

上海交通大学

计算机视觉

教师: 赵旭

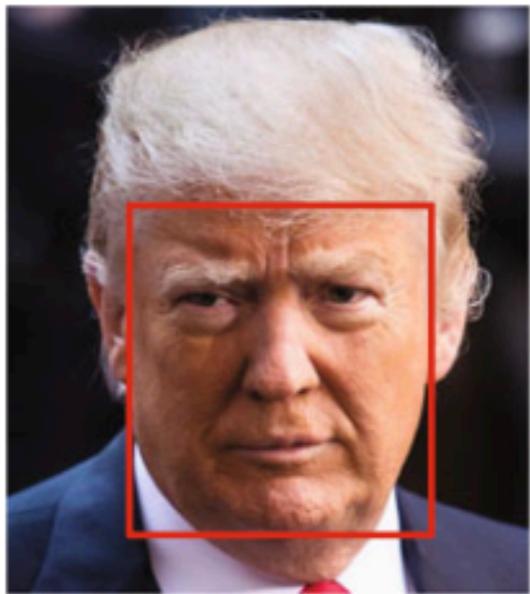
班级: AI4701

2024 春

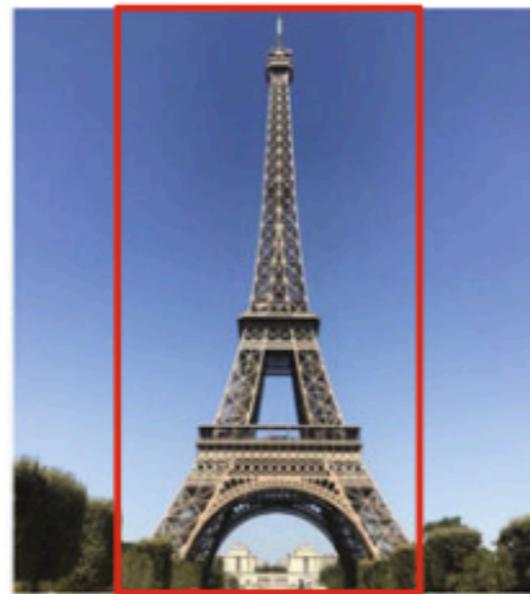
17. 物体检测

内容

- ❖ 传统方法
- ❖ 数据集与评估准则
- ❖ 基于深度学习的方法
- ❖ 关键问题



Donald Trump's face



Eiffel Tower



Mona Lisa by
Leonardo da Vinci

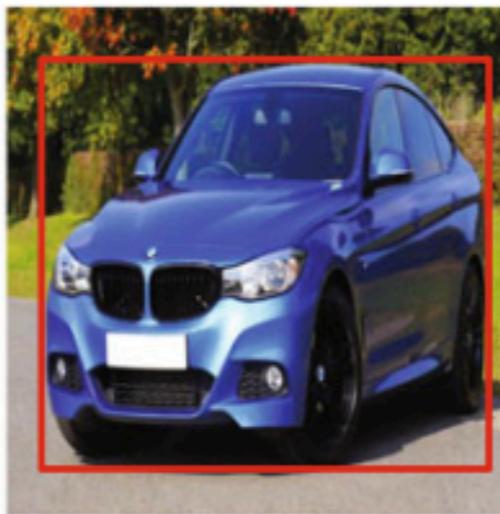


My neighbour's dog

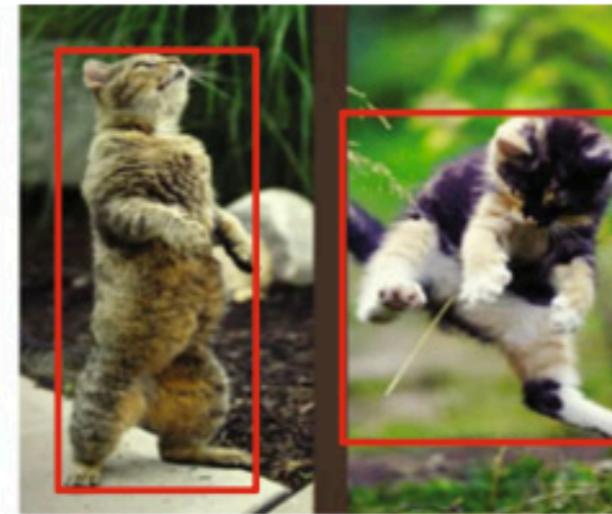
Specific Objects



Car



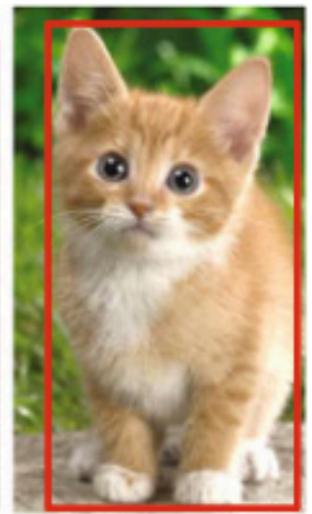
Car



Cat



Cat



Cat

Generic Object Categories

Technique Evolution

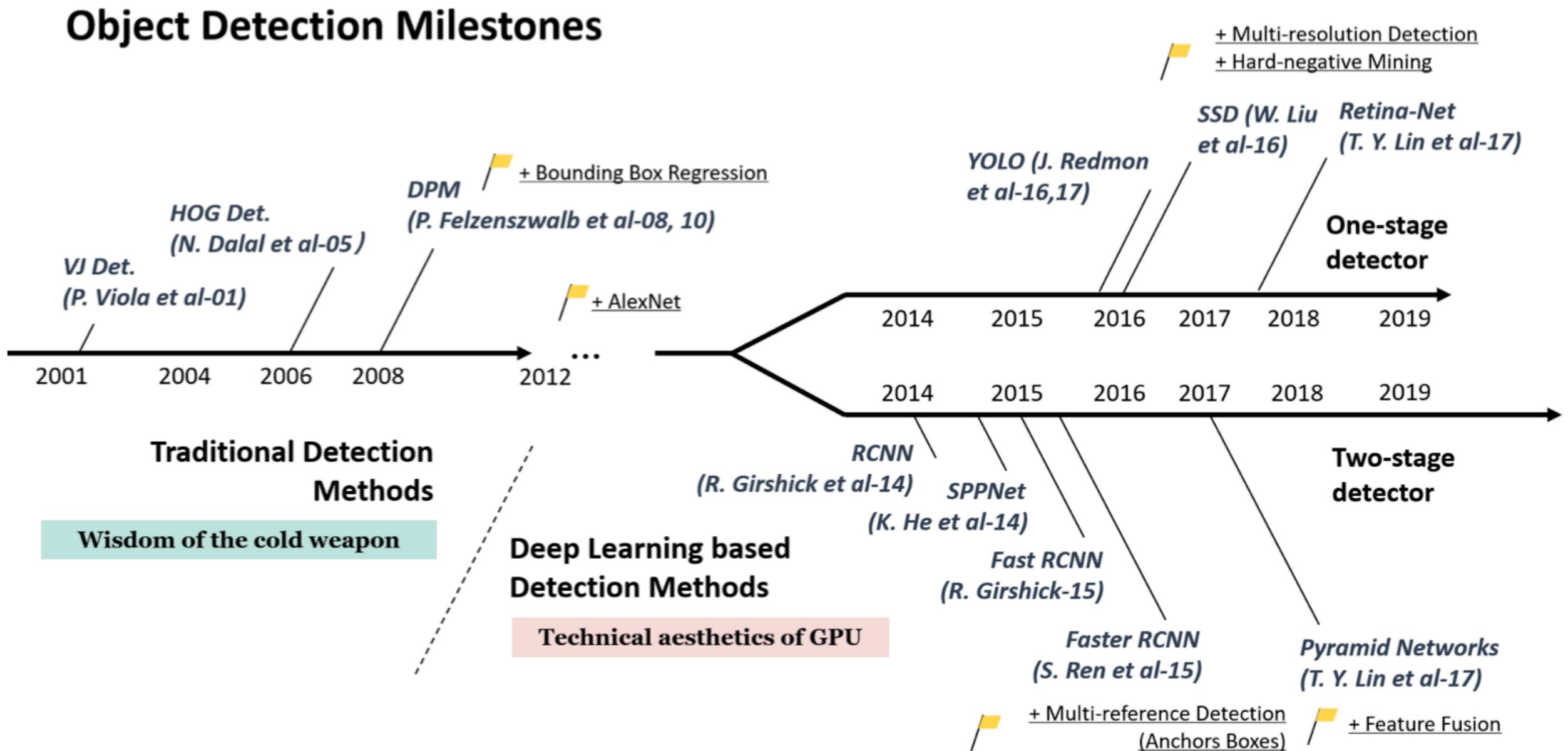
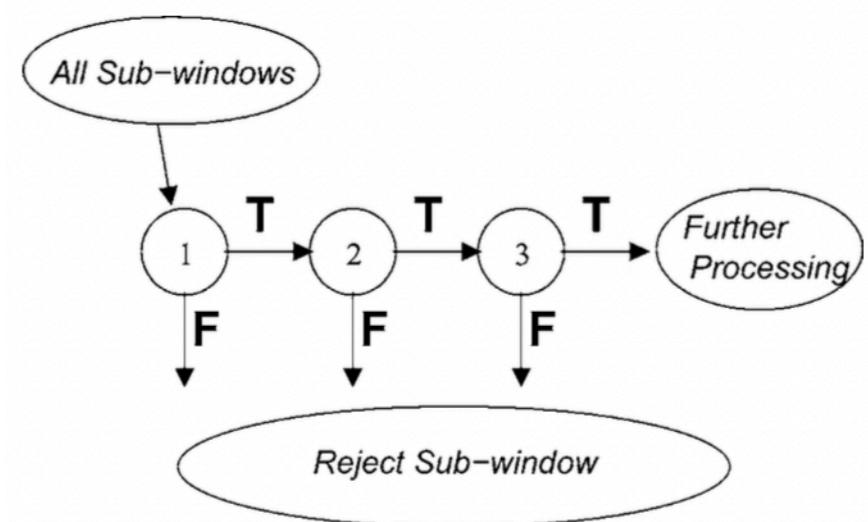
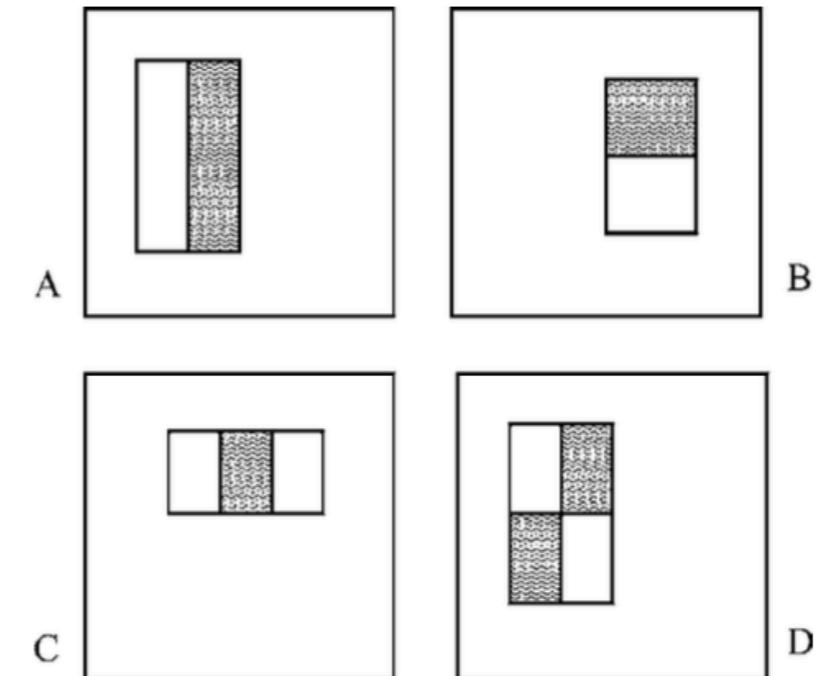


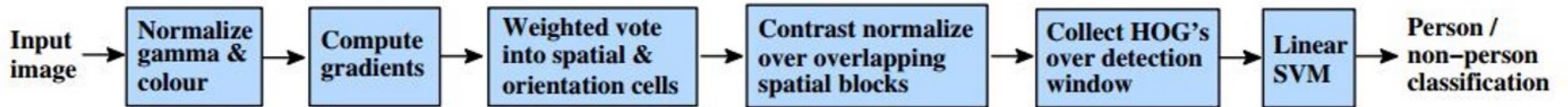
Fig. 2. A road map of object detection. Milestone detectors in this figure: VJ Det. [10, 11], HOG Det. [12], DPM [13–15], RCNN [16], SPPNet [17], Fast RCNN [18], Faster RCNN [19], YOLO [20], SSD [21], Pyramid Networks [22], Retina-Net [23].

