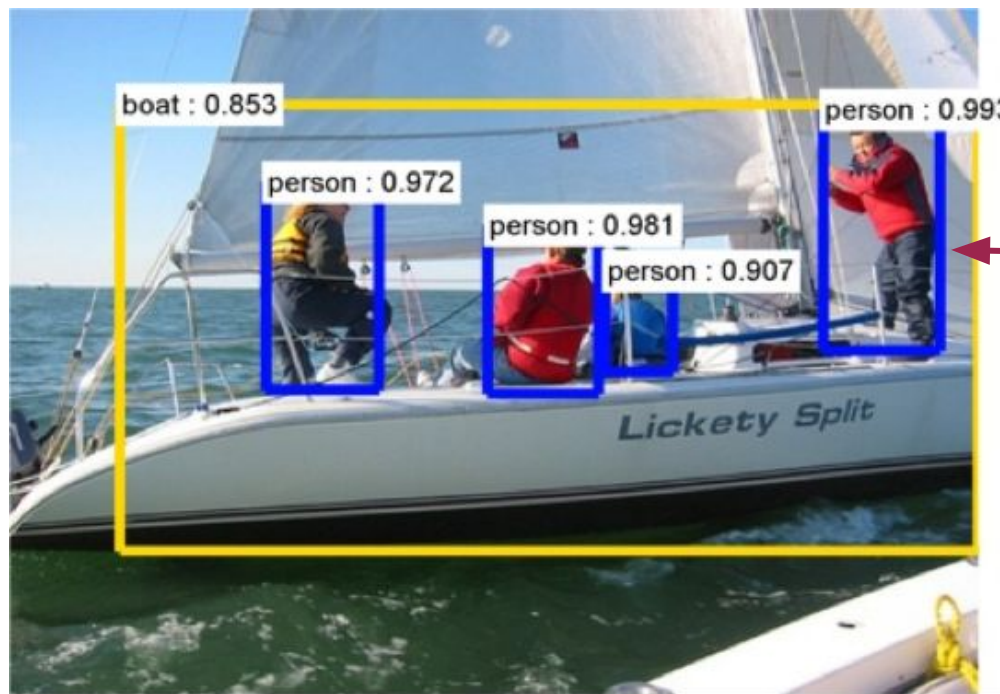


How do we represent objects

- **Bounding box**
- Instance mask
- Keypoint



Object Detection with Bounding Boxes



What? - Recognition/Classification

Where? - Localization/Regression

“Object detection”

Object Detection with Segmentation Masks



What? - Recognition

Where? - Segmentation

"Instance segmentation"

Semantic Segmentation

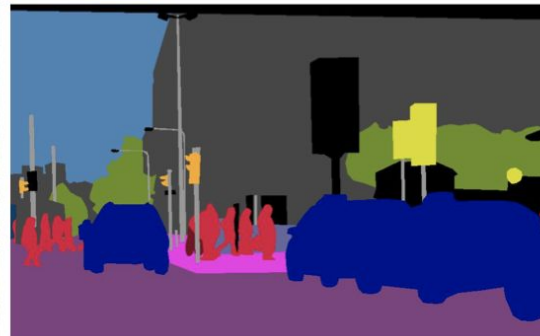
Predict a pixel-wise class label

Stuff: walls, buildings, sky, road

Things: human, cars, bikes



(a) image



(b) semantic segmentation



(c) instance segmentation



(d) panoptic segmentation

Datasets



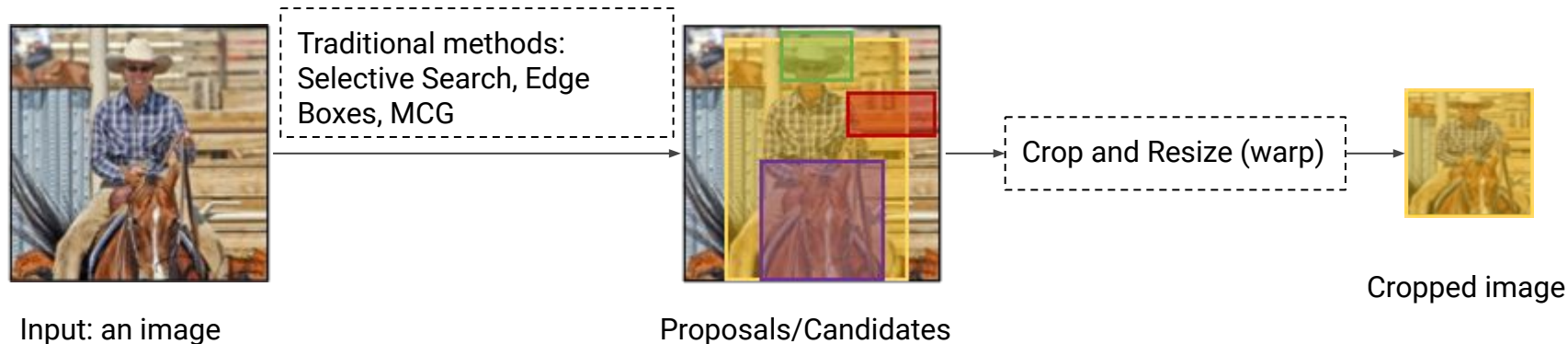
Microsoft COCO



Visual Object Classes Challenge 2012 (VOC2012)

