

## YOLOv4

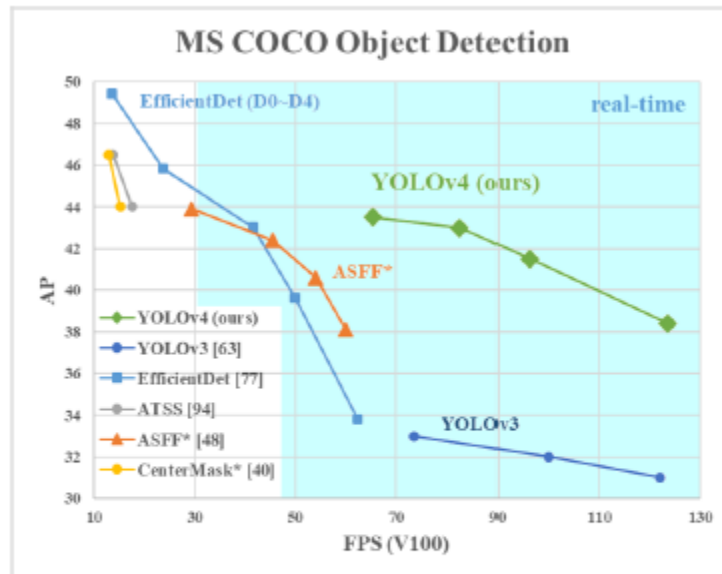


Figure 1: Comparison of the proposed YOLOv4 and other state-of-the-art object detectors. YOLOv4 runs twice faster than EfficientDet with comparable performance. Improves YOLOv3's AP and FPS by 10% and 12%, respectively.

## Object detection models:

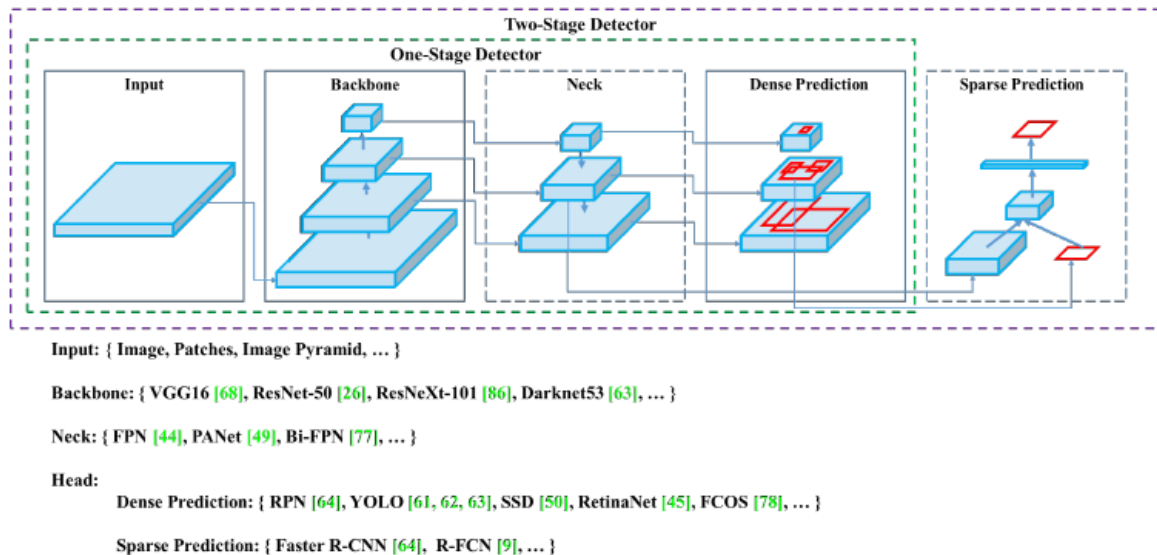


Figure 2: Object detector.

## Methodology:

The basic aim is fast operating speed of neural network, in production systems and optimization for parallel computations, rather than the low computation volume theoretical indicator (BFLOP). We present two options of real-time neural networks:

- For GPU we use a small number of groups (1 - 8) in convolutional layers: CSPResNeXt50 / CSPDarknet53
- For VPU - we use grouped-convolution, but we refrain from using Squeeze-and-excitement (SE) blocks - specifically this includes the following models: EfficientNet-lite / MixNet / GhostNet / MobileNetV3

## Selection of architecture

1) Our objective is to find the optimal balance among the input network resolution, the convolutional layer number, the parameter number (filter size<sup>2</sup> \* filters \* channel / groups), and the number of layer outputs (filters).