Traditional detectors - Viola Jones Detectors

- ❖ Sliding windows: to go through all possible locations and scales in an image to see if any window contains a human face
- ❖ Three contributions
 - ❖ Integral image
 - ❖ Feature selection
 - ❖ Detection cascades

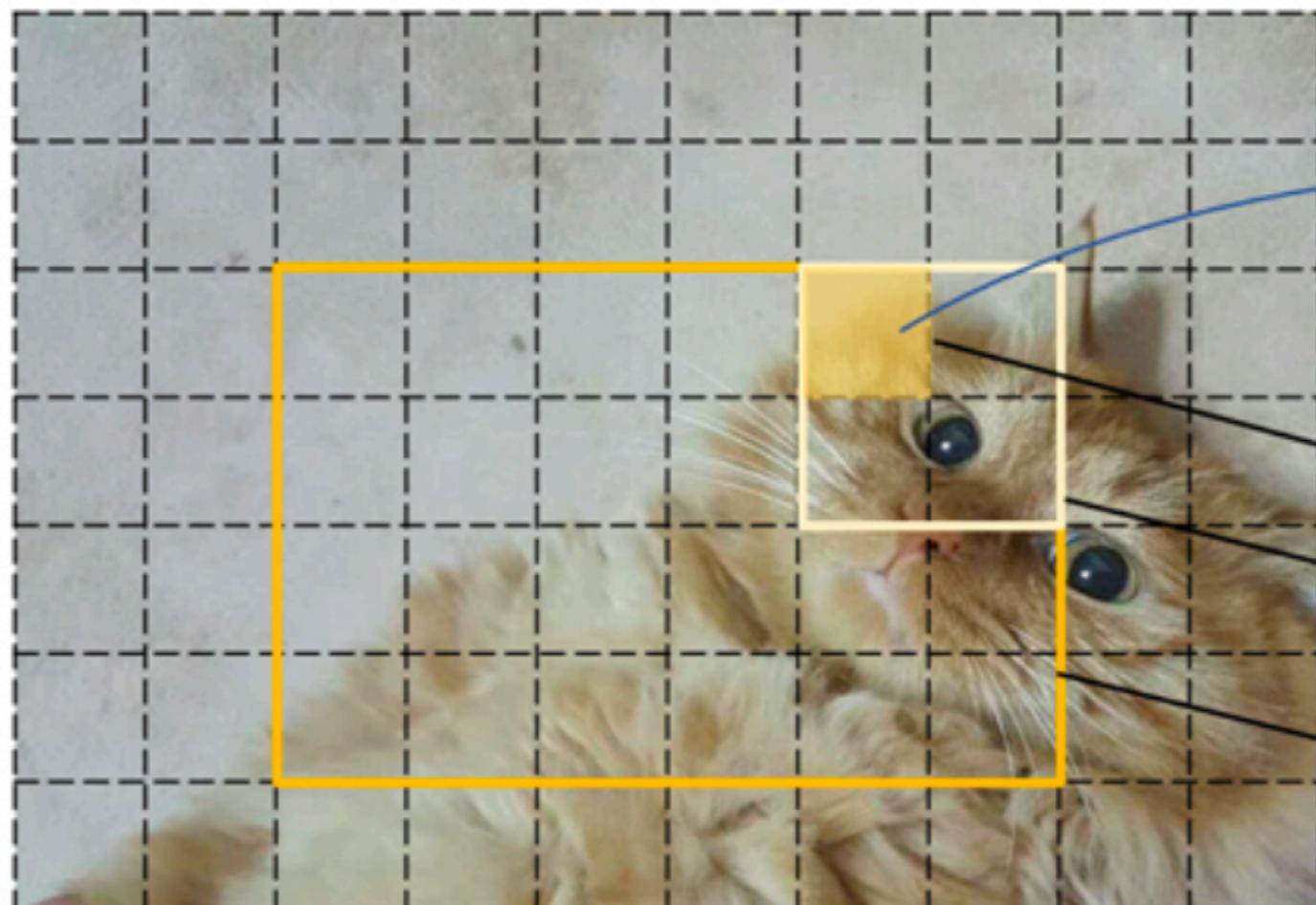


Traditional detectors - HOG

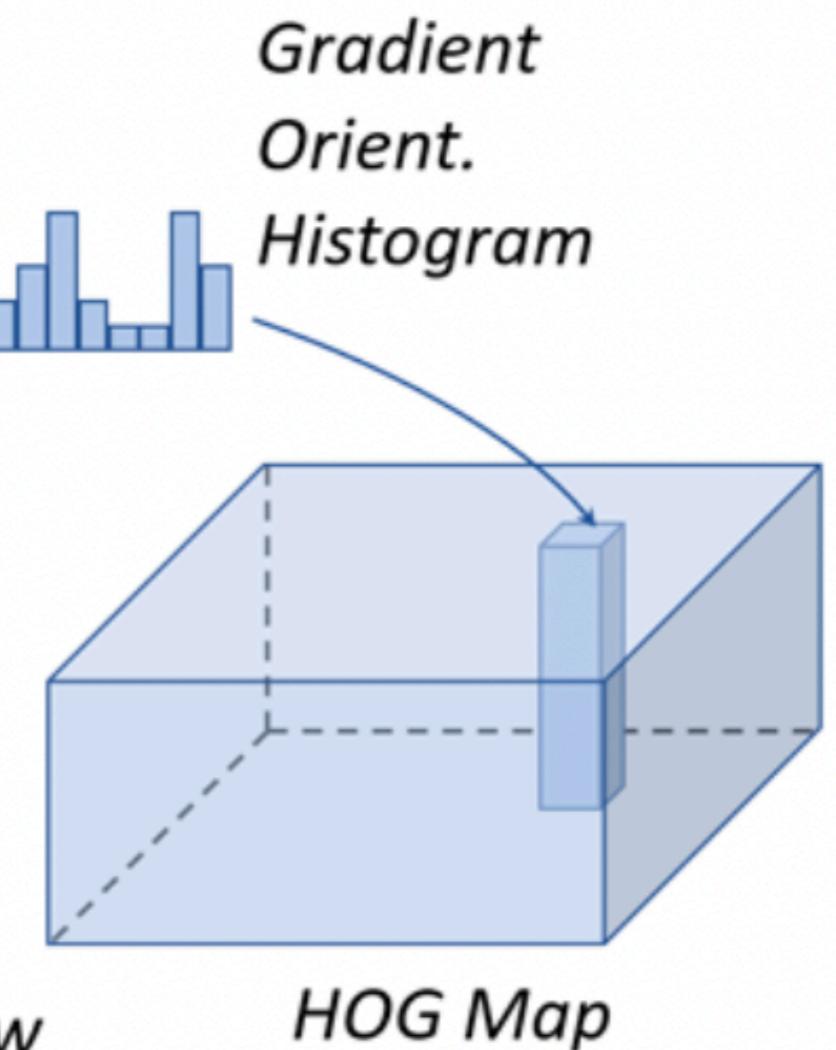


- ❖ Sliding Window (with HOG feature extraction) + Classifier (SVM)
- ❖ Designed to be computed on a dense grid of uniformly spaced cells and use overlapping local contrast normalization for improving accuracy.
- ❖ Rescale the input image for multiple times while keeping the size of a detection window unchanged.

From HOG to HOG Map

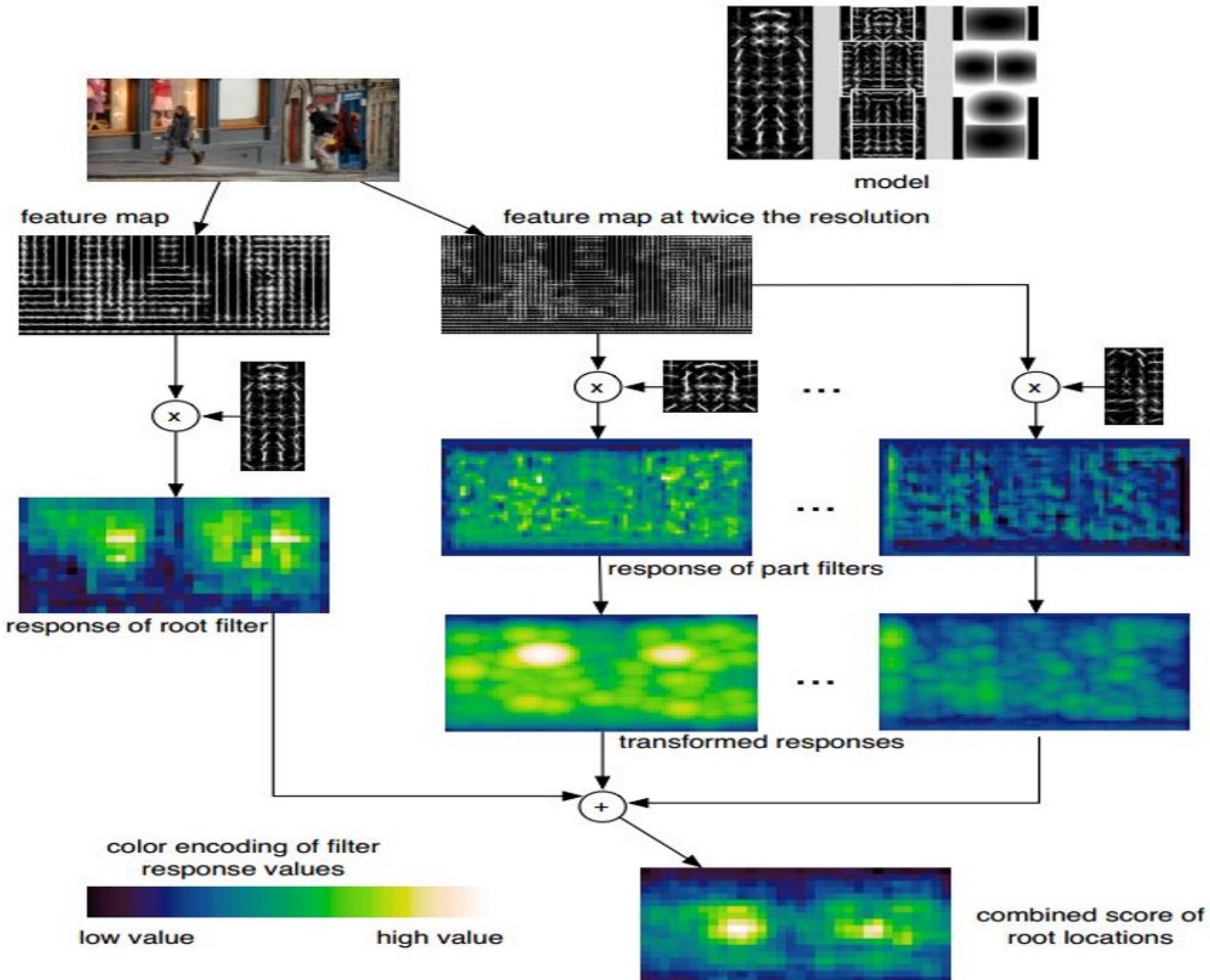


*Cell
Block
Sliding
Window*



Traditional detectors - DPM

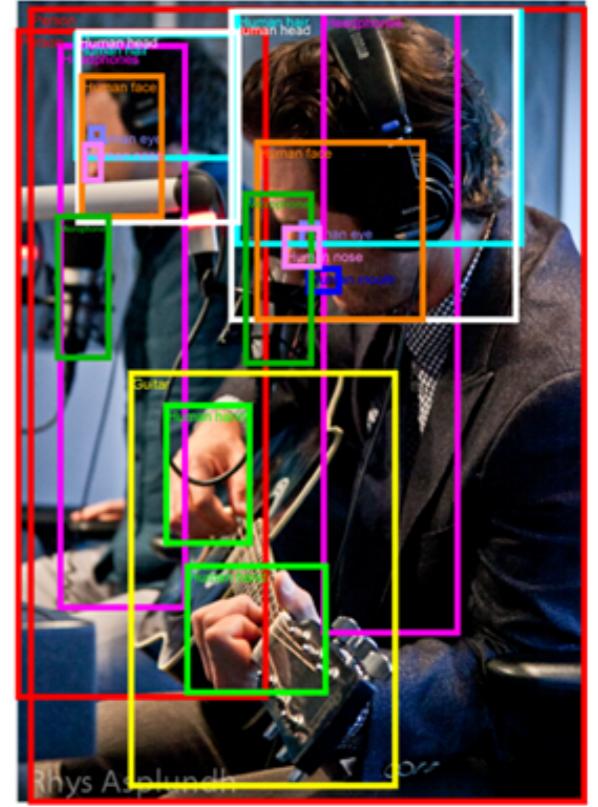
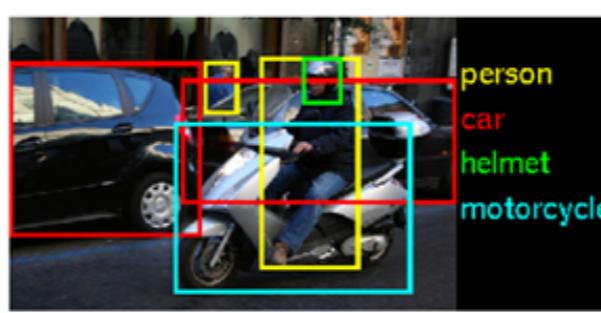
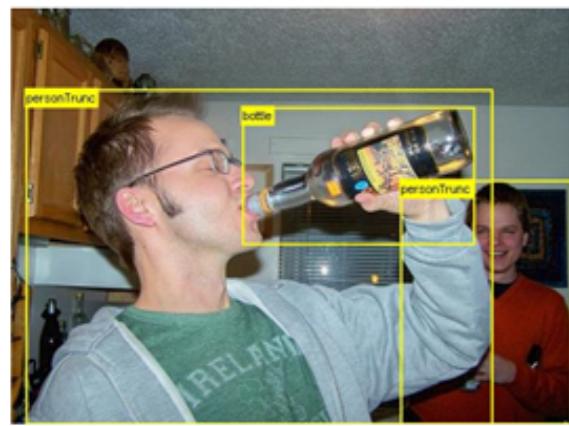
- ❖ DPM is the most successful traditional detector in recent years and is firstly proposed in 2008.
- ❖ It won the first place in Pascal VOC object detection challenge in 07,08,09.
- ❖ Learning of a proper way of decomposing an object, and the inference can be considered as an ensemble of detections on different object parts.
- ❖ A typical DPM detector consists of a root-filter and a number of part-filters.
- ❖ Propose some insightful techs: mixture models, hard negative mining, bounding box regression.



Datasets

Dataset	train		validation		trainval		test	
	images	objects	images	objects	images	objects	images	objects
VOC-2007	2,501	6,301	2,510	6,307	5,011	12,608	4,952	14,976
VOC-2012	5,717	13,609	5,823	13,841	11,540	27,450	10,991	-
ILSVRC-2014	456,567	478,807	20,121	55,502	476,688	534,309	40,152	-
ILSVRC-2017	456,567	478,807	20,121	55,502	476,688	534,309	65,500	-
MS-COCO-2015	82,783	604,907	40,504	291,875	123,287	896,782	81,434	-
MS-COCO-2018	118,287	860,001	5,000	36,781	123,287	896,782	40,670	-
OID-2018	1,743,042	14,610,229	41,620	204,621	1,784,662	14,814,850	125,436	625,282

TABLE 1
Some well-known object detection datasets and their statistics.