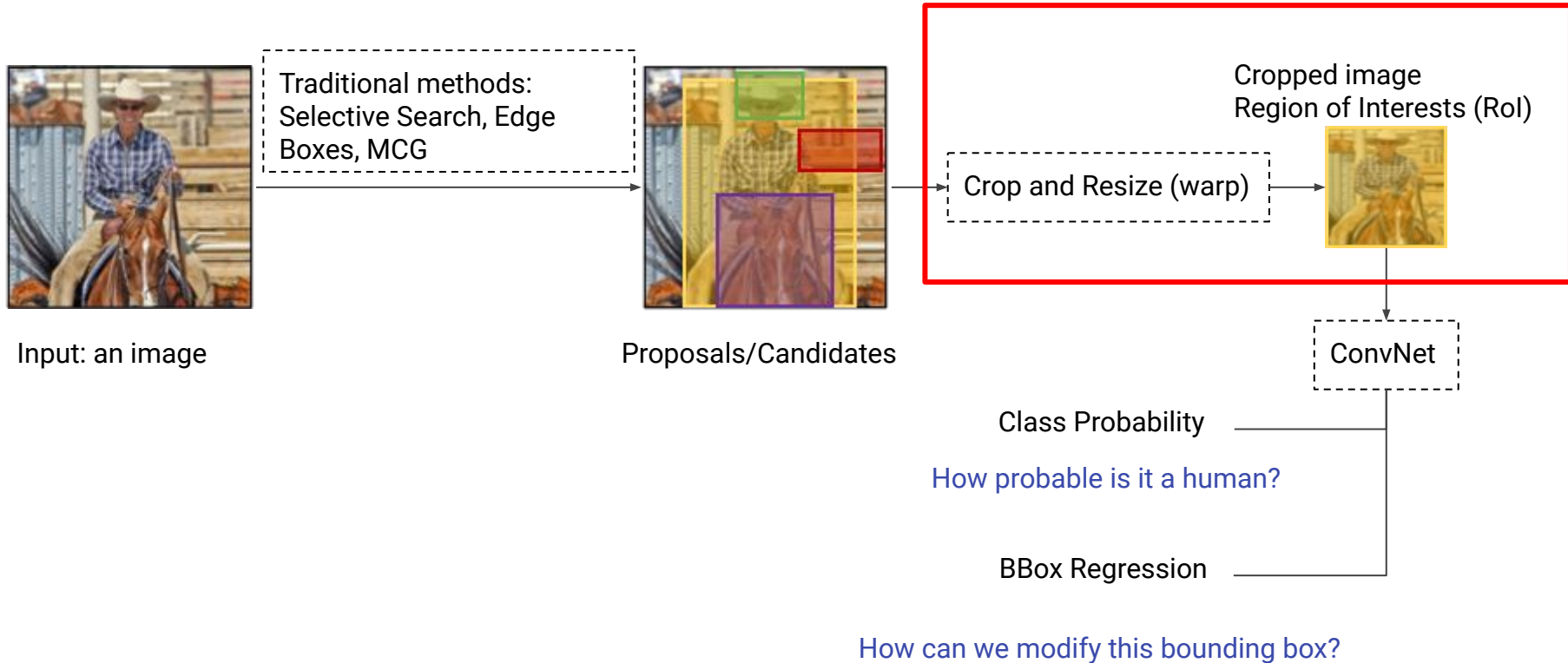
R-CNN: Region-based CNN

Object Detection → Object Classification



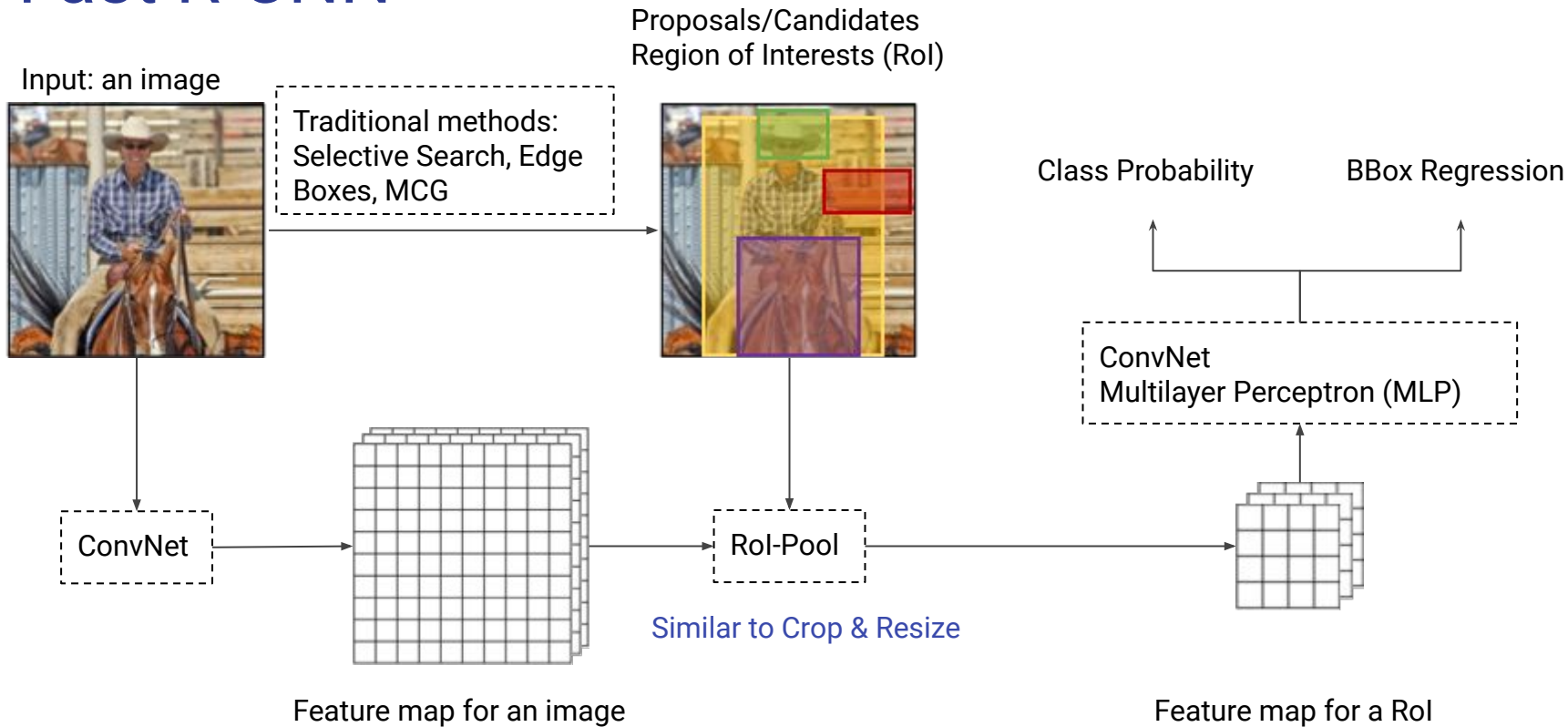
We've already reduced object detection to object classification!

R-CNN

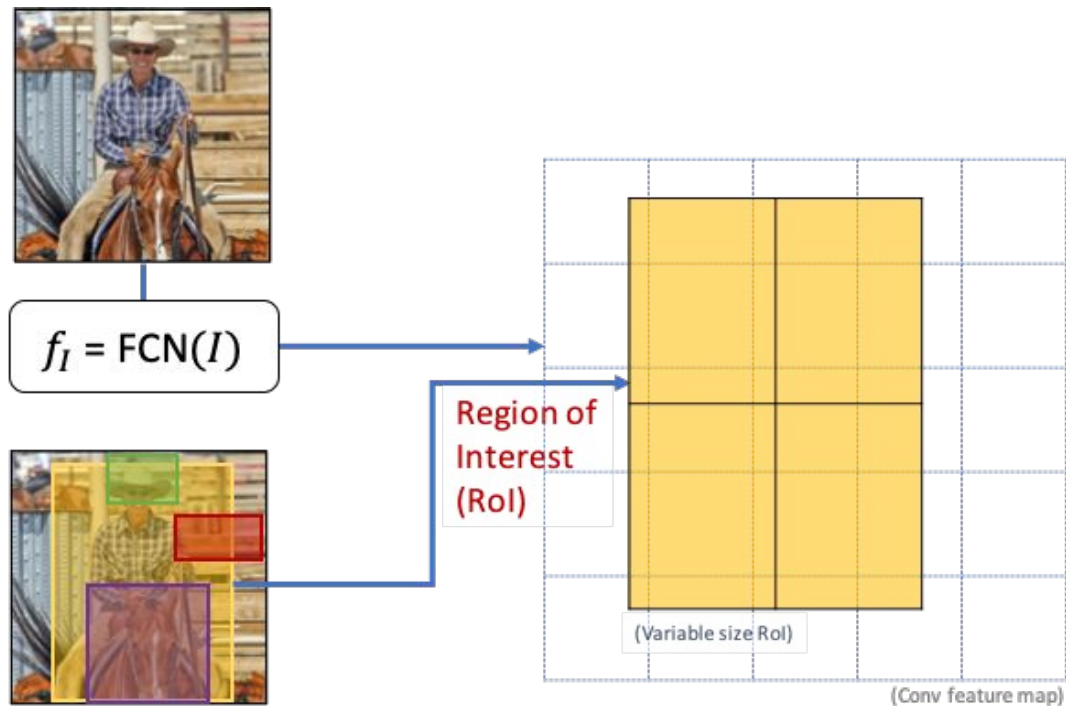


R-CNN: Crop \rightarrow CNN v.s. Fast R-CNN: CNN \rightarrow Crop

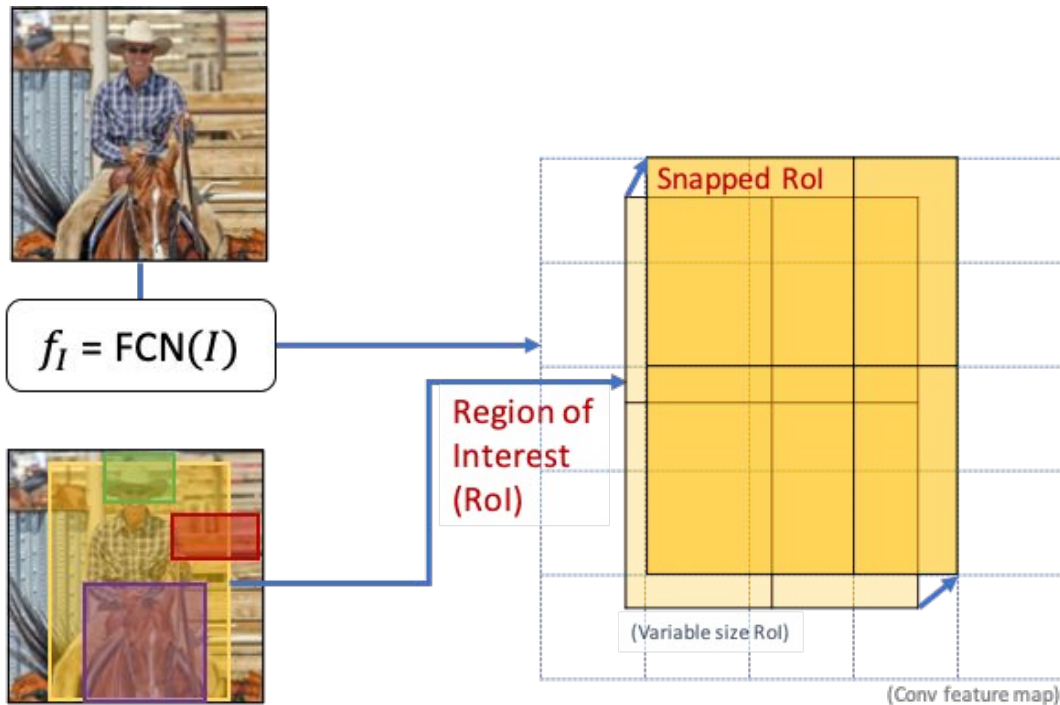
Fast R-CNN



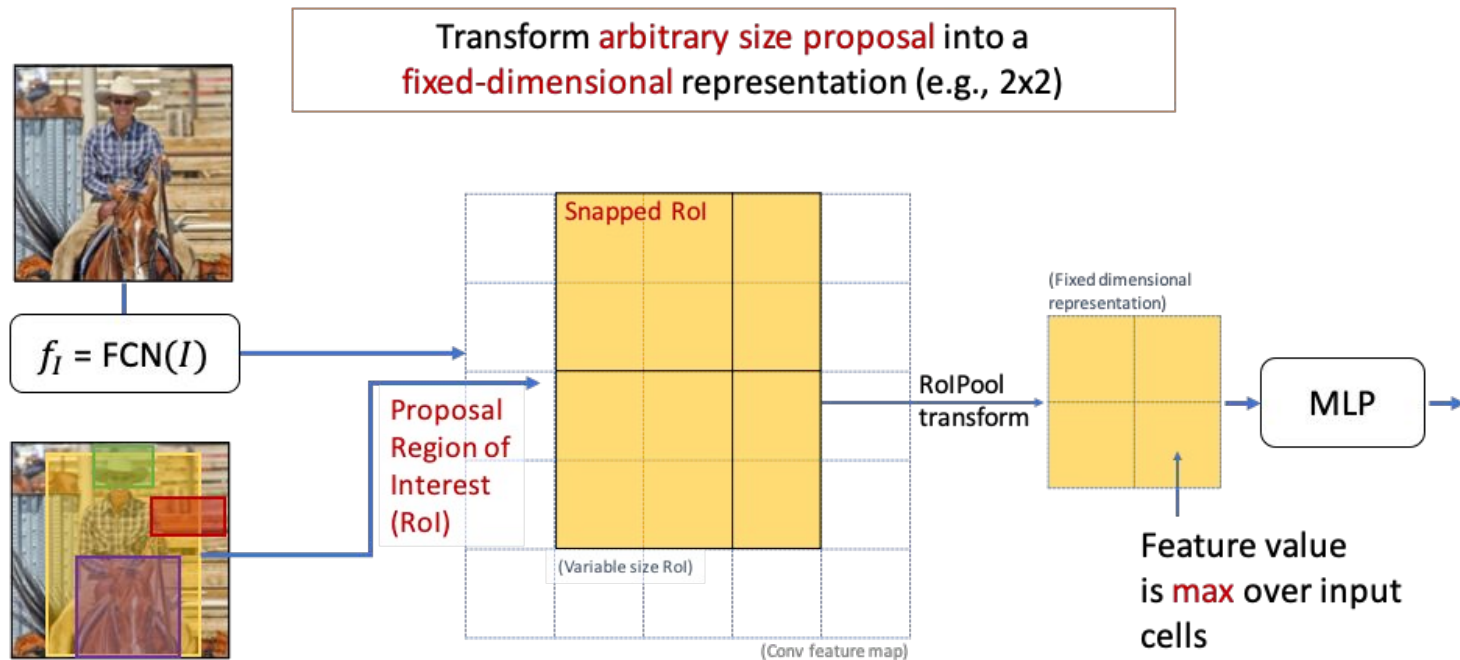
RoI Pooling (for each proposal)



