

# Lecture 6: Modern Object Detection

Gang Yu

Face++ Researcher

yugang@megvii.com

# Visual Recognition

A fundamental task in computer vision

- Classification
  - Object Detection
  - Semantic Segmentation
  - Instance Segmentation
  - Key point Detection
  - VQA
- ...



# Category-level Recognition



Category-level Recognition



Instance-level Recognition

# Representation

- Bounding-box
  - Face Detection, Human Detection, Vehicle Detection, Text Detection, general Object Detection
- Point
  - Semantic segmentation (will be discussed in next week)
- Keypoint
  - Face landmark
  - Human Keypoint

# Outline

- Detection
- Human Keypoint
- Conclusion

# Outline

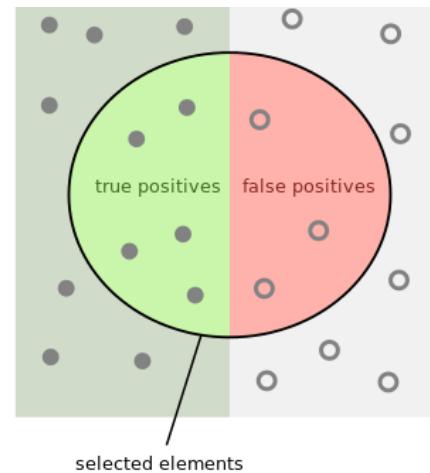
- **Detection**
- Human Keypoint
- Conclusion

# Detection - Evaluation Criteria

## Average Precision (AP) and mAP

Precision and recall are single-value metrics based on the whole list of documents returned by the system. For systems that return a ranked sequence of documents, it is desirable to also consider the order in which the returned documents are presented. By computing a precision and recall at every position in the ranked sequence of documents, one can plot a precision-recall curve, plotting precision  $p(r)$  as a function of recall  $r$ . Average precision computes the average value of  $p(r)$  over the interval from  $r = 0$  to  $r = 1$ .<sup>[9]</sup>

$$\text{AveP} = \int_0^1 p(r)dr$$



How many selected items are relevant?  
How many relevant items are selected?

$$\text{Precision} = \frac{\text{true positives}}{\text{selected elements}}$$
$$\text{Recall} = \frac{\text{true positives}}{\text{relevant items}}$$

# Detection - Evaluation Criteria

mmAP

```
Average Precision (AP):
    AP                  % AP at IoU=.50:.05:.95 (primary challenge metric)
    APIoU=.50        % AP at IoU=.50 (PASCAL VOC metric)
    APIoU=.75        % AP at IoU=.75 (strict metric)

    AP Across Scales:
        APsmall          % AP for small objects: area < 322
        APmedium         % AP for medium objects: 322 < area < 962
        APlarge          % AP for large objects: area > 962

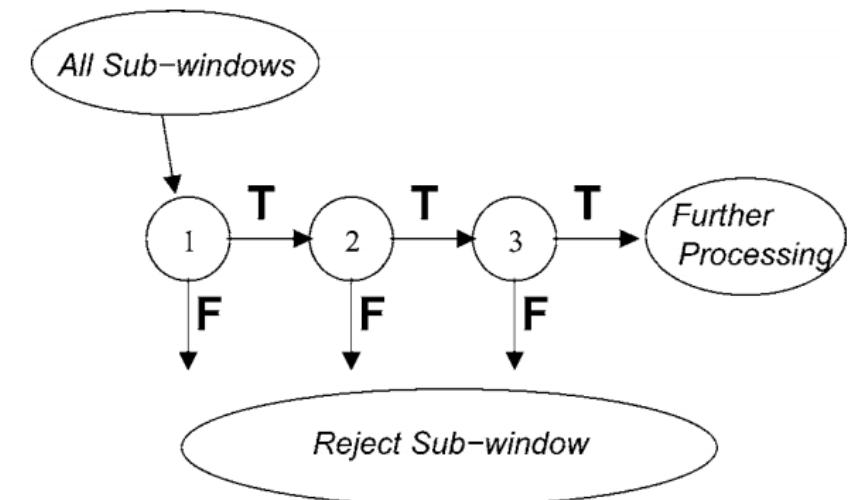
    Average Recall (AR):
        ARmax=1          % AR given 1 detection per image
        ARmax=10         % AR given 10 detections per image
        ARmax=100        % AR given 100 detections per image

    AR Across Scales:
        ARsmall          % AR for small objects: area < 322
        ARmedium         % AR for medium objects: 322 < area < 962
        ARlarge          % AR for large objects: area > 962
```

1. Unless otherwise specified, *AP* and *AR* are averaged over multiple *Intersection over Union (IoU)* values. Specifically we use 10 IoU thresholds of .50:.05:.95. This is a break from tradition, where AP is computed at a single IoU of .50 (which corresponds to our metric  $AP^{IoU=.50}$ ). Averaging over IoUs rewards detectors with better localization.
2. AP is averaged over all categories. Traditionally, this is called "mean average precision" (mAP). We make no distinction between AP and mAP (and likewise AR and mAR) and assume the difference is clear from context.
3. AP (averaged across all 10 IoU thresholds and all 80 categories) will determine the challenge winner. This should be considered the single most important metric when considering performance on COCO.
4. In COCO, there are more small objects than large objects. Specifically: approximately 41% of objects are small ( $area < 32^2$ ), 34% are medium ( $32^2 < area < 96^2$ ), and 24% are large ( $area > 96^2$ ). Area is measured as the number of pixels in the segmentation mask.
5. AR is the maximum recall given a fixed number of detections per image, averaged over categories and IoUs. AR is related to the metric of the same name used in *proposal evaluation* but is computed on a per-category basis.
6. All metrics are computed allowing for at most 100 top-scoring detections per image (across all categories).
7. The evaluation metrics for detection with bounding boxes and segmentation masks are identical in all respects except for the IoU computation (which is performed over boxes or masks, respectively).

# How to perform a detection?

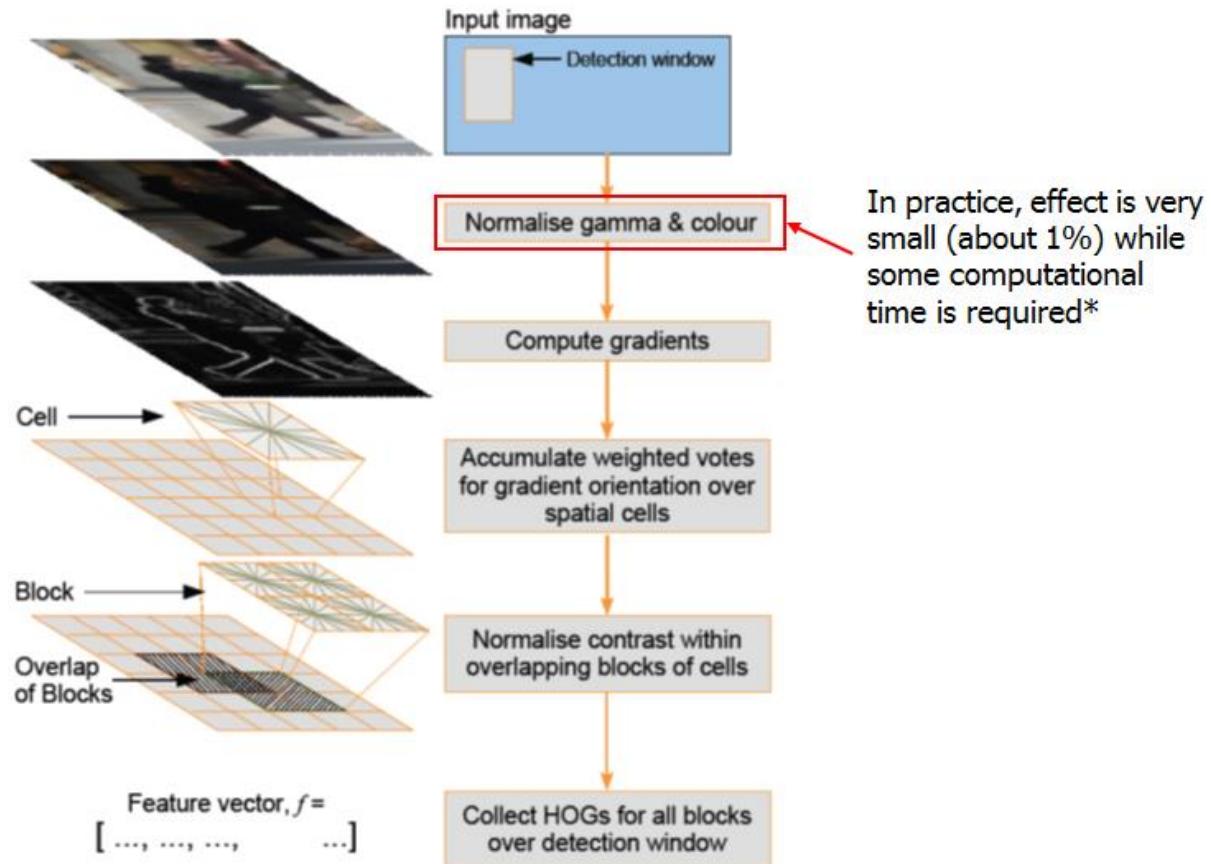
- Sliding window: enumerate all the windows (up to millions of windows)
  - VJ detector: cascade chain
- Fully Convolutional network
  - shared computation



# General Detection Before Deep Learning

- Feature + classifier
- Feature
  - Haar Feature
  - HOG (Histogram of Gradient)
  - LBP (Local Binary Pattern)
  - ACF (Aggregated Channel Feature)
  - ...
- Classifier
  - SVM
  - Bootsing
  - Random Forest

# Traditional Hand-crafted Feature: HoG



\*Navneet Dalal and Bill Triggs. Histograms of Oriented Gradients for Human Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, SanDiego, USA, June 2005. Vol. II, pp. 886-893.

# Traditional Hand-crafted Feature: HoG



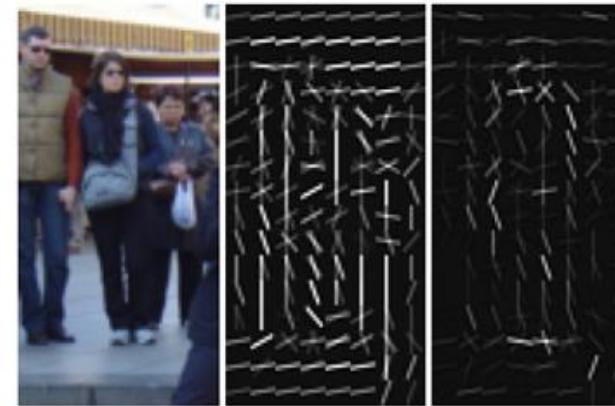
(a)



(b)



(c)



(d)

In each triplet: (1) the input image, (2) the corresponding R-HOG feature vector (only the dominant orientation of each cell is shown),  
the dominant orientations selected by the SVM (obtained by multiplying the feature vector by the corresponding weights from the linear SVM). (3)

# General Detection Before Deep Learning

## Traditional Methods

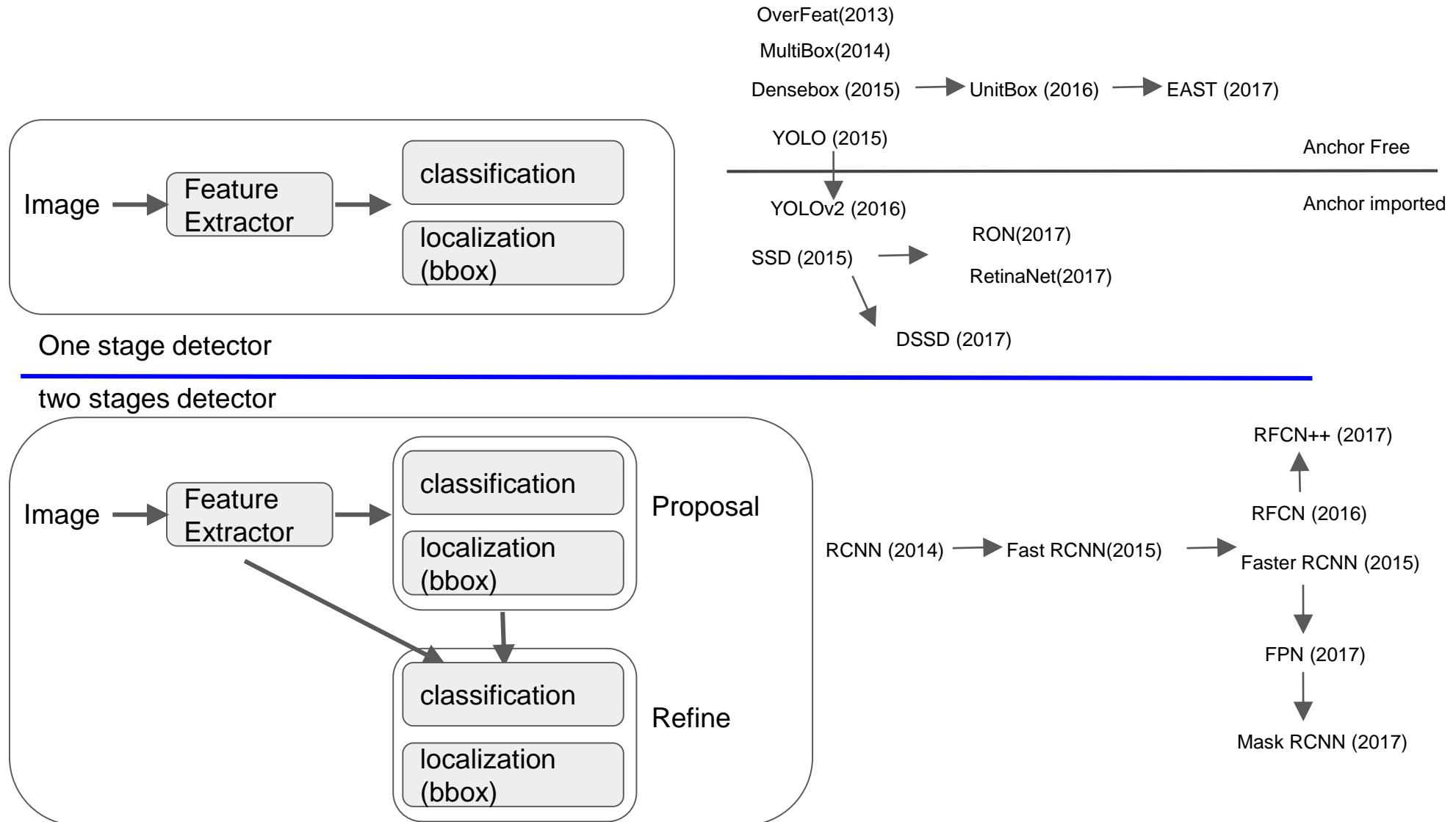
- Pros
  - Efficient to compute (e.g., HAAR, ACF) on CPU
  - Easy to debug, analyze the bad cases
  - reasonable performance on limited training data
- Cons
  - Limited performance on large dataset
  - Hard to be accelerated by GPU