(a)

(b)

(c)

(d)

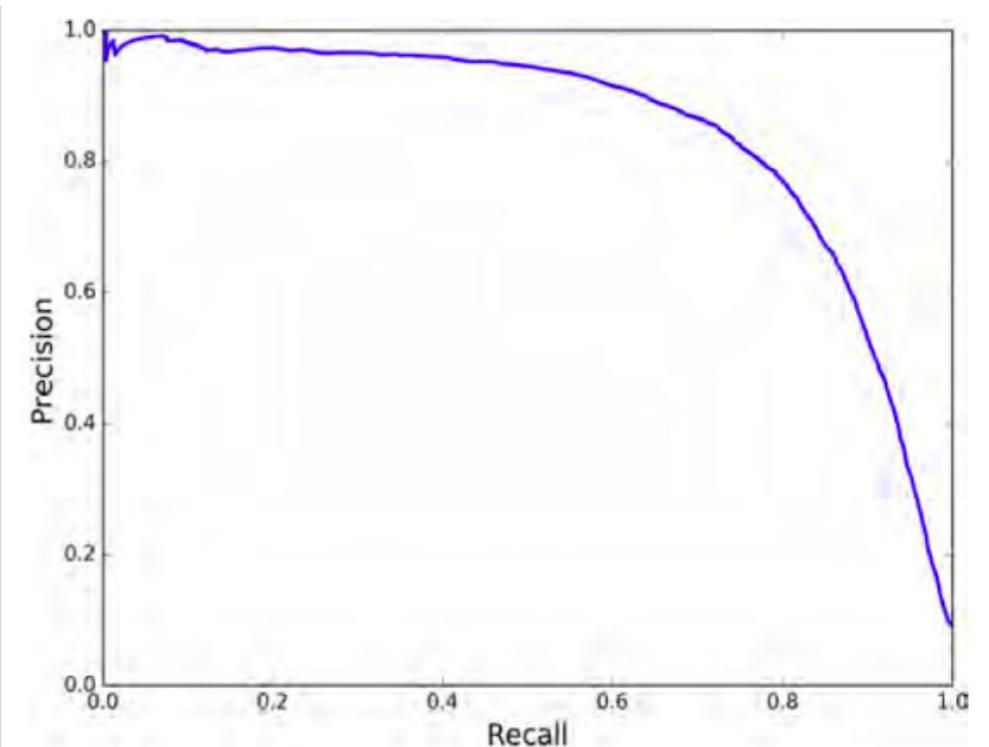
Some example images and annotations in (a) PASCAL-VOC07, (b) ILSVRC, (c) MS-COCO, and (d) Open Images.

Evaluation metrics

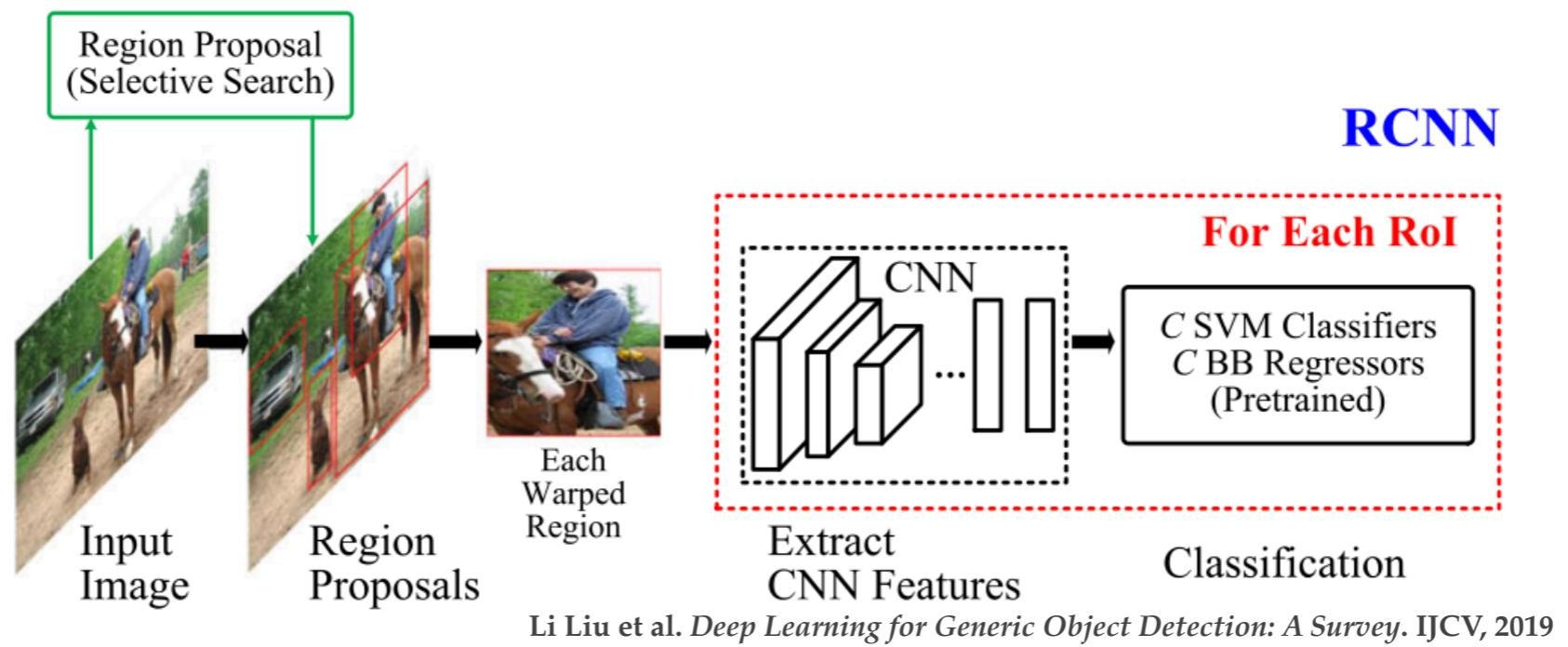
- ❖ How fast?
 - ❖ FPS - Frames Per Second
- ❖ How good?
 - ❖ Localization: **IoU** - measure **correctness** (**positiveness**) of a detection bounding box.
 - ❖ Classification: **confidence** represented with probability
 - ❖ Recall = TP / (TP+FN)
 - ❖ Precision = TP / (TP+FP)
 - ❖ AP - average detection precision under different recalls (average over different IoUs)
 - ❖ mAP - average over all the object categories

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


fendoubasaonian



CNN based two-stage detectors - RCNN

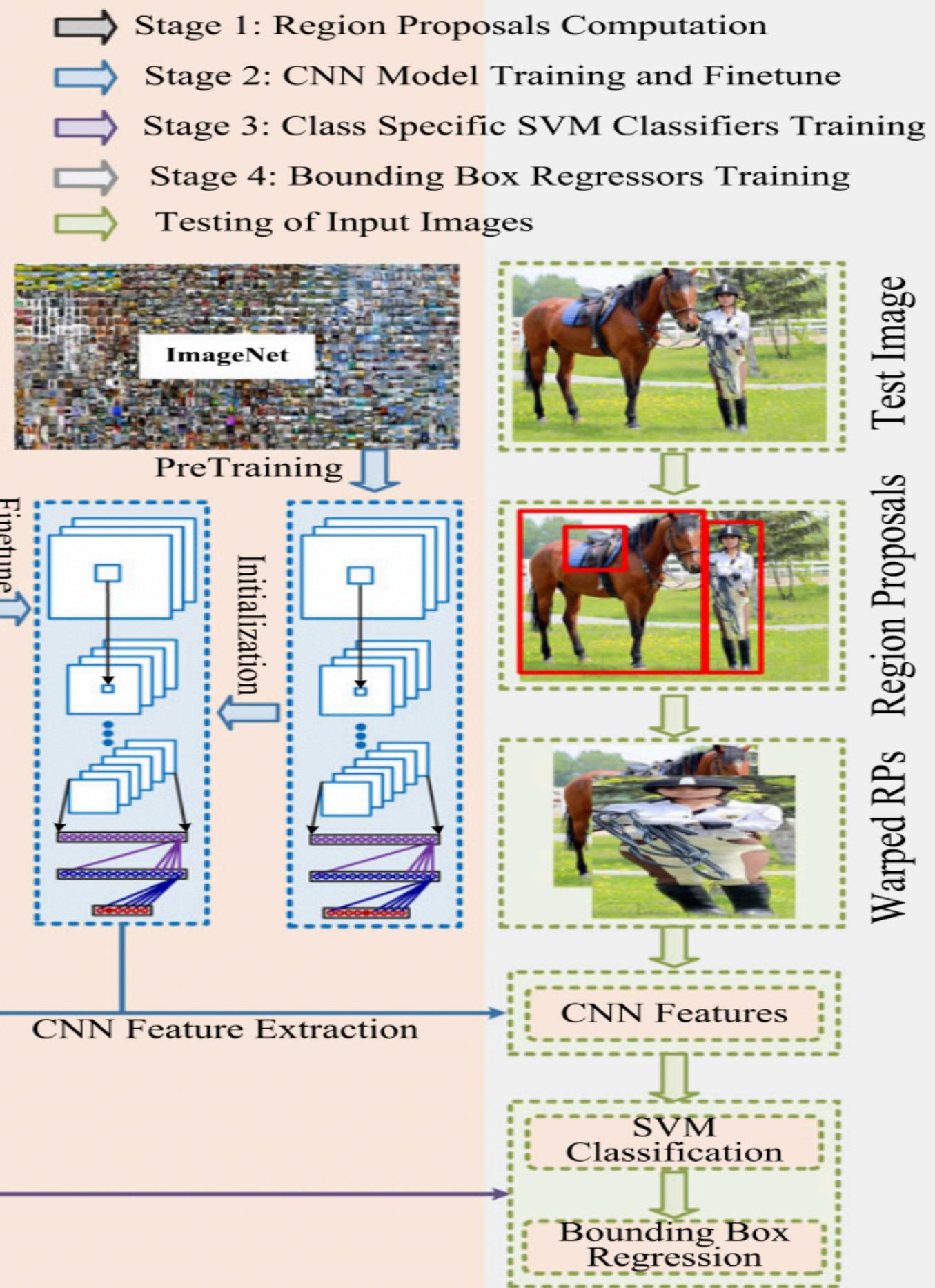


- ❖ Use the traditional selective search to obtain the proposals / RoI (Region of Interests)
- ❖ Crop the image patches within each proposal, and employ CNN to extract features.
- ❖ Use the extracted CNN features for SVM classification learning and inference.

Training



Testing



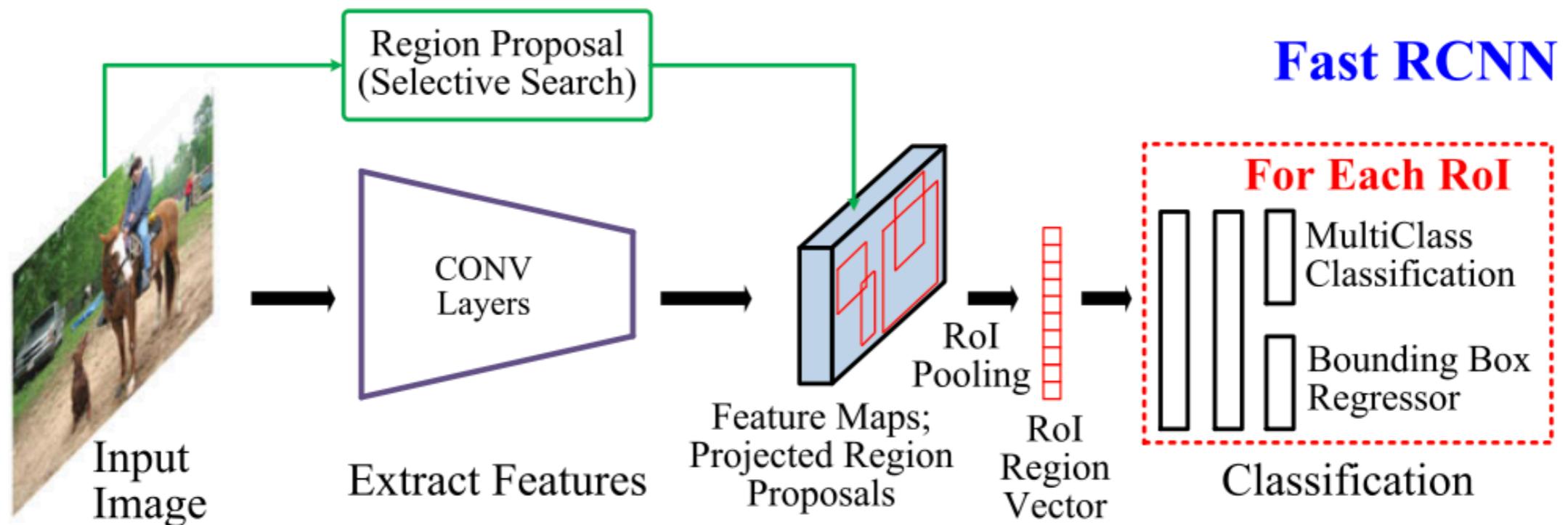
CNN based two-stage detectors - RCNN

- ❖ Pros:
 - ❖ Introducing CNN into object detection which achieves excellent performances.
- ❖ Cons:
 - ❖ Training is a multistage pipeline, slow and hard to optimize because each individual stage must be trained separately.
 - ❖ The CNN computation is redundant and time-consuming.
 - ❖ Testing is slow, since CNN features are extracted per object proposal in each test image, without shared computation.

