

# 计算机图像实体检测综述

A Survey on Computer Image Entity Detection

郑涵潇 林超超 吴浩凯

# 汇报目录

01

简单  
介绍

Introduce

02

进展  
总结

Survey

03

项目  
展示

Demo

04

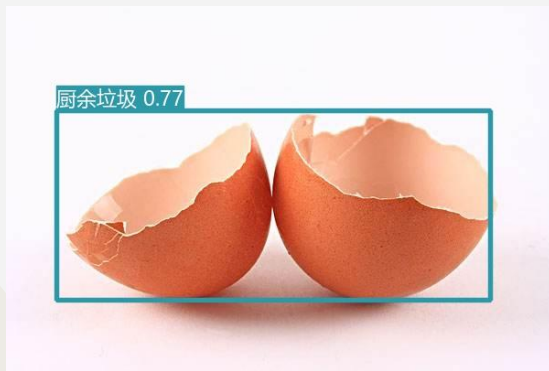
未来  
展望

Future

# 01 简单介绍

Introduce

# What is 计算机图像实体检测?



定位 + 分类

# What is 计算机图像实体检测?

Classification



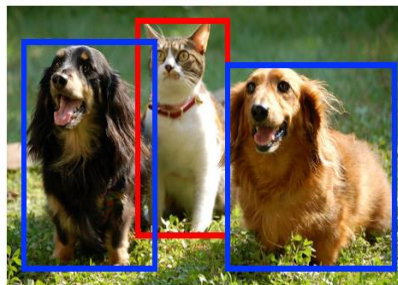
CAT

Classification  
+ Localization



CAT

Object Detection



CAT, DOG

Instance Segmentation



CAT, DOG

Single object

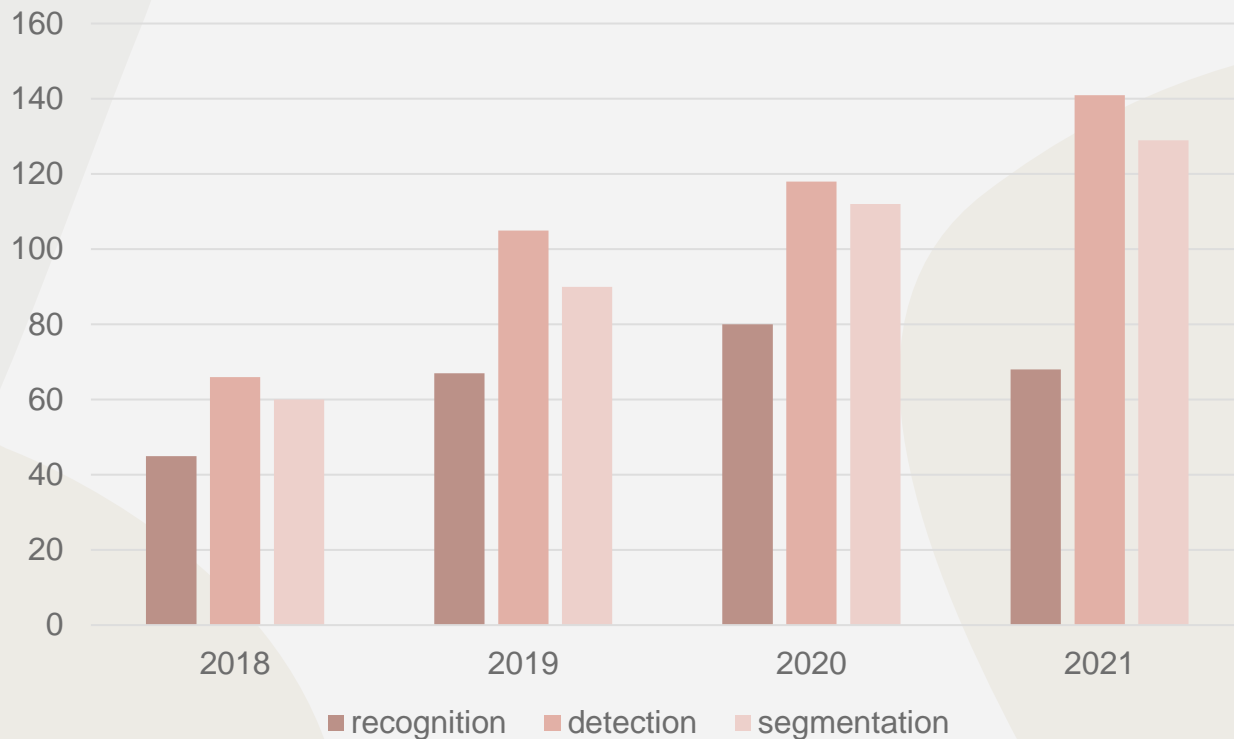
Multiple objects

# Where is 计算机图像实体检测?



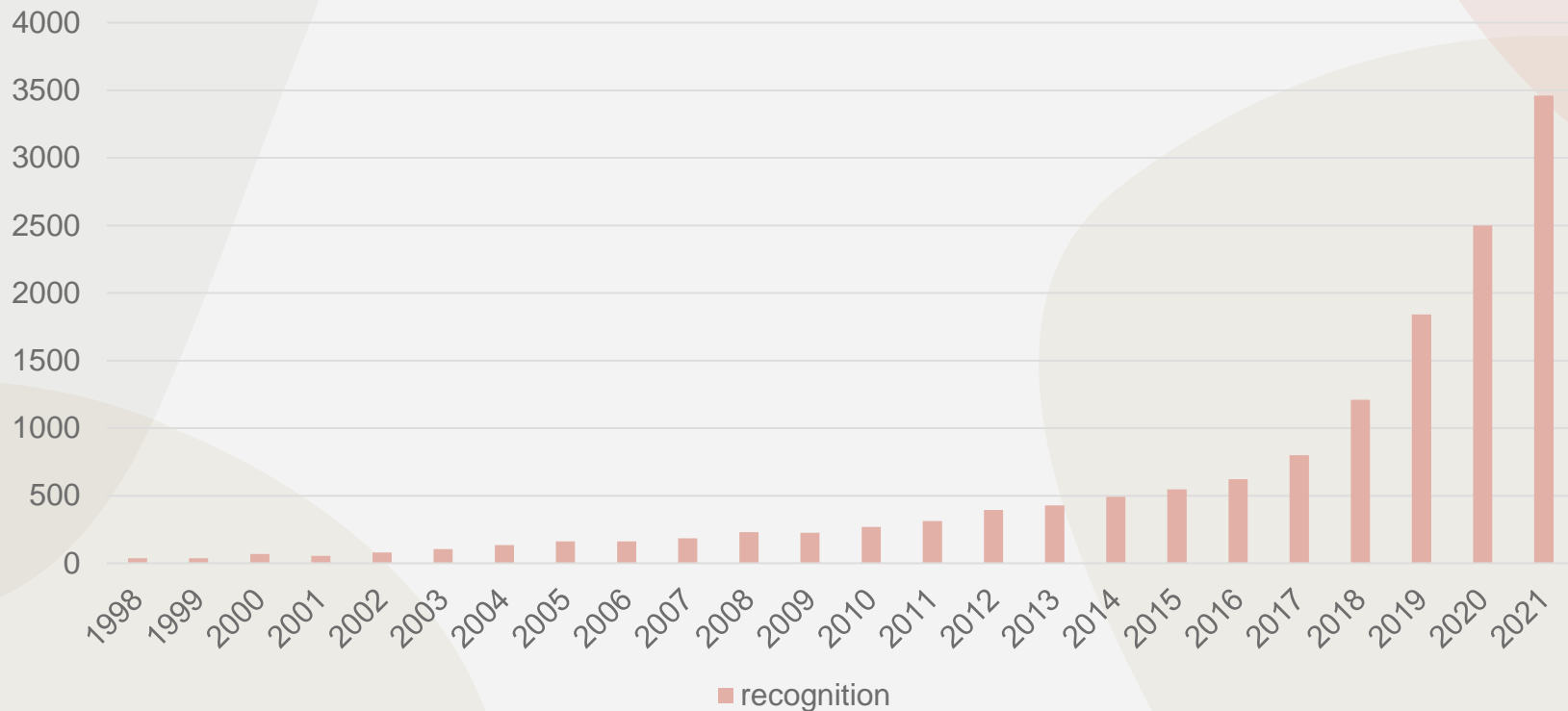
# How is 计算机图像实体检测?