RoI Pooling (for each proposal)

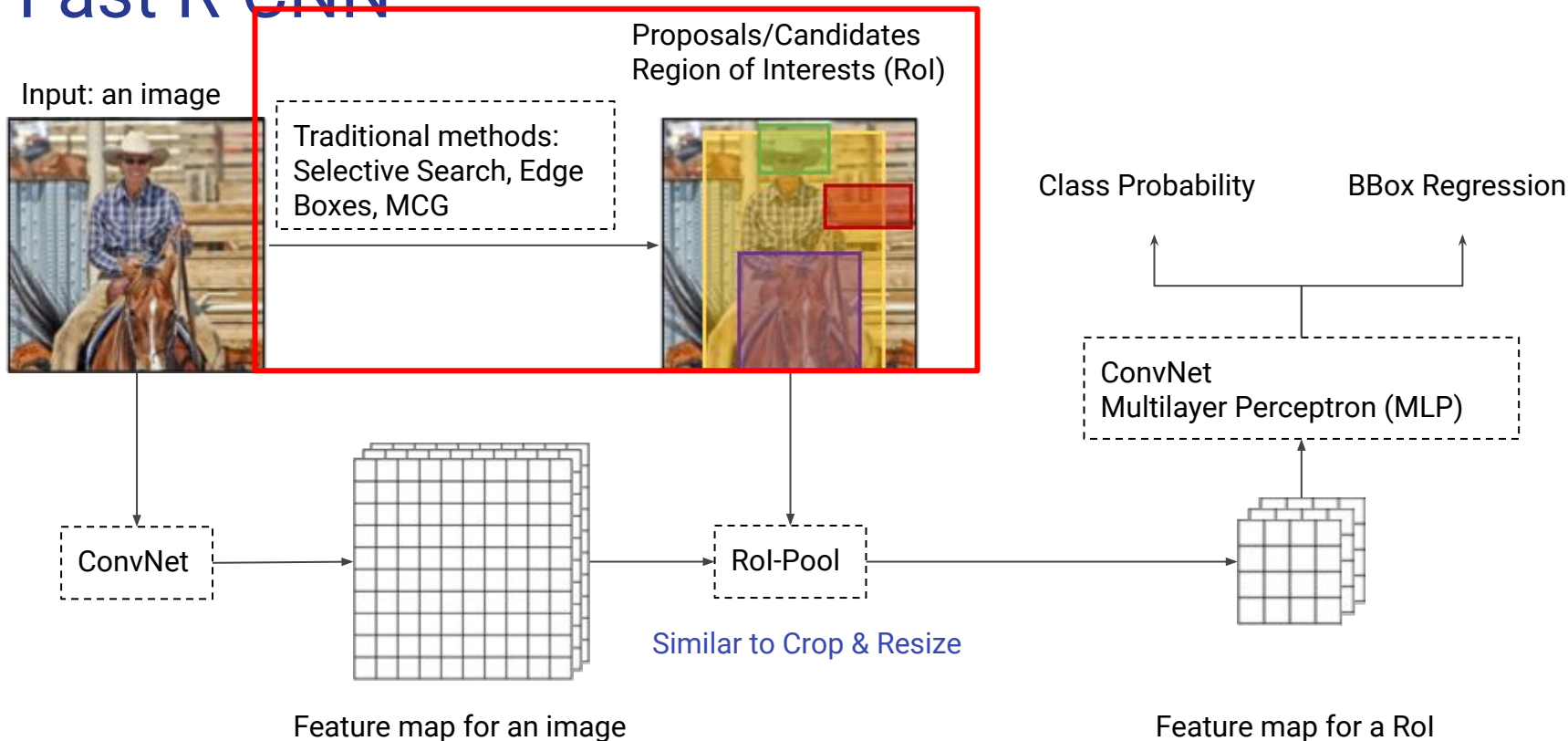


Rol Pooling (for each proposal)

