

An illustration of three people flying a kite. A woman on the left, a child in the middle, and a man on the right are holding strings attached to a kite flying in the sky. The image is overlaid with semi-transparent blue boxes and labels: 'Kite' is labeled above the kite, and 'Person' is labeled above each of the three individuals. The background is a light blue gradient.

# Object detection and recognition

Cengiz Öztireli

# Object Detection

The task of assigning a label and a bounding box to all predefined objects in the image

- Localization

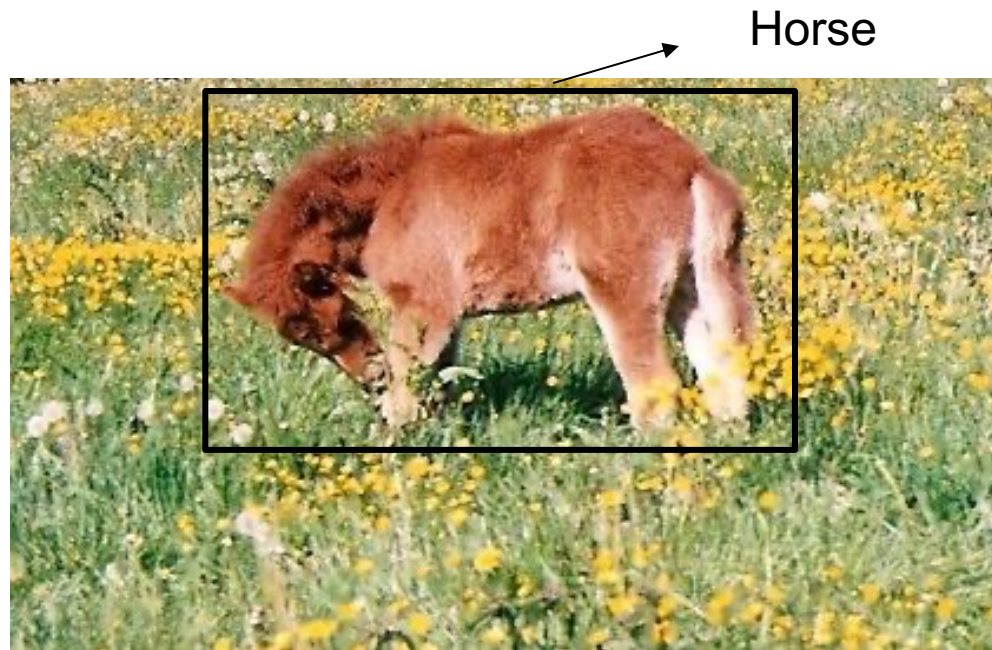
- Use a bounding box to localize the objects of interests



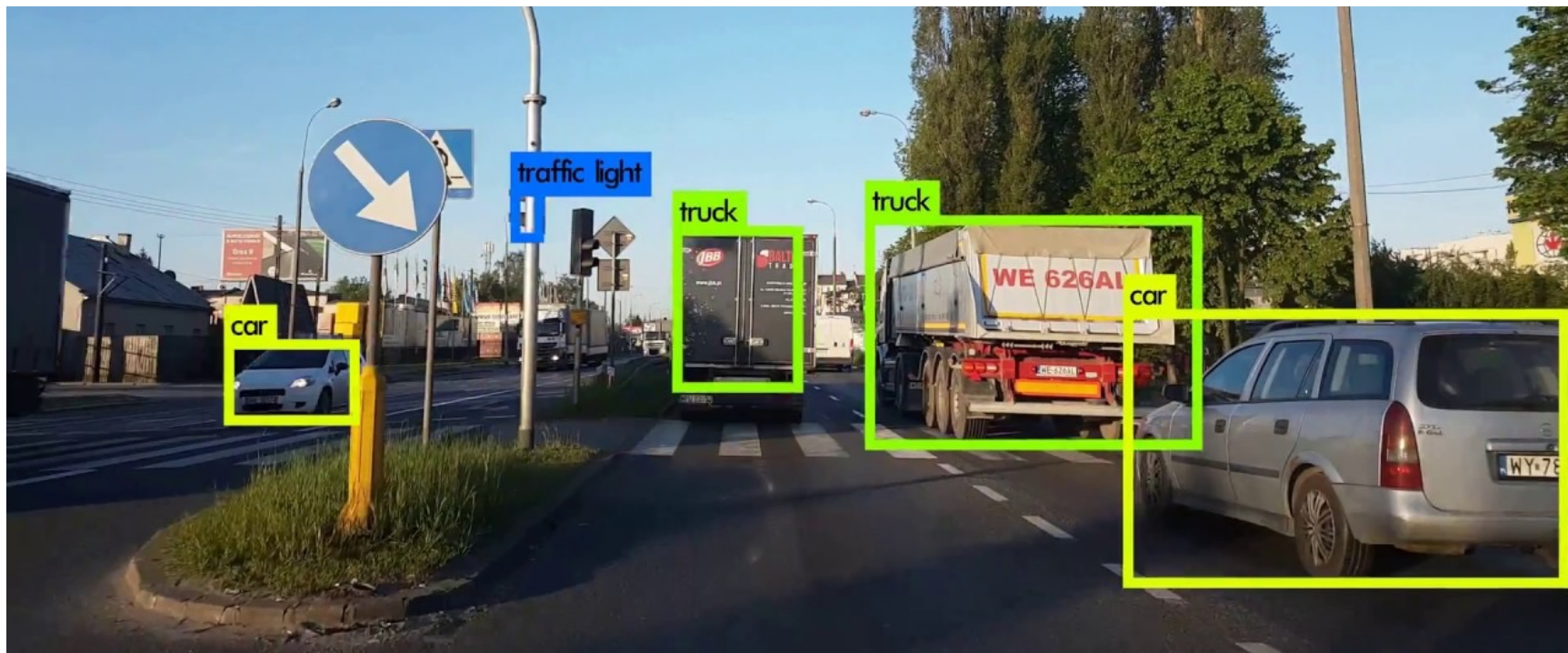
# Object Detection

The task of assigning a label and a bounding box to all predefined objects in the image.