# Deep Learning for Object Detection

Based on the whether following the “proposal and refine”

- One Stage
  - Example: Densebox, YOLO (YOLO v2), SSD, Retina Net
  - Keyword: [Anchor, Divide and conquer, loss sampling](#)
- Two Stage
  - Example: RCNN (Fast RCNN, Faster RCNN), RFCN, FPN, MaskRCNN
  - Keyword: [speed, performance](#)

# A bit of History



# One Stage Detector: Densebox

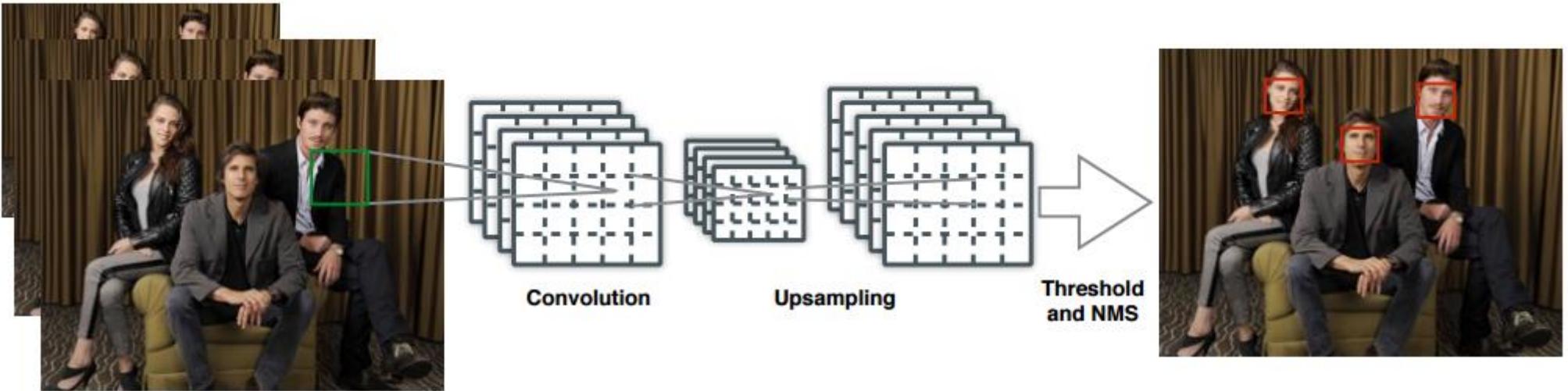


Figure 1: **The DenseBox Detection Pipeline.** 1) Image pyramid is fed to the network. 2) After several layers of convolution and pooling, upsampling feature map back and apply convolution layers to get final output. 3) Convert output feature map to bounding boxes , and apply non-maximum suppression to all bounding boxes over the threshold.

DenseBox: Unifying Landmark Localization with End to End Object Detection, Huang etc, 2015

# One Stage Detector: Densebox

- No Anchor: GT Assignment
  - A sub-circle in the GT is labeled as positive
    - fail when two GT highly overlaps
    - the size of the sub-circle matters
    - more attention (loss) will be placed to large faces
- Loss sampling
  - All pos/negative positions will be used to compute the cls loss

# One Stage Detector: Densebox

## Problems

- L2 loss is not robust to scale variation (**UnitBox**)
  - learnt features are not robust
- GT assignment issue (**SSD**)
  - Fail to handle the crowd case
- relatively large localization error (**Two stages detector**)
- more false positive (FP) (**Two stages detector**)
  - does not obviously kill the fp

# One Stage Detector: Densebox -> UnitBox

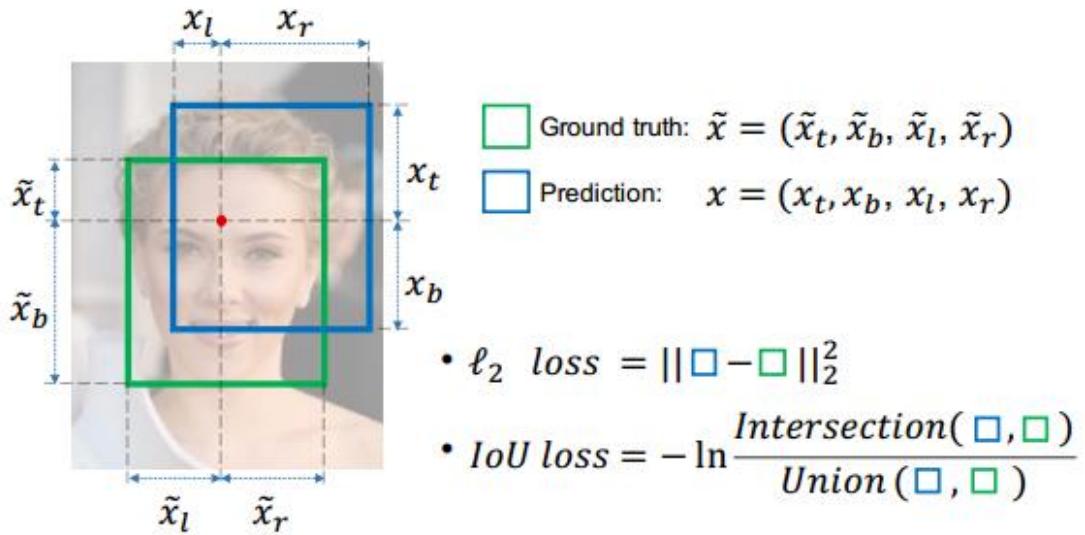


Figure 1: Illustration of  $IoU$  loss and  $\ell_2$  loss for pixel-wise bounding box prediction.

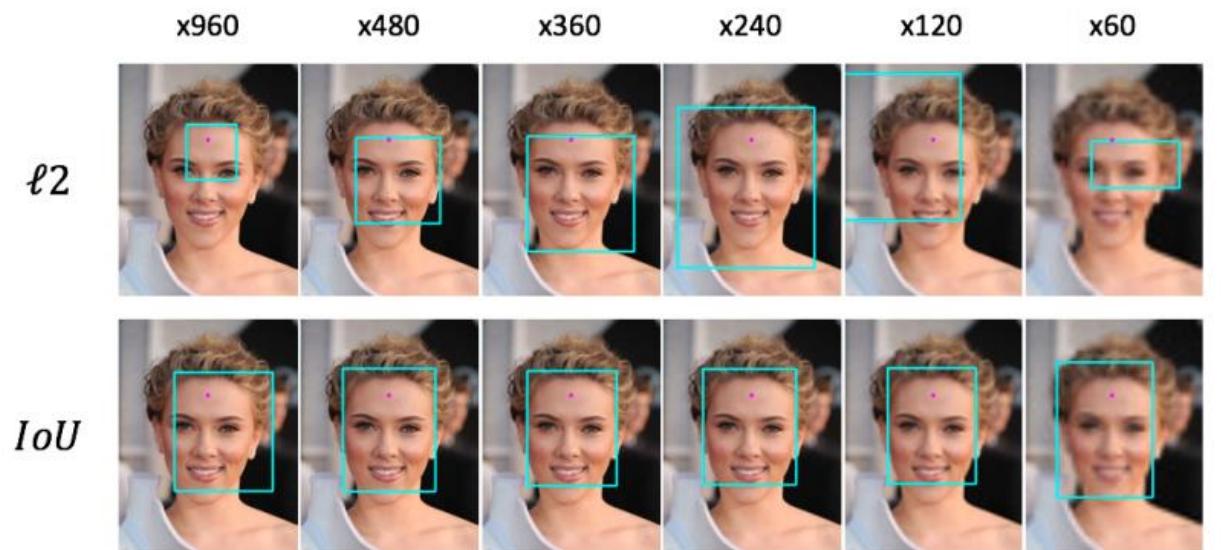


Figure 5: Compared to  $\ell_2$  loss, the  $IoU$  loss is much more robust to scale variations for bounding box prediction.