2) The next objective is to select additional blocks for increasing the receptive field and the best method of parameter aggregation from different backbone levels for different detector levels: e.g. FPN, PAN, ASFF, BiFPN.

In contrast to the classifier, the detector requires the following:

- Higher input network size (resolution) – for detecting multiple small-sized objects
- More layers – for a higher receptive field to cover the increased size of input network
- More parameters – for greater capacity of a model to detect multiple objects of different sizes in a single image

Table 1: Parameters of neural networks for image classification.

Backbone model	Input network resolution	Receptive field size	Parameters	Average size of layer output (WxHxC)	BFLOPs (512x512 network resolution)	FPS (GPU RTX 2070)
CSPResNext50	512x512	425x425	20.6 M	<b>1058 K</b>	31 (15.5 FMA)	62
CSPDarknet53	512x512	725x725	<b>27.6 M</b>	950 K	52 (26.0 FMA)	<b>66</b>
EfficientNet-B3 (ours)	512x512	<b>1311x1311</b>	12.0 M	668 K	11 (5.5 FMA)	26

The influence of the receptive field with different sizes is summarized as follows:

- Up to the object size - allows viewing the entire object
- Up to network size - allows viewing the context around the object
- Exceeding the network size - increases the number of connections between the image point and the final activation

3) We add the SPP block over the CSPDarknet53. We use PANet as the method of parameter aggregation from different backbone levels for different detector levels , instead of the FPN used in YOLOv3.

4) Finally, we choose CSPDarknet53 backbone, SPP additional module, PANet path-aggregation neck, and YOLOv3 (anchor based) head as the architecture of YOLOv4.

### **Selection of BoF(Bag of Freebies) and BoS(Bag of Specials)**

For improving the object detection training, a CNN usually uses the following:

- Activations: ReLU, leaky-ReLU, parametric-ReLU, ReLU6, SELU, Swish, or Mish
- Bounding box regression loss: MSE, IoU, GloU, CloU, DIoU
- Data augmentation: CutOut, MixUp, CutMix
- Regularization method: DropOut, DropPath [36], Spatial DropOut [79], or DropBlock
- Normalization of the network activations by their mean and variance: Batch Normalization (BN) [32], Cross-GPU Batch Normalization (CGBN or SyncBN) [93], Filter Response Normalization (FRN) [70], or Cross-Iteration Batch Normalization (CBN) [89]
- Skip-connections: Residual connections, Weighted residual connections, Multi-input weighted residual connections, or Cross stage partial connections (CSP)

In order to make the designed detector more suitable for training on single GPU, we made additional design and improvement as follows:

- We introduce a new method of data augmentation Mosaic, and Self-Adversarial Training (SAT)
- We select optimal hyper-parameters while applying genetic algorithms
- We modify some exciting methods to make our design suitable for efficient training and detection - modified SAM, modified PAN, and Cross mini-Batch Normalization (CmBN)

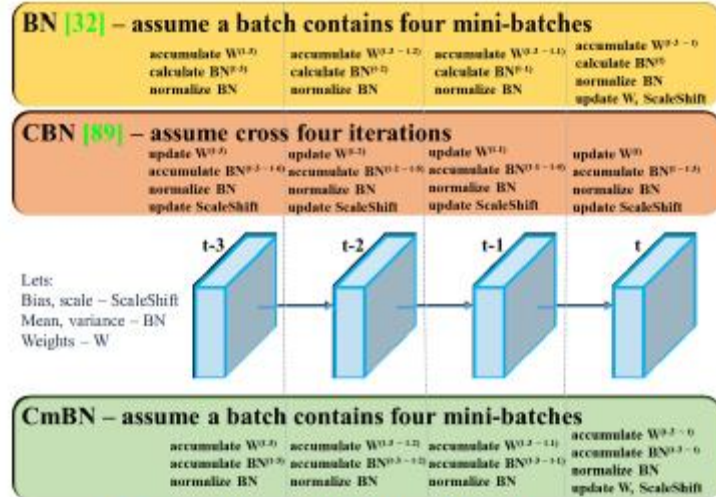


Figure 4: Cross mini-Batch Normalization.