Keypoint of CVPR 2018-2021

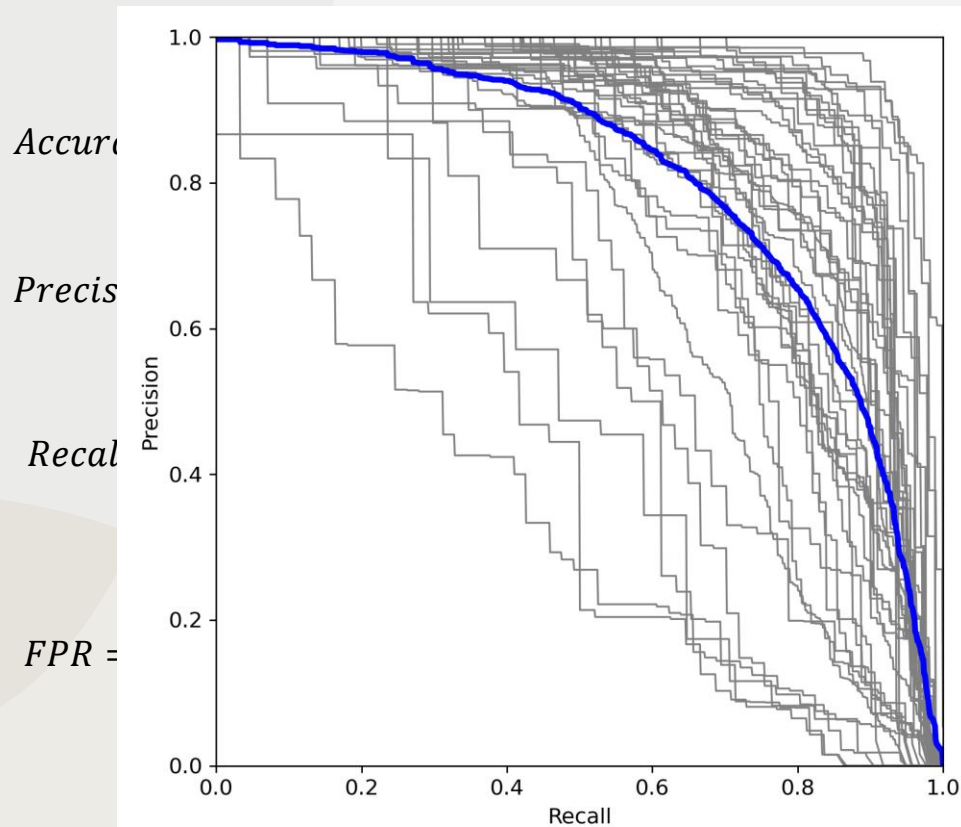


# How is 计算机图像实体检测?

“Object Detection” & “Detecting Objects” in Google Scholar



# How to 计算机图像实体检测?



# How to 计算机图像实体检测?

IMAGENET

Download

Download Images

March 11, 2011

October 10, 2011

You have been

Winter 2021

- ImageNet
- Processes
- ImageNet

People subscribe

- Descriptions
- Unsafe images
- ImageNet
- Due to security

PASCAL

Welcome, Ro

Welcome to

- LEAD
- PASC
- PASC
- PASC
- PASC
- PASC

Announcements

- There
- You can
- For ne
- You can
- anony

You are advised

应用领域	数据集	发布时间	数据集描述
	WebFace	2021	该数据集是最大的百万级公开人脸识别训练集, 包含了有噪声的 4M 身份/260M 人脸 (WebFace260M) 和清理过的 2M 身份/42M 人脸(WebFace42M) 链接: <a href="https://www.face-benchmark.org">https://www.face-benchmark.org</a>
文本检测	MSRA-TD500	2012	该数据集包含 500 张自然图像, 作为评估文本检测算法的基准。从室内(办公室和商场)和室外(街道)场景中拍摄, 由于文本的多样性和图像背景的复杂性, 数据集具有挑战性。链接: <a href="http://www.iapr-c11.org/mediawiki/index.php/MSRA_Text_Detection_500_Database_(MSRA-TD500)">http://www.iapr-c11.org/mediawiki/index.php/MSRA_Text_Detection_500_Database_(MSRA-TD500)</a>
	COCOText	2016	该数据集基于 MSCOCO 数据集, 用于自然图像中的文本检测和识别, 最新版包含 63,686 个图像, 145,859 个文本实例, 包括手写、打印、易读、难读、英文和非英文, 数据集多样化。 链接: <a href="https://vision.cornell.edu/se3/coco-text-2/">https://vision.cornell.edu/se3/coco-text-2/</a>
	Google FSNS	2017	该数据集包含超过一百万张从法国 Google 街景图像中裁剪出来的街道名称标志图像。每个图像包含相同街道名称标志的四个视图, 路牌中的文字最多可以跨越三行。 链接: <a href="https://irc.cvc.uab.es/?ch=6">https://irc.cvc.uab.es/?ch=6</a>
遥感检测	VeDAI	2015	该数据集用于遥感图像中的多类车辆检测, 包含 1210 张图像, 总共 3640 个车辆实例与 9 个类别, 空间分辨率 12.5 cm。车辆实例具有多方向、光照/阴影变化、镜面反射或遮挡等表现。 链接: <a href="https://downloads.greyc.fr/vedai/">https://downloads.greyc.fr/vedai/</a>
	DOTA	2017	该数据集来源 GoogleEarth 和中国卫星, 是目前最大的光学遥感图像数据集, 包含 2806 张遥感图像(图片尺寸从 800*800 到 4000*4000), 一共 188,282 个实例, 分为 15 个类别。 链接: <a href="https://captain-whu.github.io/DOTA/dataset.html">https://captain-whu.github.io/DOTA/dataset.html</a>