# One Stage Detector: Densebox -> UnitBox->EAST

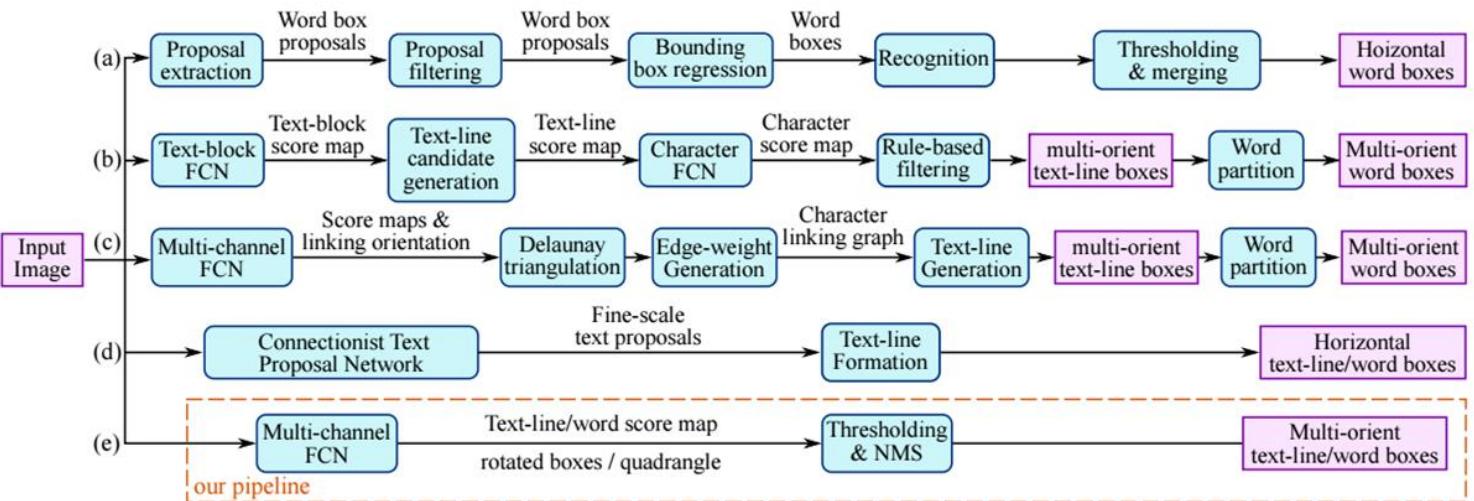
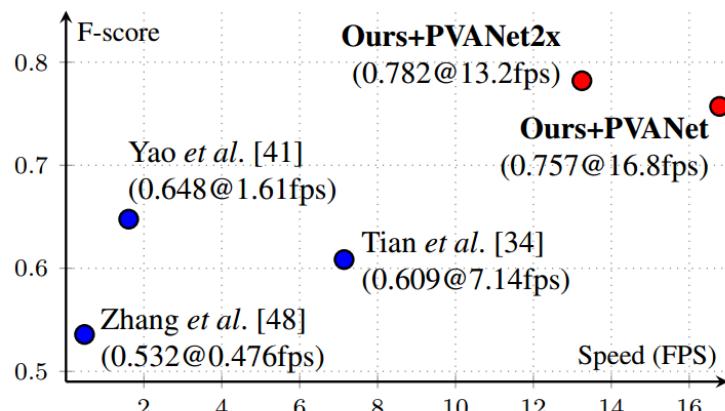
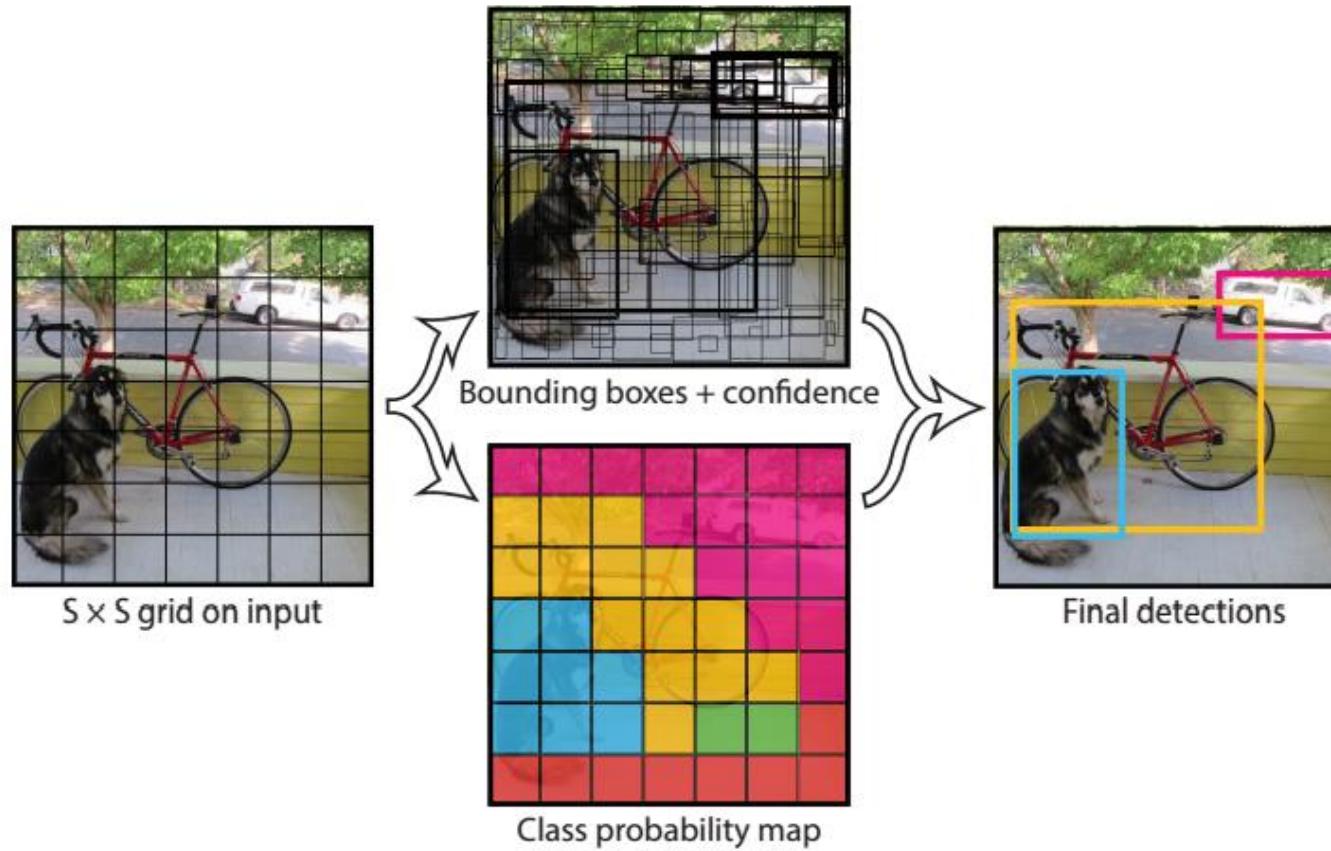


Figure 2. Comparison of pipelines of several recent works on scene text detection: (a) Horizontal word detection and recognition pipeline proposed by Jaderberg *et al.* [12]; (b) Multi-orient text detection pipeline proposed by Zhang *et al.* [48]; (c) Multi-orient text detection pipeline proposed by Yao *et al.* [41]; (d) Horizontal text detection using CTPN, proposed by Tian *et al.* [34]; (e) Our pipeline, which eliminates most intermediate steps, consists of only two stages and is much simpler than previous solutions.

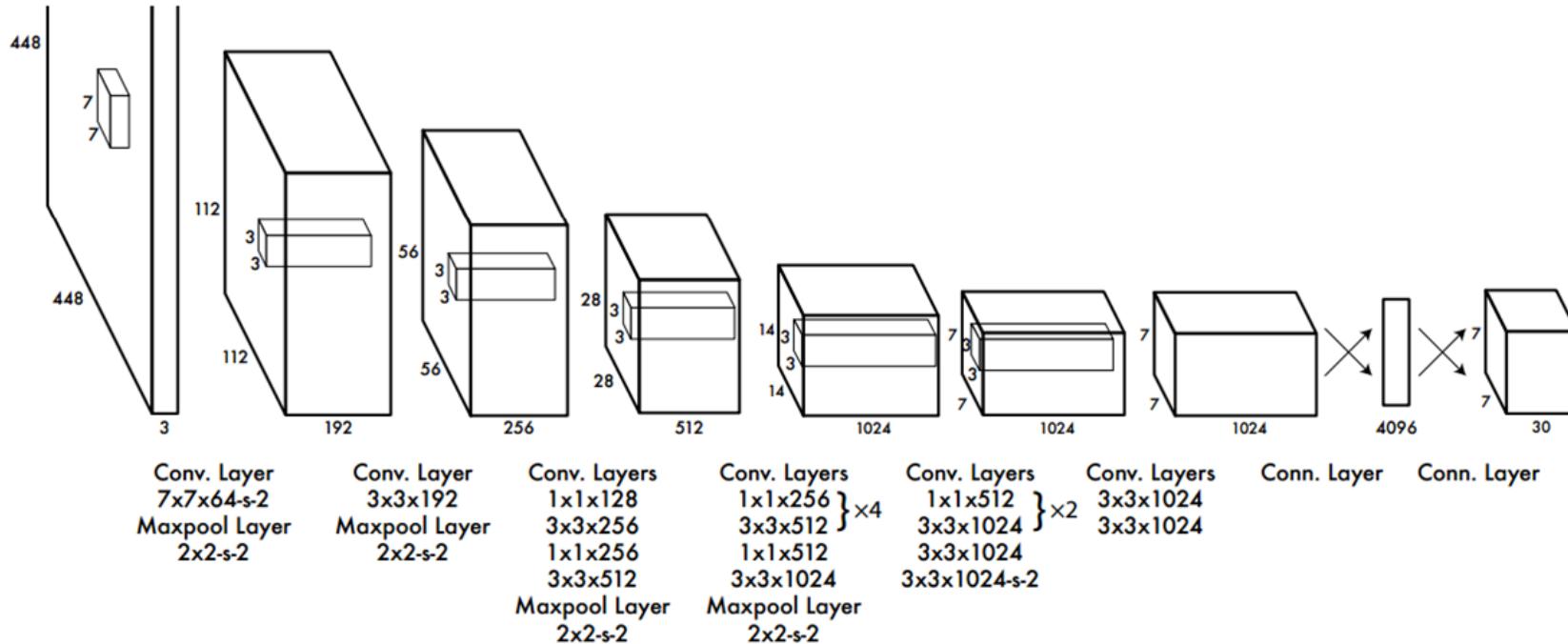
EAST: An Efficient and Accurate Scene Text Detector, Zhou etc, CVPR 2017

# One Stage Detector: YOLO



You Only Look Once: Unified, Real-Time Object Detection, Redmon etc, CVPR 2016

# One Stage Detector: YOLO



**Figure 3: The Architecture.** Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating  $1 \times 1$  convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution ( $224 \times 224$  input image) and then double the resolution for detection.

# One Stage Detector: YOLO

- No Anchor
  - GT assignment is based on the cells (7x7)
- Loss sampling
  - all pos/neg predictions are evaluated (but more **sparse** than densebox)

# One Stage Detector: YOLO

## Discussion

- fc reshape (4096-> 7x7x30)
  - more context
  - but not fully convolutional
- One cell can output up to two boxes in one category
  - fail to work on the crowd case
- Fast speed
  - small imagenet base model
  - small input size (448x448)

# One Stage Detector: YOLO

Experiments on general detection

Method	VOC 2007 test	VOC 2012 test	COCO	time
YOLO	57.9/NA	52.7/63.4	NA	fps: 45/155

# One Stage Detector: YOLO -> YOLOv2

	YOLO							YOLOv2
batch norm?		✓	✓	✓	✓	✓	✓	✓
hi-res classifier?			✓	✓	✓	✓	✓	✓
convolutional?				✓	✓	✓	✓	✓
anchor boxes?					✓	✓		
new network?						✓	✓	✓
dimension priors?							✓	✓
location prediction?							✓	✓
passthrough?							✓	✓
multi-scale?								✓
hi-res detector?								✓
VOC2007 mAP	63.4	65.8	69.5	69.2	69.6	74.4	75.4	76.8
								<b>78.6</b>

**Table 2: The path from YOLO to YOLOv2.** Most of the listed design decisions lead to significant increases in mAP. Two exceptions are switching to a fully convolutional network with anchor boxes and using the new network. Switching to the anchor box style approach increased recall without changing mAP while using the new network cut computation by 33%.

# One Stage Detector: YOLO -> YOLOv2

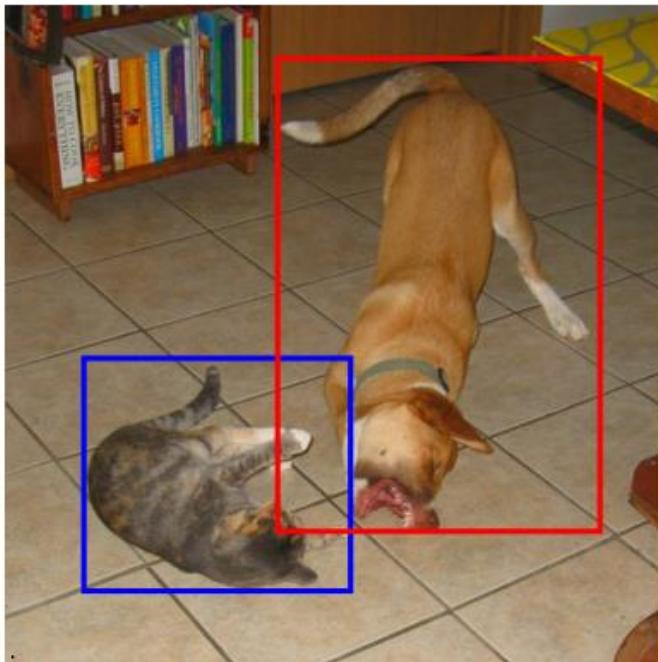
Experiments:

Method	VOC 2007 test	VOC 2012 test	COCO	time
YOLO	52.7/63.4	57.9/NA	NA	fps: 45/155
YOLOv2	78.6	73.4	21.6	fps: 40