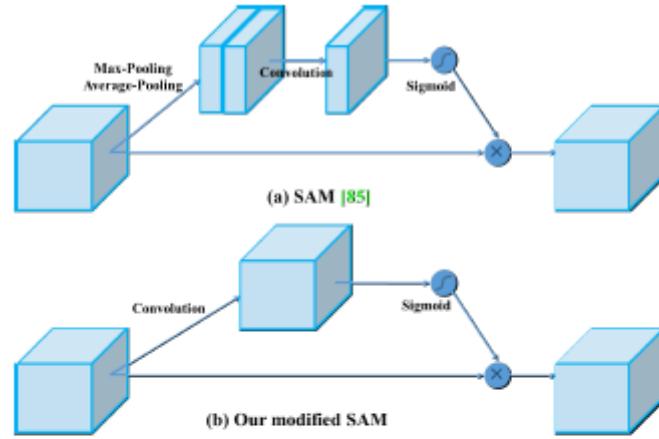


Figure 5: Modified SAM.

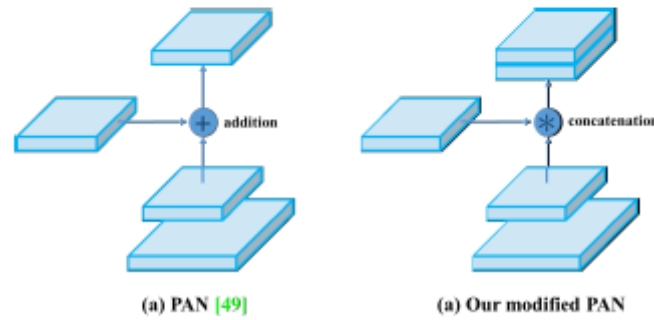


Figure 6: Modified PAN.

## **YOLOv4:**

YOLOv4 consists of:

- Backbone: CSPDarknet53
- Neck: SPP, PAN
- Head: YOLOv3

YOLO v4 uses:

- Bag of Freebies (BoF) for backbone: CutMix and Mosaic data augmentation, DropBlock regularization, Class label smoothing
- Bag of Specials (BoS) for backbone: Mish activation, Cross-stage partial connections (CSP), Multiinput weighted residual connections (MiWRC)
- Bag of Freebies (BoF) for detector: CloU-loss, CmBN, DropBlock regularization, Mosaic data augmentation, Self-Adversarial Training, Eliminate grid sensitivity, Using multiple anchors for a single ground truth, Cosine annealing scheduler [52], Optimal hyperparameters, Random training shapes
- Bag of Specials (BoS) for detector: Mish activation, SPP-block, SAM-block, PAN path-aggregation block, DIoU-NMS

## **Experiment**

Dataset:

test the influence of different training improvement techniques on accuracy of the classifier on ImageNet (ILSVRC 2012 val) dataset.

on the accuracy of the detector on MS COCO (test-dev 2017) dataset.

Table 2: Influence of BoF and Mish on the CSPResNeXt-50 classifier accuracy.

MixUp	CutMix	Mosaic	Blurring	Label Smoothing	Swish	Mish	Top-1	Top-5
							77.9%	94.0%
✓							77.2%	<b>94.0%</b>
	✓						<b>78.0%</b>	<b>94.3%</b>
		✓					<b>78.1%</b>	<b>94.5%</b>
			✓				77.5%	93.8%
				✓			<b>78.1%</b>	<b>94.4%</b>
					✓		64.5%	86.0%
						✓	<b>78.9%</b>	<b>94.5%</b>
	✓	✓		✓			<b>78.5%</b>	<b>94.8%</b>
	✓	✓		✓		✓	<b>79.8%</b>	<b>95.2%</b>

Table 3: Influence of BoF and Mish on the CSPDarknet-53 classifier accuracy.

MixUp	CutMix	Mosaic	Blurring	Label Smoothing	Swish	Mish	Top-1	Top-5
							77.2%	93.6%
	✓	✓		✓			<b>77.8%</b>	<b>94.4%</b>
	✓	✓		✓		✓	<b>78.7%</b>	<b>94.8%</b>

Table 4: Ablation Studies of Bag-of-Freebies. (CSPResNeXt50-PANet-SPP, 512x512).