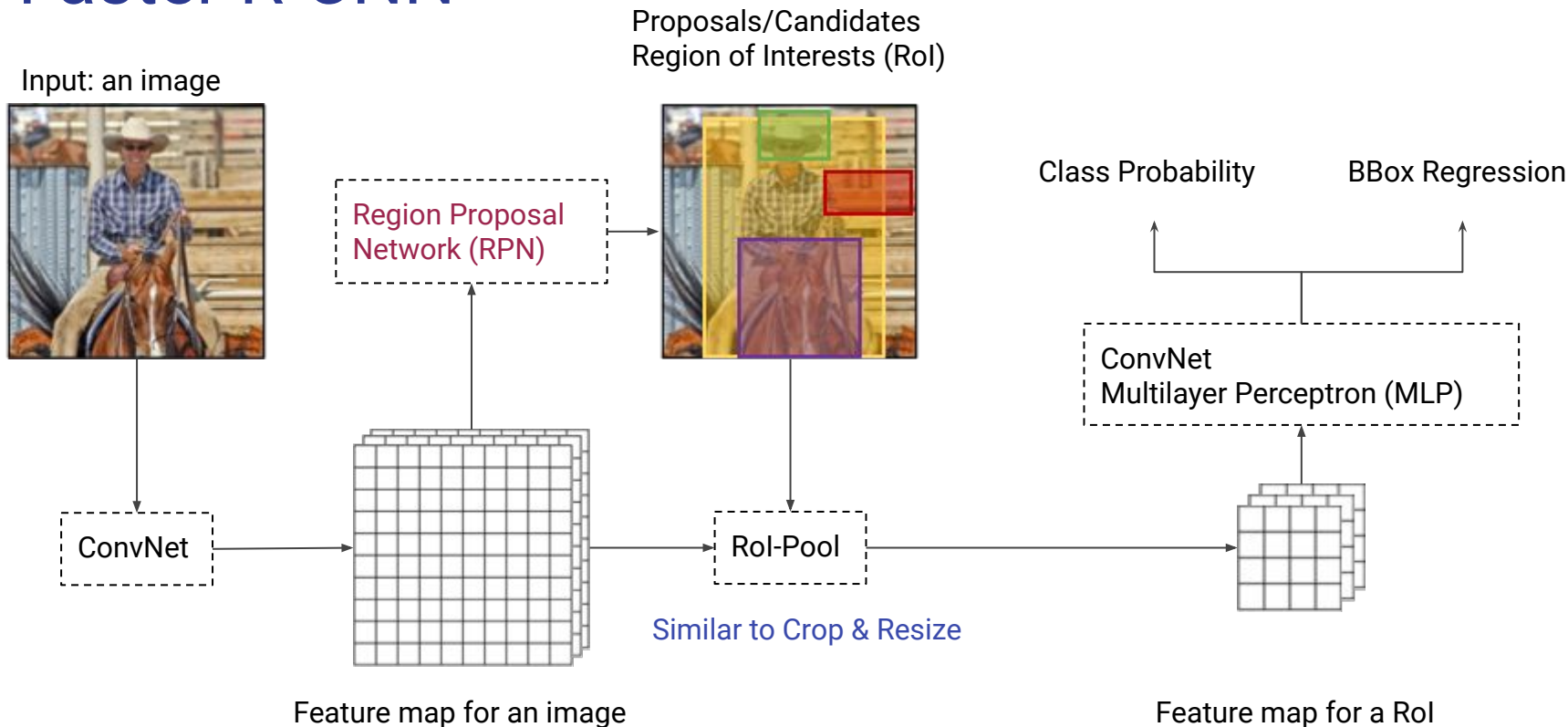


Fast R-CNN

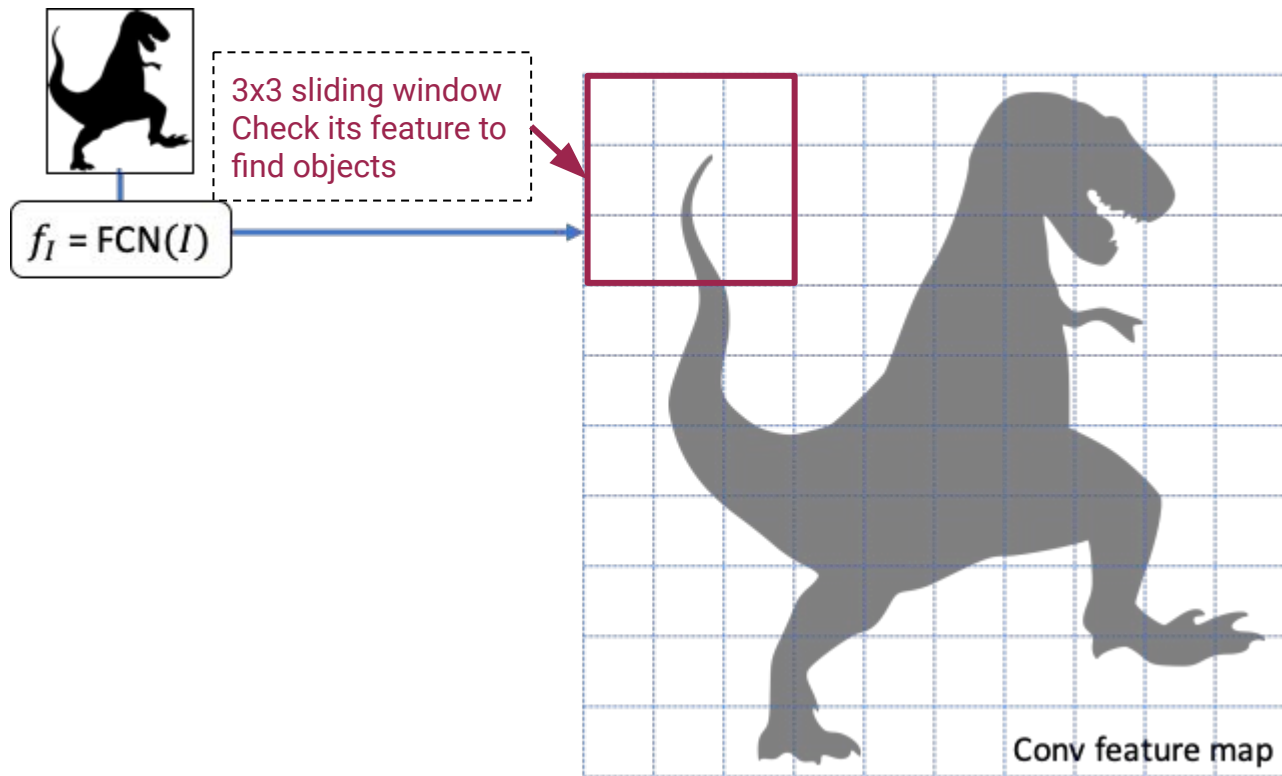
Not good enough



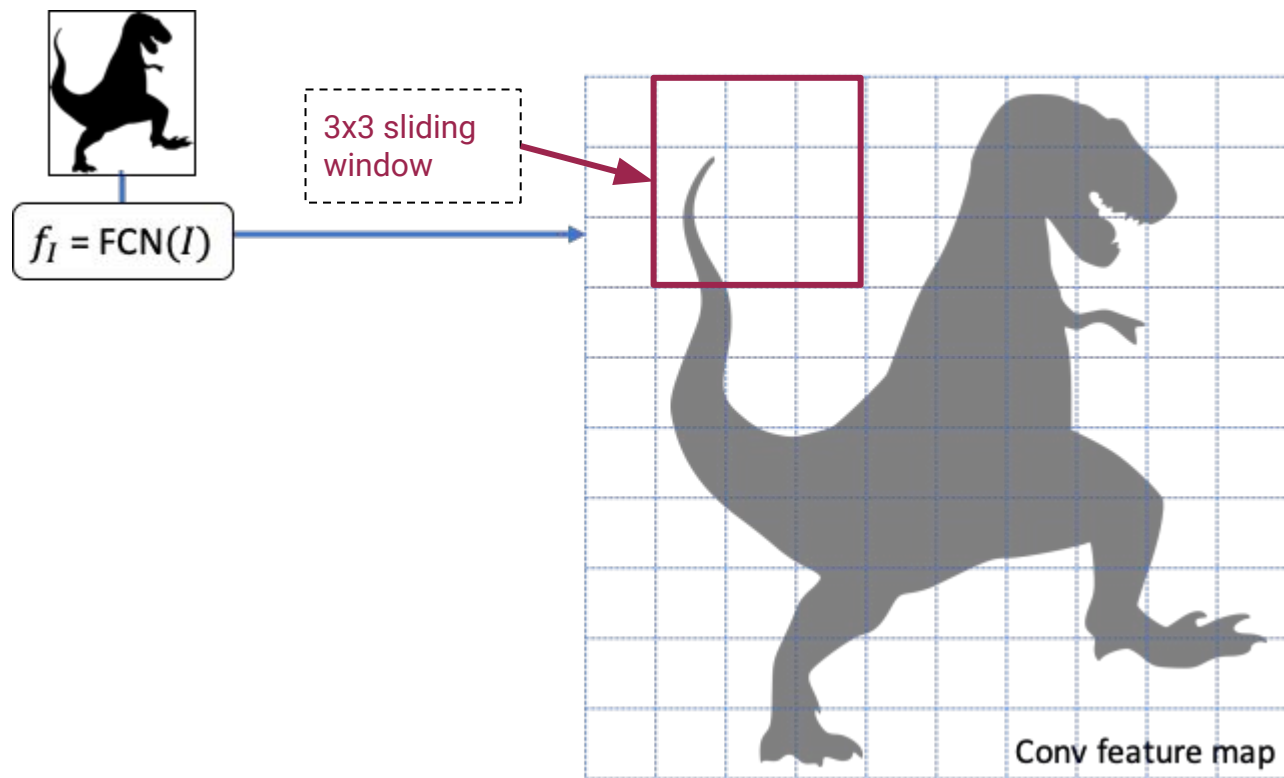
Faster R-CNN



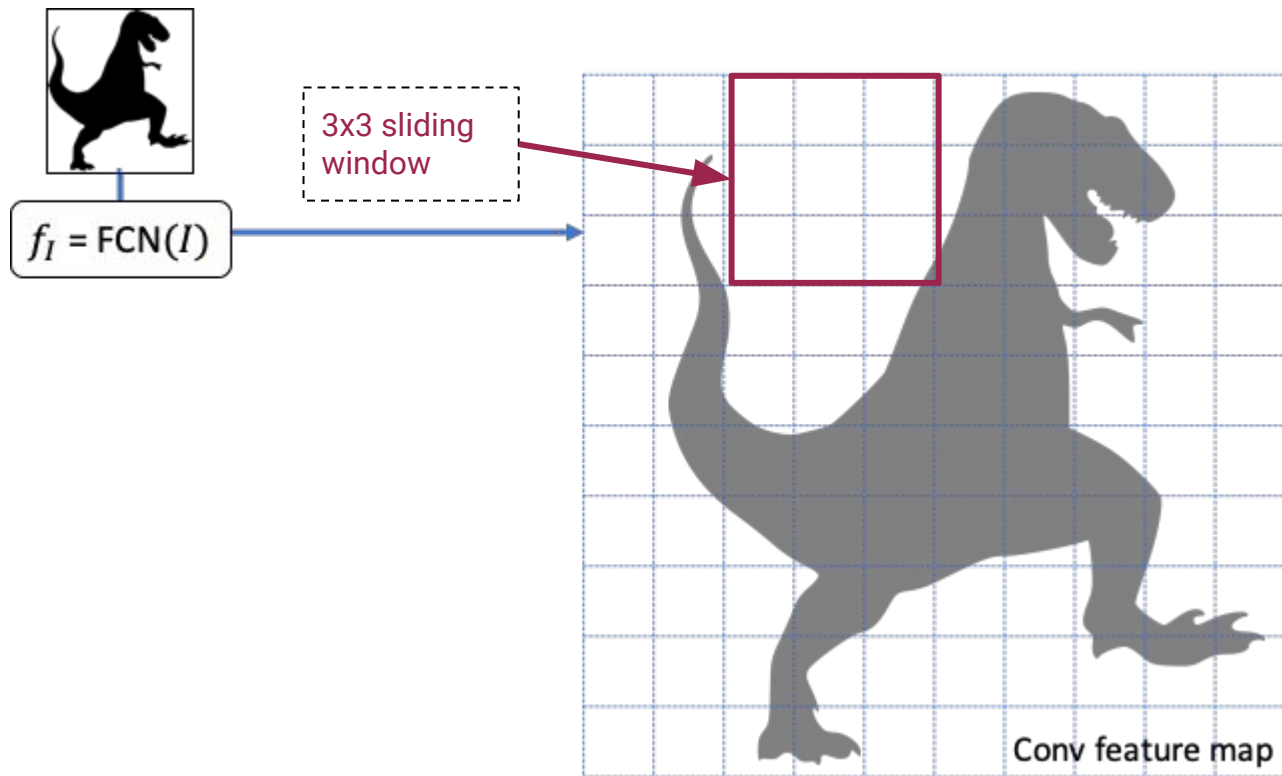
RPN: Region Proposal Network