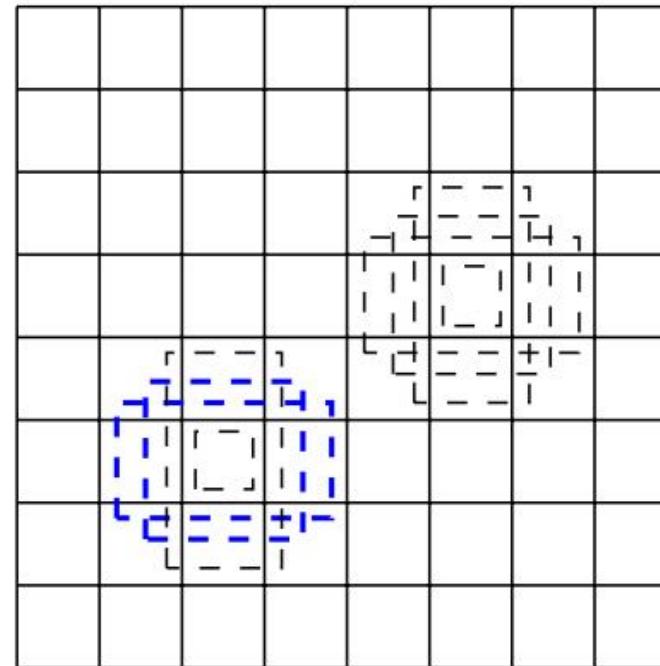
# One Stage Detector: YOLO -> YOLOv2

Video demo: <https://pjreddie.com/darknet/yolo/>

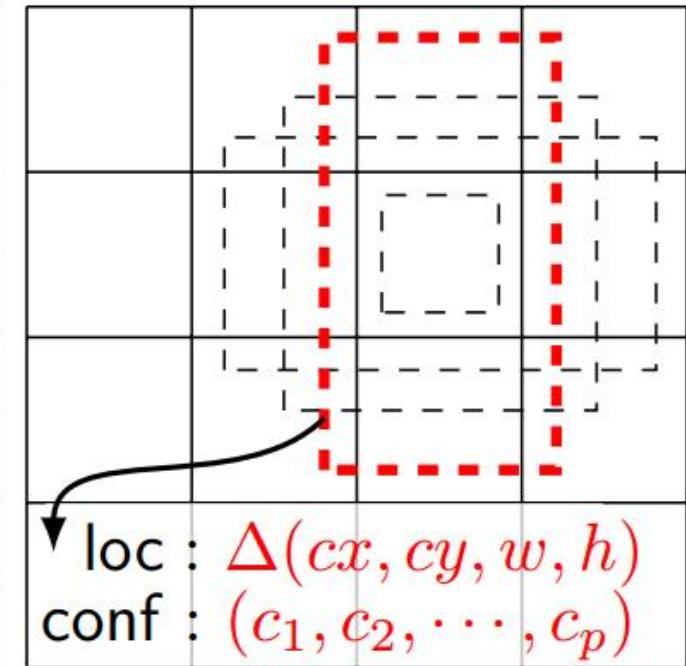
# One Stage Detector: SSD



(a) Image with GT boxes



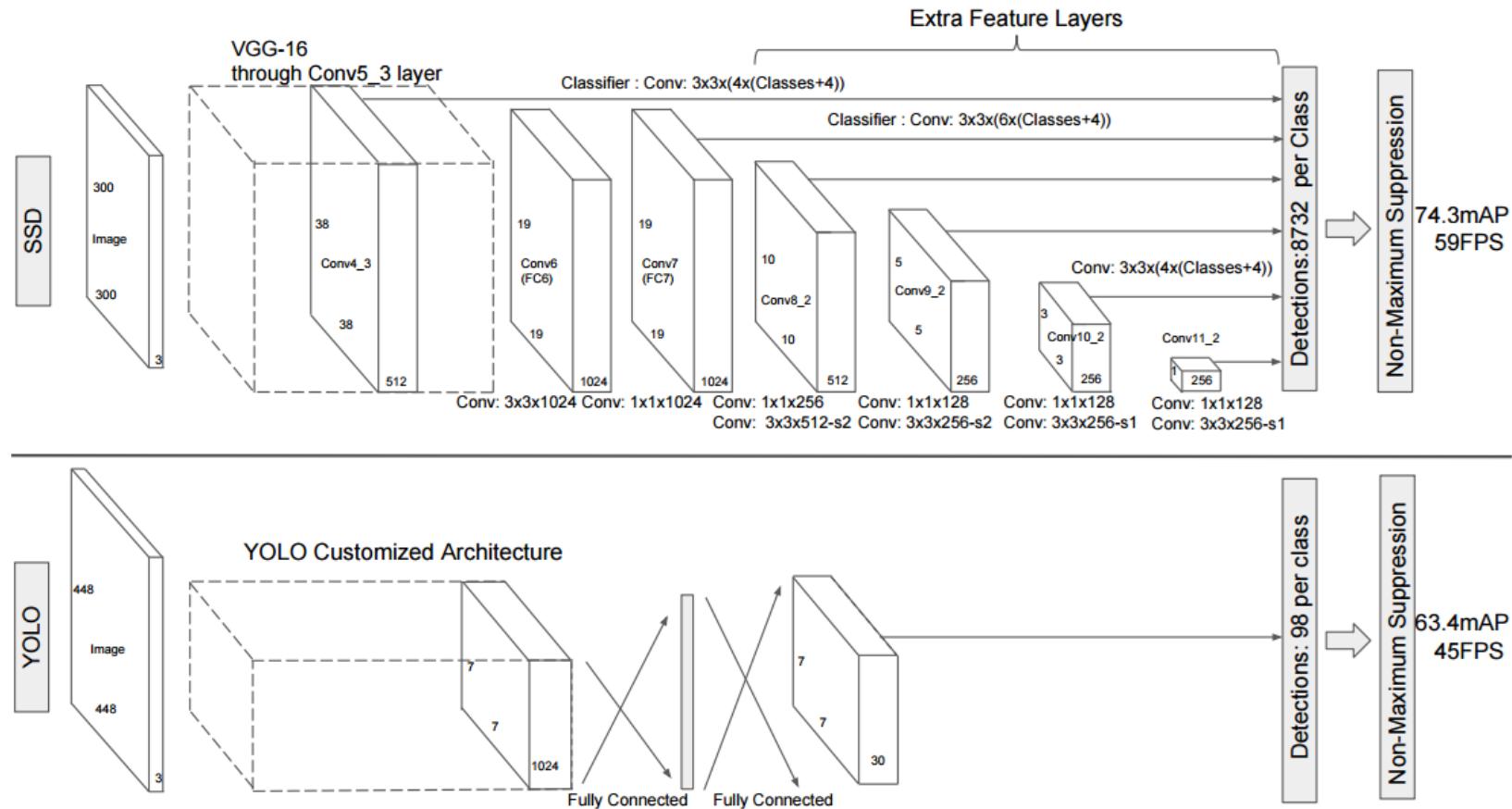
(b)  $8 \times 8$  feature map



loc :  $\Delta(cx, cy, w, h)$   
conf :  $(c_1, c_2, \dots, c_p)$

(c)  $4 \times 4$  feature map

# One Stage Detector: SSD



# One Stage Detector: SSD

- Anchor
  - GT-anchor assignment
    - GT is predicted by one best matched (IOU) anchor or matched with an anchor with  $\text{IOU} > 0.5$ 
      - better recall
    - dense or sparse anchor?
- Divide and Conquer
  - Different layers handle the objects with different scales
    - Assume small objects can be predicted in earlier layers (not very strong semantics)
- Loss sampling
  - OHEM: negative positions are sampled (not balanced pos/neg ratio)
  - negative:pos is at most 3:1

# One Stage Detector: SSD

## Discussion:

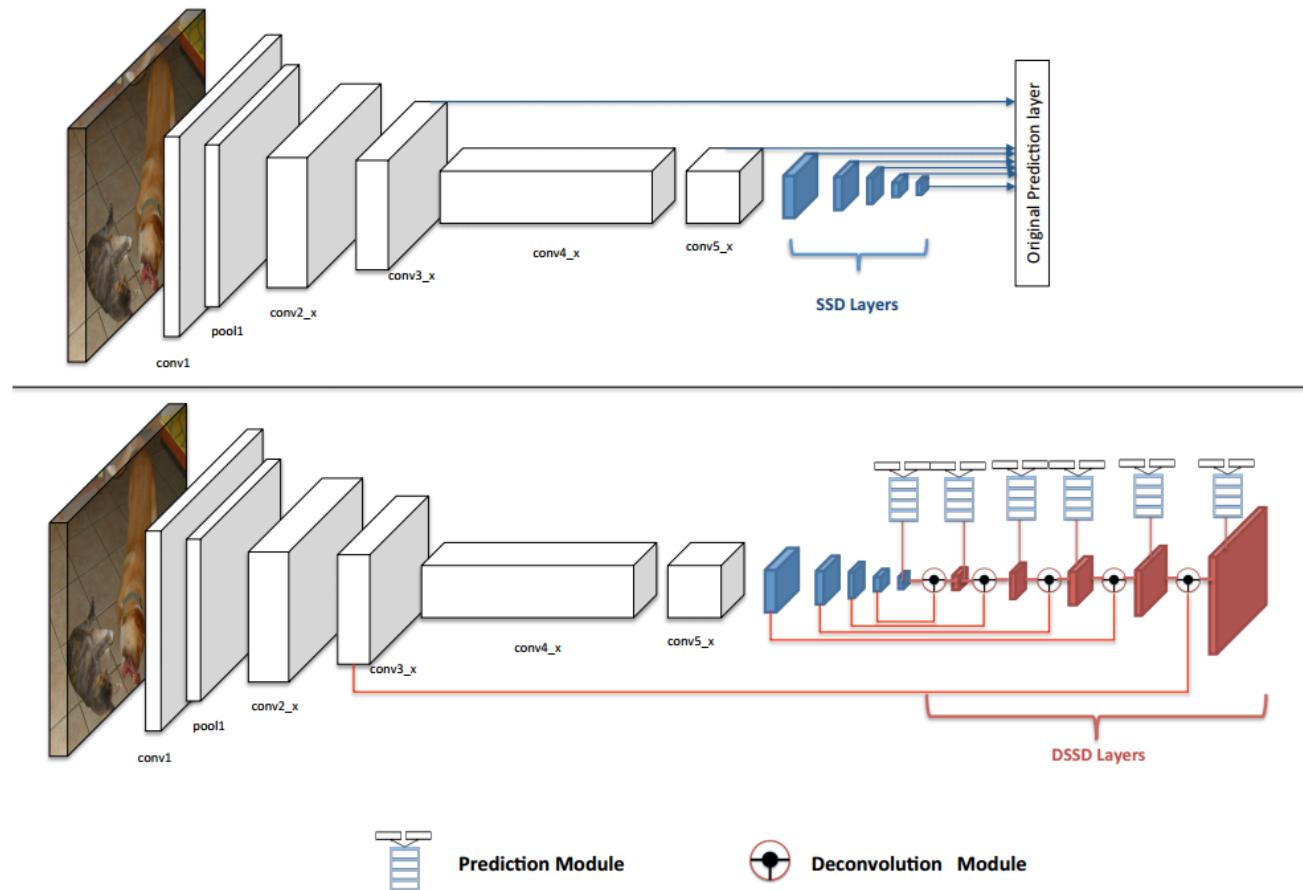
- Assume small objects can be predicted in earlier layers (not very strong semantics) (DSSD, RON, RetinaNet)
- strong data augmentation
- VGG model (Replace by resnet in DSSD)
  - cannot be easily adapted to other models
  - a lot of hacks
- A long tail (Large computation)

# One Stage Detector: SSD

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19

# One Stage Detector: SSD -> DSSD



# One Stage Detector: DSSD

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5

# One Stage Detector: SSD -> RON

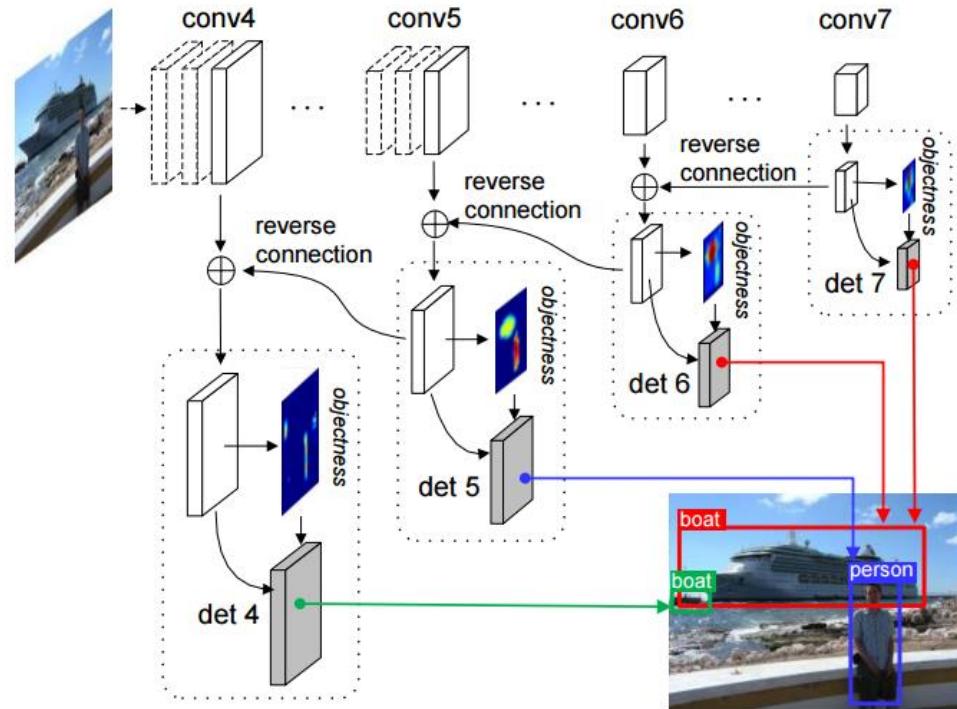
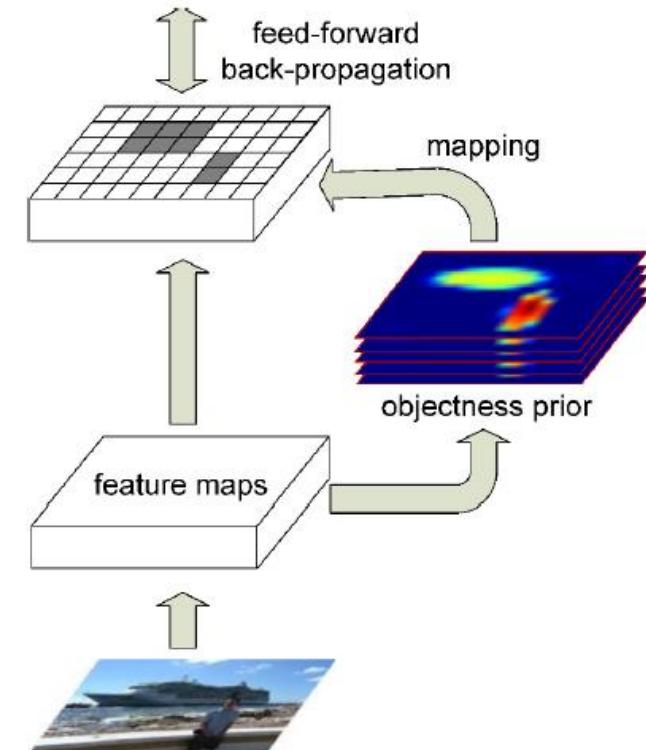


Figure 2. RON object detection overview. Given an input image, the network firstly computes features of the backbone network. Then at each detection scale: (a) adds reverse connection; (b) generates objectness prior; (c) detects object on its corresponding CNN scales and locations. Finally, all detection results are fused and selected with non-maximum suppression.