CNN based two-stage detectors - SPPNet

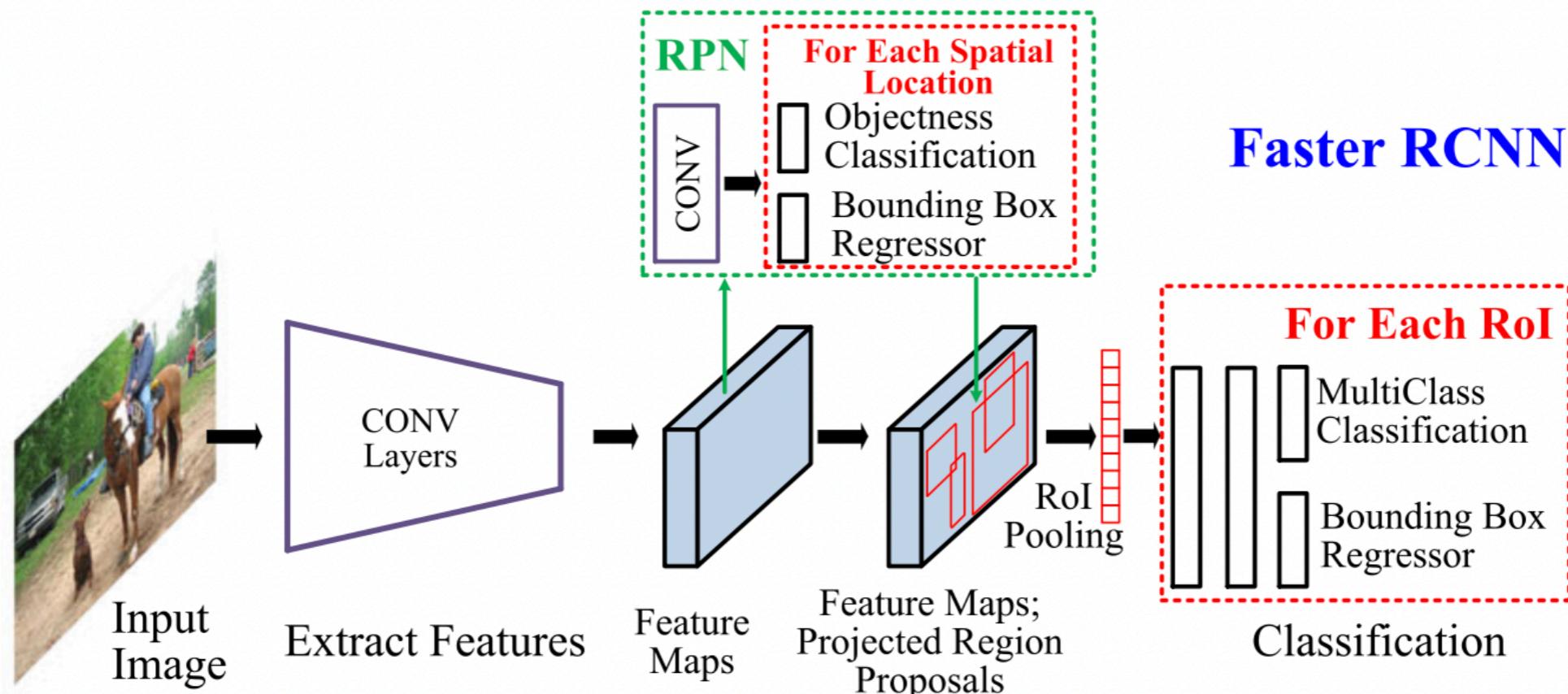
- ❖ Pros:
 - ❖ Enables a CNN to generate a fixed-length representation regardless of the size of image/region of interest.
 - ❖ The feature maps can be computed from the entire image only once
 - ❖ More than 20 times faster than R-CNN without sacrificing any detection accuracy (VOC07 mAP=59.2%).
- ❖ Cons:
 - ❖ Training still is multistage
 - ❖ Only fine-tunes its fully connected layers while simply ignores all previous layers

CNN based two-stage detectors - FastRCNN



- ❖ Pros:
 - ❖ Enables simultaneously training of a detector and a bounding box regressor under the same network configurations.
 - ❖ Fast RCNN increased the mAP from 58.5% (RCNN) to 70.0% while with a detection speed over 200 times faster than R-CNN.
- ❖ Cons:
 - ❖ Selective search still is used to generate proposals, so still slow.

CNN based two-stage detectors - FasterRCNN



❖ Pros:

- ❖ The selective search is replaced by a CNN in producing region proposals. An efficient and accurate Region Proposal Network (RPN) for generating region proposals.
- ❖ **Propose the anchor mechanism**, set anchor boxes with various scales and aspect ratios at each location on the feature map and use Network to learn the classification and regression
- ❖ ~5 FPS (including all stages) on a GPU, achieving good accuracy on PASCAL VOC 2007 using 300 proposals per image

❖ Cons

- ❖ Still slow with some fully connected layers

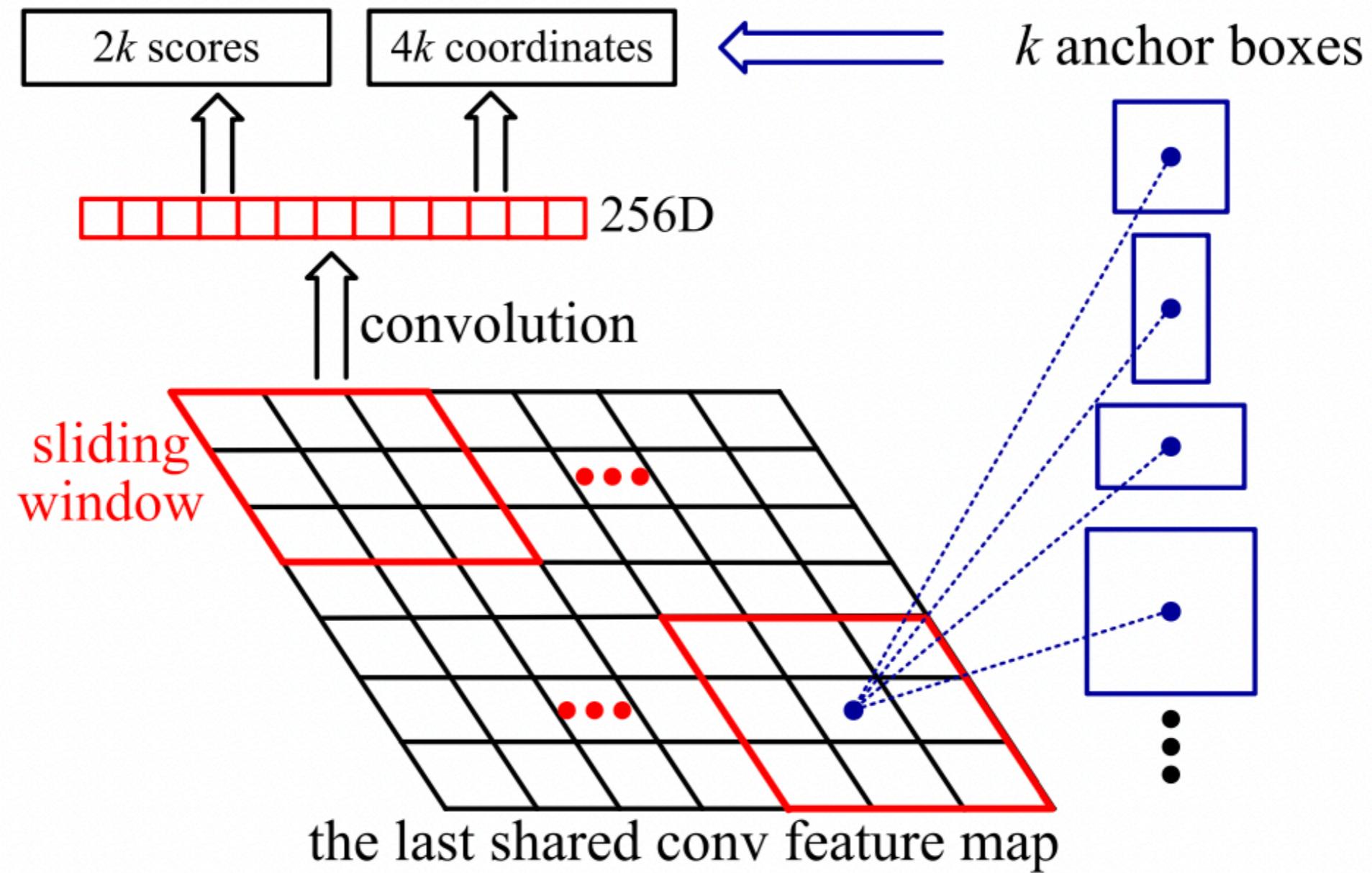
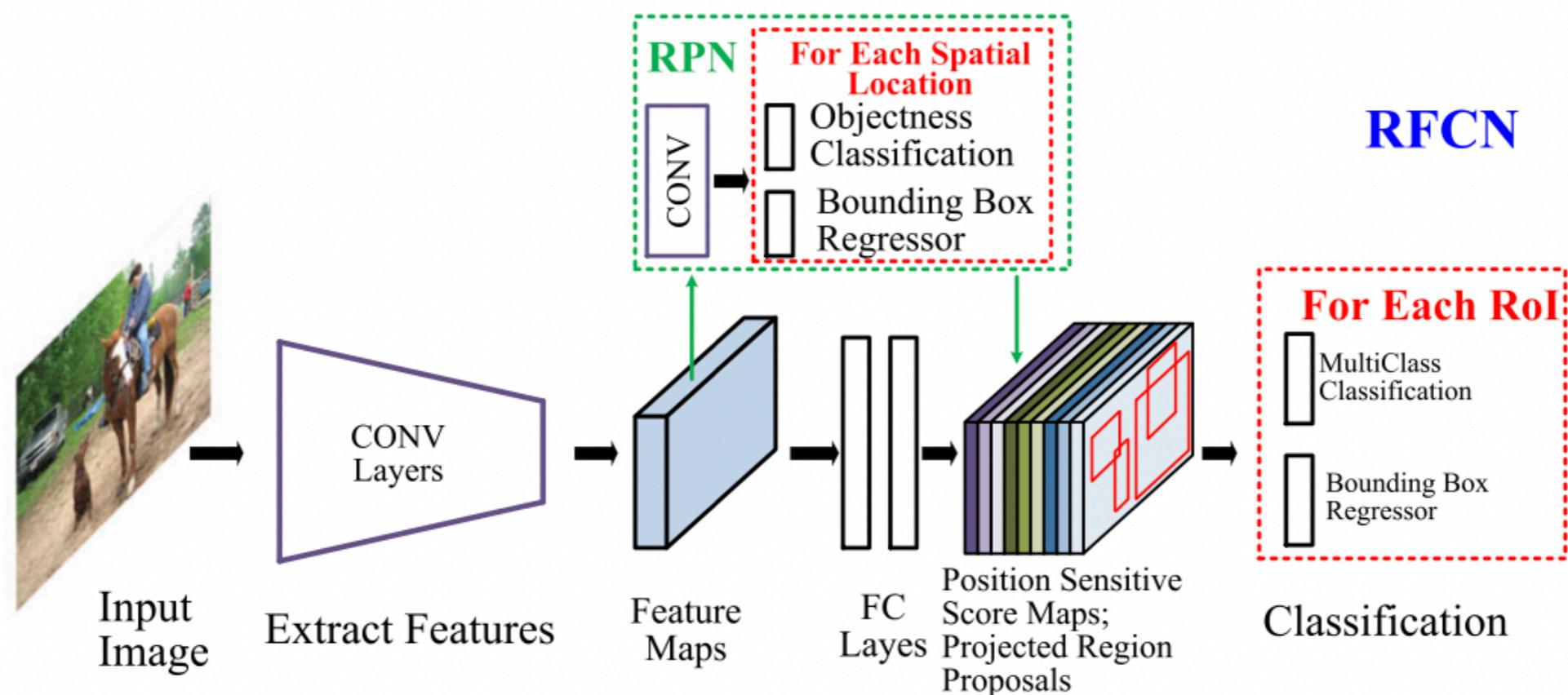


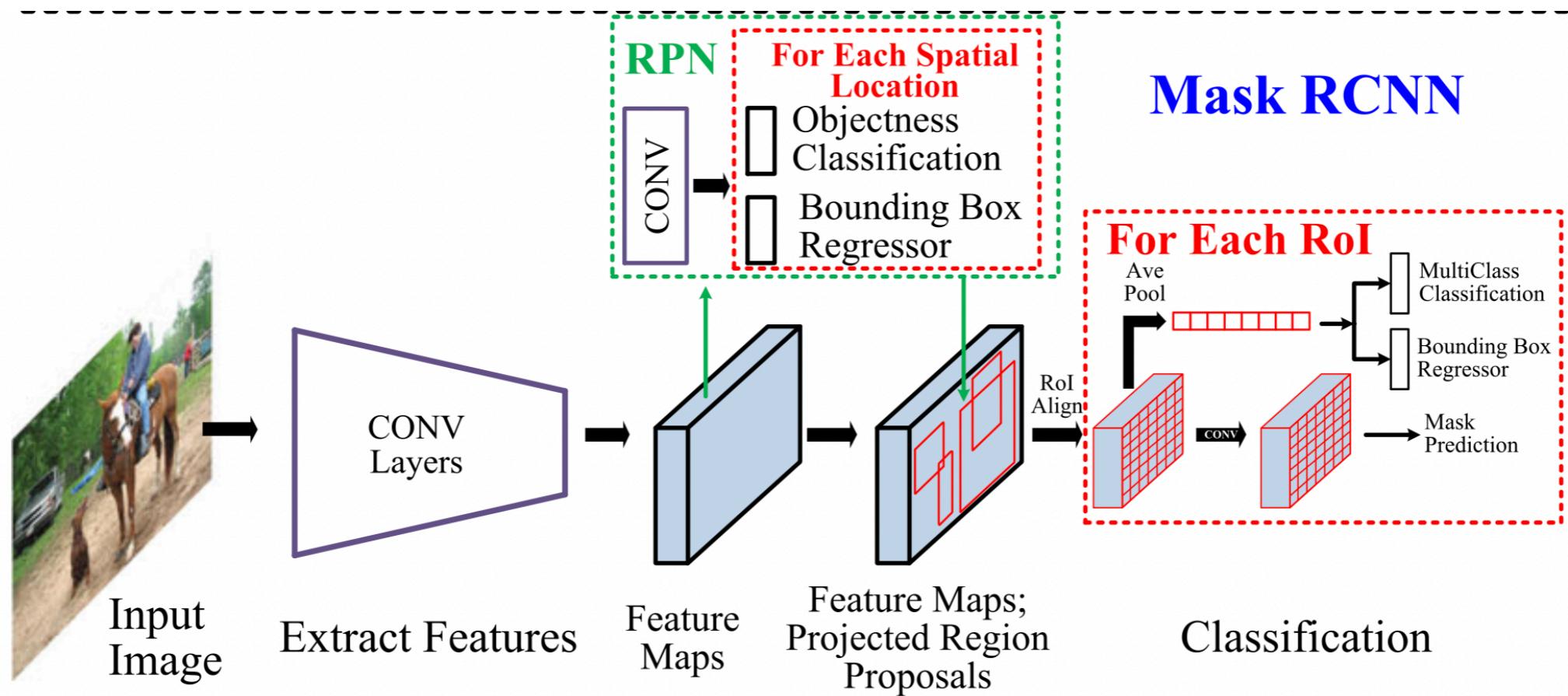
Illustration of the region proposal network (RPN)