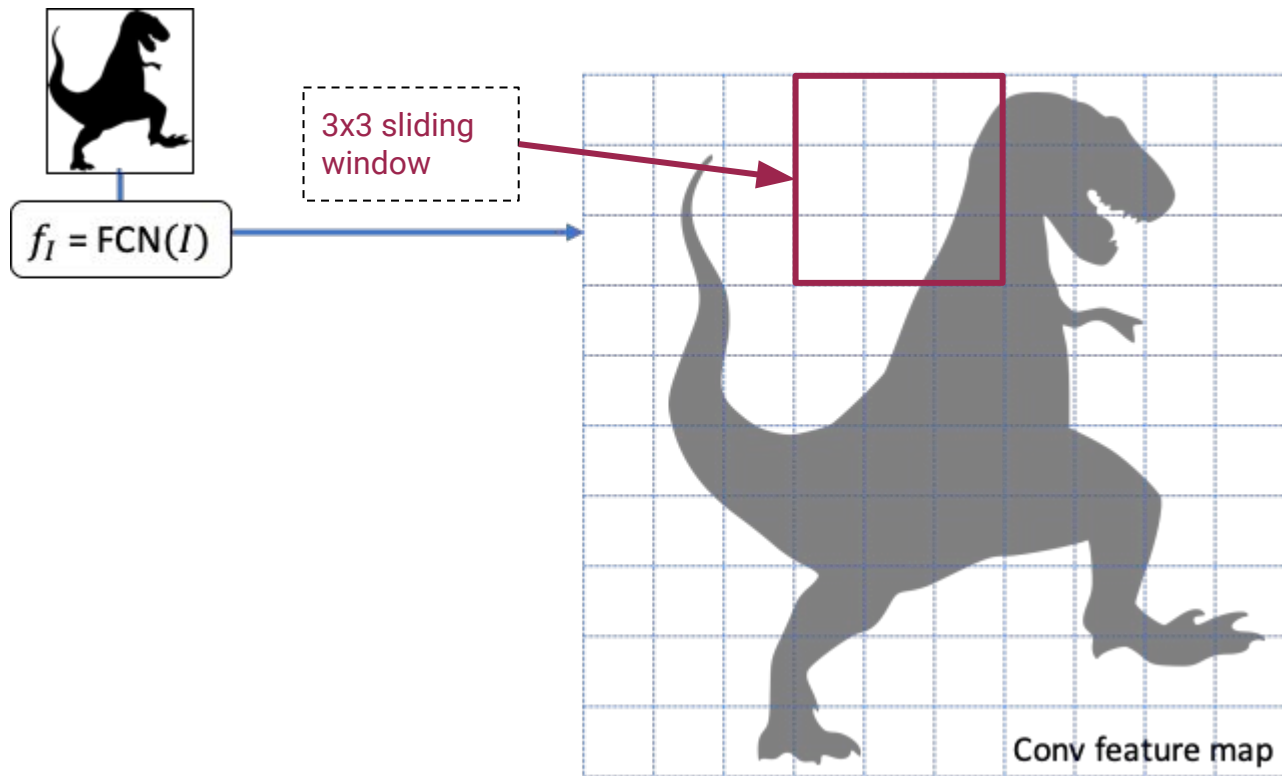
RPN: Region Proposal Network



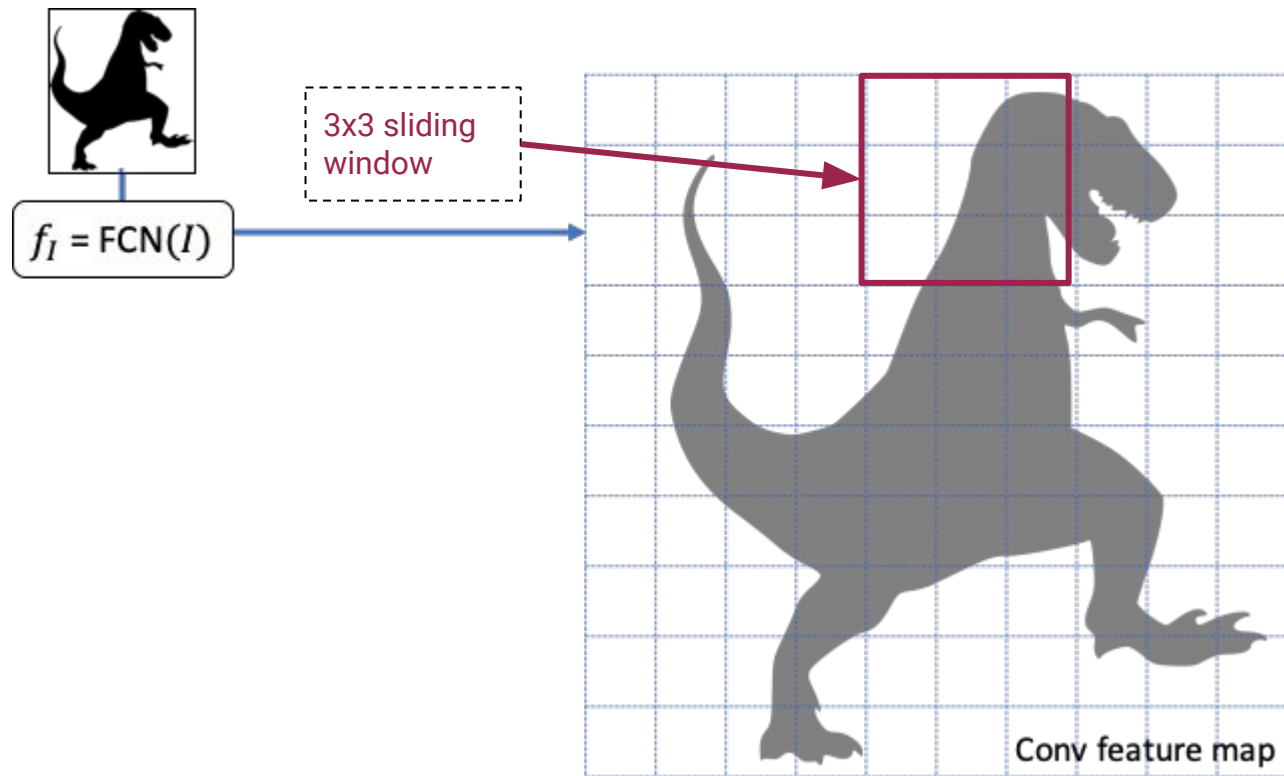
RPN: Region Proposal Network



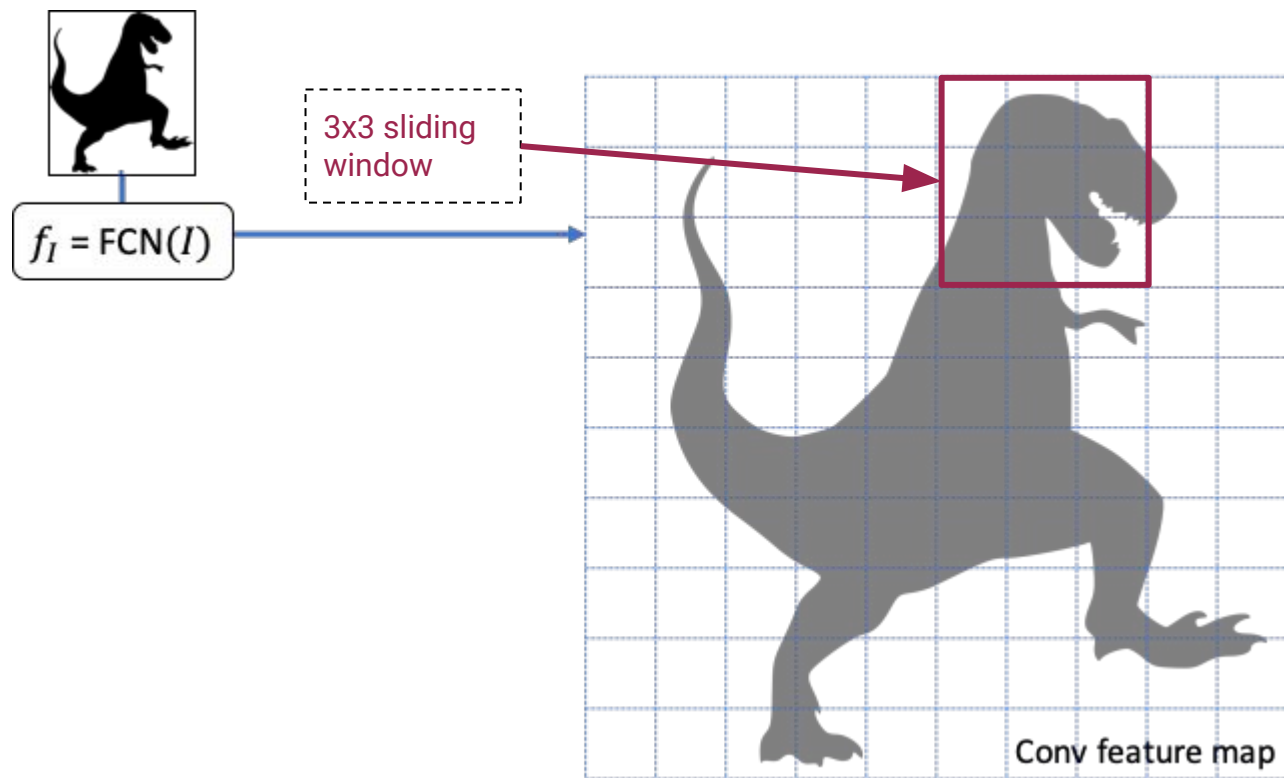
RPN: Region Proposal Network



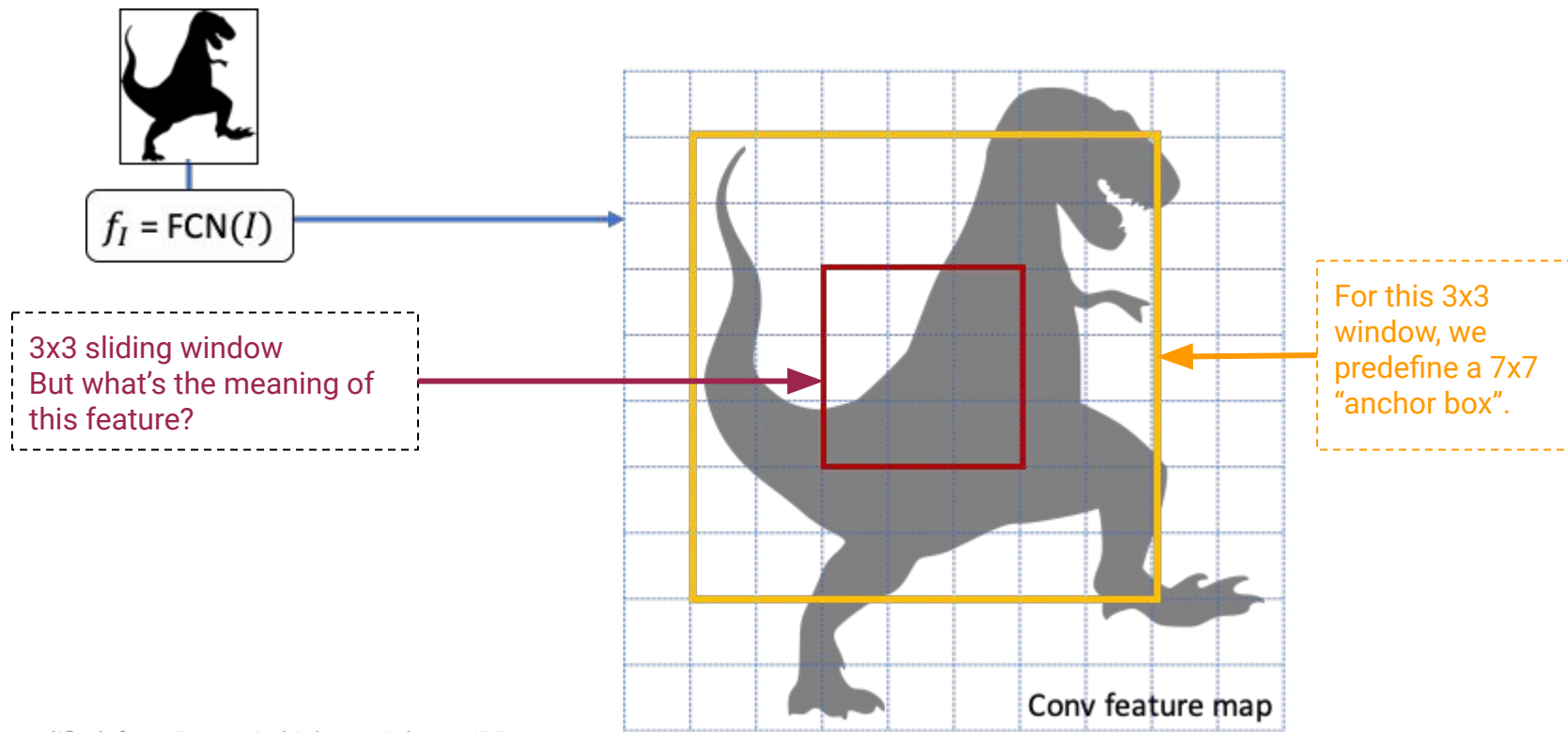