# One Stage Detector: RON

- Anchor
- Divide and conquer
  - Reverse Connect (similar to FPN)
- Loss Sampling
  - Objectness prior
    - pos/neg unbalanced issue
    - split to 1) binary cls 2) multi-class cls



# One Stage Detector: RON

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15

# One Stage Detector: SSD -> RetinaNet

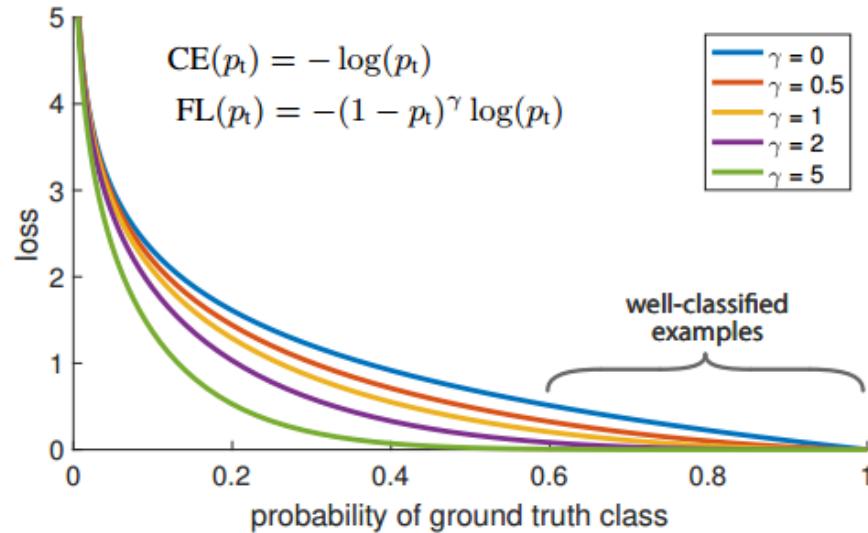


Figure 1. We propose a novel loss we term the *Focal Loss* that adds a factor  $(1 - p_t)^\gamma$  to the standard cross entropy criterion. Setting  $\gamma > 0$  reduces the relative loss for well-classified examples ( $p_t > .5$ ), putting more focus on hard, misclassified examples. As our experiments will demonstrate, the proposed focal loss enables training highly accurate dense object detectors in the presence of vast numbers of easy background examples.

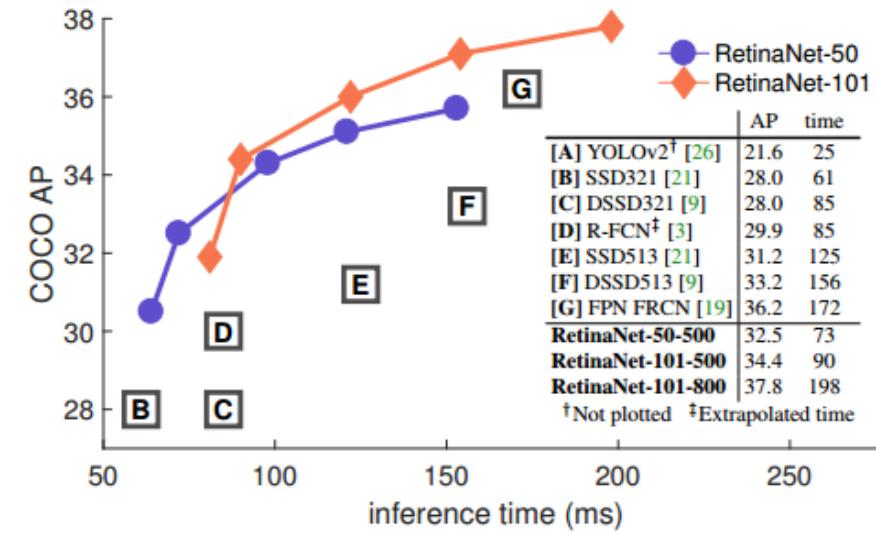


Figure 2. Speed (ms) versus accuracy (AP) on COCO test-dev. Enabled by the focal loss, our simple one-stage *RetinaNet* detector outperforms all previous one-stage and two-stage detectors, including the best reported Faster R-CNN [27] system from [19]. We show variants of RetinaNet with ResNet-50-FPN (blue circles) and ResNet-101-FPN (orange diamonds) at five scales (400-800 pixels). Ignoring the low-accuracy regime ( $AP < 25$ ), RetinaNet forms an upper envelope of all current detectors, and a variant trained for longer (not shown) achieves 39.1 AP. Details are given in §5.

# One Stage Detector: SSD -> RetinaNet

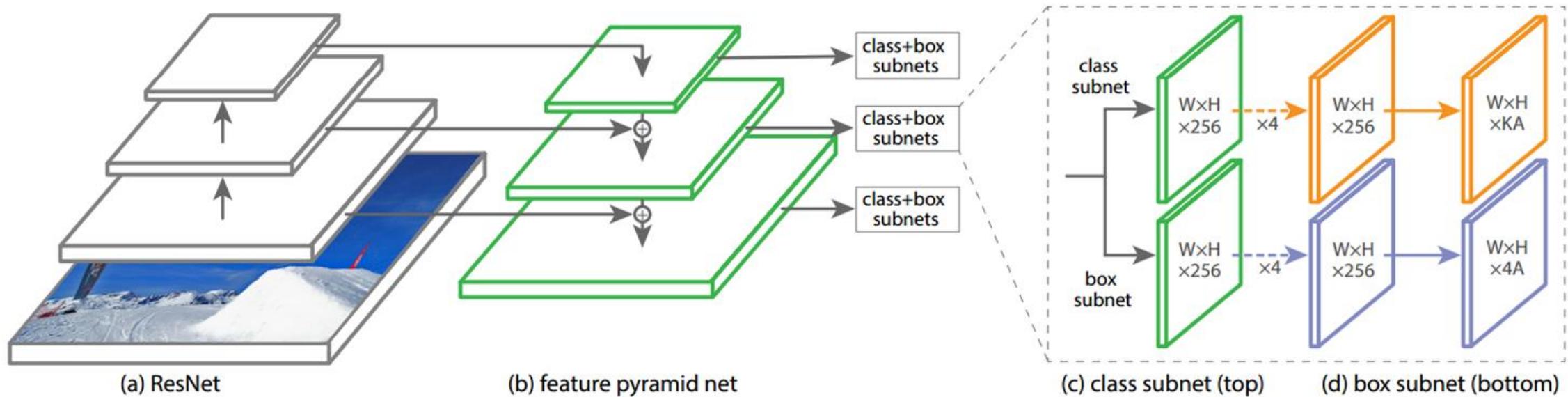


Figure 3. The one-stage **RetinaNet** network architecture uses a Feature Pyramid Network (FPN) [19] backbone on top of a feedforward ResNet architecture [15] (a) to generate a rich, multi-scale convolutional feature pyramid (b). To this backbone RetinaNet attaches two subnetworks, one for classifying anchor boxes (c) and one for regressing from anchor boxes to ground-truth object boxes (d). The network design is intentionally simple, which enables this work to focus on a novel focal loss function that eliminates the accuracy gap between our one-stage detector and state-of-the-art two-stage detectors like Faster R-CNN with FPN [19] while running at faster speeds.

# One Stage Detector: RetinaNet

- Anchor
- Divide and Conquer
  - FPN
- Loss Sampling
  - Focal loss
    - pos/neg unbalanced issue
    - new setting (e.g., more anchor)

# One Stage Detector: RetinaNet

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5

# One Stage Detector: Summary

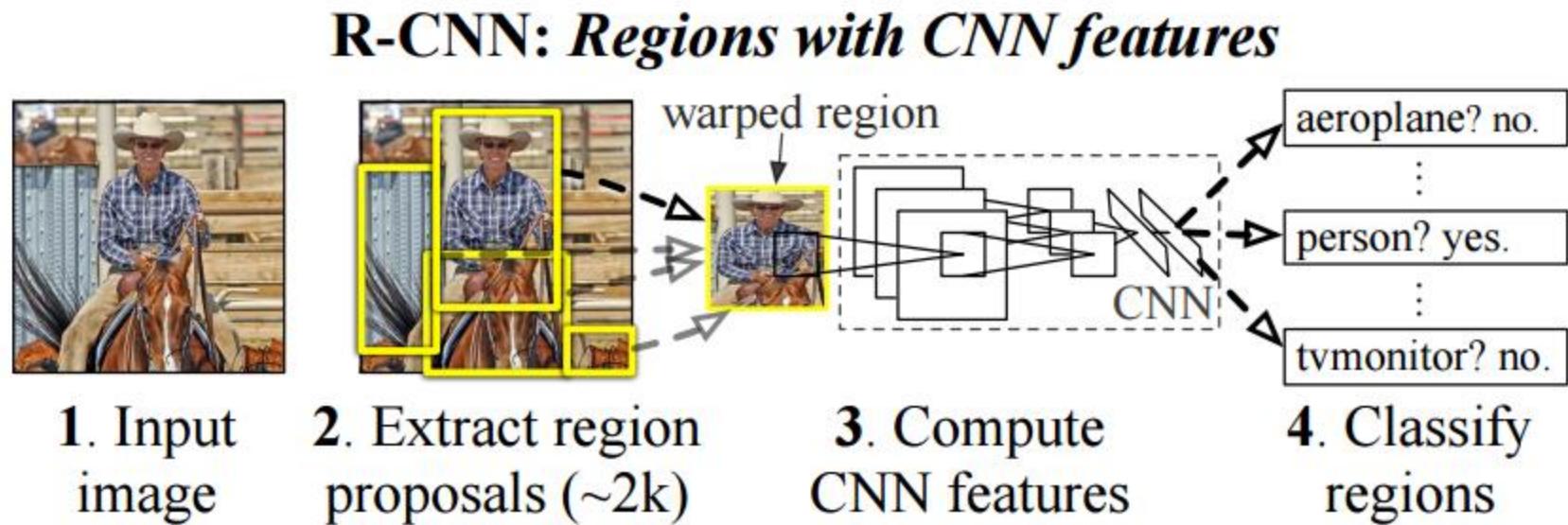
- Anchor
  - No anchor: YOLO, densebox/unitbox/east
  - Anchor: YOLOv2, SSD, DSSD, RON, RetinaNet
- Divide and conquer
  - SSD, DSSD, RON, RetinaNet
- loss sample
  - all sample: densebox
  - OHEM: SSD
  - focal loss: RetinaNet

# One Stage Detector: Discussion

Anchor (YOLO v2, SSD, RetinaNet) or Without Anchor (Densebox, YOLO)

- Model Complexity
  - Difference on the extremely small model (< 30M flops on 224x224 input)
- Sampling
- Application
  - No Anchor: Face
  - With Anchor: Human, General Detection
- Problem for one stage detector
  - Unbalanced pos/neg data
  - Pool localization precision

# Two Stages Detector: RCNN



# Two Stages Detector: RCNN

## Discussion

- Extremely slow speed
  - selective search proposal (CPU)/warp
- not end-to-end optimized
- Good for small objects

# Two Stages Detector: RCNN

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s

# Two Stages Detector: RCNN -> Fast RCNN

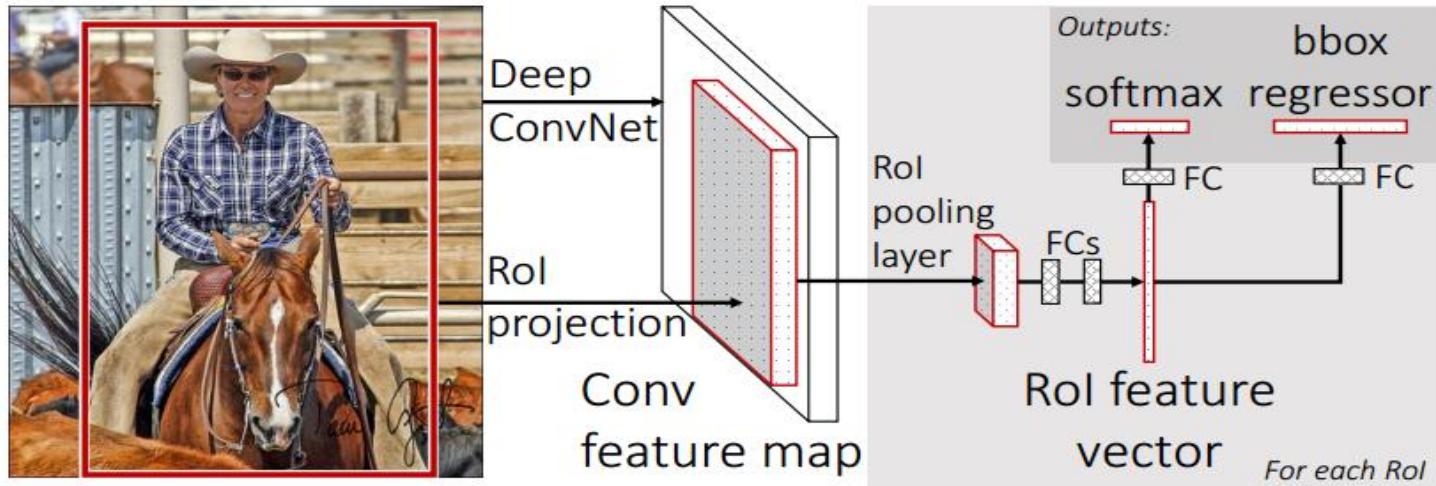


Figure 1. Fast R-CNN architecture. An input image and multiple regions of interest (RoIs) are input into a fully convolutional network. Each ROI is pooled into a fixed-size feature map and then mapped to a feature vector by fully connected layers (FCs). The network has two output vectors per ROI: softmax probabilities and per-class bounding-box regression offsets. The architecture is trained end-to-end with a multi-task loss.