- Classification
  - Box-level classification



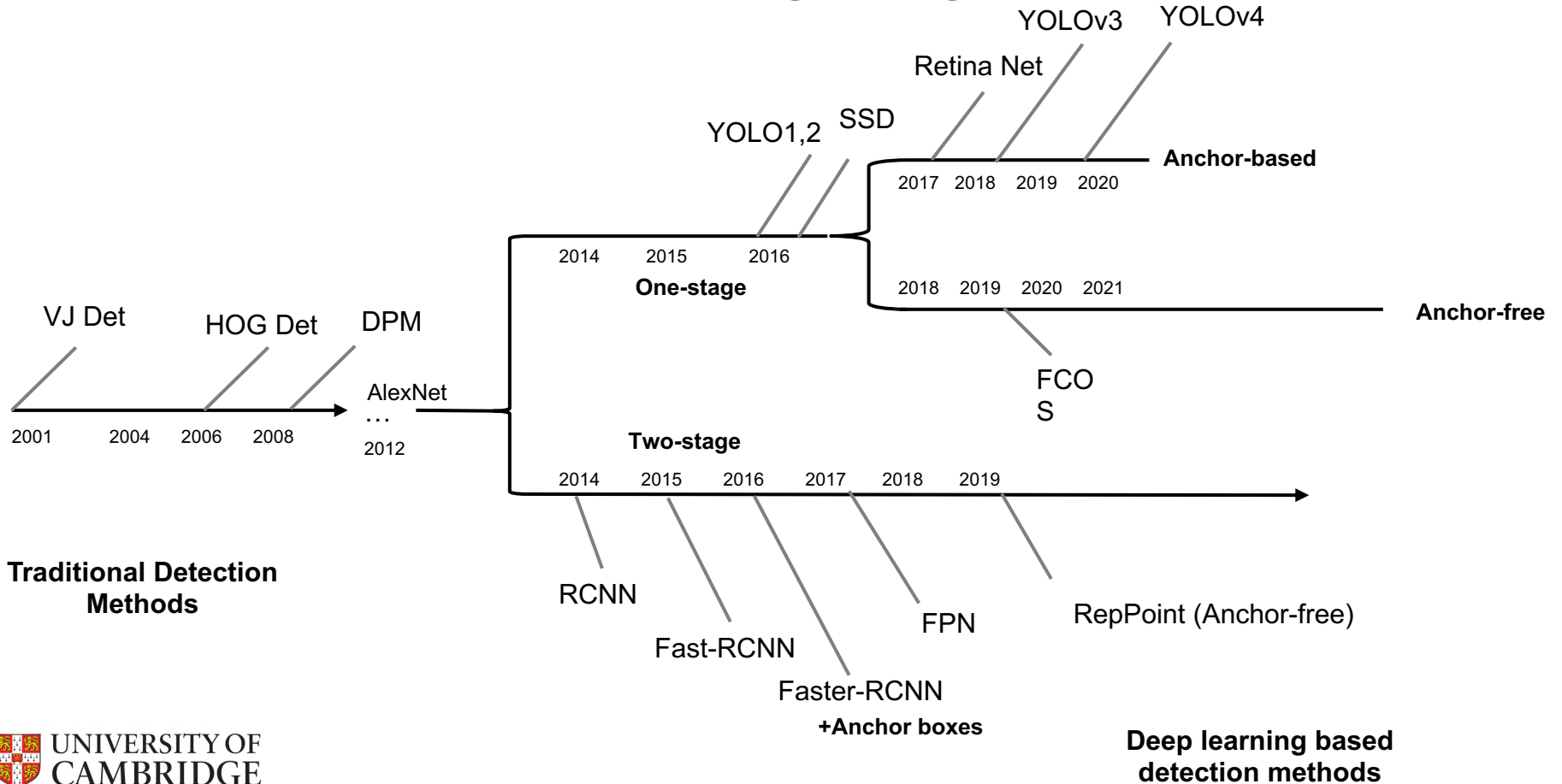
# Applications



# Applications



# Timeline



# Traditional Methods

## Sliding windows: Classification

Class: Dog, Cat, Human, Background



Dog: Yes

Cat : No

Human: No

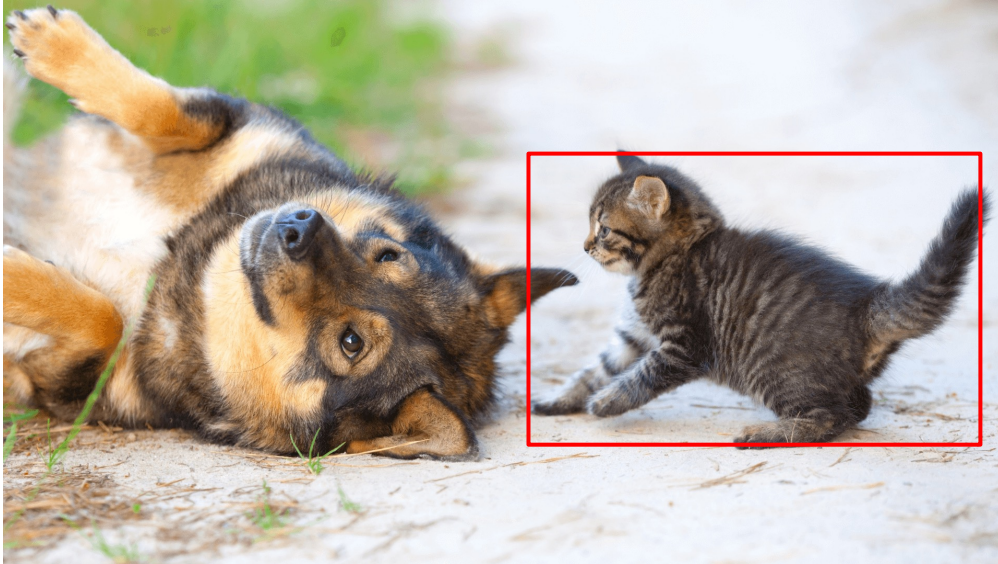
Background: No



# Traditional Methods

## Sliding windows: Classification

Class: Dog, Cat, Human, Background



Dog: No

Cat : Yes

Human: No

Background: No



# Traditional Methods

## Sliding windows: Classification

Class: Dog, Cat, Human, Background



Dog: No

Cat : No

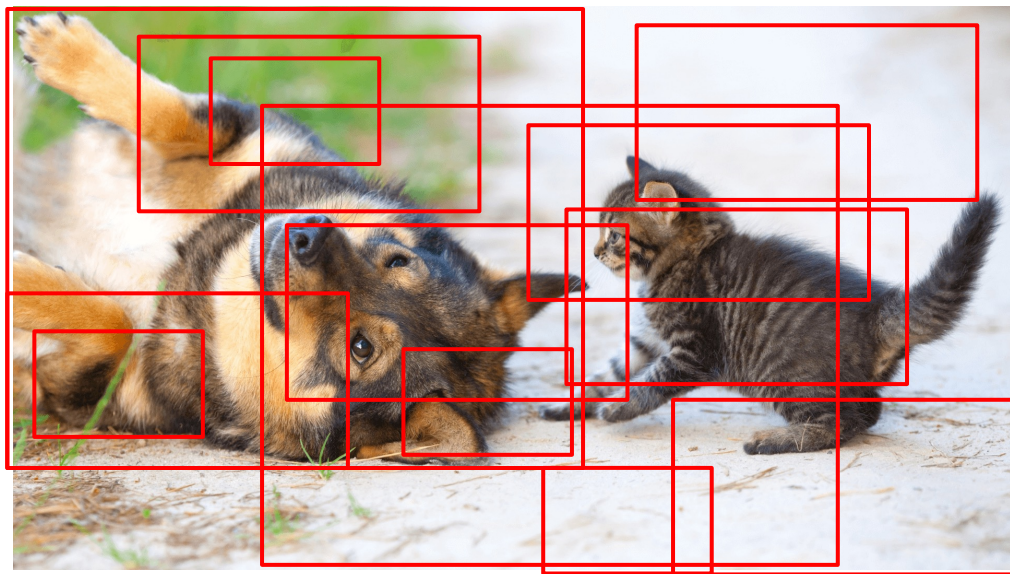
Human: No

Background: Yes

# Traditional Methods

## Sliding windows: Classification

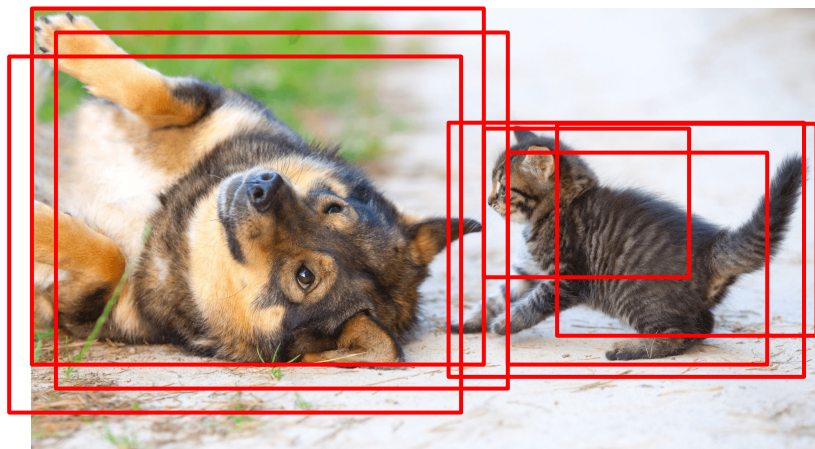
Class: Dog, Cat, Human, Background



1. Compute features on multiple resolutions
2. Scoring every sliding windows
3. Applying Non-maxima suppression

# Non-maxima Suppression (NMS)

- **Input:** A list of proposal boxes  $B$ , corresponding confidence scores  $S$  and overlap threshold  $N$ .
- **Output:** A list of filtered proposals  $D$ .



Before NMS

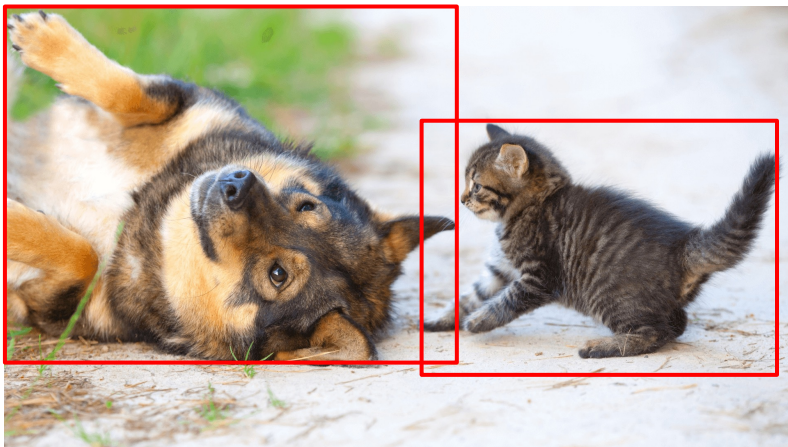
```
D = []
while B is not empty:
    i = Argmax(S)
    D.append[Bi]
    B.delete[Bi]

    for Br in B:
        if iou[Br,Bi] > th:
            B.delete[Br]

Return D
```

# Non-maxima Suppression (NMS)

- **Input:** A list of proposal boxes  $B$ , corresponding confidence scores  $S$  and overlap threshold  $N$ .
- **Output:** A list of filtered proposals  $D$ .