RPN: Region Proposal Network



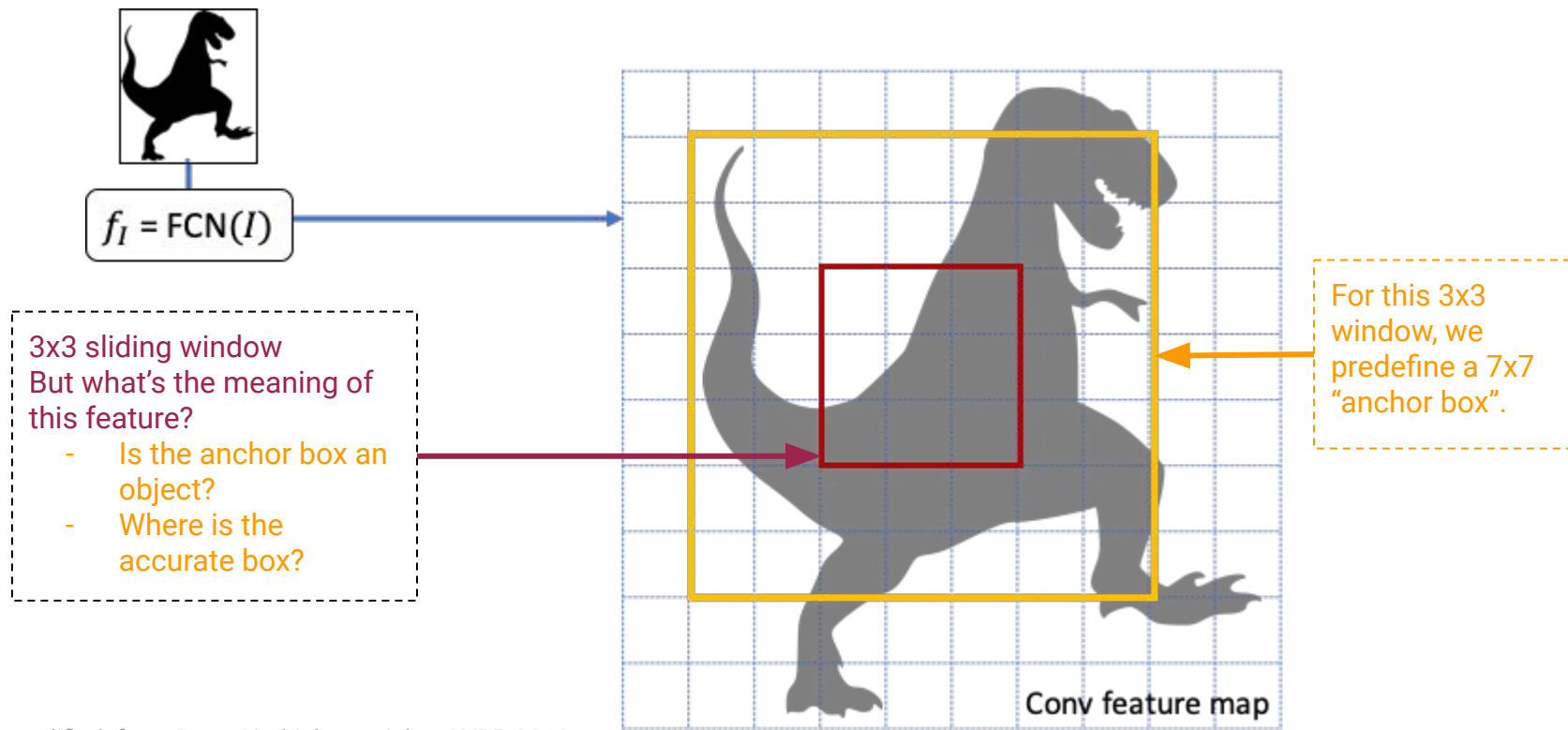
RPN: Region Proposal Network



RPN: Anchor Box

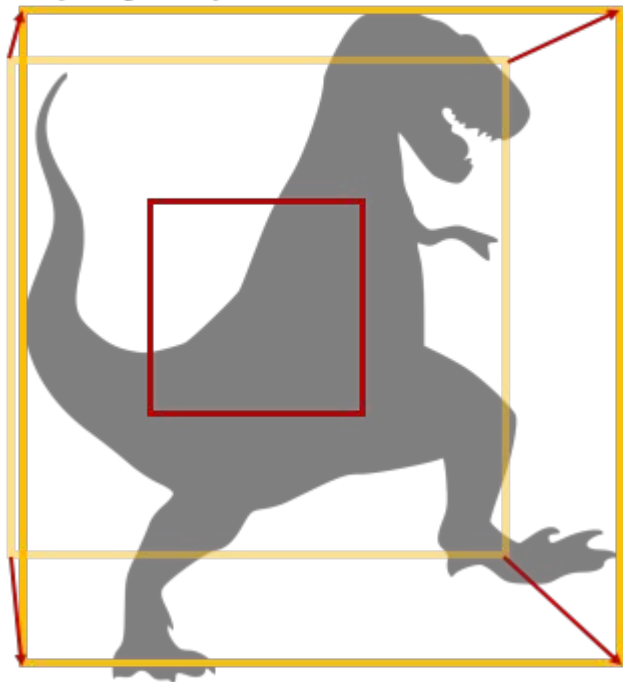


RPN: Anchor Box



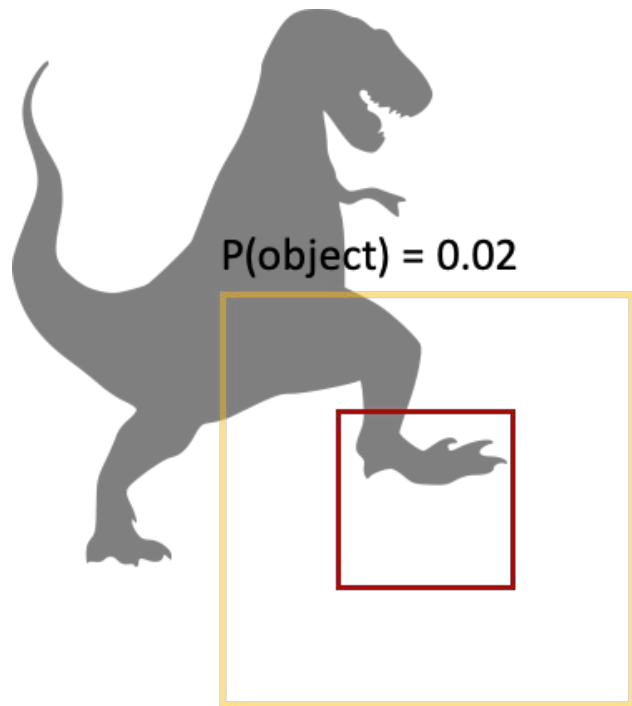
RPN: Prediction (on object)