# Two Stages Detector: Fast RCNN

## Discussion

- slow speed
  - selective search proposal (CPU)
- not end-to-end optimized
- ROI pooling
  - alignment issue
  - sampling
  - aspect ratio changes

# Two Stages Detector: Fast RCNN

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s
Fast RCNN	77.0	82.3 (wth coco data)	NA	0.5s

# Two Stages Detector: RCNN -> Fast RCNN -> FasterRCNN

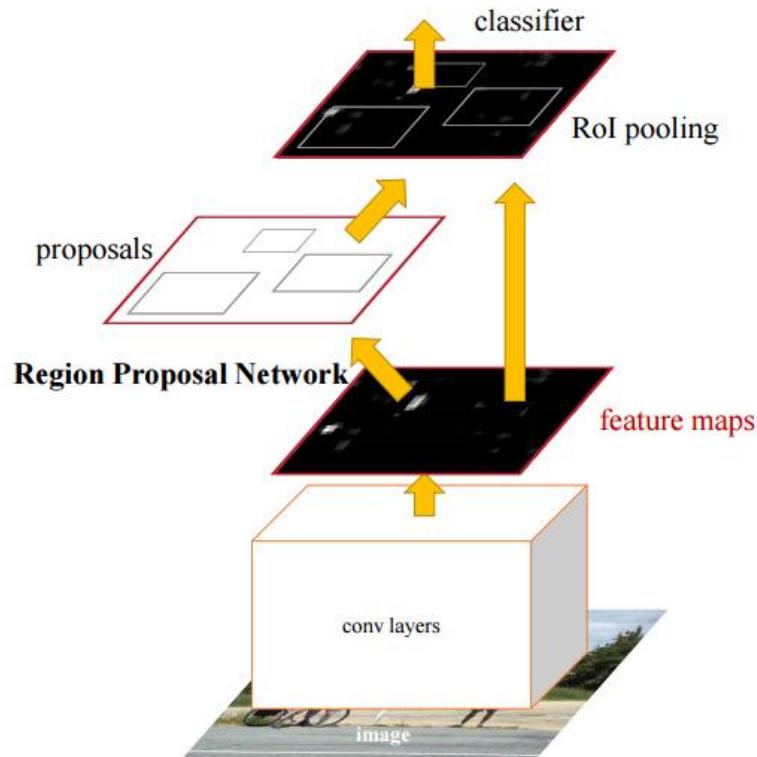


Figure 2: Faster R-CNN is a single, unified network for object detection. The RPN module serves as the ‘attention’ of this unified network.

# Two Stages Detector: Faster RCNN

## Discussion

- speed
  - selective search proposal (CPU) -> RPN
- alternative optimization/end-to-end optimization
- Recall issue due to two stages detector

# Two Stages Detector: Faster RCNN

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s
Fast RCNN	77.0	82.3 (wth coco data)	NA	0.5s
Faster RCNN	73.2	70.4	NA	5

# Two Stages Detector: RCNN -> Fast RCNN -> FasterRCNN -> RFCN

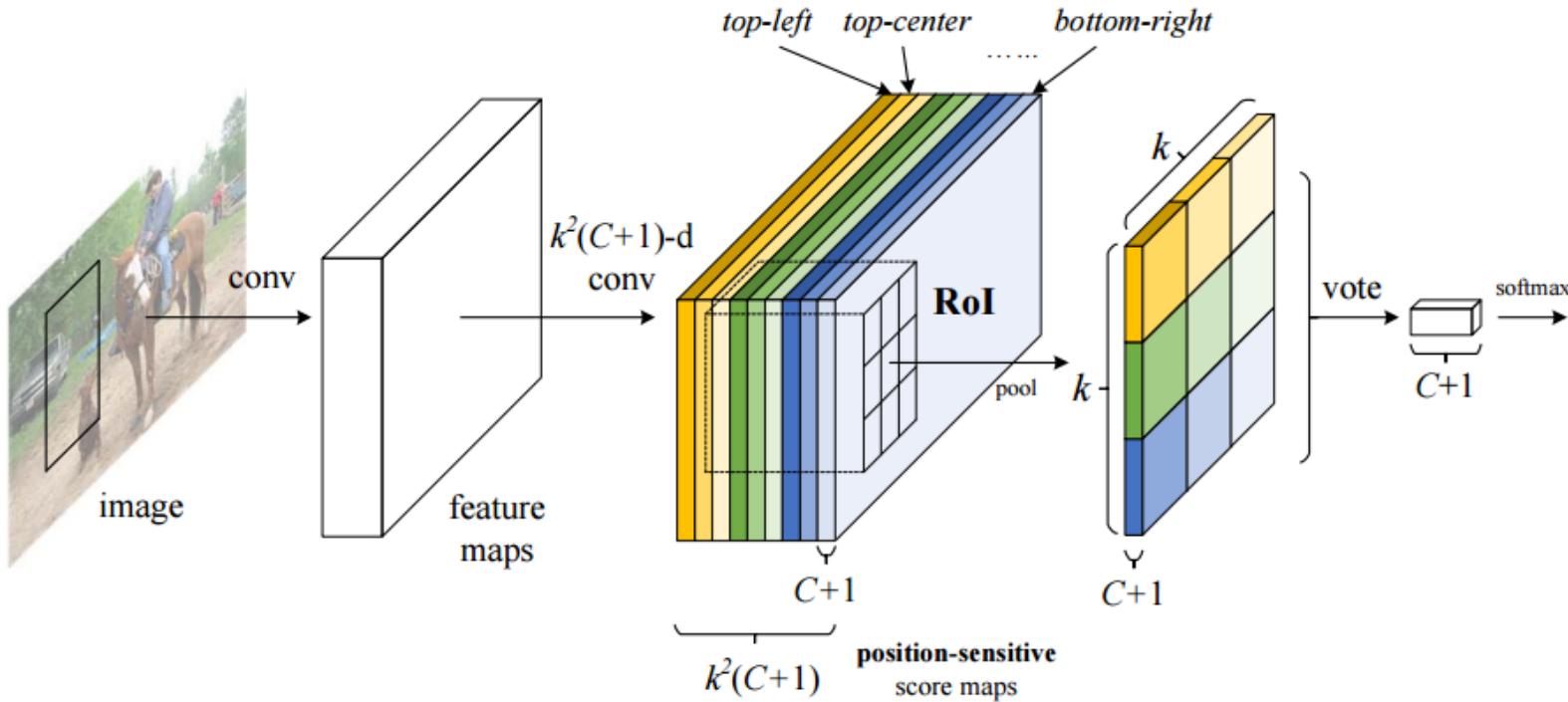


Figure 1: Key idea of **R-FCN** for object detection. In this illustration, there are  $k \times k = 3 \times 3$  position-sensitive score maps generated by a fully convolutional network. For each of the  $k \times k$  bins in an ROI, pooling is only performed on one of the  $k^2$  maps (marked by different colors).

# Two Stages Detector: RFCN

## Discussion

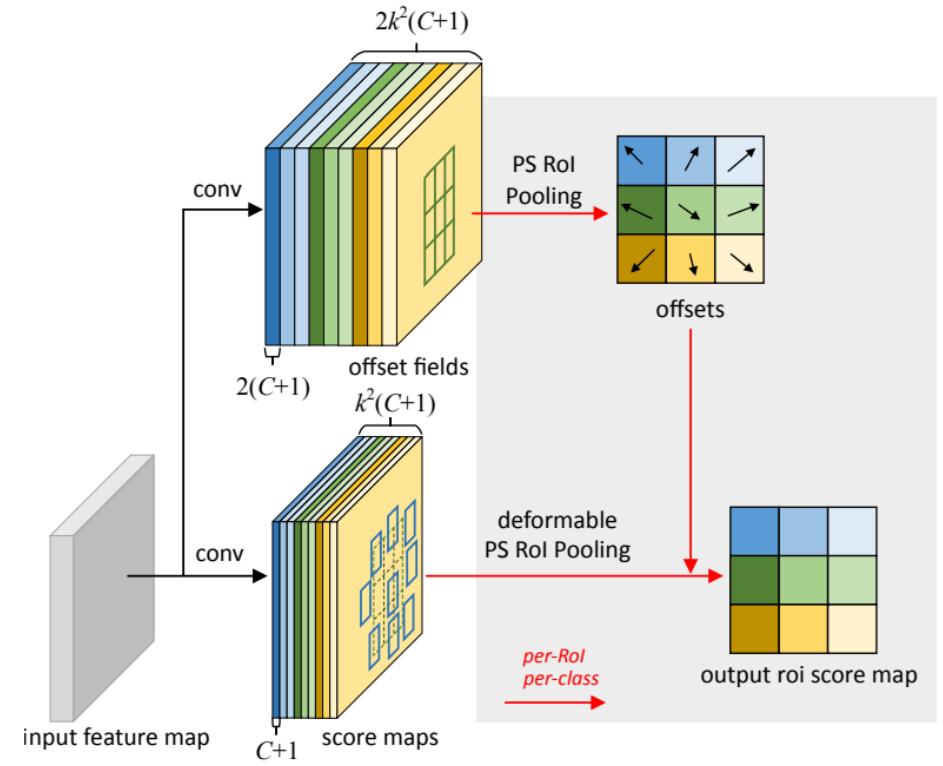
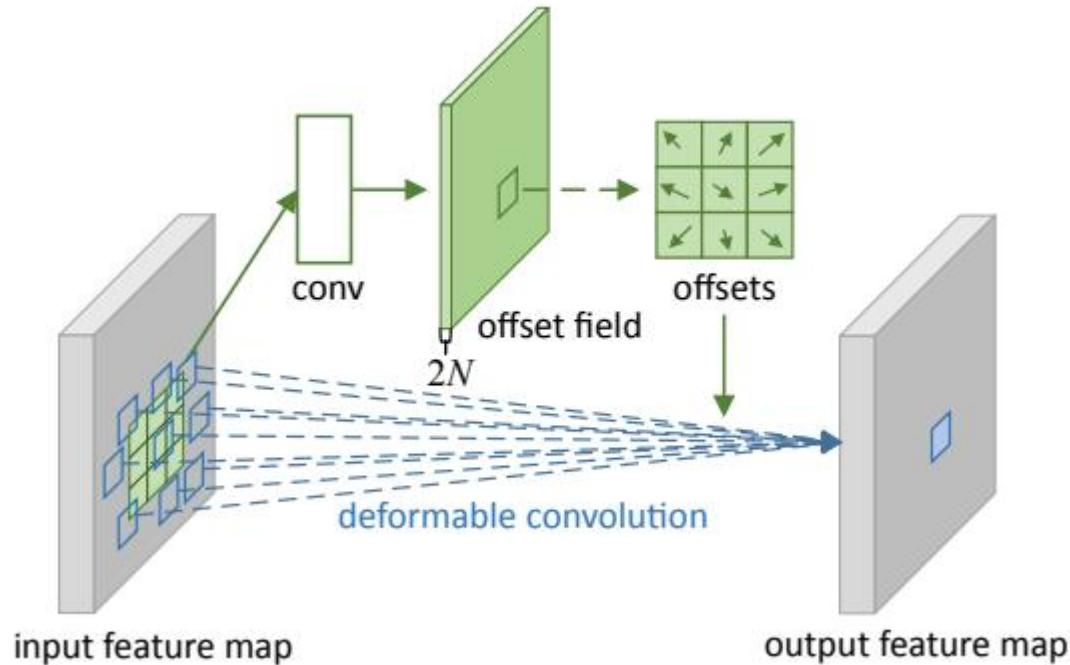
- Share convolution
  - fasterRCNN: shared Res1-4 (RPN), not shared Res5 (RCNN)
  - RFCN: shared Res1-5 (both RPN and RCNN)
- PSPooling
  - a large number of channels:  $(7 \times 7 \times C) \times W \times H$
  - Problems in ROI Pooling also exist
- Fully connected vs Convolution
  - fc: global context
  - conv: can be shared but the context is relative small
  - trade-off: large kernel

# Two Stages Detector: RFCN

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s
Fast RCNN	77.0	82.3 (wth coco data)	NA	0.5s
Faster RCNN	73.2	70.4	NA	200ms
RFCN	79.5	77.6	29.9	170ms

# Two Stages Detector: RFCN -> Deformable Convolutional Networks



# Two Stages Detector: RFCN -> Deformable Convolutional Networks



Figure 6: Each image triplet shows the sampling locations ( $9^3 = 729$  red points in each image) in three levels of  $3 \times 3$  deformable filters (see Figure 5 as a reference) for three activation units (green points) on the background (left), a small object (middle), and a large object (right), respectively.