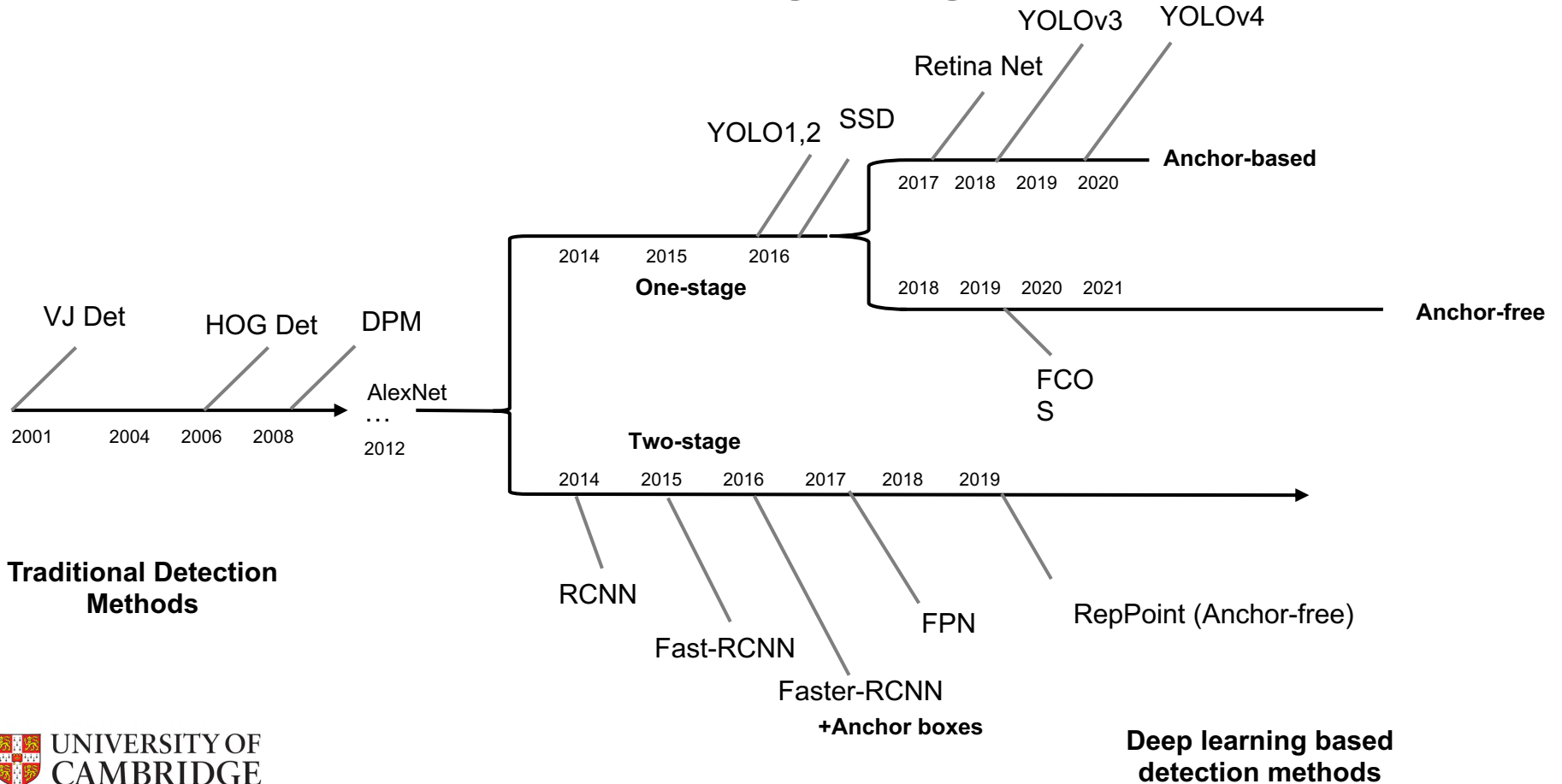
After NMS

```
D = []
while B is not empty:
    i = Argmax(S)
    D.append[Bi]
    B.delete[Bi]

    for Br in B:
        if iou[Br,Bi] > th:
            B.delete[Br]

Return D
```

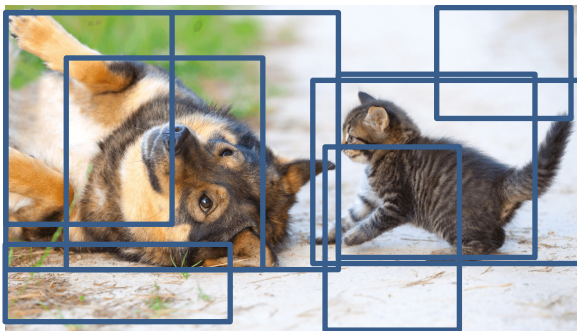
# Timeline



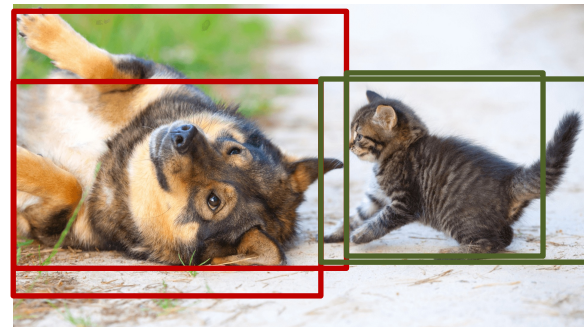
# Two stage vs. One stage



Input image



Proposals



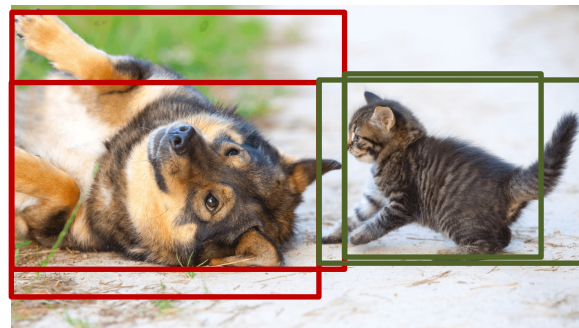
Classify the proposals

Two stage methods

# Two stage vs. One stage



Input image



Predict classified boxes

One stage methods

# RCNN

Step 1: Selective Search  $\rightarrow$   $\sim$ 2k proposals

- find image regions that likely contain objects

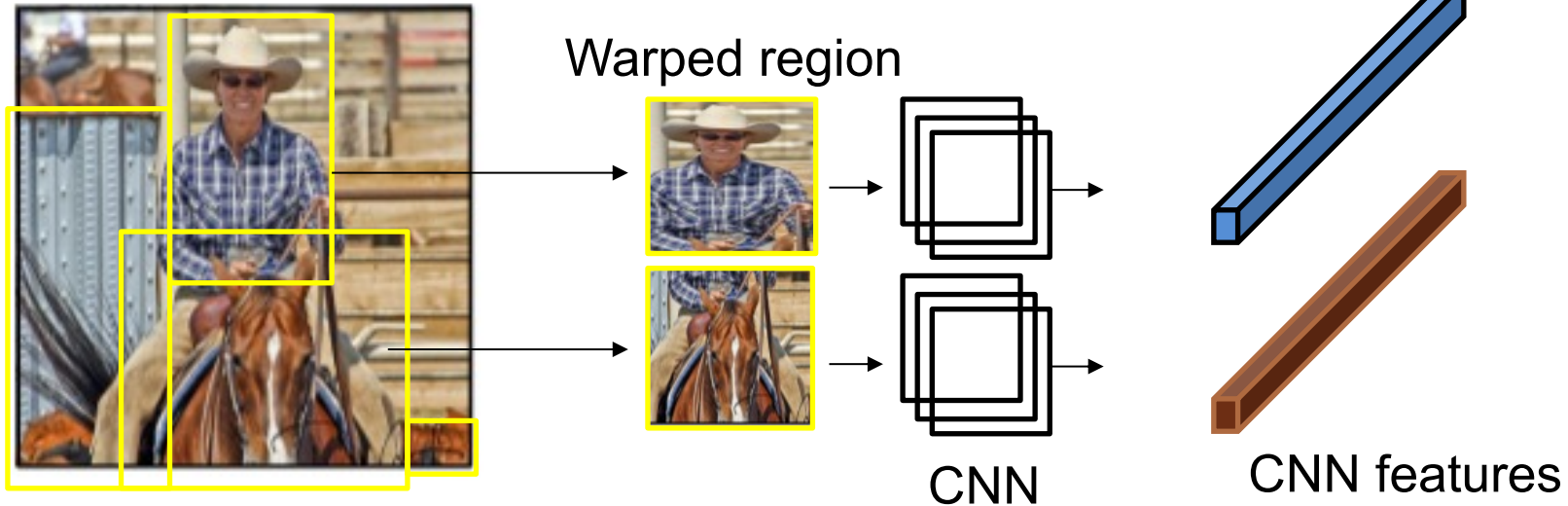




# RCNN

## Step 2: CNN extract features

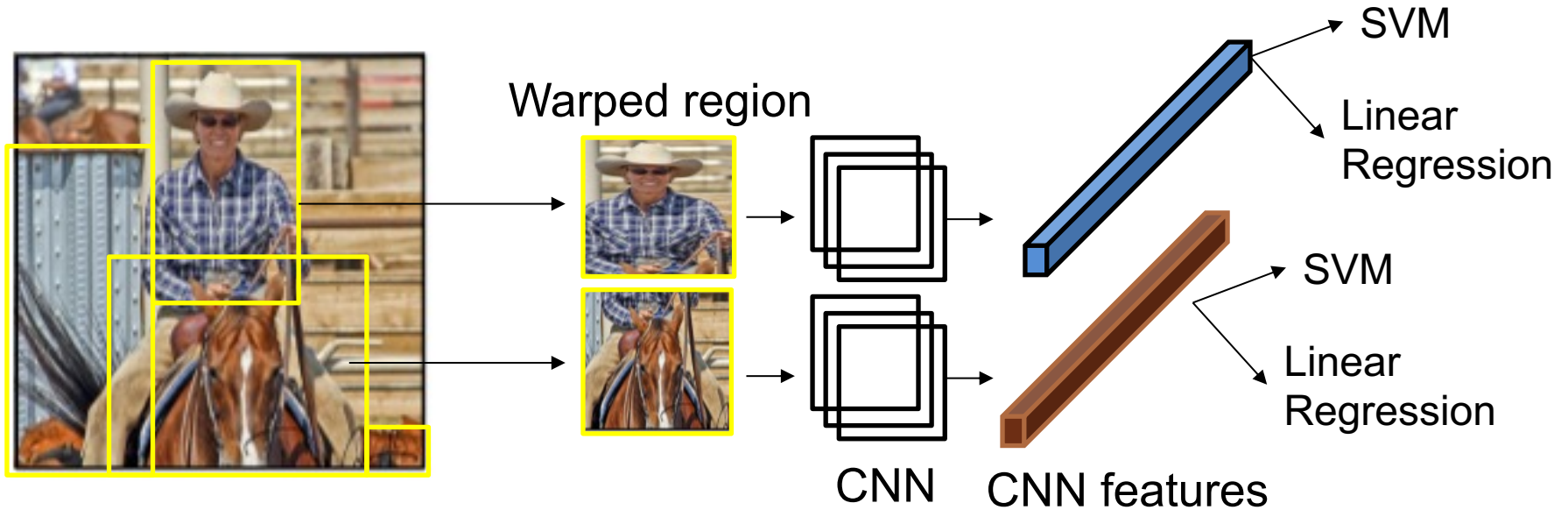
- Affine image warping: Get a fixed input size