CNN based two-stage detectors - RFCN



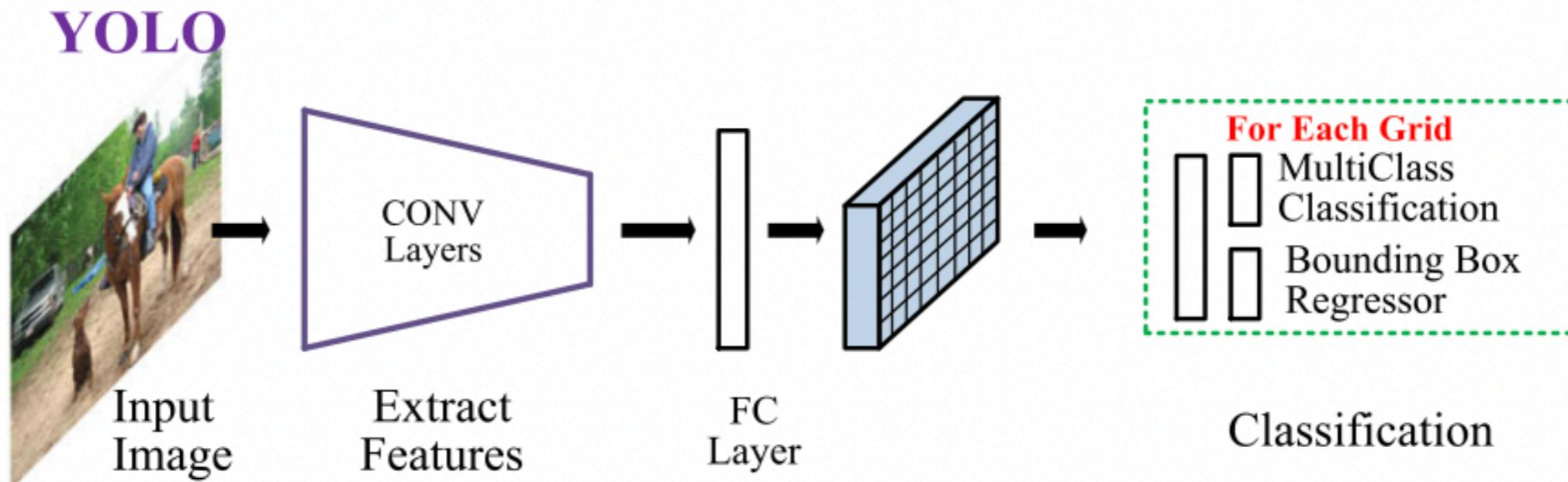
- ❖ Pros:
 - ❖ No hidden FC layers
 - ❖ All CONV layers to construct a shared ROI sub-network

CNN based two-stage detectors - MaskRCNN



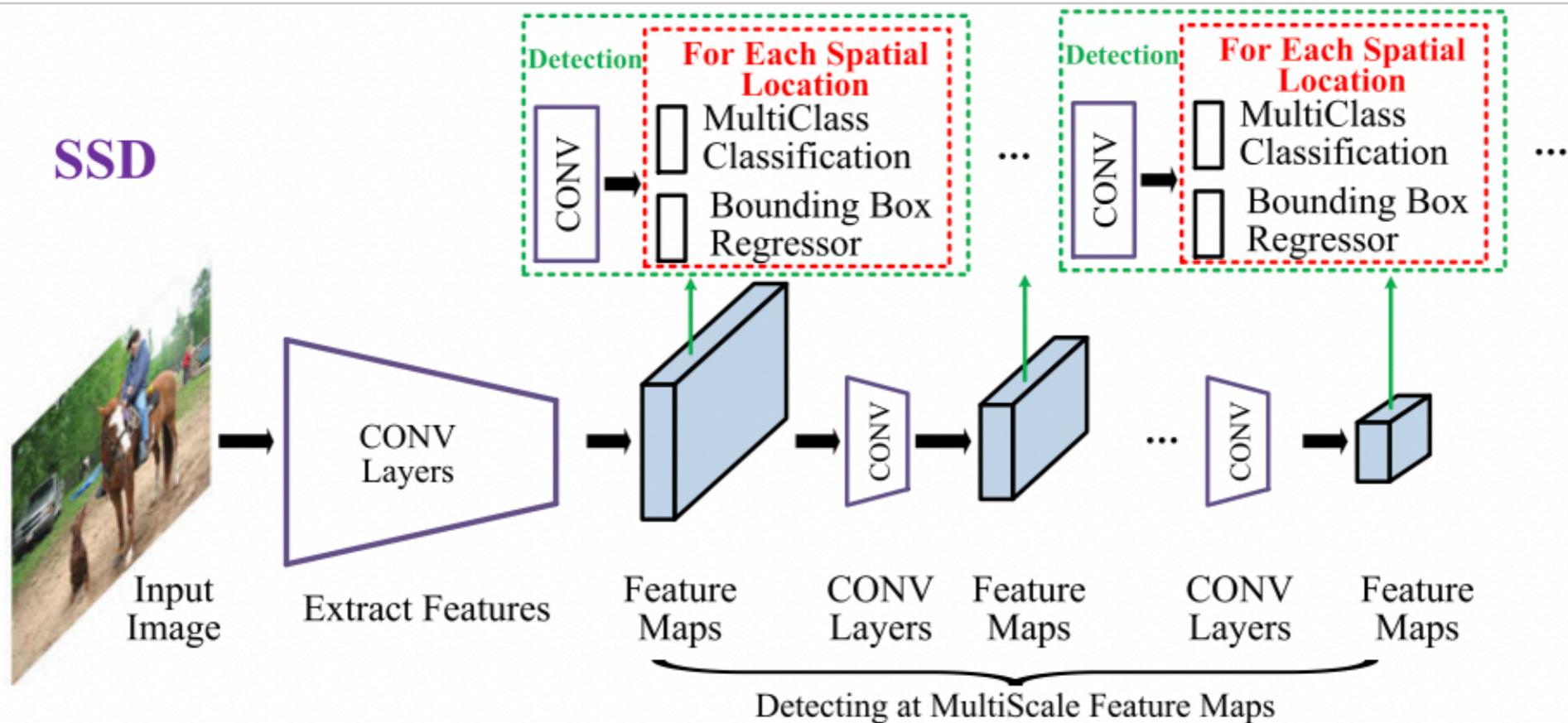
- ❖ Predict the spatial segmentation mask along with classification and bounding box regression
- ❖ Segmentation loss, classification loss and bounding box regression loss are optimized together.
- ❖ RoI Align is proposed to use interpolation operation to obtain the precise RoI pooling feature value.
- ❖ 5FPS

CNN based one-stage detectors - YOLO



- ❖ The first one-stage detector in deep learning era
- ❖ Directly predicts detections using a small set of candidate regions
- ❖ Divides an image into an $S \times S$ grid, each predicting C class probabilities, B bounding box locations, and confidence scores.
- ❖ 45FPS-155FPS
- ❖ YOLO-v2, v3,...,v5, v8
- ❖ YOLO9000 (allows weakly supervised detection)
- ❖ Cons: coarse localization

CNN based one-stage detectors - SSD

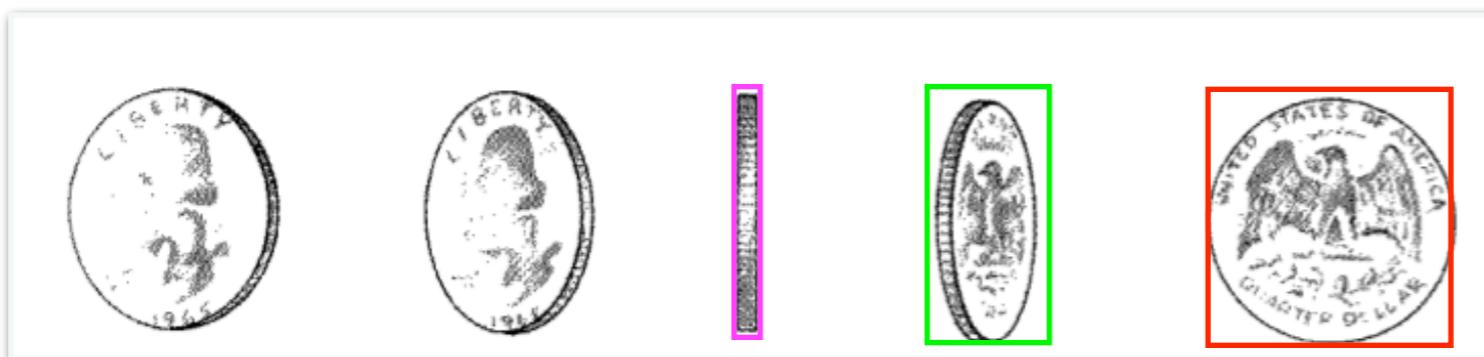


- ❖ Performs detection over multiple scales by operating on multiple CONV feature maps, each of which predicts category scores and box offsets for bounding boxes of appropriate sizes.
- ❖ For a 300×300 input, SSD achieves 74.3% mAP on the VOC2007 test at 59 FPS versus Faster RCNN 7 FPS / mAP 73.2% or YOLO 45 FPS / mAP 63.4%.

Popular DCNN architectures

No.	DCNN architecture	#Paras ($\times 10^6$)	#Layers (CONV+FC)	Test error (Top 5)	First used in	Highlights
1	AlexNet (Krizhevsky et al. 2012b)	57	5 + 2	15.3%	Girshick et al. (2014)	The first DCNN found effective for ImageNet classification; the historical turning point from hand-crafted features to CNN; Winning the ILSVRC2012 Image classification competition
2	ZFNet (fast) (Zeiler and Fergus 2014)	58	5 + 2	14.8%	He et al. (2014)	Similar to AlexNet, different in stride for convolution, filter size, and number of filters for some layers
3	OverFeat (Sermanet et al. 2014)	140	6 + 2	13.6%	Sermanet et al. (2014)	Similar to AlexNet, different in stride for convolution, filter size, and number of filters for some layers
4	VGGNet (Simonyan and Zisserman 2015)	134	13 + 2	6.8%	Girshick (2015)	Increasing network depth significantly by stacking 3×3 convolution filters and increasing the network depth step by step
5	GoogLeNet (Szegedy et al. 2015)	6	22	6.7%	Szegedy et al. (2015)	Use Inception module, which uses multiple branches of convolutional layers with different filter sizes and then concatenates feature maps produced by these branches. The first inclusion of bottleneck structure and global average pooling
6	Inception v2 (Ioffe and Szegedy 2015)	12	31	4.8%	Howard et al. (2017)	Faster training with the introduce of batch normalization
7	Inception v3 (Szegedy et al. 2016)	22	47	3.6%		Inclusion of separable convolution and spatial resolution reduction
8	YOLONet (Redmon et al. 2016)	64	24 + 1	—	Redmon et al. (2016)	A network inspired by GoogLeNet used in YOLO detector
9	ResNet50 (He et al. 2016)	23.4	49	3.6% (ResNets)	He et al. (2016)	With identity mapping, substantially deeper networks can be learned

Important Issues - Scale



Important Issues - Scale

- ❖ Scale problems: “size” and “aspect ratio”
- ❖ Feature pyramids and sliding windows (2004-2014)
 - ❖ VJ\HOG\DPN
 - ❖ mixture model for variations in aspect ratio
- ❖ Object proposals: generate boxes at varies size and aspect ratio (2010-2015)
- ❖ Deep regression: directly predict the coordinates of a bounding box based on the deep features (2013-2016)
- ❖ Multi-reference/-resolution detection (after 2015)
 - ❖ Anchor box
 - ❖ Faster RCNN, SSD