S	M	IT	GA	LS	CBN	CA	DM	OA	loss	AP	AP <sub>50</sub>	AP <sub>75</sub>
									MSE	38.0%	60.0%	40.8%
✓									MSE	37.7%	59.9%	40.5%
	✓								MSE	<b>39.1%</b>	<b>61.8%</b>	<b>42.0%</b>
		✓							MSE	36.9%	59.7%	39.4%
			✓						MSE	<b>38.9%</b>	<b>61.7%</b>	<b>41.9%</b>
				✓					MSE	33.0%	55.4%	35.4%
					✓				MSE	<b>38.4%</b>	<b>60.7%</b>	<b>41.3%</b>
						✓			MSE	<b>38.7%</b>	<b>60.7%</b>	<b>41.9%</b>
							✓		MSE	35.3%	57.2%	38.0%
✓									GIoU	<b>39.4%</b>	59.4%	<b>42.5%</b>
✓									DIoU	<b>39.1%</b>	58.8%	<b>42.1%</b>
✓									CIoU	<b>39.6%</b>	59.2%	<b>42.6%</b>
✓	✓	✓	✓						CIoU	<b>41.5%</b>	<b>64.0%</b>	<b>44.8%</b>
	✓		✓						CIoU	36.1%	56.5%	38.4%
✓	✓	✓	✓					✓	MSE	<b>40.3%</b>	<b>64.0%</b>	<b>43.1%</b>
✓	✓	✓	✓					✓	GIoU	<b>42.4%</b>	<b>64.4%</b>	<b>45.9%</b>
✓	✓	✓	✓					✓	CIoU	<b>42.4%</b>	<b>64.4%</b>	<b>45.9%</b>

Table 5: Ablation Studies of Bag-of-Specials. (Size 512x512).

Model	AP	AP <sub>50</sub>	AP <sub>75</sub>
CSPResNeXt50-PANet-SPP	42.4%	64.4%	45.9%
CSPResNeXt50-PANet-SPP-RFB	41.8%	62.7%	45.1%
CSPResNeXt50-PANet-SPP-SAM	<b>42.7%</b>	<b>64.6%</b>	<b>46.3%</b>
CSPResNeXt50-PANet-SPP-SAM-G	41.6%	62.7%	45.0%
CSPResNeXt50-PANet-SPP-ASFF-RFB	41.1%	62.6%	44.4%

Table 6: Using different classifier pre-trained weightings for detector training (all other training parameters are similar in all models) .

Model (with optimal setting)	Size	AP	AP <sub>50</sub>	AP <sub>75</sub>
<b>CSPResNeXt50-PANet-SPP</b>	512x512	42.4	64.4	45.9
<b>CSPResNeXt50-PANet-SPP</b> (BoF-backbone)	512x512	42.3	64.3	45.7
<b>CSPResNeXt50-PANet-SPP</b> (BoF-backbone + Mish)	512x512	42.3	64.2	45.8
<b>CSPDarknet53-PANet-SPP</b> (BoF-backbone)	512x512	42.4	64.5	46.0
<b>CSPDarknet53-PANet-SPP</b> (BoF-backbone + Mish)	512x512	43.0	64.9	46.5

Table 7: Using different mini-batch size for detector training.

Model (without OA)	Size	AP	AP <sub>50</sub>	AP <sub>75</sub>
<b>CSPResNeXt50-PANet-SPP</b> (without BoF/BoS, mini-batch 4)	608	37.1	59.2	39.9
<b>CSPResNeXt50-PANet-SPP</b> (without BoF/BoS, mini-batch 8)	608	38.4	60.6	41.6
<b>CSPDarknet53-PANet-SPP</b> (with BoF/BoS, mini-batch 4)	512	41.6	64.1	45.0
<b>CSPDarknet53-PANet-SPP</b> (with BoF/BoS, mini-batch 8)	512	41.7	64.2	45.2