# RCNN

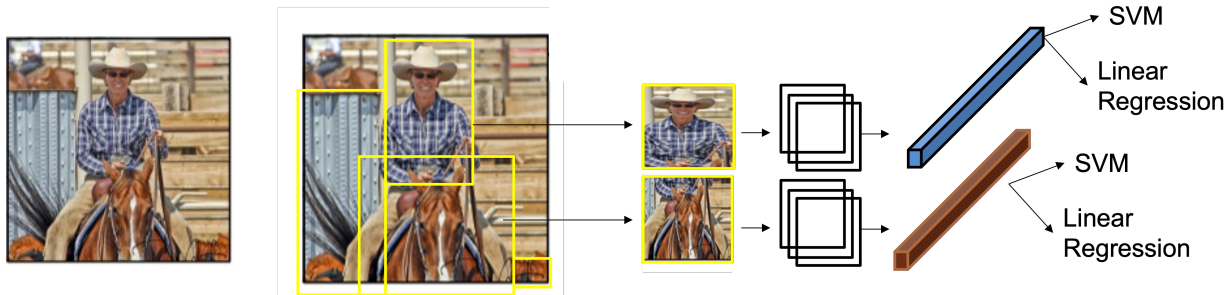
## Step 3: Classification and regression

- Classify with a linear SVM, linear regression for the bounding box offset



# Fast-RCNN

RCNN



Fast-RCNN

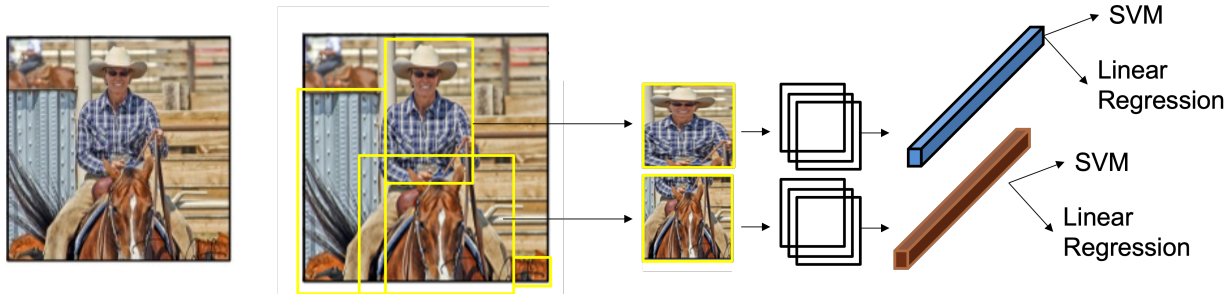


1. Selective search

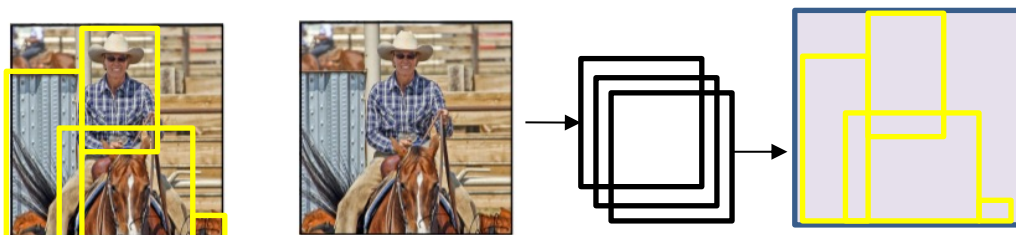
2. Input image to CNN

# Fast-RCNN

RCNN



Fast-RCNN

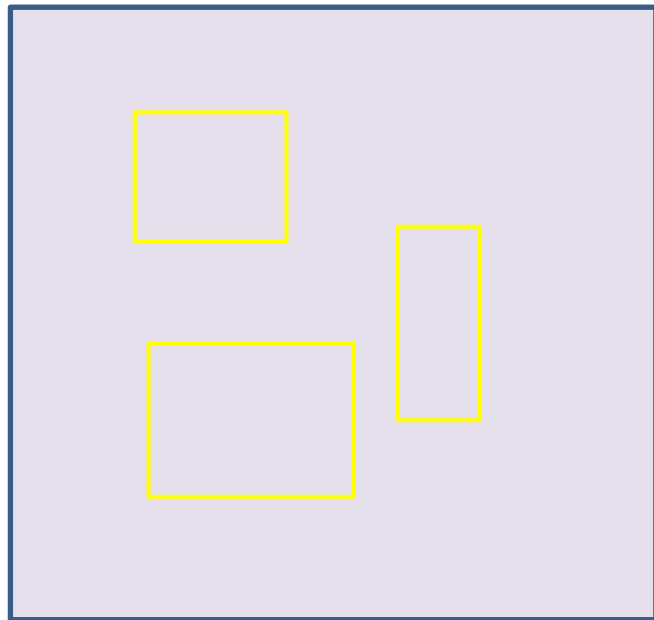


1. Selective search

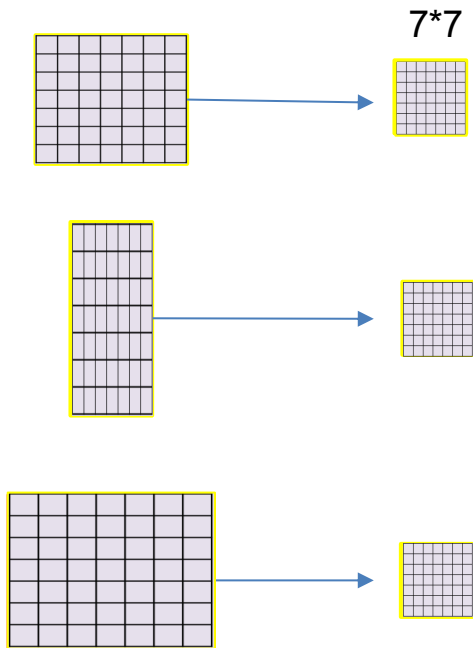
2. Input image to CNN

3. ROI pooling

# ROI Pooling



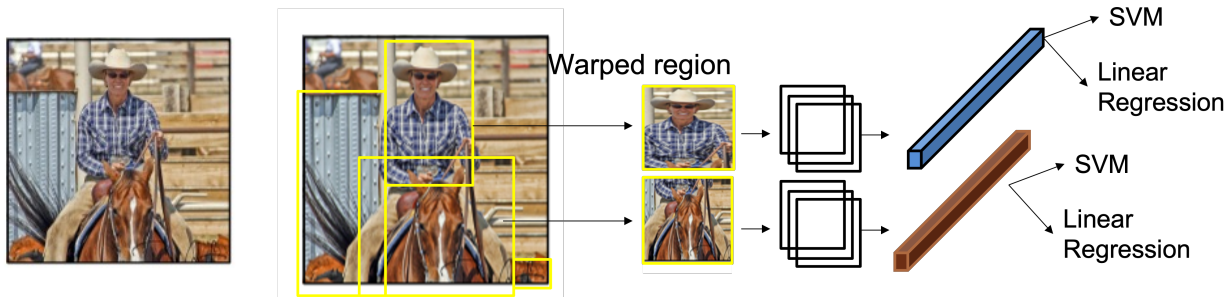
Feature map



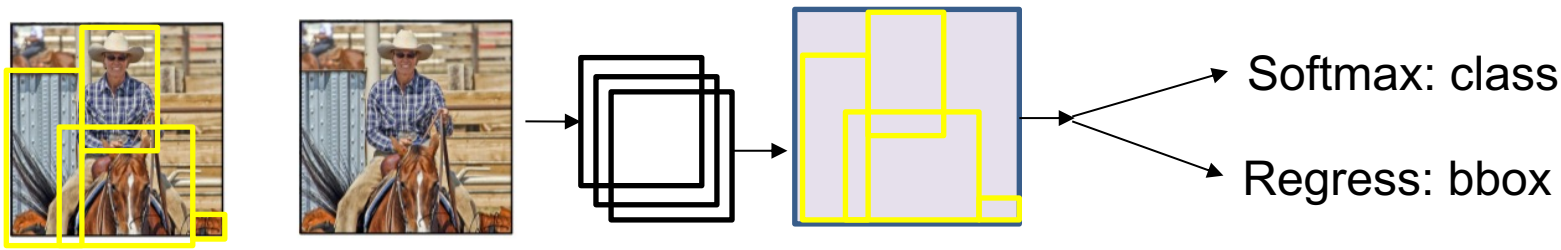
ROI Pooling: max pooling for each grid

# Fast-RCNN

RCNN



Fast-RCNN



1. Selective search

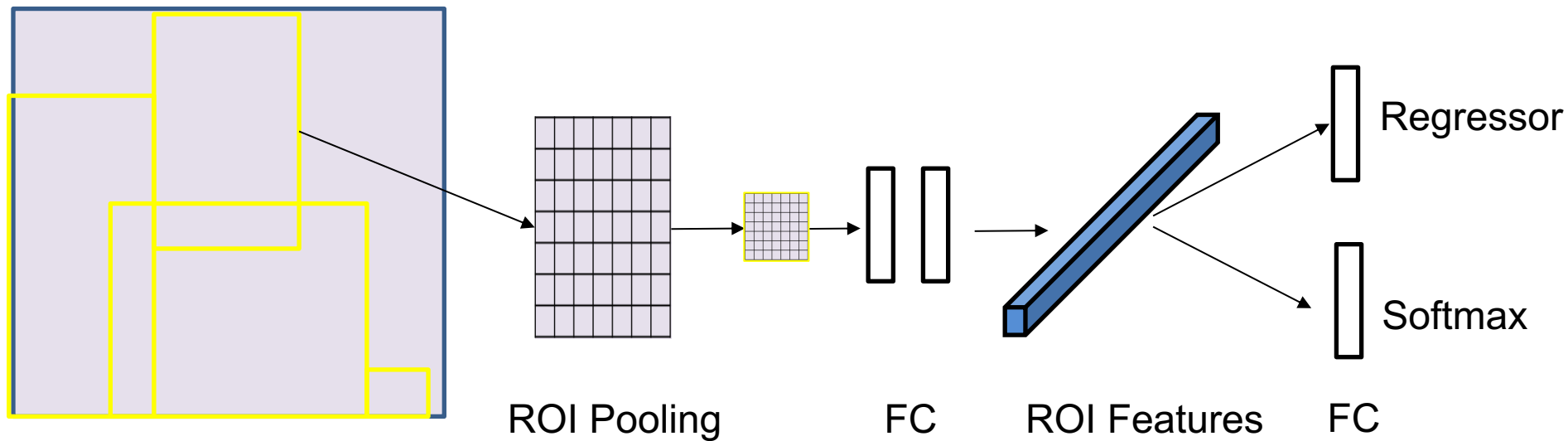
2. Input image to CNN

3. ROI pooling

4. CNN heads

# Fast-RCNN

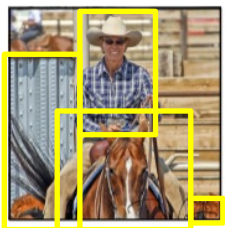