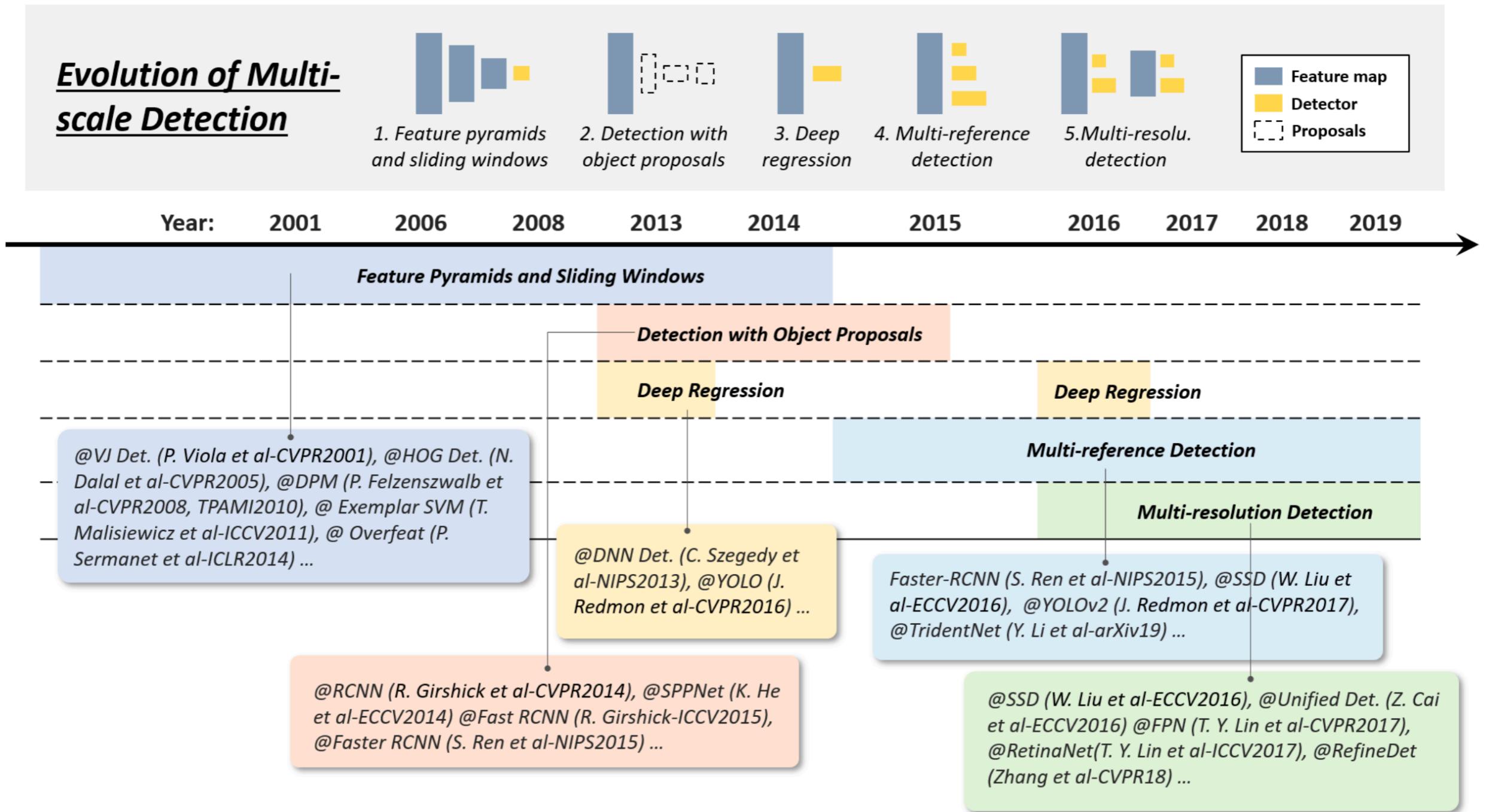
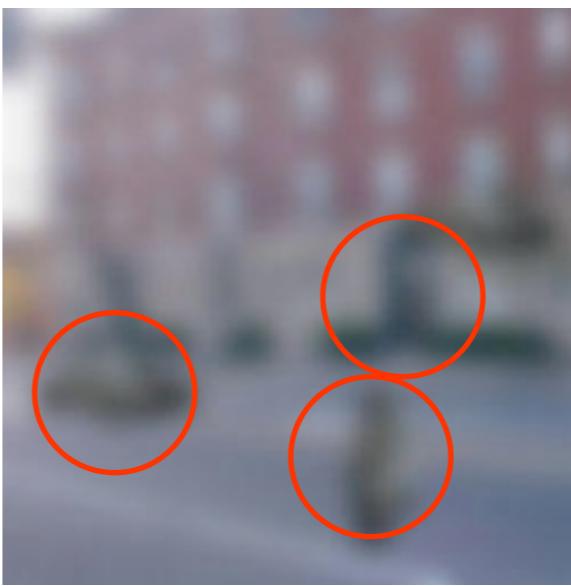
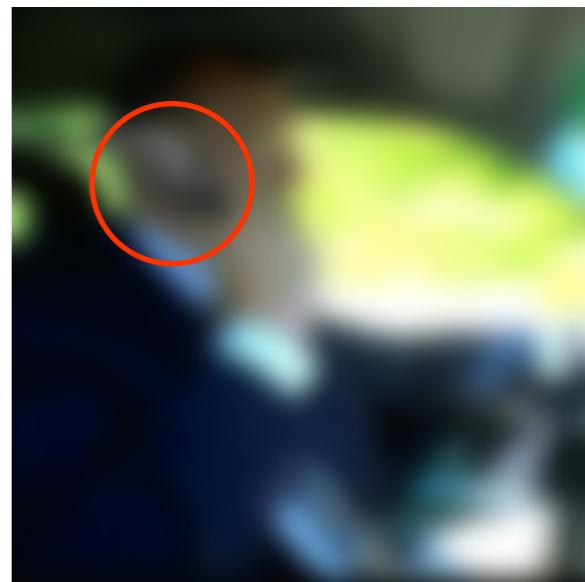
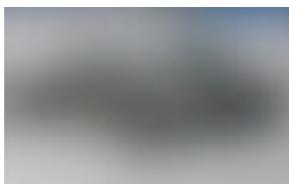


Fig. 6. Evolution of multi-scale detection techniques in object detection from 2001 to 2019: 1) feature pyramids and sliding windows, 2) detection with object proposals, 3) deep regression, 4) multi-reference detection, and 5) multi-resolution detection. Detectors in this figure: VJ Det. [10], HOG Det. [12], DPM [13, 15], Exemplar SVM [36], Overfeat [103], RCNN [16], SPPNet [17], Fast RCNN [18], Faster RCNN [19], DNN Det. [104], YOLO [20], YOLO-v2 [48], SSD [21], Unified Det. [105], FPN [22], RetinaNet [23], RefineDet [55], TridentNet [56].

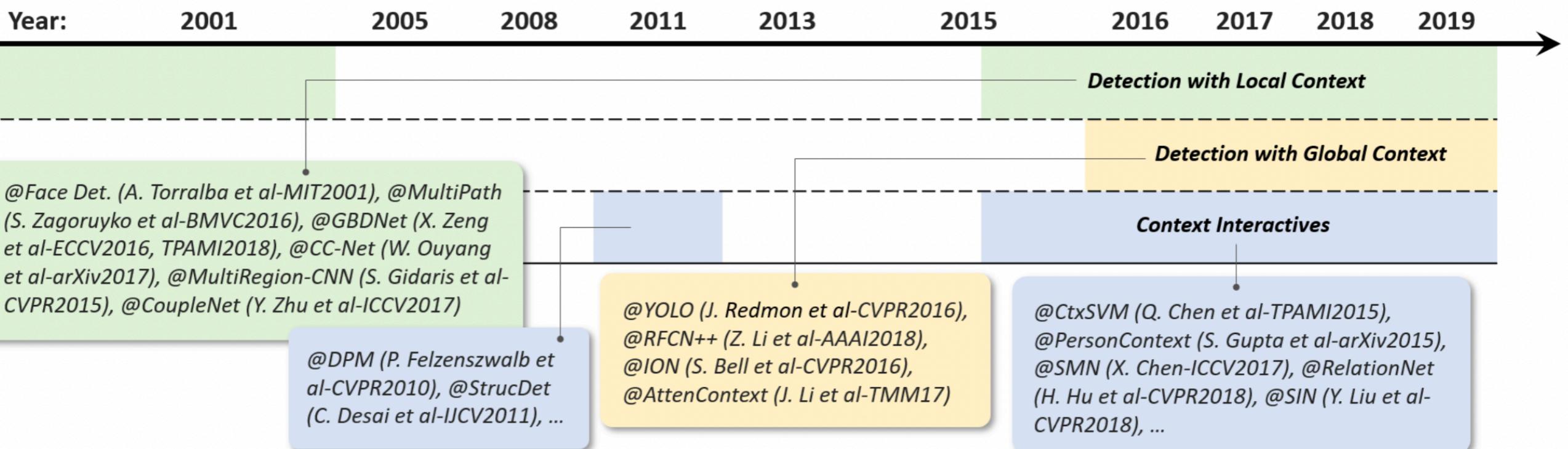
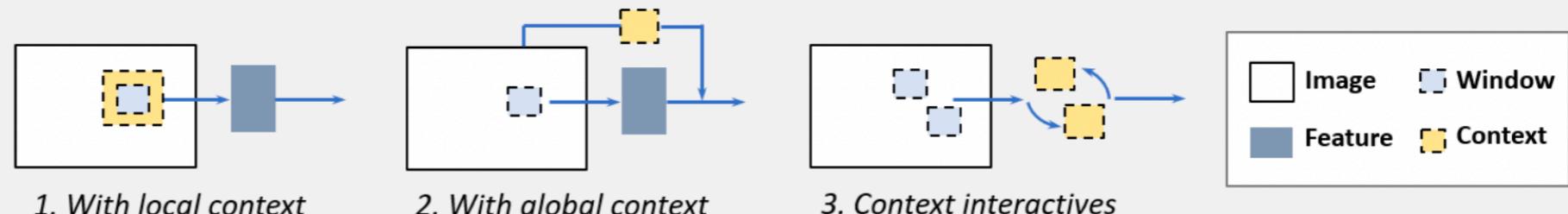
Important Issues - Context



Important Issues - Context

- ❖ Local context: Local context refers to the visual information in the area that surrounds the object to detect
 - ❖ enlarge the networks' receptive field or the size of object proposals
 - ❖ incorporate a small amount of background information
 - ❖ inclusion of local contextual regions
- ❖ Global context: exploits scene configuration as an additional source of information for object detection
 - ❖ integrate a statistical summary of the scene
 - ❖ large receptive field or global pooling operation
 - ❖ think of the global context as a kind of sequential information and to learn it with the recurrent neural networks
- ❖ Context interactive: object-object, object-scene

Evolution of Context Priming in Object Detection



Important Issues - NMS

- ❖ NMS - Non-Maximum Suppression
 - ❖ post-processing step to remove the replicated boxes
- ❖ Greedy selection: the bounding box with the maximum detection score is selected while its neighboring boxes are removed according to a predefined overlap threshold