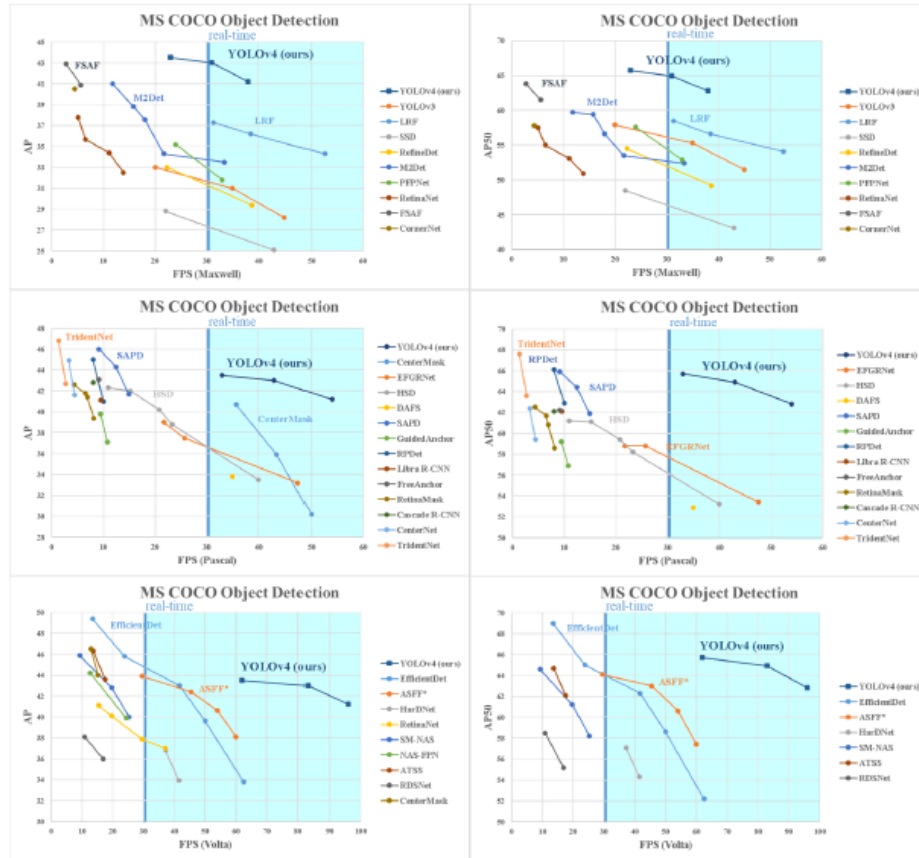


Figure 8: Comparison of the speed and accuracy of different object detectors. (Some articles stated the FPS of their detectors for only one of the GPUs: Maxwell/Pascal/Volta)

Table 8: Comparison of the speed and accuracy of different object detectors on the MS COCO dataset (test-dev 2017). (Real-time detectors with FPS 30 or higher are highlighted here. We compare the results with batch=1 without using tensorRT.)

Method	Backbone	Size	FPS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<b>YOLOv4: Optimal Speed and Accuracy of Object Detection</b>									
YOLOv4	CSPDarknet-53	416	38 (M)	41.2%	62.8%	44.3%	20.4%	44.4%	56.0%
YOLOv4	CSPDarknet-53	512	31 (M)	<b>43.0%</b>	<b>64.9%</b>	<b>46.5%</b>	<b>24.3%</b>	<b>46.1%</b>	<b>55.2%</b>
YOLOv4	CSPDarknet-53	608	23 (M)	43.5%	65.7%	47.3%	26.7%	46.7%	53.3%
<b>Learning Rich Features at High-Speed for Single-Shot Object Detection [84]</b>									
LRF	VGG-16	300	76.9 (M)	32.0%	51.5%	33.8%	12.6%	34.9%	47.0%
LRF	ResNet-101	300	52.6 (M)	34.3%	54.1%	36.6%	13.2%	38.2%	50.7%
LRF	VGG-16	512	38.5 (M)	36.2%	56.6%	38.7%	19.0%	39.9%	48.8%
LRF	ResNet-101	512	31.3 (M)	37.3%	58.5%	39.7%	19.7%	42.8%	50.1%
<b>Receptive Field Block Net for Accurate and Fast Object Detection [47]</b>									
RFBNet	VGG-16	300	66.7 (M)	30.3%	49.3%	31.8%	11.8%	31.9%	45.9%
RFBNet	VGG-16	512	33.3 (M)	33.8%	54.2%	35.9%	16.2%	37.1%	47.4%
RFBNet-E	VGG-16	512	30.3 (M)	34.4%	55.7%	36.4%	17.6%	37.0%	47.6%
<b>YOLOv3: An incremental improvement [63]</b>									
YOLOv3	Darknet-53	320	45 (M)	28.2%	51.5%	29.7%	11.9%	30.6%	43.4%
YOLOv3	Darknet-53	416	35 (M)	31.0%	55.3%	32.3%	15.2%	33.2%	42.8%
YOLOv3	Darknet-53	608	20 (M)	33.0%	57.9%	34.4%	18.3%	35.4%	41.9%