# 02 进展总结

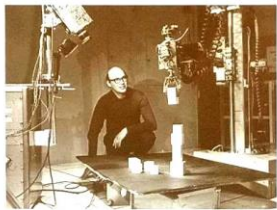
Survey

# Life of 计算机图像实体检测

## 早期探索阶段(1966-1998)

1966

Marvin Minsky  
&  
Gerald Sussman



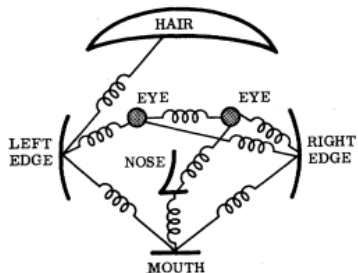
Marvin Minsky  
1927 - 2016

"Spend the summer linking a camera to a computer  
and getting the computer to describe what it saw."

Marvin Minsky to his undergrad student Gerald Sussman in 1966

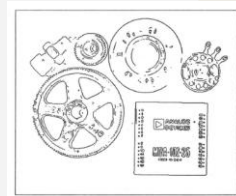
1973

图画结构  
Pictorial Structure



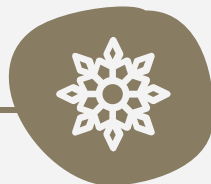
1980s

人工智能寒冬  
各种滤波器诞生



1998

神经网络  
人脸检测



# Life of 计算机图像实体检测

## 传统计算阶段(1999-2011)

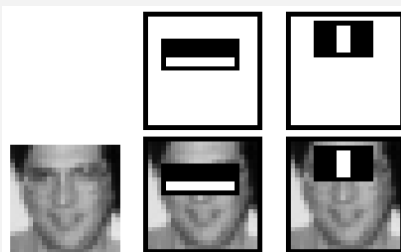
1999

SIFT  
尺度不变特征变换



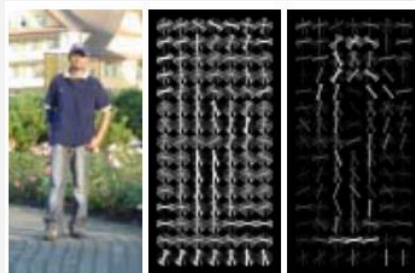
2001

VJ检测器  
Harr-like特征



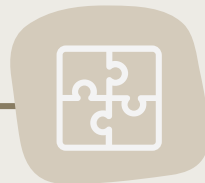
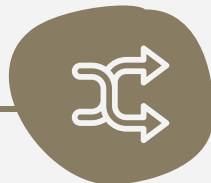
2005