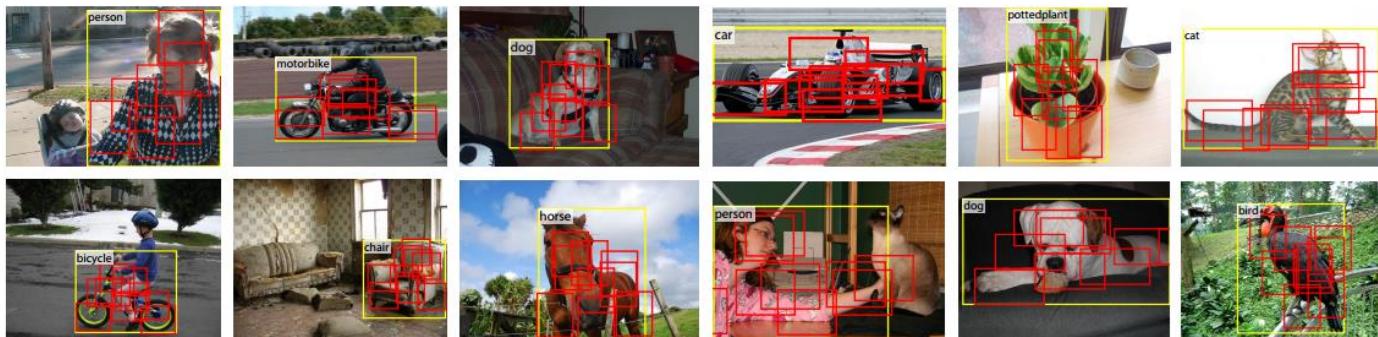


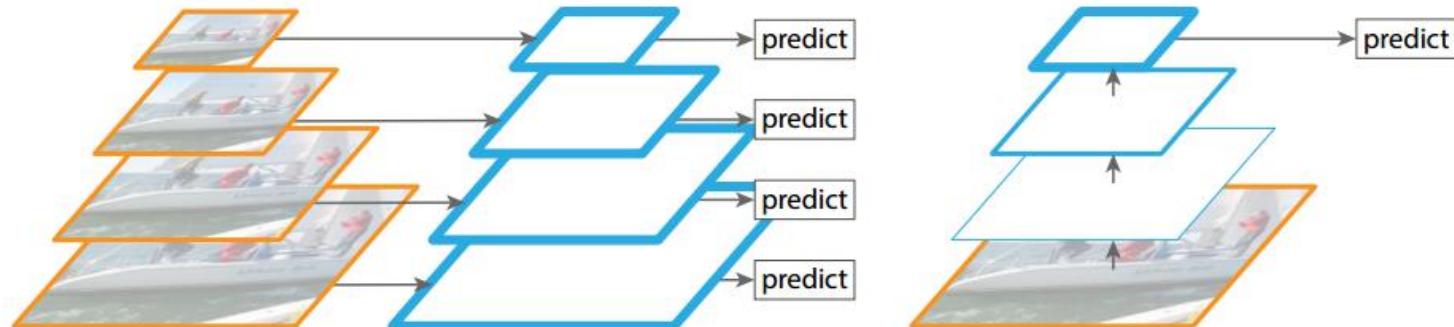
Figure 7: Illustration of offset parts in deformable (positive sensitive) ROI pooling in R-FCN [7] and  $3 \times 3$  bins (red) for an input ROI (yellow). Note how the parts are offset to cover the non-rigid objects.

# Two Stages Detector: RFCN -> Deformable Convolutional Networks

## Discussion

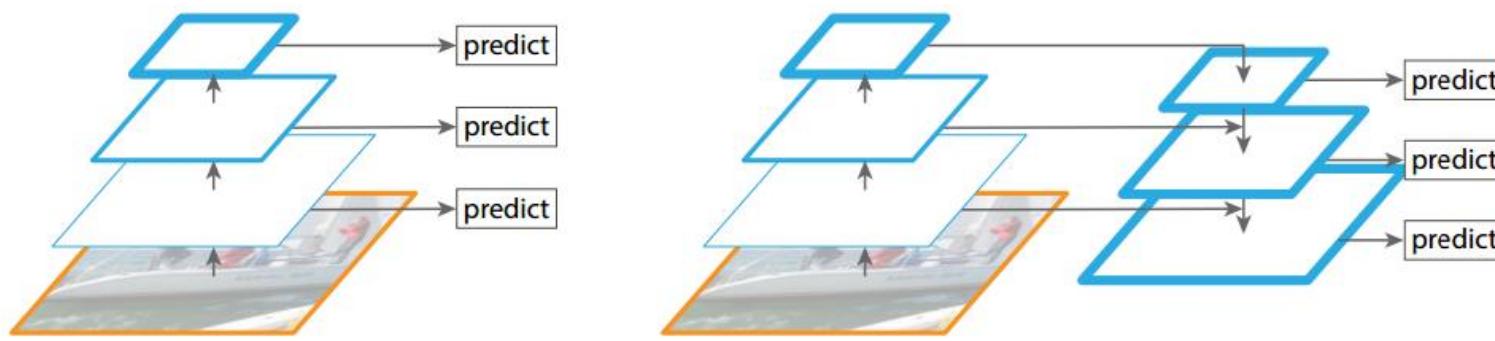
- Deformable pool is similar to ROIAlign (in Mask RCNN)
- Deformable conv
  - flexible to learn the non-rigid objects

# Two Stages Detector: RCNN -> Fast RCNN -> FasterRCNN -> FPN



(a) Featurized image pyramid

(b) Single feature map



(c) Pyramidal feature hierarchy

(d) Feature Pyramid Network

# Two Stages Detector: FPN

## Discussion

- FasterRCNN reproduced (setting)
- Deeply supervised (better feature)

# Two Stages Detector: FPN

## Experiments

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s
Fast RCNN	77.0	82.3 (wth coco data)	NA	0.5s
Faster RCNN	73.2	70.4	NA	200ms
RFCN	79.5	77.6	29.9	170ms
FPN	NA	NA	36.2	6

# Two Stages Detector: RCNN -> Fast RCNN -> FasterRCNN -> FPN -> MaskRCNN

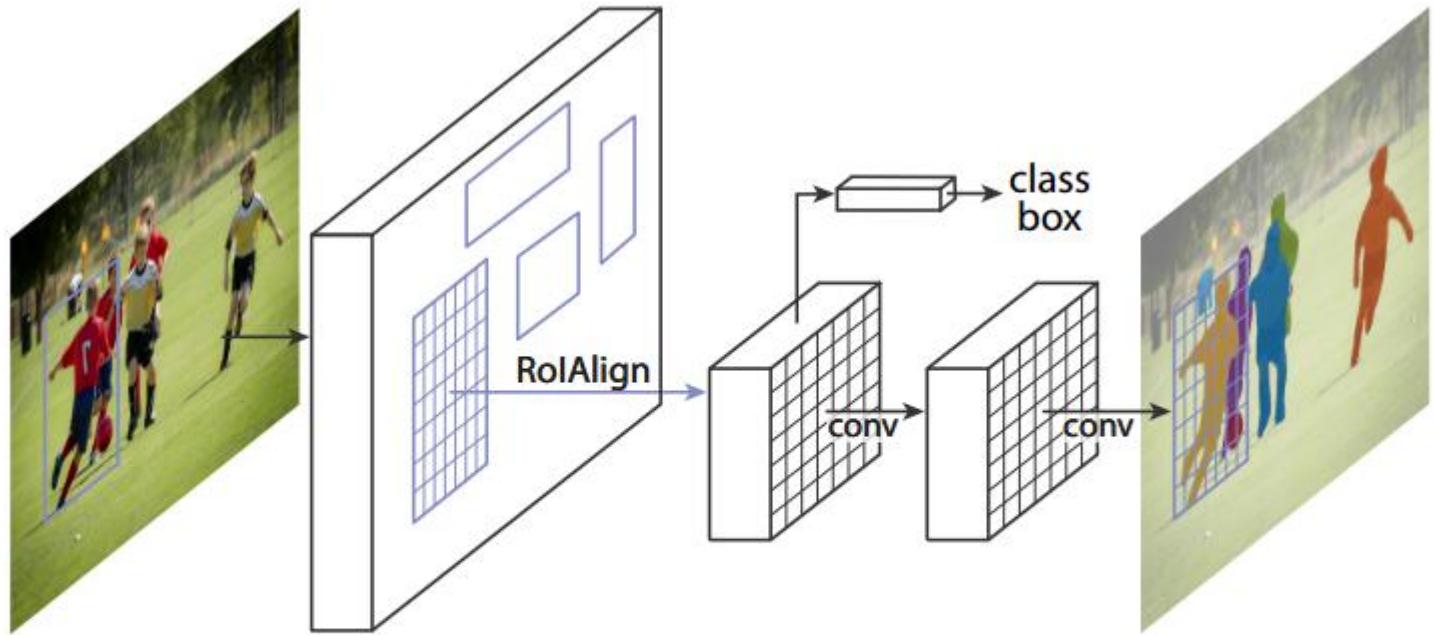


Figure 1. The **Mask R-CNN** framework for instance segmentation.

# Two Stages Detector:

RCNN -> Fast RCNN -> FasterRCNN -> FPN ->  
MaskRCNN

	align?	bilinear?	agg.	AP	AP <sub>50</sub>	AP <sub>75</sub>
<i>RoIPool</i> [12]			max	26.9	48.8	26.4
<i>RoIWarp</i> [10]		✓	max	27.2	49.2	27.1
		✓	ave	27.1	48.9	27.1
<i>RoIAlign</i>	✓	✓	max	<b>30.2</b>	<b>51.0</b>	<b>31.8</b>
	✓	✓	ave	<b>30.3</b>	<b>51.2</b>	<b>31.5</b>

(c) **RoIAlign** (ResNet-50-C4): Mask results with various RoI layers. Our RoIAlign layer improves AP by ~3 points and AP<sub>75</sub> by ~5 points. Using proper alignment is the only factor that contributes to the large gap between RoI layers.

# Two Stages Detector: Mask RCNN

## Discussion

- Alignment issue in ROI Pooling -> ROI Align
- Multi-task learning: detection & mask

# Experiments Two Stages Detector: Mask RCNN

Method	VOC 2007 test	VOC 2012 test	COCO	time (fps)
YOLO	52.7/63.4	57.9/NA	NA	45/155
YOLOv2	78.6	73.4	21.6	40
SSD	77.2/79.8	75.8/78.5	25.1/28.8	46/19
DSSD	81.5	80.0	33.2	5.5
RON	81.3	80.7	27.4	15
RetinaNet	NA	N	39.1	5
RCNN	66	NA	NA	47s
Fast RCNN	77.0	82.3 (wth coco data)	NA	0.5s
Faster RCNN	73.2	70.4	NA	200ms
RFCN	79.5	77.6	29.9	170ms
FPN	NA	NA	36.2	6
Mask RCNN	NA	NA	38.2	2.5

# Two Stages Detector: Summary

- Speed
  - RCNN -> Fast RCNN -> Faster RCNN -> RFCN
- performance
  - Divide and conquer
    - FPN
  - Deformable Pool/ROIAlign
  - Deformable Conv
  - Multi-task learning

# Two Stages Detector: Discussion

FasterRCNN vs RFCN

One stage vs two Stage

# MegDetection

Introduction & Demo Video

# Open Problem in Detection

- FP
- NMS (detection in crowd)
- GT assignment issue
- Detection in video
  - detect & track in a network