YOLOv3-SPP	Darknet-53	608	20 (M)	36.2%	60.6%	38.2%	20.6%	37.4%	46.1%
<b>SSD: Single shot multibox detector [50]</b>									
SSD	VGG-16	300	43 (M)	25.1%	43.1%	25.8%	6.6%	25.9%	41.4%
SSD	VGG-16	512	22 (M)	28.8%	48.5%	30.3%	10.9%	31.8%	43.5%
<b>Single-shot refinement neural network for object detection [95]</b>									
RefineDet	VGG-16	320	38.7 (M)	29.4%	49.2%	31.3%	10.0%	32.0%	44.4%
RefineDet	VGG-16	512	22.3 (M)	33.0%	54.5%	35.5%	16.3%	36.3%	44.3%
<b>M2det: A single-shot object detector based on multi-level feature pyramid network [98]</b>									
M2det	VGG-16	320	33.4 (M)	33.5%	52.4%	35.6%	14.4%	37.6%	47.6%
M2det	ResNet-101	320	21.7 (M)	34.3%	53.5%	36.5%	14.8%	38.8%	47.9%
M2det	VGG-16	512	18 (M)	37.6%	56.6%	40.5%	18.4%	43.4%	51.2%
M2det	ResNet-101	512	15.8 (M)	38.8%	59.4%	41.7%	20.5%	43.9%	53.4%
M2det	VGG-16	800	11.8 (M)	41.0%	59.7%	45.0%	22.1%	46.5%	53.8%
<b>Parallel Feature Pyramid Network for Object Detection [34]</b>									
PFPNet-R	VGG-16	320	33 (M)	31.8%	52.9%	33.6%	12%	35.5%	46.1%
PFPNet-R	VGG-16	512	24 (M)	35.2%	57.6%	37.9%	18.7%	38.6%	45.9%
<b>Focal Loss for Dense Object Detection [45]</b>									
RetinaNet	ResNet-50	500	13.9 (M)	32.5%	50.9%	34.8%	13.9%	35.8%	46.7%
RetinaNet	ResNet-101	500	11.1 (M)	34.4%	53.1%	36.8%	14.7%	38.5%	49.1%
RetinaNet	ResNet-50	800	6.5 (M)	35.7%	55.0%	38.5%	18.9%	38.9%	46.3%
RetinaNet	ResNet-101	800	5.1 (M)	37.8%	57.5%	40.8%	20.2%	41.1%	49.2%
<b>Feature Selective Anchor-Free Module for Single-Shot Object Detection [102]</b>									
AB+FSAF	ResNet-101	800	5.6 (M)	40.9%	61.5%	44.0%	24.0%	44.2%	51.3%
AB+FSAF	ResNeXt-101	800	2.8 (M)	42.9%	63.8%	46.3%	26.6%	46.2%	52.7%
<b>CornerNet: Detecting objects as paired keypoints [37]</b>									
CornerNet	Hourglass	512	4.4 (M)	40.5%	57.8%	45.3%	20.8%	44.8%	56.7%