$P(\text{object}) = 0.94$

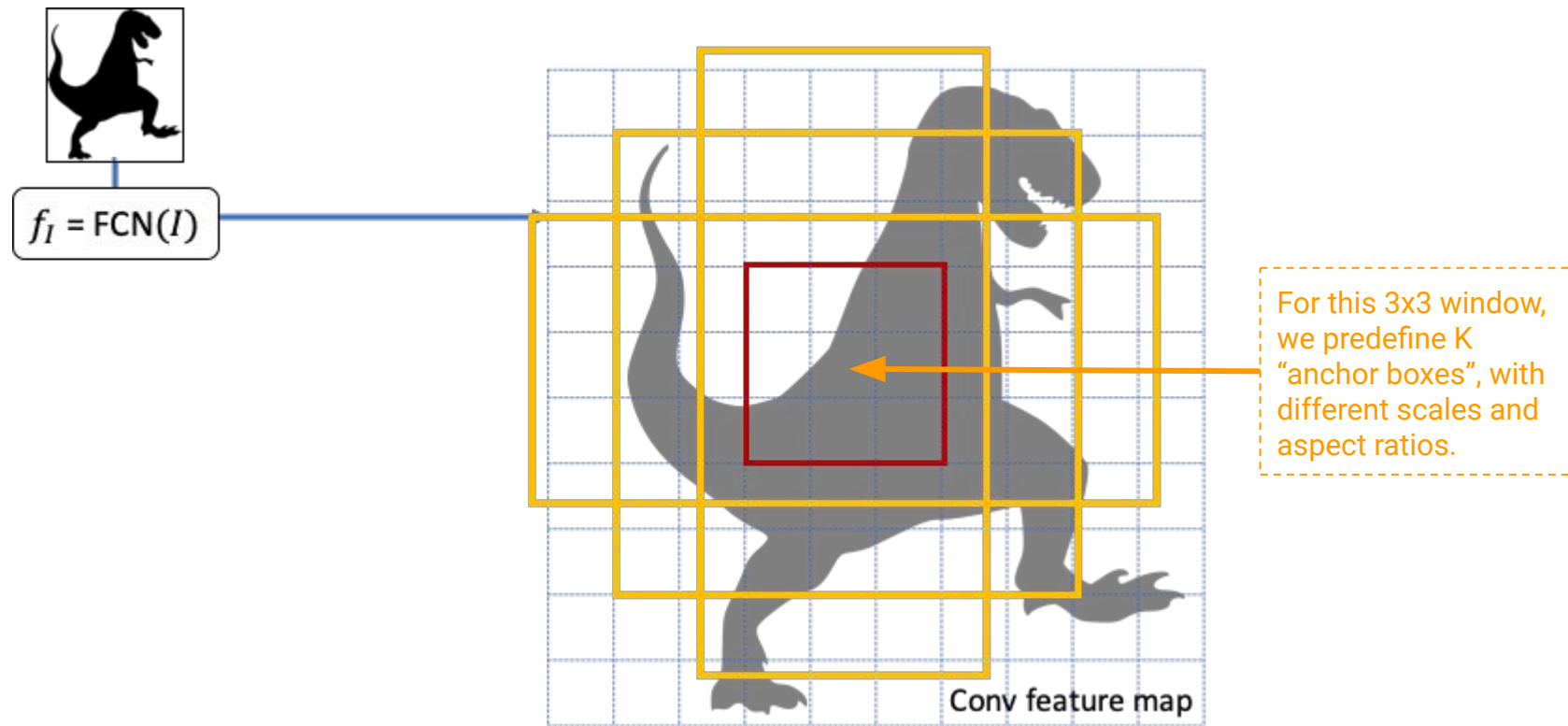


→ Direction to the accurate box

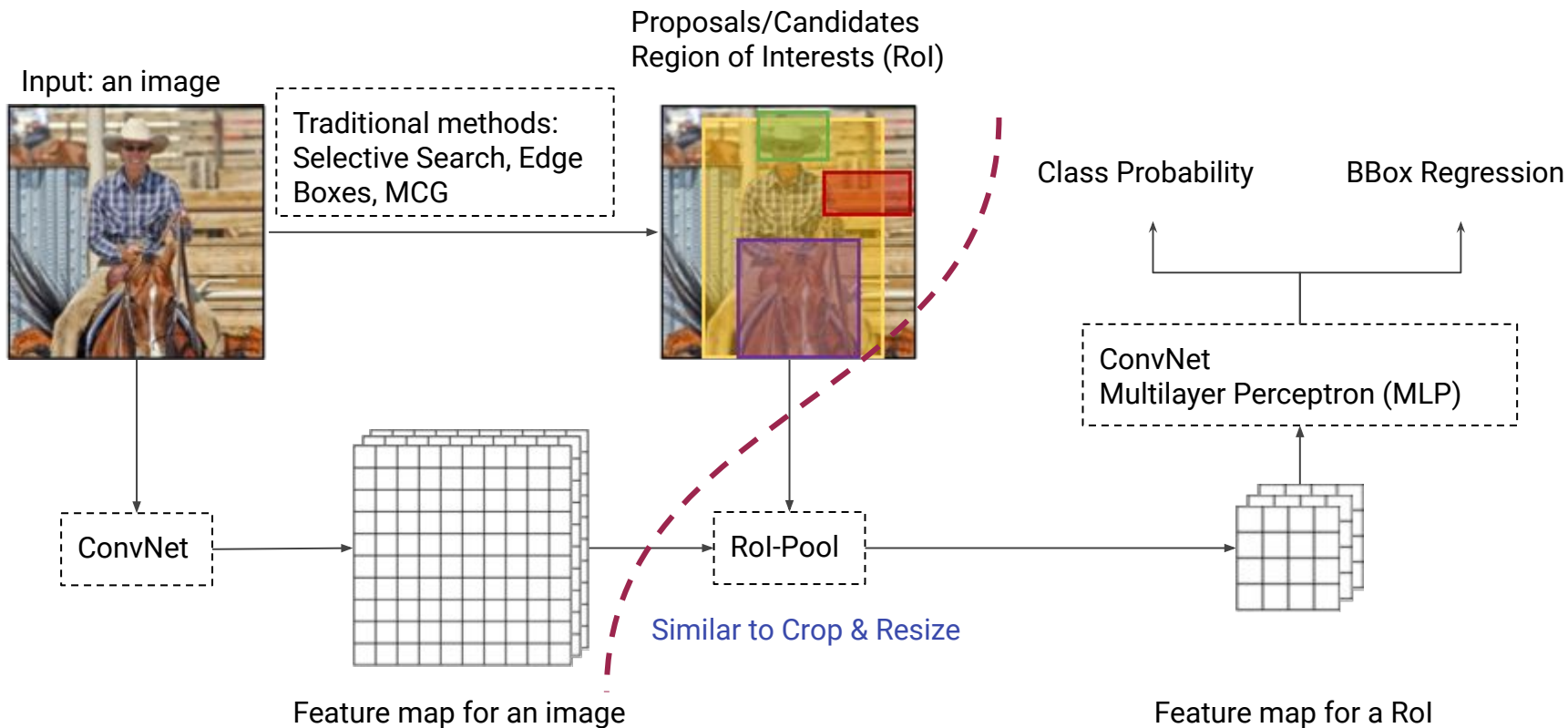
RPN: Prediction (off object)



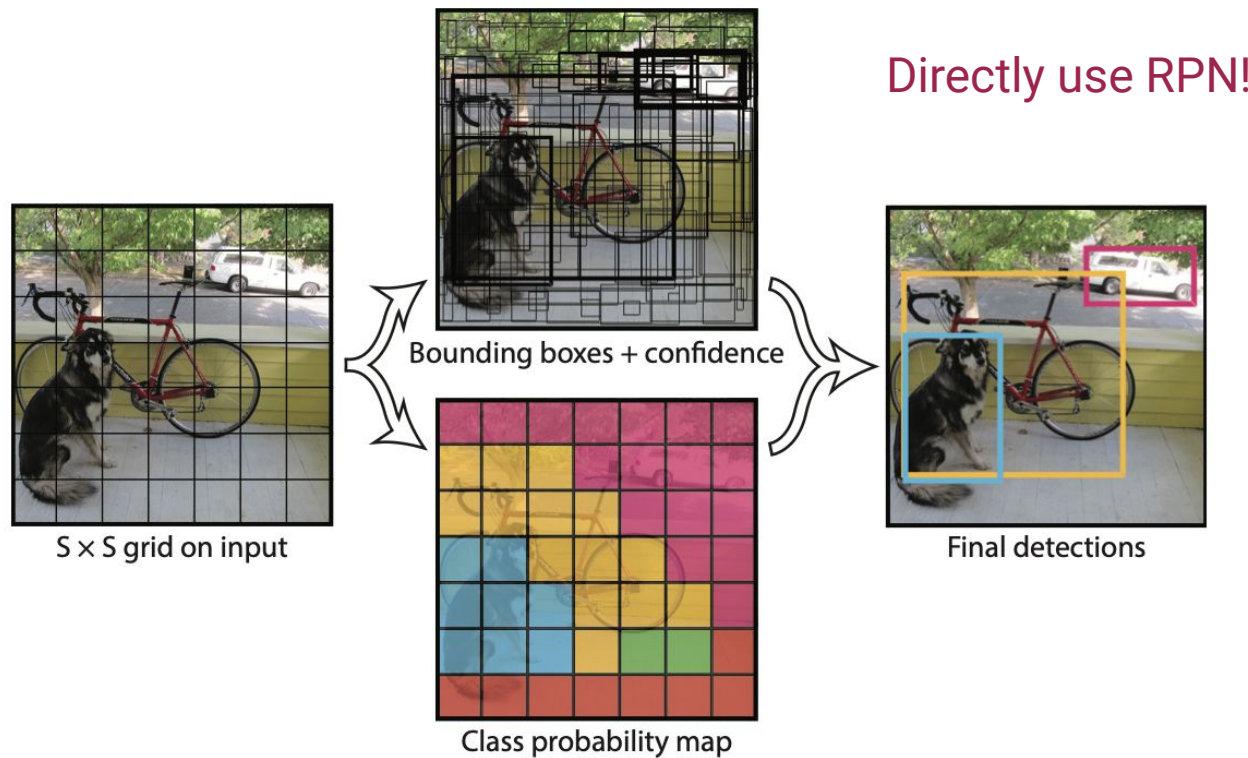
RPN: Multiple Anchors



Two stages or one stage?



You Look Only Once (YOLO)



Other Methods

http://web.stanford.edu/class/cs231a/lectures/lecture12_2D_detection.pdf

- VJ Face
- Deformable Part Model
- Implicit Shape Model



U-Net

Semantic vs. Instance Segmentation

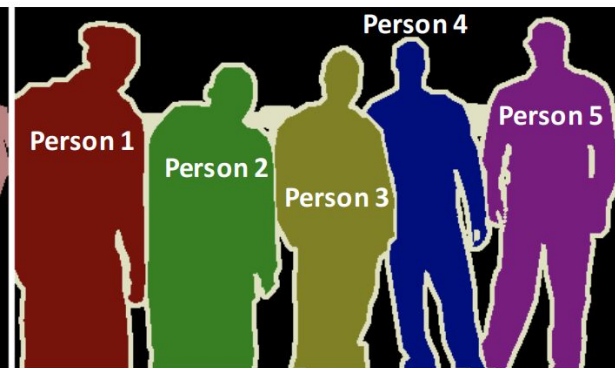
Object detection



Semantic segmentation



Instance segmentation

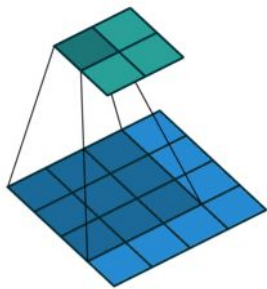


How shall we modify CNN?

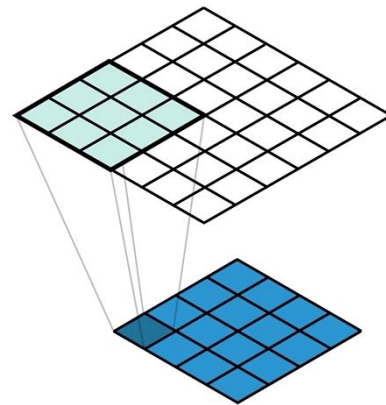
- We need to predict a label for each pixel, but CNN has downsampled our input to a very small scale (e.g. from 224×224 to 7×7).
- Thus, we need a layer to upsample feature maps to the original size.