## Details of the CNN head



# Faster-RCNN

Main difference: Use a CNN to generate region proposals instead of selective search

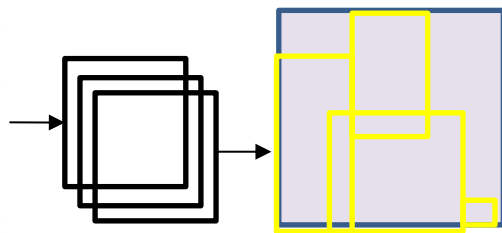
Fast-RCNN



1. Selective search



2. Input image into CNN



3. ROI pooling

Softmax: class

Regress: bbox

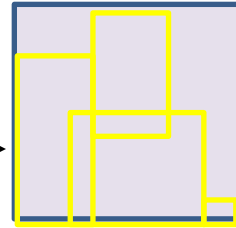
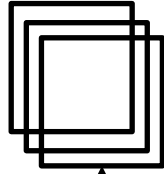
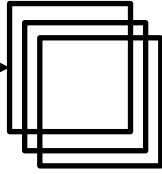
4. CNN heads



# Faster-RCNN

Main difference: Use a CNN to generate region proposals instead of selective search

Fast-RCNN



Softmax: class

Regress: bbox

1. Input image into CNN

2. Region Proposal Network

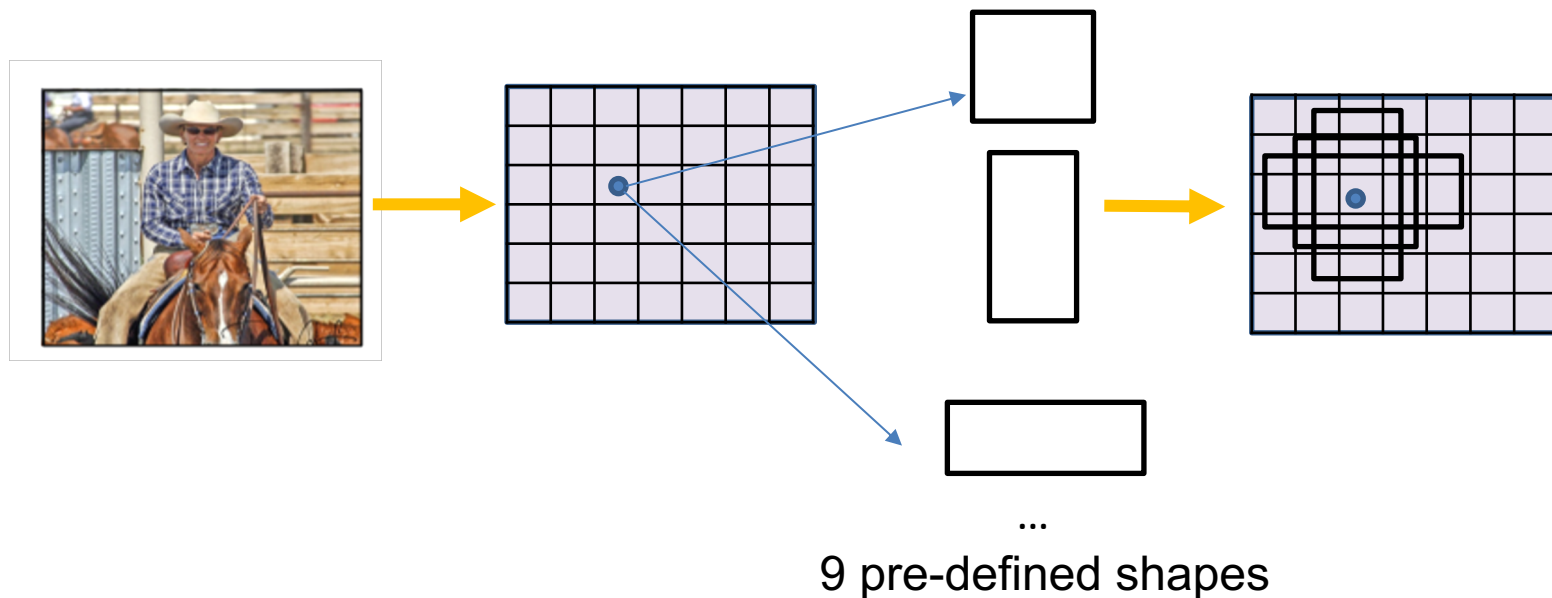
RPN

3. ROI pooling

4. CNN heads

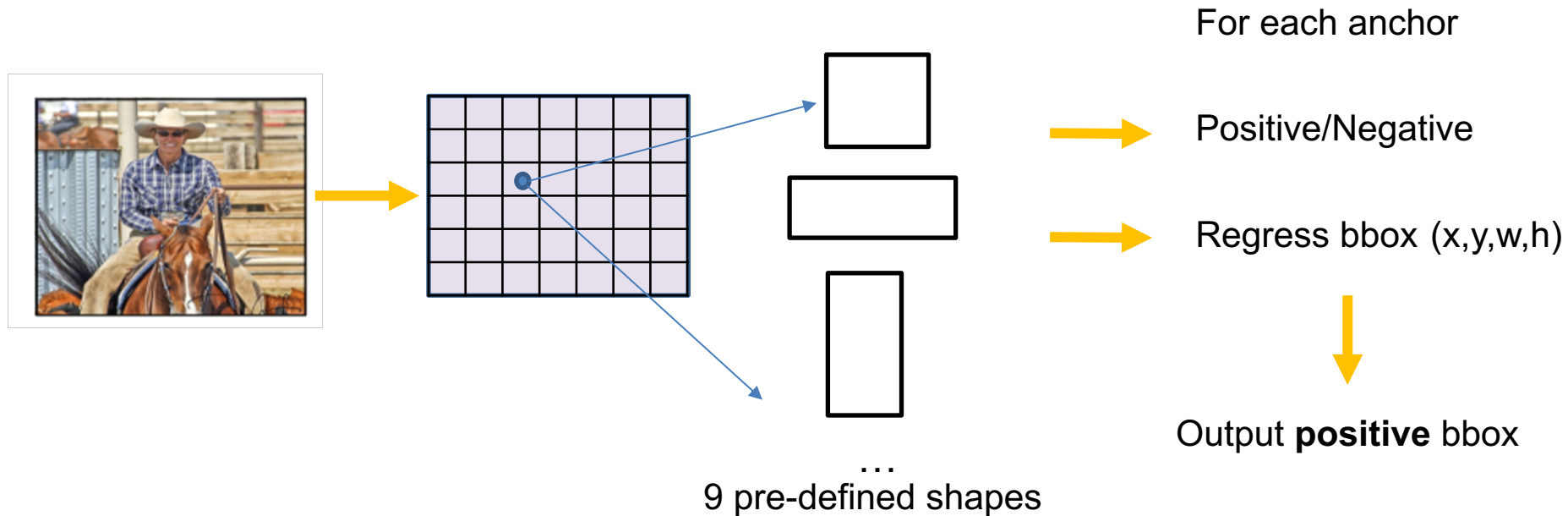
# Faster-RCNN: RPN

Anchors: pre-defined boxes



# Faster-RCNN: RPN

Anchors: pre-defined boxes



# Test Performance

	<b>RCNN</b>	<b>Fast-RCNN</b>	<b>Faster RCNN</b>
Time	50 s	2 s	0.2 s
mAP on Pascal VOC	66.0	66.9	66.9