Table 9: Comparison of the speed and accuracy of different object detectors on the MS COCO dataset (test-dev 2017). (Real-time detectors with FPS 30 or higher are highlighted here. We compare the results with batch=1 without using tensorRT.)

Method	Backbone	Size	FPS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<b>YOLOv4: Optimal Speed and Accuracy of Object Detection</b>									
YOLOv4	CSPDarknet-53	416	54 (P)	41.2%	62.8%	44.3%	20.4%	44.4%	56.0%
YOLOv4	CSPDarknet-53	512	43 (P)	43.0%	64.9%	46.5%	24.3%	46.1%	55.2%
YOLOv4	CSPDarknet-53	608	33 (P)	43.5%	65.7%	47.3%	26.7%	46.7%	53.3%
<b>CenterMask: Real-Time Anchor-Free Instance Segmentation [40]</b>									
CenterMask-Lite	MobileNetV2-FPN	600×	50.0 (P)	30.2%	-	-	14.2%	31.9%	40.9%
CenterMask-Lite	VoVNet-19-FPN	600×	43.5 (P)	35.9%	-	-	19.6%	38.0%	45.9%
CenterMask-Lite	VoVNet-39-FPN	600×	35.7 (P)	40.7%	-	-	22.4%	43.2%	53.5%
<b>Enriched Feature Guided Refinement Network for Object Detection [57]</b>									
EFGRNet	VGG-16	320	47.6 (P)	33.2%	53.4%	35.4%	13.4%	37.1%	47.9%
EFGRNet	VG-G16	512	25.7 (P)	37.5%	58.8%	40.4%	19.7%	41.6%	49.4%
EFGRNet	ResNet-101	512	21.7 (P)	39.0%	58.8%	42.3%	17.8%	43.6%	54.5%
<b>Hierarchical Shot Detector [3]</b>									
HSD	VGG-16	320	40 (P)	33.5%	53.2%	36.1%	15.0%	35.0%	47.8%
HSD	VGG-16	512	23.3 (P)	38.8%	58.2%	42.5%	21.8%	41.9%	50.2%
HSD	ResNet-101	512	20.8 (P)	40.2%	59.4%	44.0%	20.0%	44.4%	54.9%
HSD	ResNeXt-101	512	15.2 (P)	41.9%	61.1%	46.2%	21.8%	46.6%	57.0%
HSD	ResNet-101	768	10.9 (P)	42.3%	61.2%	46.9%	22.8%	47.3%	55.9%
<b>Dynamic anchor feature selection for single-shot object detection [41]</b>									
DAFS	VGG16	512	35 (P)	33.8%	52.9%	36.9%	14.6%	37.0%	47.7%
<b>Soft Anchor-Point Object Detection [101]</b>									
SAPD	ResNet-50	-	14.9 (P)	41.7%	61.9%	44.6%	24.1%	44.6%	51.6%
SAPD	ResNet-50-DCN	-	12.4 (P)	44.3%	64.4%	47.7%	25.5%	47.3%	57.0%
SAPD	ResNet-101-DCN	-	9.1 (P)	46.0%	65.9%	49.6%	26.3%	49.2%	59.6%
<b>Region proposal by guided anchoring [82]</b>									
RetinaNet	ResNet-50	-	10.8 (P)	37.1%	56.9%	40.0%	20.1%	40.1%	48.0%
Faster R-CNN	ResNet-50	-	9.4 (P)	39.8%	59.2%	43.5%	21.8%	42.6%	50.7%
<b>RepPoints: Point set representation for object detection [87]</b>									
RPDet	ResNet-101	-	10 (P)	41.0%	62.9%	44.3%	23.6%	44.1%	51.7%
RPDet	ResNet-101-DCN	-	8 (P)	45.0%	66.1%	49.0%	26.6%	48.6%	57.5%
<b>Libra R-CNN: Towards balanced learning for object detection [58]</b>									
Libra R-CNN	ResNet-101	-	9.5 (P)	41.1%	62.1%	44.7%	23.4%	43.7%	52.5%
<b>FreeAnchor: Learning to match anchors for visual object detection [96]</b>									
FreeAnchor	ResNet-101	-	9.1 (P)	43.1%	62.2%	46.4%	24.5%	46.1%	54.8%
<b>RetinaMask: Learning to Predict Masks Improves State-of-The-Art Single-Shot Detection for Free [14]</b>									
RetinaMask	ResNet-50-FPN	800×	8.1 (P)	39.4%	58.6%	42.3%	21.9%	42.0%	51.0%
RetinaMask	ResNet-101-FPN	800×	6.9 (P)	41.4%	60.8%	44.6%	23.0%	44.5%	53.5%
RetinaMask	ResNet-101-FPN-GN	800×	6.5 (P)	41.7%	61.7%	45.0%	23.5%	44.7%	52.8%
RetinaMask	ResNeXt-101-FPN-GN	800×	4.3 (P)	42.6%	62.5%	46.0%	24.8%	45.6%	53.8%
<b>Cascade R-CNN: Delving into high quality object detection [2]</b>									
Cascade R-CNN	ResNet-101	-	8 (P)	42.8%	62.1%	46.3%	23.7%	45.5%	55.2%
<b>CenterNet: Object detection with keypoint triplets [13]</b>									
CenterNet	Hourglass-52	-	4.4 (P)	41.6%	59.4%	44.2%	22.5%	43.1%	54.1%
CenterNet	Hourglass-104	-	3.3 (P)	44.9%	62.4%	48.1%	25.6%	47.4%	57.4%
<b>Scale-Aware Trident Networks for Object Detection [42]</b>									
TridentNet	ResNet-101	-	2.7 (P)	42.7%	63.6%	46.5%	23.9%	46.6%	56.6%
TridentNet	ResNet-101-DCN	-	1.3 (P)	46.8%	67.6%	51.5%	28.0%	51.2%	60.5%

Table 10: Comparison of the speed and accuracy of different object detectors on the MS COCO dataset (test-dev 2017). (Real-time detectors with FPS 30 or higher are highlighted here. We compare the results with batch=1 without using tensorRT.)