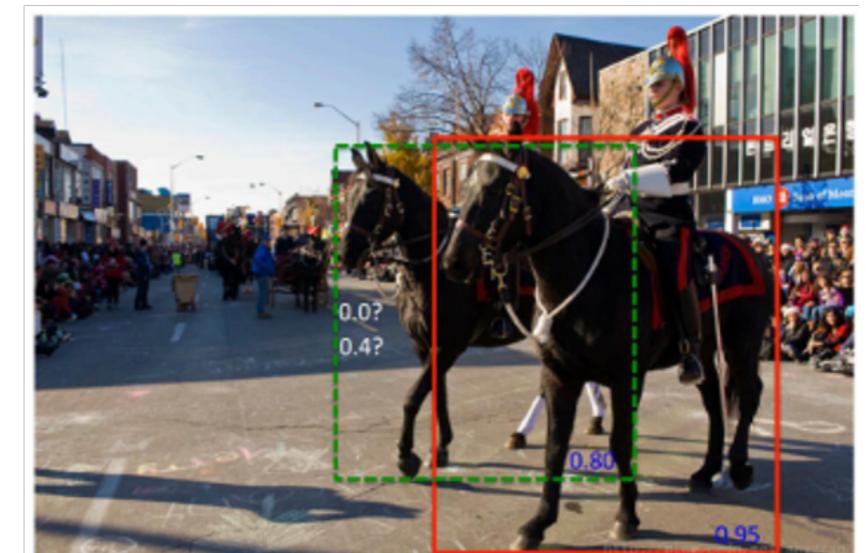
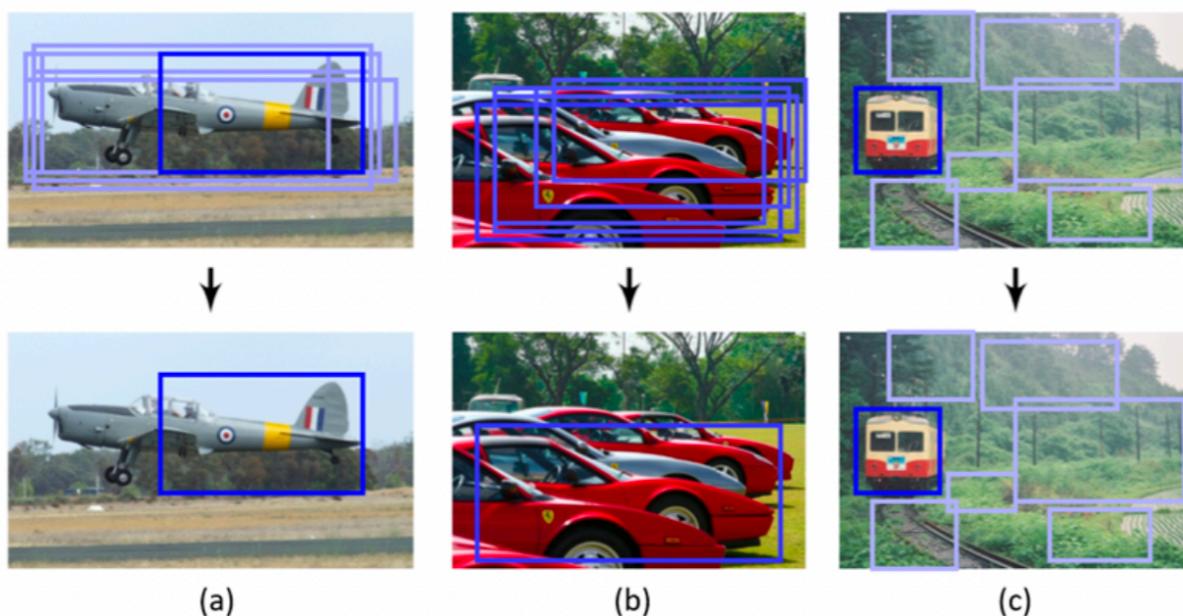
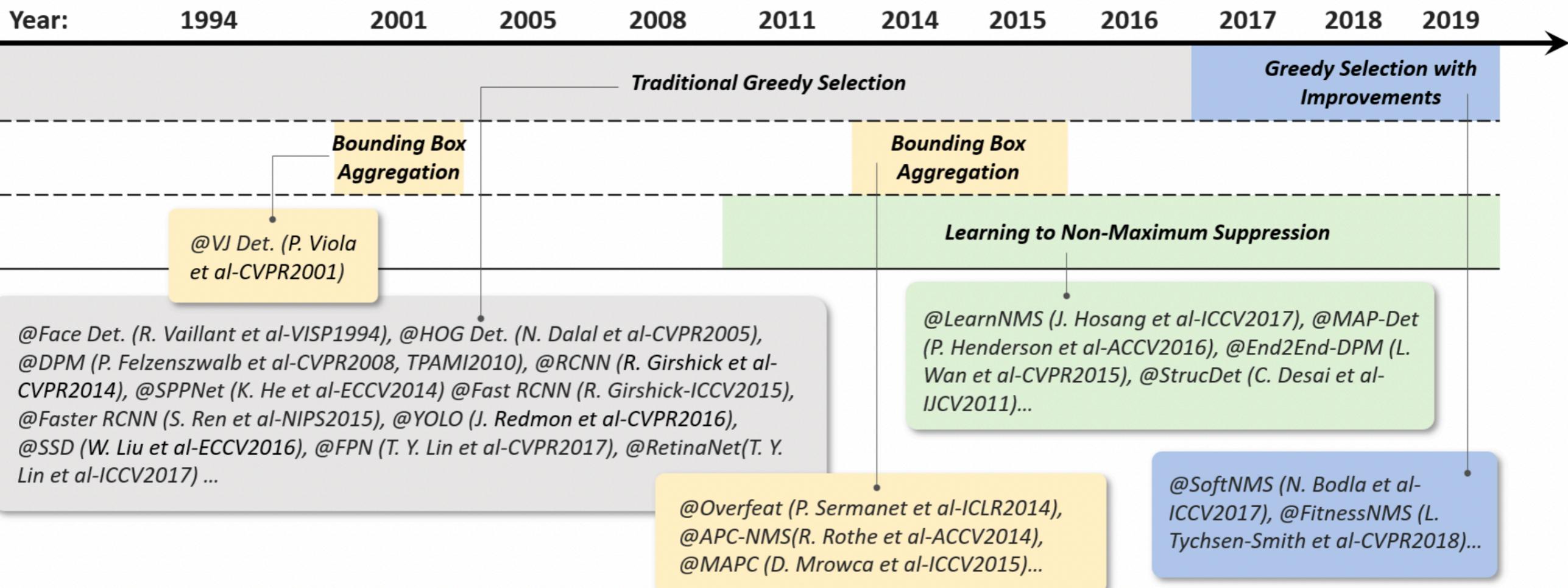
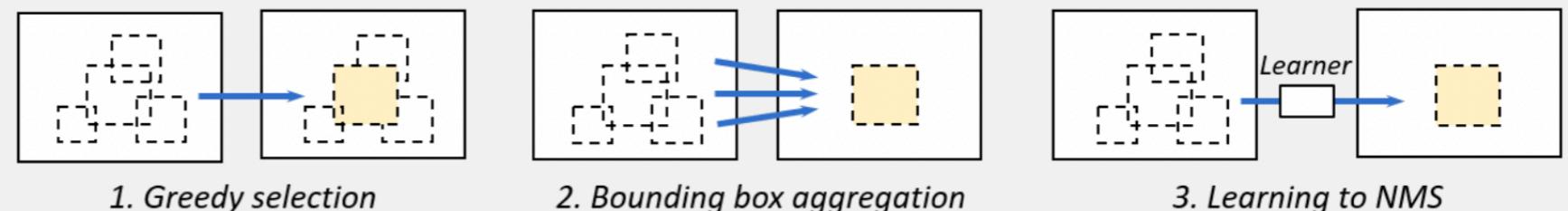


Fig. 11. Examples of possible failures when using a standard greedy selection based non-max suppression: (a) the top-scoring box may not be the best fit, (b) it may suppress nearby objects, and (c) it does not suppress false positives. Images from *R. Rothe et al. ACCV2014* [156].

Important Issues - NMS

- ❖ NMS - Non-Maximum Suppression
 - ❖ post-processing step to remove the replicated boxes
- ❖ BB aggregation:
 - ❖ combining or clustering multiple overlapped bounding boxes into one final detection
- ❖ Learning to NMS:
 - ❖ think of NMS as a filter to re-score all raw detections and to train the NMS as part of a network in an end-to-end fashion

Evolution of Non-Max Suppression



Important Issues - HNM

- ❖ Aims to deal with the problem of imbalanced data during training
 - ❖ Easy negatives: background samples, too many
 - ❖ Bootstrap: the training starts with a small part of background samples and then iteratively add new misclassified backgrounds during the training process
 - ❖ In DNN:
 - ❖ Weight-balancing
 - ❖ Redesign the loss function, so that more focus on hard, misclassified examples
 - ❖ Anchor refinement module

Evolution of Hard Negative Mining

Year:	1994	2001	2005	2008	2014	2015	2016	2017	2018	2019	
Method	Bootstrap				Without Hard Negative Mining		Bootstrap + New Loss Functions				
Remarks	Bootstrap was widely used to deal with the insufficient computing resources of early time				By simply balancing the weights between object and background classes		Focusing on hard examples. Computing power is no longer a problem.				
@Face Det. (H. A. Rowley et al-CMUTechRep1995), @Haar Det. (C. P. Papageorgiou et al-ICCV1998), @VJ Det. (P. Viola et al-CVPR2001), @HOG Det. (N. Dalal et al-CVPR2005), @DPM (P. Felzenszwalb et al-CVPR2008, TPAMI2010)...					@RCNN (R. Girshick et al-CVPR2014), @SPPNet (K. He et al-ECCV2014) @Fast RCNN (R. Girshick-ICCV2015), @Faster RCNN (S. Ren et al-NIPS2015), @YOLO (J. Redmon et al-CVPR2016)...					@SSD (W. Liu et al-ECCV2016), @FasterPed (L. Zhang et al-ECCV2016), @OHEM (A. Shrivastava et al-CVPR2016), @RetinaNet (T. Y. Lin et al-ICCV2017), @RefineDet (Zhang et al-CVPR18)...	

Important Issues - BB Regression

- ❖ Bounding Box regression: refine the location of a predicted bounding box based on the initial proposal or the anchor box
 - ❖ Without BB regression: VJ
 - ❖ From BB to BB: DPM
 - ❖ From feature to BB

Evolution of Bounding Box Regression

weak ← *Invariance of translation and scale* → strong

Year:	2001	2006	2008	2013	2014	2015	2016	2017	2018	2019	
Method	<i>Without Bounding Box Regression</i>	<i>From Bounding Box to Bounding Box</i>				<i>From Feature to Bounding Box</i>					
Remarks			Icing on the cake, optional				Essential, integrated with the model				
	@VJ Det. (P. Viola et al-CVPR2001), @HOG Det. (N. Dalal et al-CVPR2005), @ Exemplar SVM (T. Malisiewicz et al-ICCV2011) ...	@DPM (P. Felzenszwalb et al-CVPR2008, TPAMI2010)				@Overfeat (P. Sermanet et al-ICLR2014), @RCNN (R. Girshick et al-CVPR2014), @SPPNet (K. He et al-ECCV2014) @Fast RCNN (R. Girshick-ICCV2015), @Faster RCNN (S. Ren et al-NIPS2015), @YOLO (J. Redmon et al-CVPR2016), @SSD (W. Liu et al-ECCV2016), @YOLOv2 (J. Redmon et al-CVPR2017), @Unified Det. (Z. Cai et al-ECCV2016) @FPN (T. Y. Lin et al-CVPR2017), @RetinaNet(T. Y. Lin et al-ICCV2017), @RefineDet (Zhang et al-CVPR18), @TridentNet (Y. Li et al-arXiv19) ...					

Important Issues - Backbones

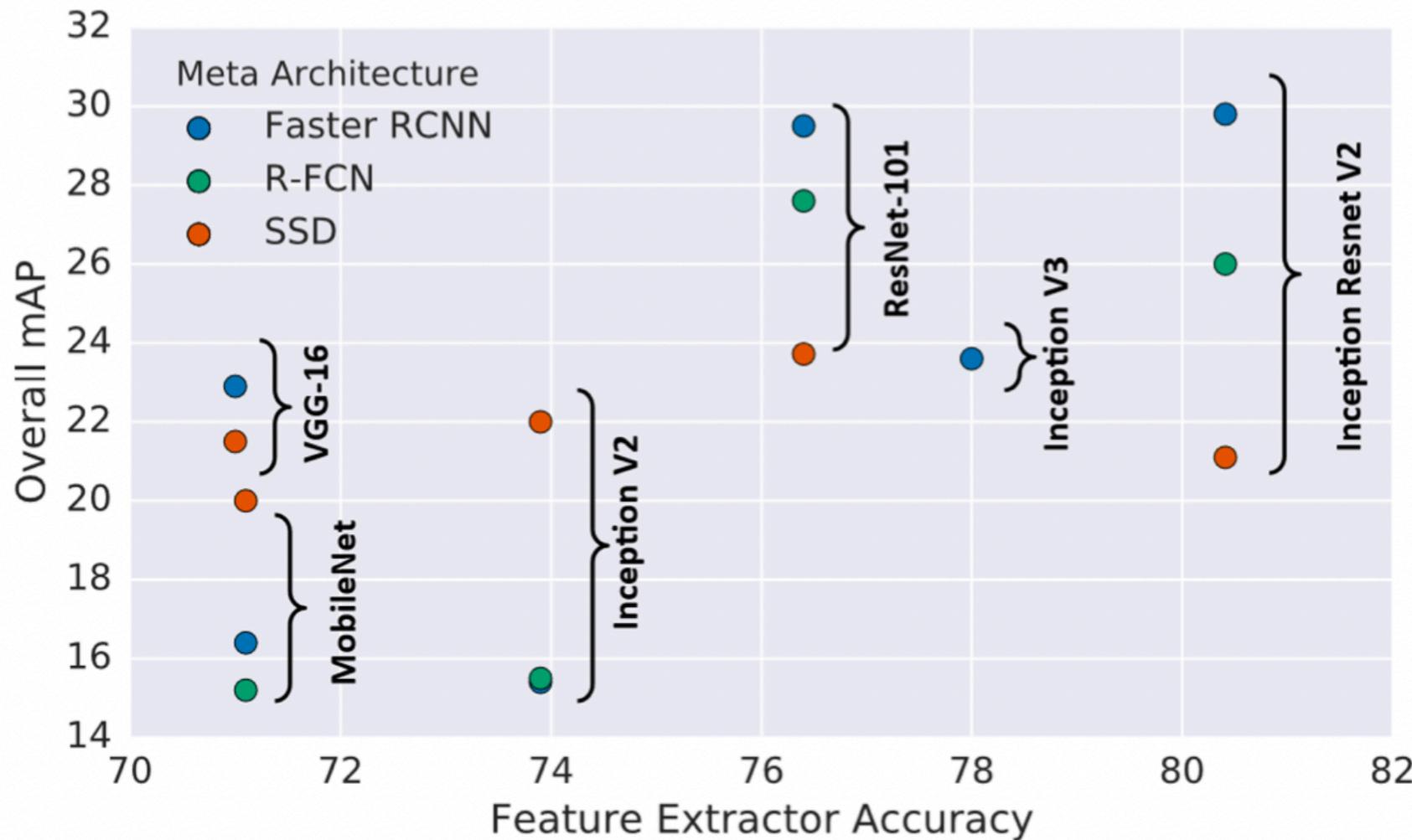


Fig. 17. A comparison of detection accuracy of three detectors: Faster RCNN [19], R-FCN [46] and SSD [21] on MS-COCO dataset with different detection engines. Image from J. Huang et al. CVPR2017 [27].