# Mean Average Precision

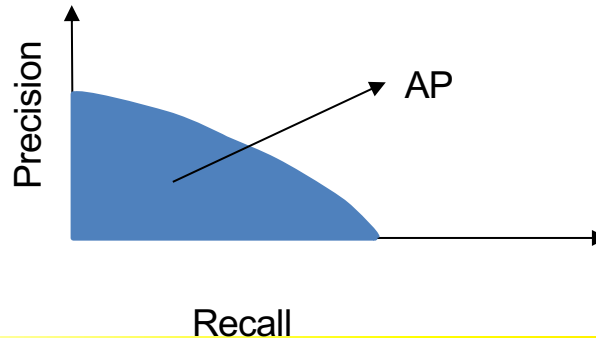
$$\textit{Precision} = \frac{TP}{TP + FP}$$

$$\textit{Recall} = \frac{TP}{TP + FN}$$

*TP = True Positives (Predicted as positive as was correct)*

*FN = False Negatives (Failed to predict an object)*

*FP = False Positives (Predicted as positive but was incorrect)*



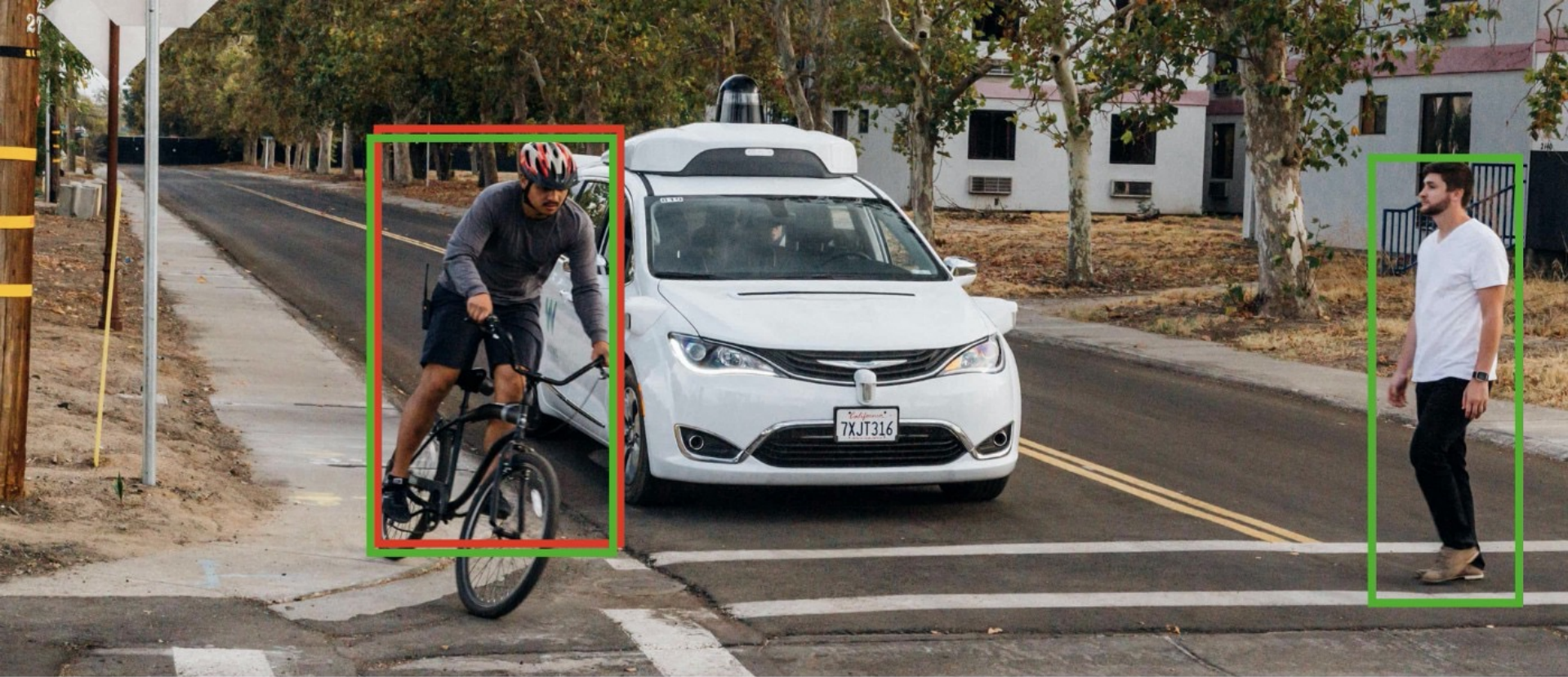


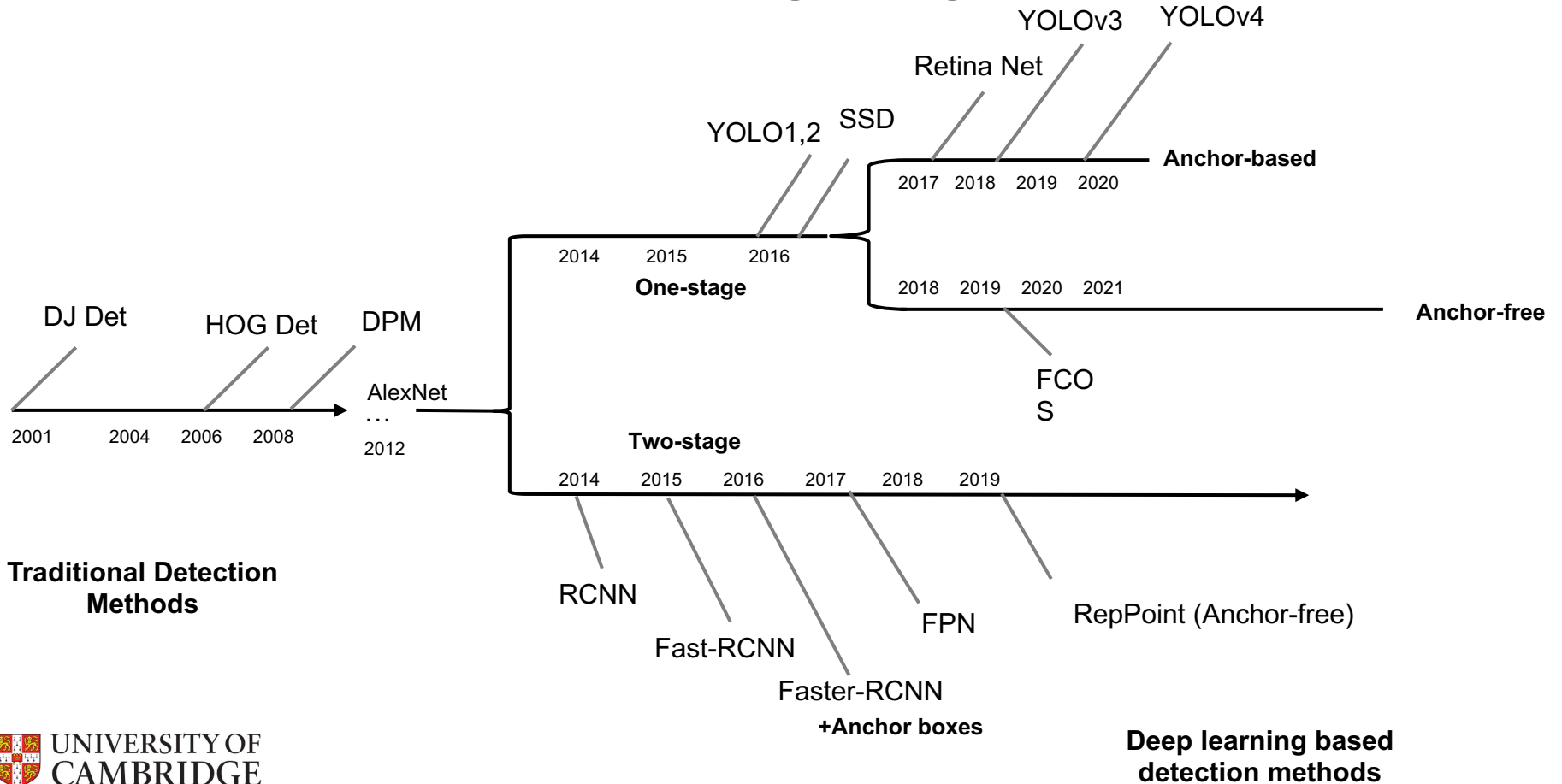
image: Waymo

 = Predicted Bounding Box

 = Ground Truth Bounding Box

TP=1, FP=0, FN=1

# Timeline



# Two-stage vs One-stage

---

One stage: No region proposal networks

YOLO, SSD

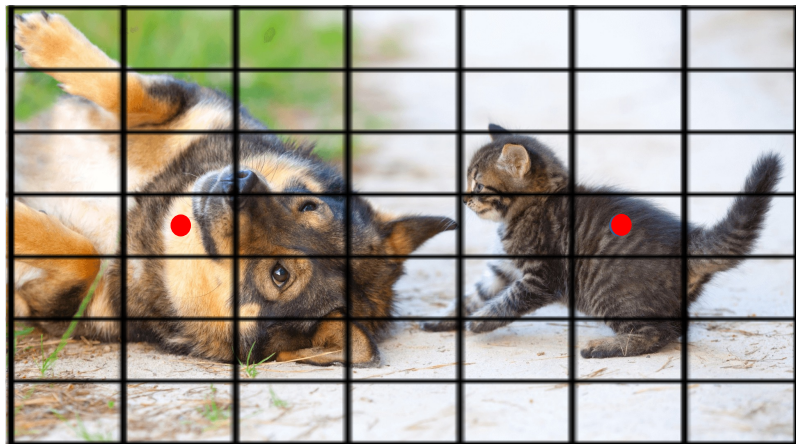
Directly regress the bbox (dx, dy, dh, dw, confidence)