Transposed Convolution



3x3 Convolution



3x3 Transposed
Convolution

Encoder-Decoder

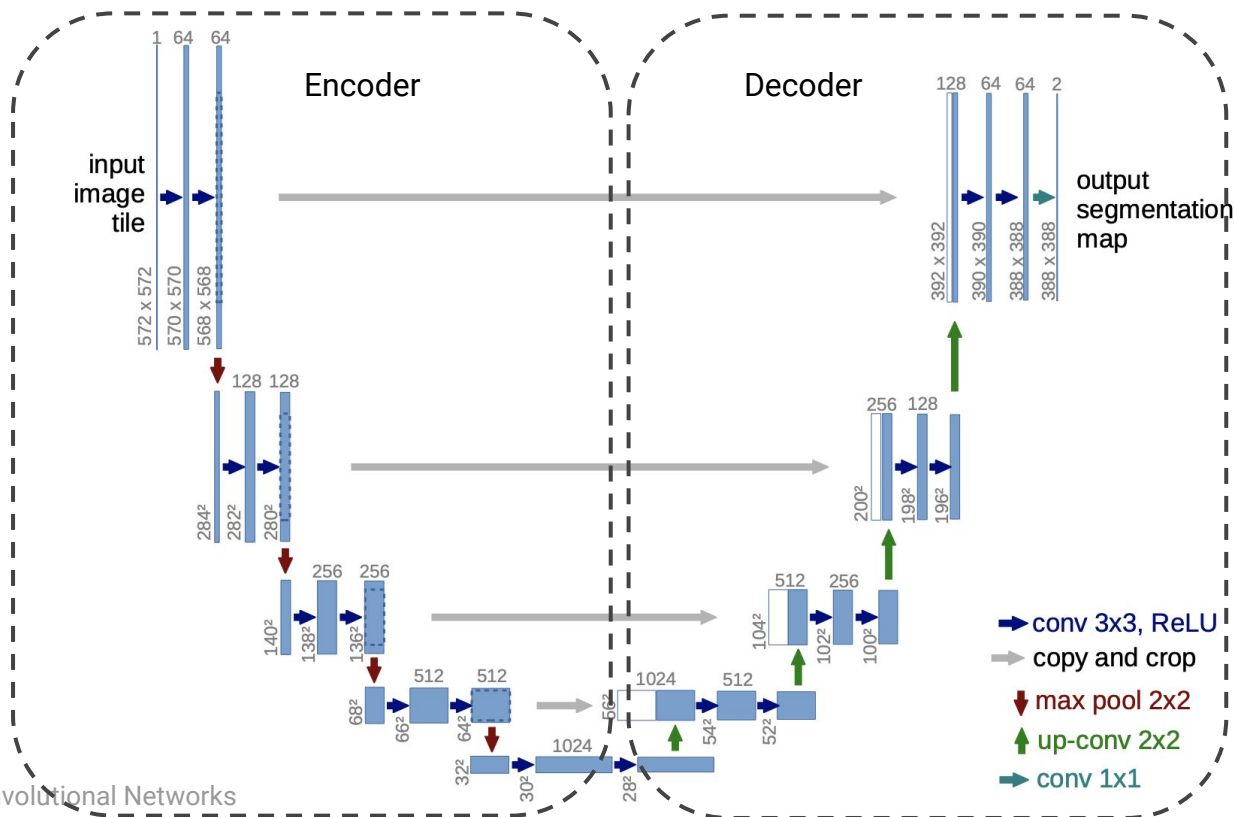


Figure from U-Net: Convolutional Networks for Biomedical Image Segmentation

Transposed Convolution

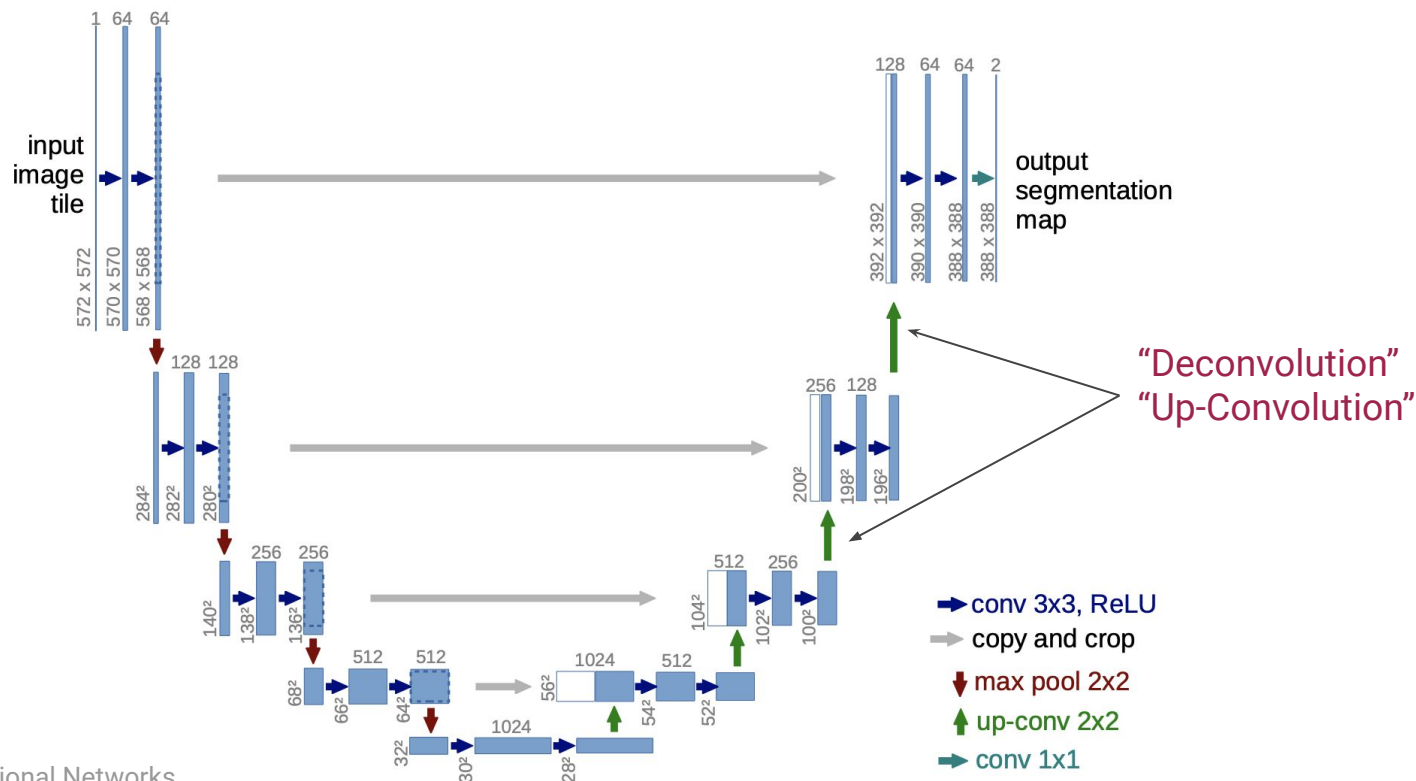


Figure from U-Net: Convolutional Networks for Biomedical Image Segmentation

Skip Layer

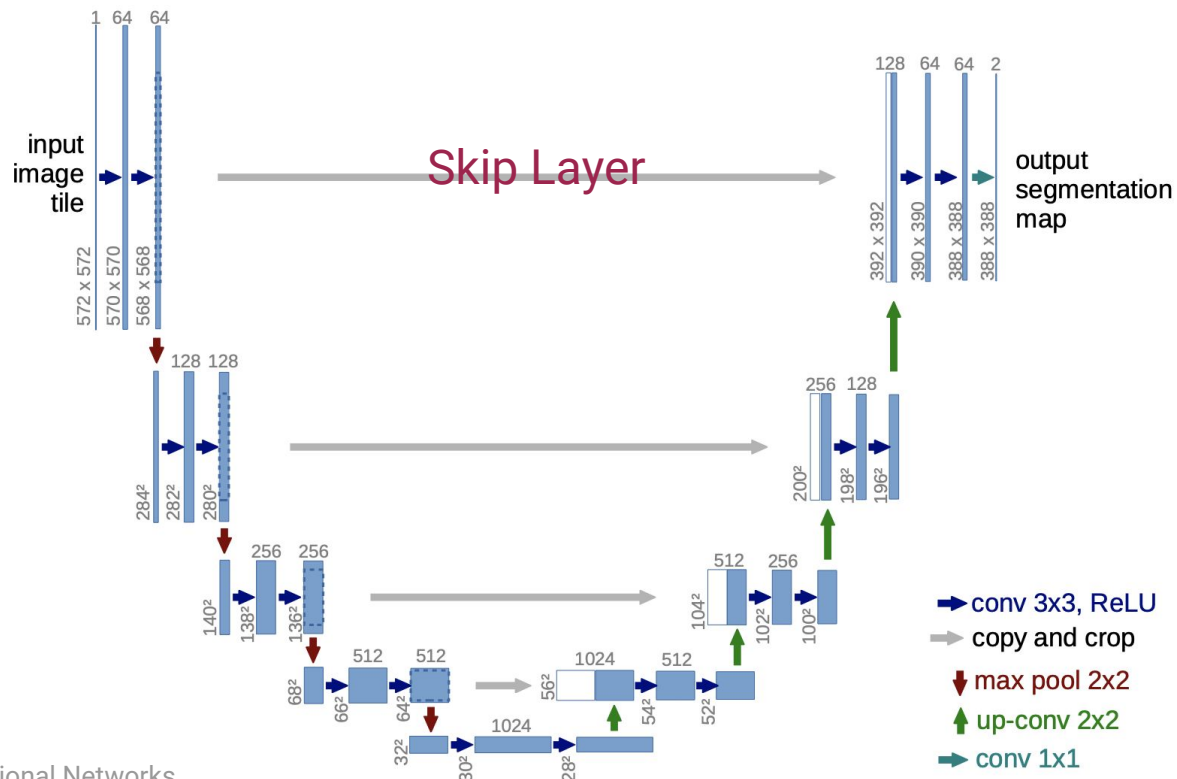


Figure from U-Net: Convolutional Networks for Biomedical Image Segmentation

References

- [1] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks,” pp. 1–9, 2015.
- [2] R. Girshick, “Fast R-CNN,” *Proc. IEEE Int. Conf. Comput. Vis.*, vol. 11-18-Dec, pp. 1440–1448, 2016.
- [3] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, pp. 580–587, 2014.
- [4] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical image computing and computer-assisted intervention*, 2015, pp. 234–241.



The background is a solid pink color. In the top right corner, there is a decorative pattern of overlapping geometric shapes, including triangles and squares, in various shades of pink and magenta.

Thanks