Important Issues - Feature Fusion

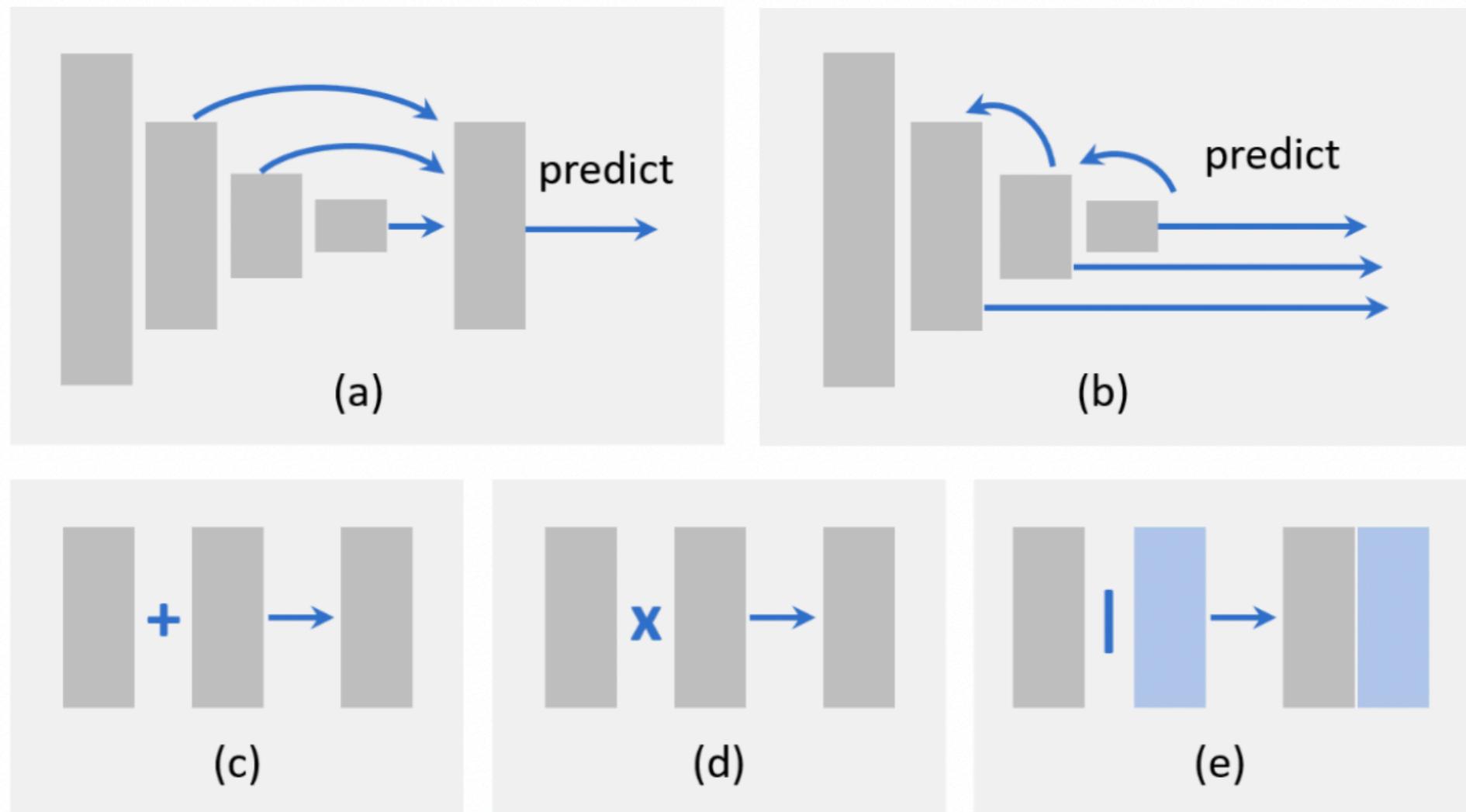
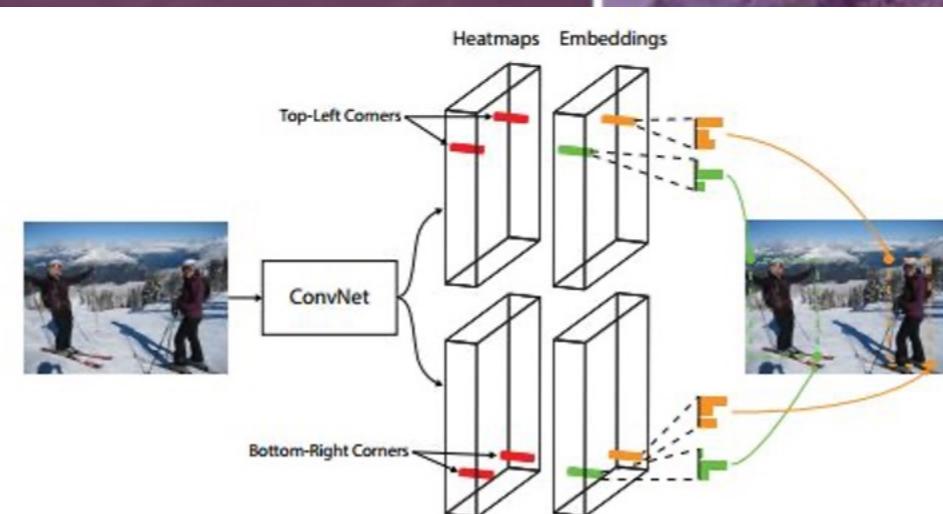
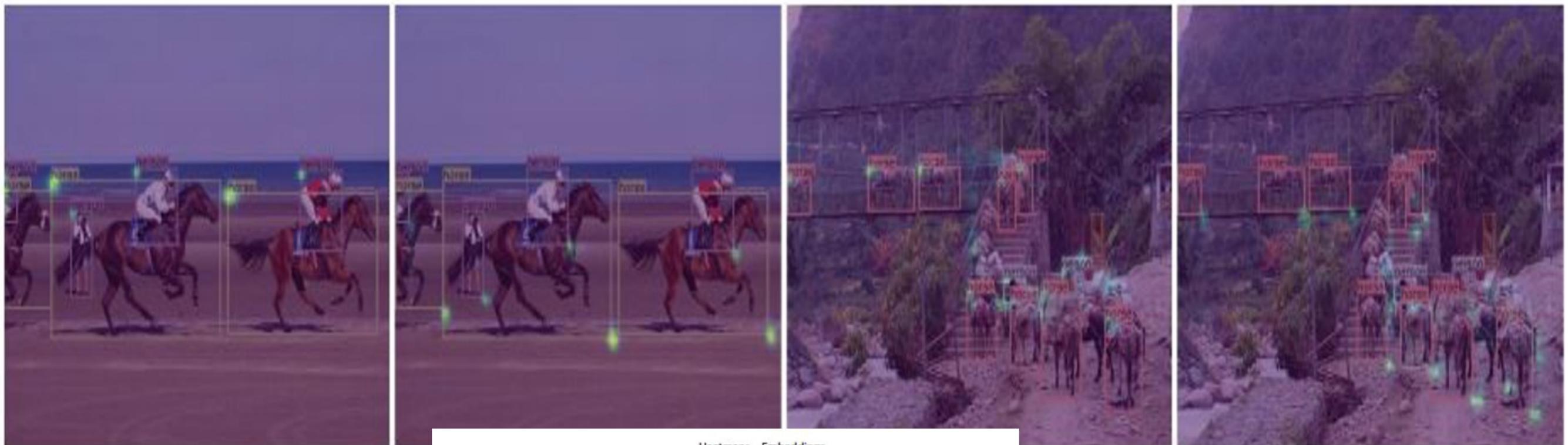


Fig. 18. An illustration of different feature fusion methods: (a) bottom-up fusion, (b) top-down fusion, (c) element-wise sum, (d) element-wise product, and (e) concatenation.

Important Issues - Detection Paradigm

- ❖ Detection as key points localization



H. Law and J. Deng, *CornerNet: Detecting objects as paired keypoints*, ECCV, 2018.

Important Issues - Rotation

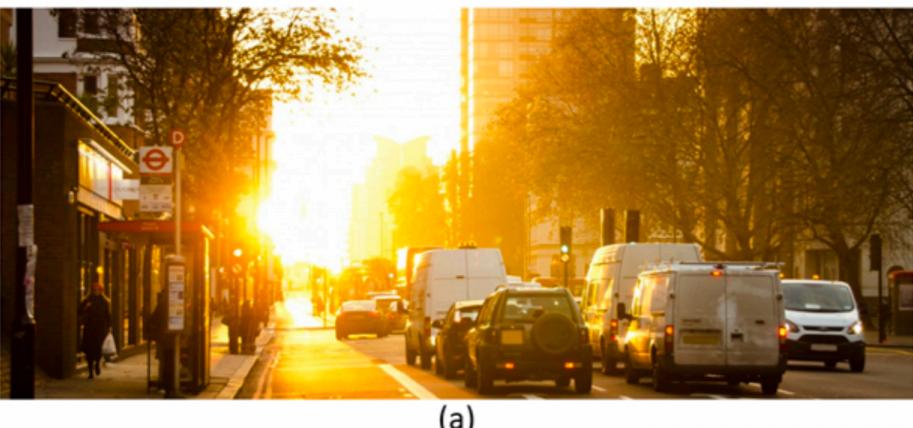
- ❖ How to achieve rotation invariant object detection
 - ❖ Data augmentation
 - ❖ Rotation invariant loss functions: introduced a constraint on the original detection loss function so that to make the features of rotated objects unchanged
 - ❖ Rotation calibration: make geometric transformations. Spatial Transformer Networks (STN)
 - ❖ Rotation RoI Pooling: from Cartesian coordinates to polar coordinates

Important Issues - Adversarial training

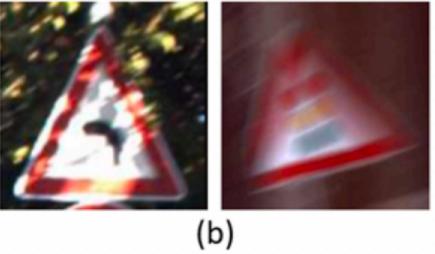
- ❖ Small object detection: GAN has been used to enhance the detection on small objects by narrowing the representations between small and large ones.
- ❖ Occlusion handling: generate occlusion masks by using adversarial training
- ❖ Adversarial attack

Applications

- ❖ Pedestrian detection
- ❖ Face detection
- ❖ Text detection
- ❖ Traffic Sign and Traffic Light Detection
- ❖ Remote Sensing Target Detection



(a)



(b)



(c)

Fig. 23. Challenges in traffic sign detection and traffic light detection: (a) Illumination changes. Image from *pxhere* (free of copyrights). (b) Motion blur. Image from GTSRB Dataset [81]. (c) Detection under bad weather. Image from *Flickr* and *Max Pixel* (free of copyrights).



(a)



(b)



(c)

Fig. 22. Challenges in text detection and recognition: (a) Variation of fonts, colors and languages. Image from *maxpixel* (free of copyrights). (b) Text rotation and perspective distortion. Image from Y. Liu et al. CVPR2017 [336]. (c) Densely arranged text localization. Image from Y. Wu et al. ICCV2017 [337].



(a)

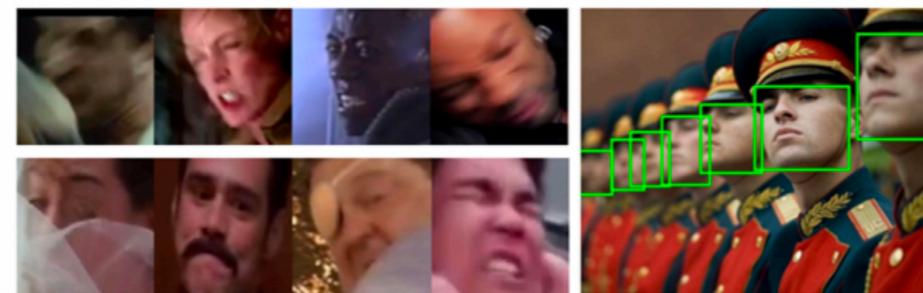


(b)



(c)

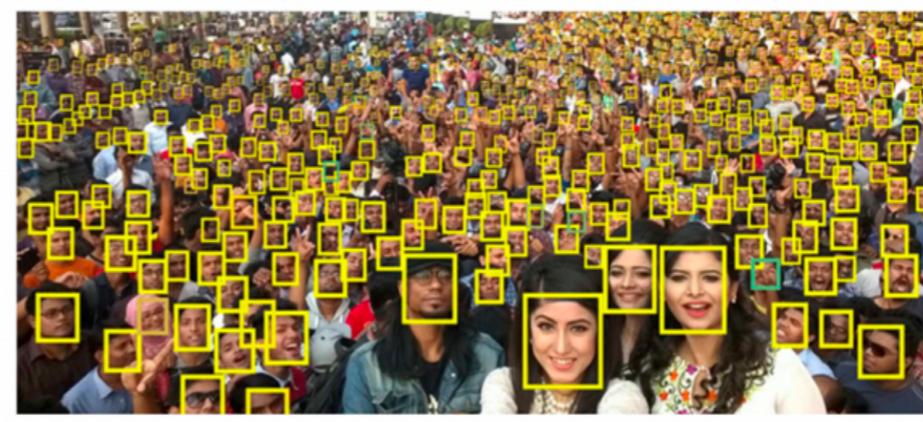
Fig. 20. Some hard examples of pedestrian detection from Caltech dataset [59, 60]: (a) small pedestrians, (b) hard negatives, and (c) dense and occluded pedestrians.



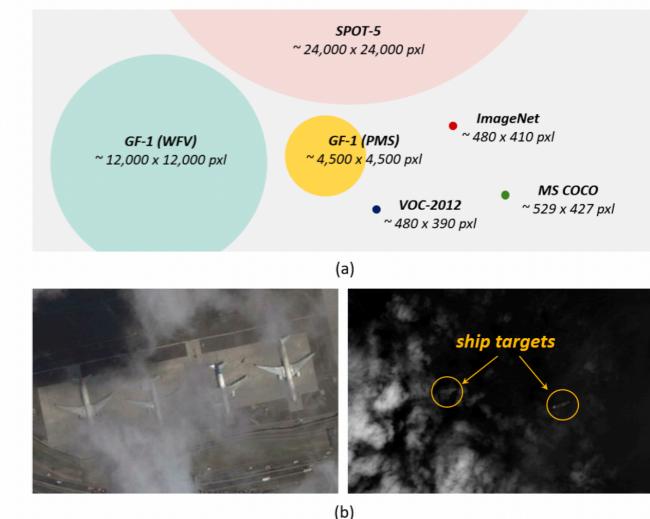
(a)



(b)



(c)



(a)



(b)

Fig. 21. Challenges in face detection: (a) Intra-class variation, image from WildestFaces Dataset [70]. (b) Face occlusion, image from UFDD Dataset [69]. (c) Multi-scale face detection. Image from P. Hu et al. /PR2017 [322].

Fig. 24. Challenges in remote sensing target detection: (a) Detection in “big data”: data volume comparison between a single-view remote sensing imagery and an average image size of VOC, ImageNet, and MS-COCO. (b) Targets occluded by cloud. Images from S. Qiu et al. JSTARS2017 [380] and Z. Zou et al. TGRS2016 [381].

Object detection: future direction

- ❖ Small object detection
- ❖ Weakly supervised detection
- ❖ Fusion of multiple modalities
- ❖ Lightweight object detection
- ❖ Real-time object detection/tracking
- ❖ Few / Zero Shot Object Detection
- ❖ ...