# Two-stage vs One-stage

One stage: No region proposal networks

YOLO, SSD



7\*7 grids

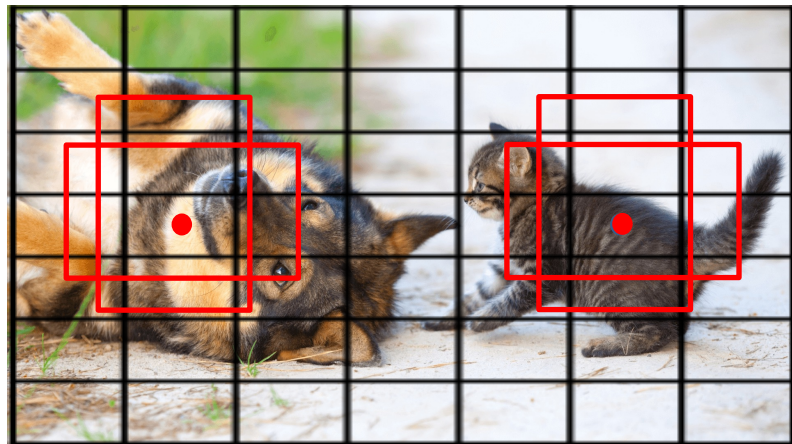
Directly regress the bbox (dx, dy, dh, dw, confidence)

YOLO V1: only one estimation per pre-defined region → low recall

# Two-stage vs One-stage

One stage: No region proposal networks

YOLO, SSD



7\*7 grids

Directly regress the bbox (dx, dy, dh, dw, confidence)

YOLO V1: only one estimation per pre-defined region → low recall

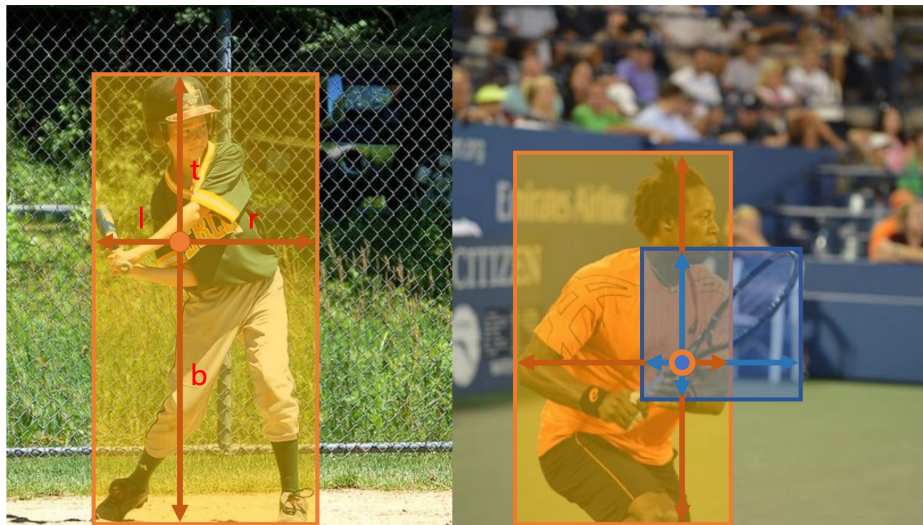
YOLOV2, 3: Use anchors → improve recall

# Anchor-free vs Anchor-based

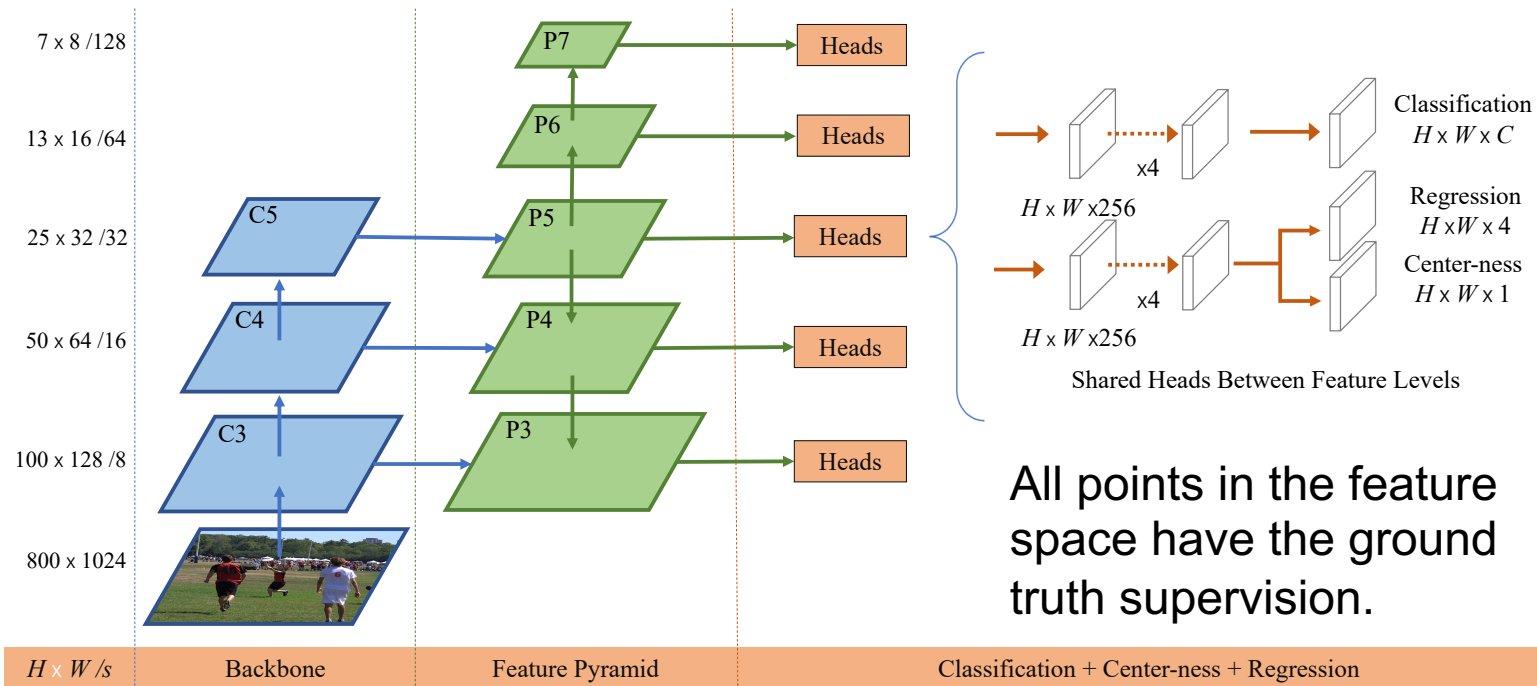
Anchor-free: do not rely on anchors

YOLOv1: only one estimation per pre-defined region

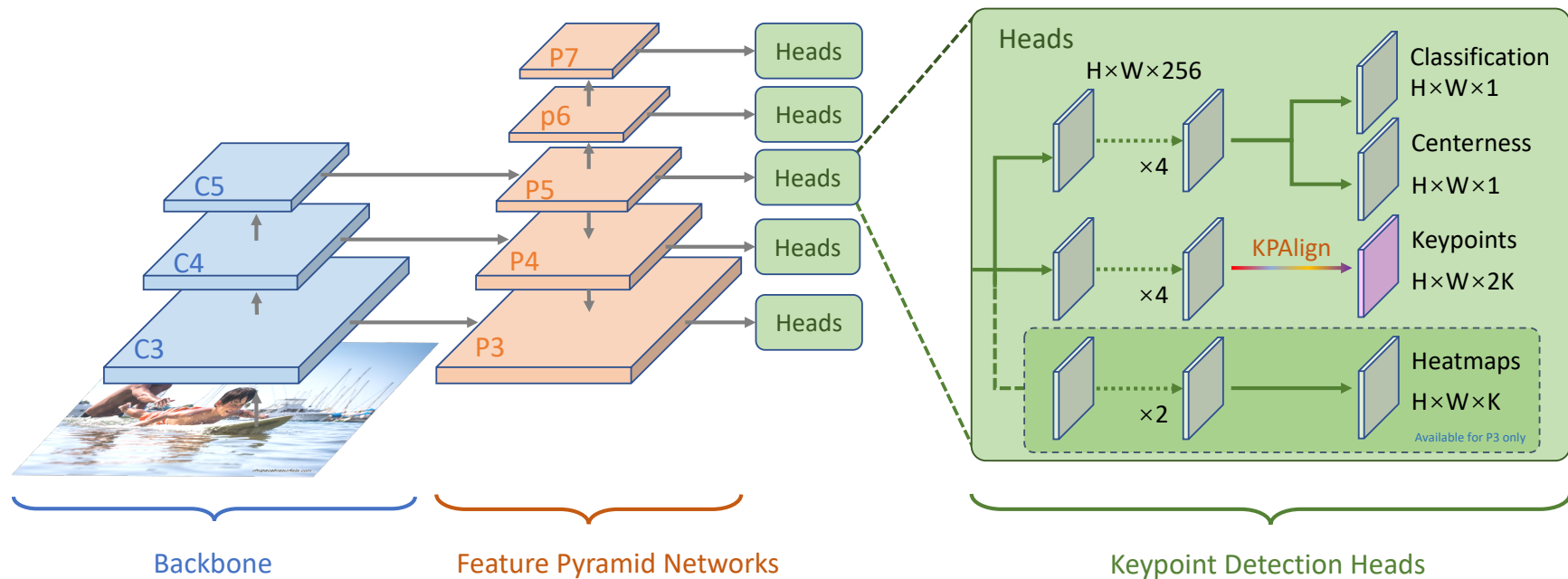
FCOS: for each point on the image plane, regress to the distances to the bounding box edges → high recall