Method	Backbone	Size	FPS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<b>YOLOv4: Optimal Speed and Accuracy of Object Detection</b>									
<b>YOLOv4</b>	CSPDarknet-53	416	96 (V)	41.2%	62.8%	44.3%	20.4%	44.4%	56.0%
<b>YOLOv4</b>	CSPDarknet-53	512	83 (V)	43.0%	64.9%	46.5%	24.3%	46.1%	55.2%
<b>YOLOv4</b>	CSPDarknet-53	608	62 (V)	<b>43.5%</b>	<b>65.7%</b>	47.3%	<b>26.7%</b>	46.7%	53.3%
<b>EfficientDet: Scalable and Efficient Object Detection [77]</b>									
EfficientDet-D0	Efficient-B0	512	62.5 (V)	33.8%	52.2%	35.8%	12.0%	38.3%	51.2%
EfficientDet-D1	Efficient-B1	640	50.0 (V)	39.6%	58.6%	42.3%	17.9%	44.3%	56.0%
EfficientDet-D2	Efficient-B2	768	41.7 (V)	43.0%	62.3%	46.2%	22.5%	<b>47.0%</b>	<b>58.4%</b>
EfficientDet-D3	Efficient-B3	896	23.8 (V)	45.8%	65.0%	49.3%	26.6%	49.4%	59.8%
<b>Learning Spatial Fusion for Single-Shot Object Detection [48]</b>									
YOLOv3 + ASFF*	Darknet-53	320	60 (V)	38.1%	57.4%	42.1%	16.1%	41.6%	53.6%
YOLOv3 + ASFF*	Darknet-53	416	54 (V)	40.6%	60.6%	45.1%	20.3%	44.2%	54.1%
YOLOv3 + ASFF*	Darknet-53	608×	45.5 (V)	42.4%	63.0%	<b>47.4%</b>	25.5%	45.7%	52.3%
YOLOv3 + ASFF*	Darknet-53	800×	29.4 (V)	43.9%	64.1%	49.2%	27.0%	46.6%	53.4%
<b>HardNet: A Low Memory Traffic Network [4]</b>									
RFBNet	HardNet68	512	41.5 (V)	33.9%	54.3%	36.2%	14.7%	36.6%	50.5%
RFBNet	HardNet85	512	37.1 (V)	36.8%	57.1%	39.5%	16.9%	40.5%	52.9%
<b>Focal Loss for Dense Object Detection [45]</b>									
RetinaNet	ResNet-50	640	37 (V)	37.0%	-	-	-	-	-
RetinaNet	ResNet-101	640	29.4 (V)	37.9%	-	-	-	-	-
RetinaNet	ResNet-50	1024	19.6 (V)	40.1%	-	-	-	-	-
RetinaNet	ResNet-101	1024	15.4 (V)	41.1%	-	-	-	-	-
<b>SM-NAS: Structural-to-Modular Neural Architecture Search for Object Detection [88]</b>									
SM-NAS: E2	-	800×600	25.3 (V)	40.0%	58.2%	43.4%	21.1%	42.4%	51.7%
SM-NAS: E3	-	800×600	19.7 (V)	42.8%	61.2%	46.5%	23.5%	45.5%	55.6%
SM-NAS: E5	-	1333×800	9.3 (V)	45.9%	64.6%	49.6%	27.1%	49.0%	58.0%
<b>NAS-FPN: Learning scalable feature pyramid architecture for object detection [17]</b>									
NAS-FPN	ResNet-50	640	24.4 (V)	39.9%	-	-	-	-	-
NAS-FPN	ResNet-50	1024	12.7 (V)	44.2%	-	-	-	-	-
<b>Bridging the Gap Between Anchor-based and Anchor-free Detection via Adaptive Training Sample Selection [94]</b>									
ATSS	ResNet-101	800×	17.5 (V)	43.6%	62.1%	47.4%	26.1%	47.0%	53.6%
ATSS	ResNet-101-DCN	800×	13.7 (V)	46.3%	64.7%	50.4%	27.7%	49.8%	58.4%
<b>RDSNet: A New Deep Architecture for Reciprocal Object Detection and Instance Segmentation [83]</b>									
RDSNet	ResNet-101	600	16.8 (V)	36.0%	55.2%	38.7%	17.4%	39.6%	49.7%
RDSNet	ResNet-101	800	10.9 (V)	38.1%	58.5%	40.8%	21.2%	41.5%	48.2%
<b>CenterMask: Real-Time Anchor-Free Instance Segmentation [40]</b>									
CenterMask	ResNet-101-FPN	800×	15.2 (V)	44.0%	-	-	25.8%	46.8%	54.9%
CenterMask	VoVNet-99-FPN	800×	12.9 (V)	46.5%	-	-	28.7%	48.9%	57.2%

Paperswithcode link: <https://paperswithcode.com/paper/yolov4-optimal-speed-and-accuracy-of-object>

Github link(official): <https://github.com/AlexeyAB/darknet>