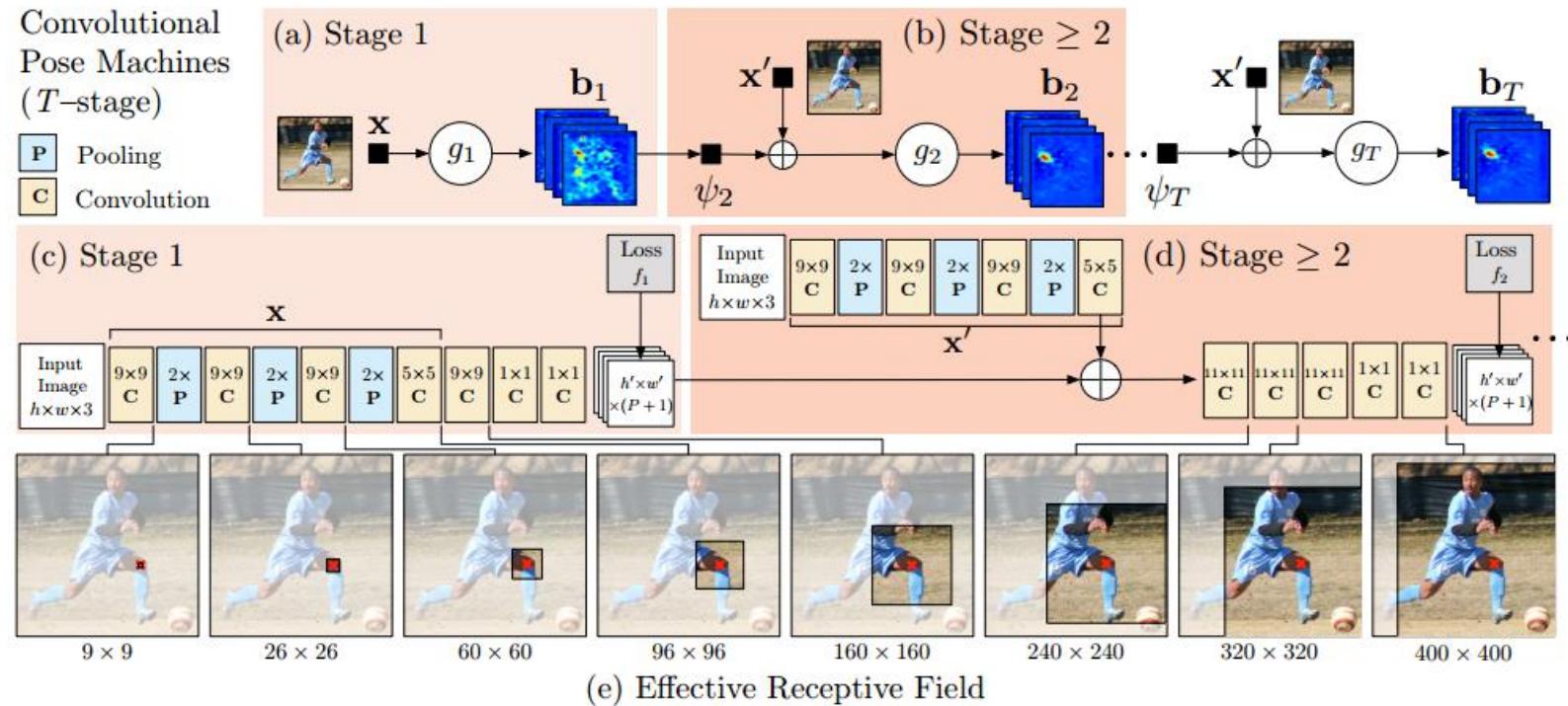
# Outline

- Detection
- **Human Keypoint**
- Conclusion

# Human Keypoint Task

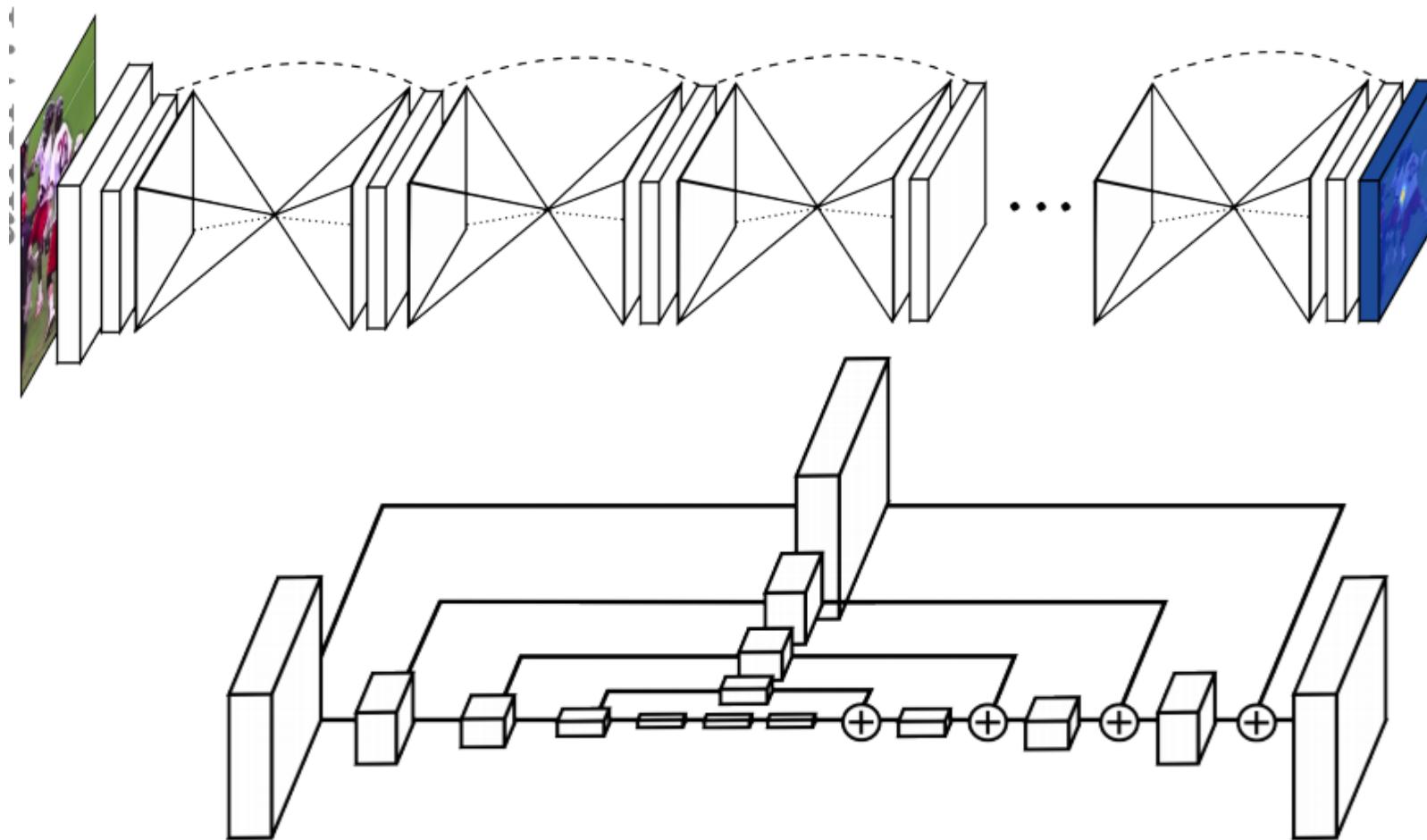
- Single Person Skeleton
  - Cropped RGB image -> 2d key points / 3d key points
  - Keyword: [inter-middle loss](#), [large receptive field](#), [context](#)
- Multiple-Person Skeleton
  - RGB image -> human localization & human Keypoint for each person

# Single Person Skeleton: CPM



**Figure 2: Architecture and receptive fields of CPMs.** We show a convolutional architecture and receptive fields across layers for a CPM with any  $T$  stages. The pose machine [29] is shown in insets (a) and (b), and the corresponding convolutional networks are shown in insets (c) and (d). Insets (a) and (c) show the architecture that operates only on image evidence in the first stage. Insets (b) and (d) shows the architecture for subsequent stages, which operate both on image evidence as well as belief maps from preceding stages. The architectures in (b) and (d) are repeated for all subsequent stages (2 to  $T$ ). The network is locally supervised after each stage using an intermediate loss layer that prevents vanishing gradients during training. Below in inset (e) we show the effective receptive field on an image (centered at left knee) of the architecture, where the large receptive field enables the model to capture long-range spatial dependencies such as those between head and knees. (Best viewed in color.)

# Single Person Skeleton: Hourglass



# Multiple-Person Skeleton

- Top Down
  - Detect -> Single person skeleton
- Bottom Up
  - Deep/Deeper Cut
  - OpenPose
  - Associative Embedding

# Multiple-Person Skeleton: OpenPose

CPM + PAF

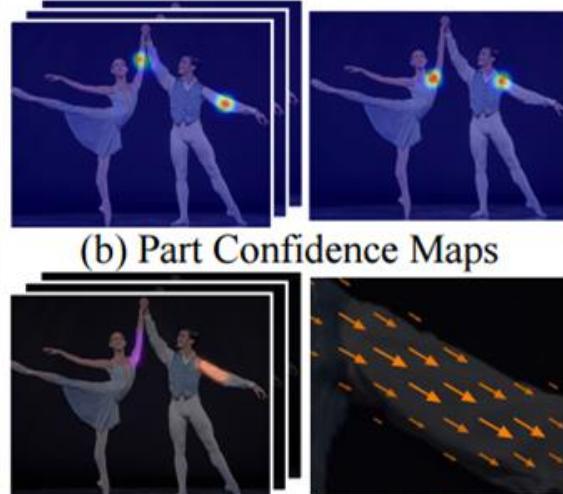


Figure 1. **Top:** Multi-person pose estimation. Body parts belonging to the same person are linked. **Bottom left:** Part Affinity Fields (PAFs) corresponding to the limb connecting right elbow and right wrist. The color encodes orientation. **Bottom right:** A zoomed in view of the predicted PAFs. At each pixel in the field, a 2D vector encodes the position and orientation of the limbs.

# Multiple-Person Skeleton: OpenPose



(a) Input Image



(b) Part Confidence Maps



(c) Part Affinity Fields



(d) Bipartite Matching



(e) Parsing Results

Figure 2. Overall pipeline. Our method takes the entire image as the input for a two-branch CNN to jointly predict confidence maps for body part detection, shown in (b), and part affinity fields for parts association, shown in (c). The parsing step performs a set of bipartite matchings to associate body parts candidates (d). We finally assemble them into full body poses for all people in the image (e).

<https://github.com/CMU-Perceptual-Computing-Lab/openpose>

# Multiple-Person Skeleton: Associative Embedding

Hourglass + AE



Figure 1. Both multi-person pose estimation and instance segmentation are examples of computer vision tasks that require detection of visual elements (joints of the body or pixels belonging to a semantic class) and grouping of these elements (as poses or individual object instances).

# Multiple-Person Skeleton: Associative Embedding

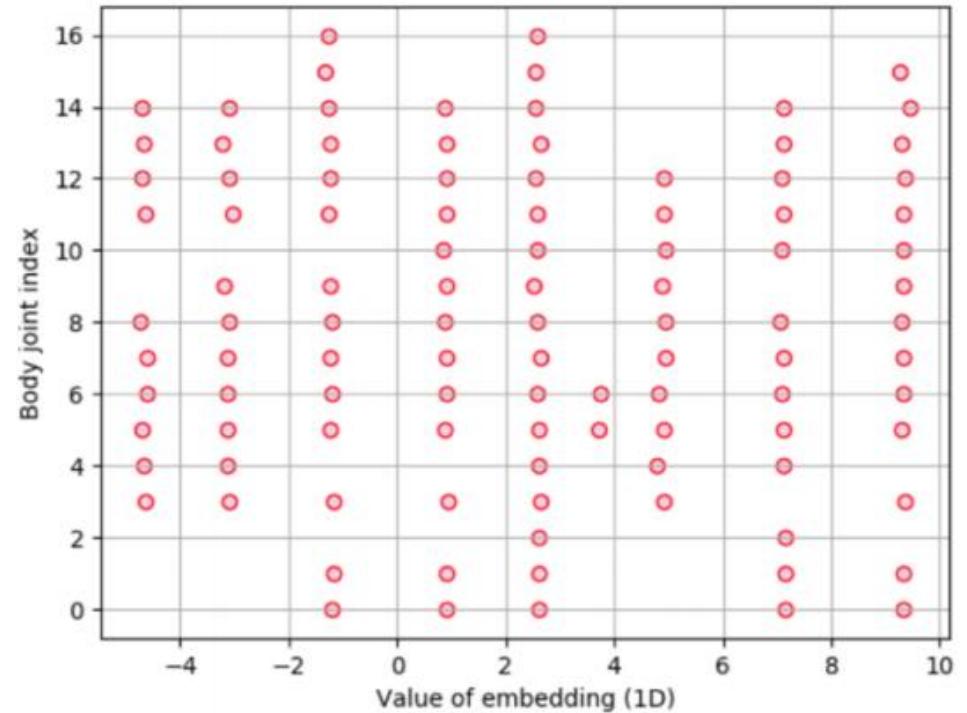


Figure 4. Tags produced by our network on a held-out validation image from the MS-COCO training set. The tag values are already well separated and decoding the groups is straightforward.

# Multiple-Person Skeleton: Discussion

- Top Down:
  - Depends on the detector
    - Fail in the crowd case
    - Fail with partial observation
    - can detect the small-scale human
  - More computation
  - Better localization when the input-size of single person skeleton is large
- Bottom up:
  - Fast computational speed
  - good at localizing the human with partial observation
  - Hard to assemble human

# Challenges in Skeleton

- combine top-down approaches with bottom-up approaches
- perform pose track
- handle the crowd case

# MegSkeleton

Introduction and demo Video

# Outline

- Detection
- Human Keypoint
- Conclusion

# Conclusion

- Detection
  - One stage: Densebox, YOLO, SSD, RetinaNet
  - Two Stage: RCNN, Fast RCNN, FasterRCNN, RFCN, FPN, Mask RCNN
- Skeleton
  - Single Person Skeleton: CPM, Hourglass
  - Multi-person Skeleton
    - Top Down
    - Bottom up: Openpose, Associative Embedding

# Thanks