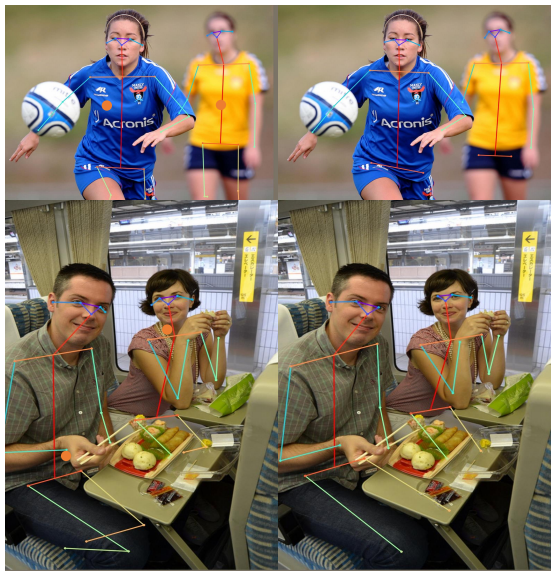
# FCOS



# Aside: Key-point Detection



# Aside: Keypoint Detection



Estimated    Ground Truth



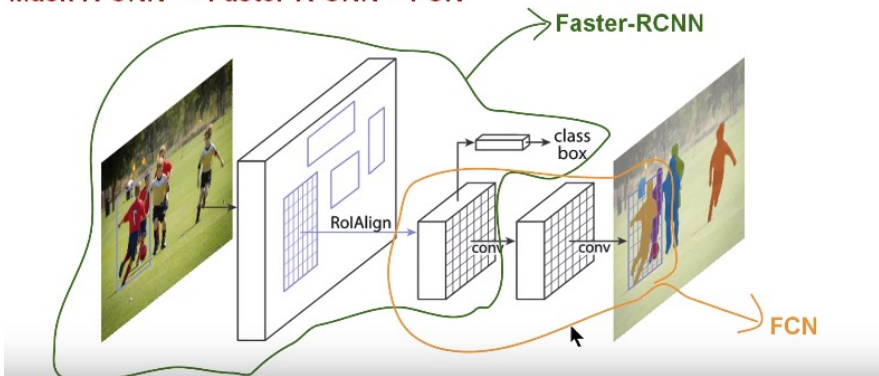
Estimated    Ground Truth

# Aside: Instance Segmentation

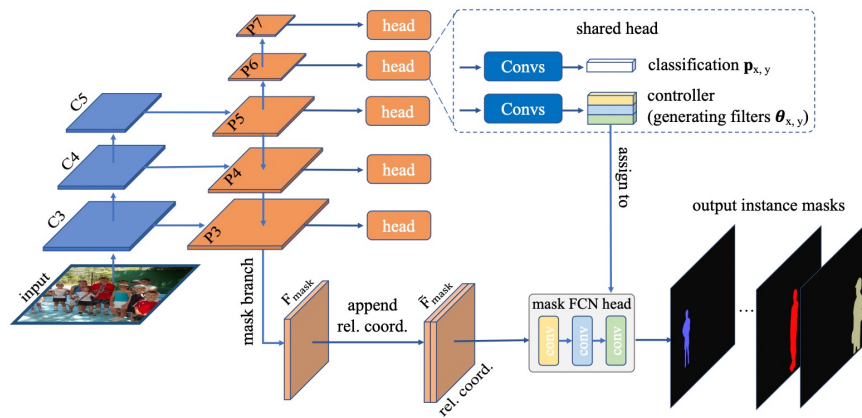


# Aside: Instance Segmentation

Mask R-CNN → Faster R-CNN + FCN

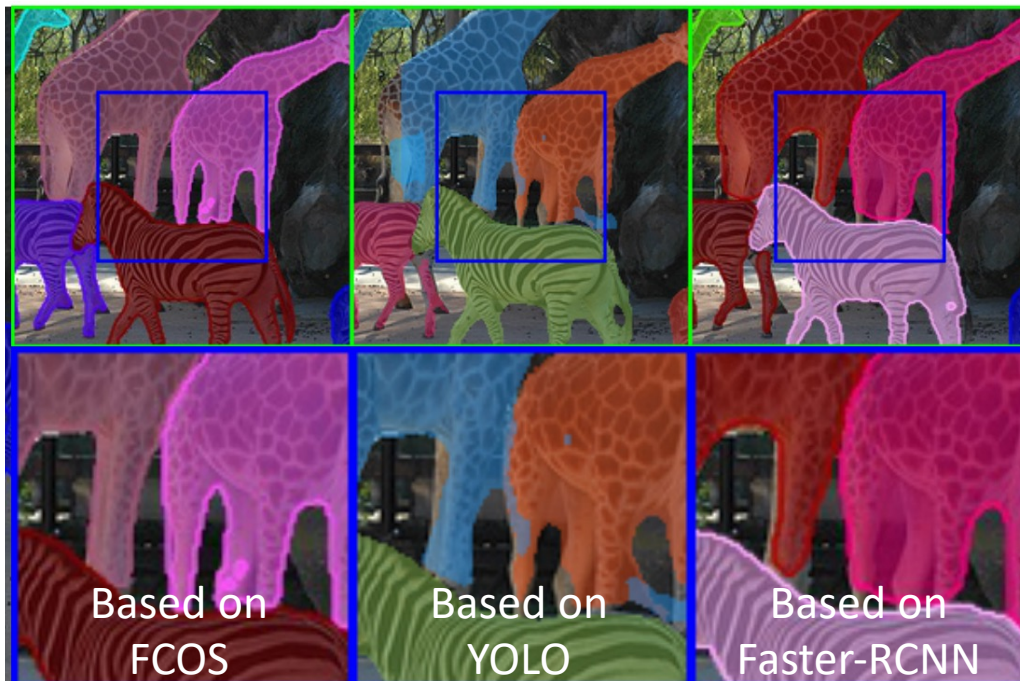


Based on Faster-RCNN



Based on FCOS





# Multiple Variations

- **Input:** Image, Patch, Image Pyramid
- **Backbones:**
  - VGG16, ResNet-50, SpineNet, EfficientNet-B0/B7, CSPResNeXt50, CSPDarknet53
- **Neck:**
  - **Additional block:** SPP, ASPP, RFB, SAM
  - **Feature Fusion:** FPN, PAN, NAS-FPN, Fully-connected FPN, BiFPN, ASFF, SFAM
- **Head:**
  - **One-stage:**
    - RPN, SSD, YOLO, RetinaNet (anchor based)
    - CornerNet, CenterNet, MatrixNet, FCOS (anchor free)
  - **Two-stage:**
    - Faster R-CNN, R-FCN, Mask RCNN (anchor based)
    - RepPoints (anchor free)

# New Framework based on Transformers

- Examples:
  - DETR: End-to-End Object Detection with Transformers
  - Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

COCO dataset

	<b>FasterRCNN</b>	<b>FasterRCNN-fpn</b>	<b>FCOS</b>	<b>Retinanet</b>	<b>DETR</b>	<b>SWIN-L</b>
mAP	41.1	42.0	43.1	40.4	44.9	58

# Reference List

- 1) Viola, Paul, and Michael Jones. "Robust real-time object detection." *International journal of computer vision* 4, no. 34-47 (2001): 4.
- 2) Dalal, Navneet, and Bill Triggs. "Histograms of oriented gradients for human detection." In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, vol. 1, pp. 886-893. Ieee, 2005.
- 3) Felzenszwalb, Pedro F., Ross B. Girshick, David McAllester, and Deva Ramanan. "Object detection with discriminatively trained part-based models." *IEEE transactions on pattern analysis and machine intelligence* 32, no. 9 (2009): 1627-1645.
- 4) He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Spatial pyramid pooling in deep convolutional networks for visual recognition." *IEEE transactions on pattern analysis and machine intelligence* 37, no. 9 (2015): 1904-1916.
- 5) Girshick R, Donahue J, Darrell T, et al. Rich feature hierarchies for accurate object detection and semantic segmentation[C]//*Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014: 580-587.
- 6) Girshick, Ross. "Fast r-cnn." In *Proceedings of the IEEE international conference on computer vision*, pp. 1440-1448. 2015.
- 7) Ren, Shaoqing, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." *Advances in neural information processing systems* 28 (2015): 91-99.
- 8) Lin, Tsung-Yi, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. "Feature pyramid networks for object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2117-2125. 2017.
- 9) Lin, Tsung-Yi, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. "Focal loss for dense object detection." In *Proceedings of the IEEE international conference on computer vision*, pp. 2980-2988. 2017.
- 10) Yang, Ze, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. "Reppoints: Point set representation for object detection." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9657-9666. 2019.
- 11) Redmon, Joseph, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779-788. 2016.
- 12) Redmon, Joseph, and Ali Farhadi. "YOLO9000: better, faster, stronger." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7263-7271. 2017.
- 13) Redmon, Joseph, and Ali Farhadi. "Yolov3: An incremental improvement." *arXiv preprint arXiv:1804.02767* (2018).
- 14) Bochkovskiy, Alexey, Chien-Yao Wang, and Hong-Yuan Mark Liao. "Yolov4: Optimal speed and accuracy of object detection." *arXiv preprint arXiv:2004.10934* (2020).