HOG检测器  
方向梯度直方图



2008

DPM  
可变形部件模型



# Life of 计算机图像实体检测

## 传统计算阶段(1999-2011)

- 选取候选区域

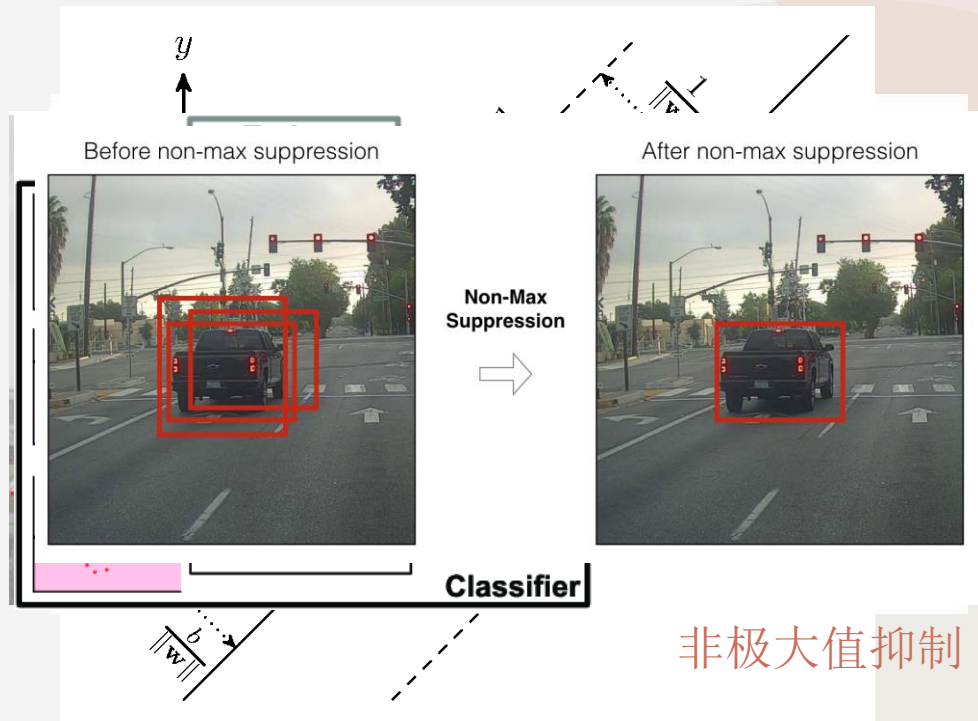
滑动窗口法    选择性搜索法    ...

- 特征提取

SIFT                      HOG                      ...

- 特征向量分类

SVM                      AdaBoost                      ...



非极大值抑制

# Life of 计算机图像实体检测

## 深度学习阶段(2012-至今)

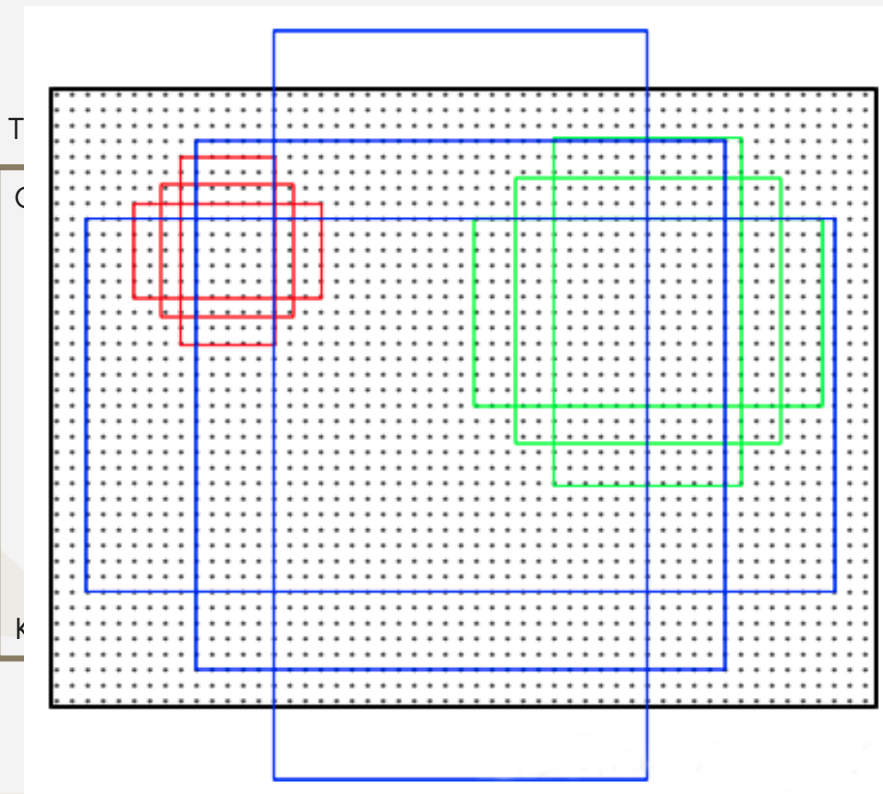
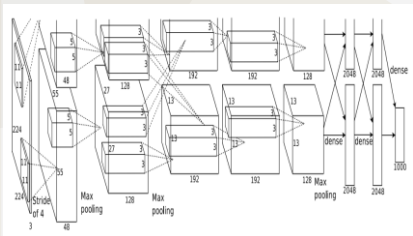
2012

深度卷积神经网络  
AlexNet



Anchor-based

Anchor-free



2018

Cascade RCNN

2019

Grid RCNN

2020

YOLOv4

2021

YOLOv5

2019

CornerNet-Lite ResPoint  
ExtremeNet

2021

SAPD

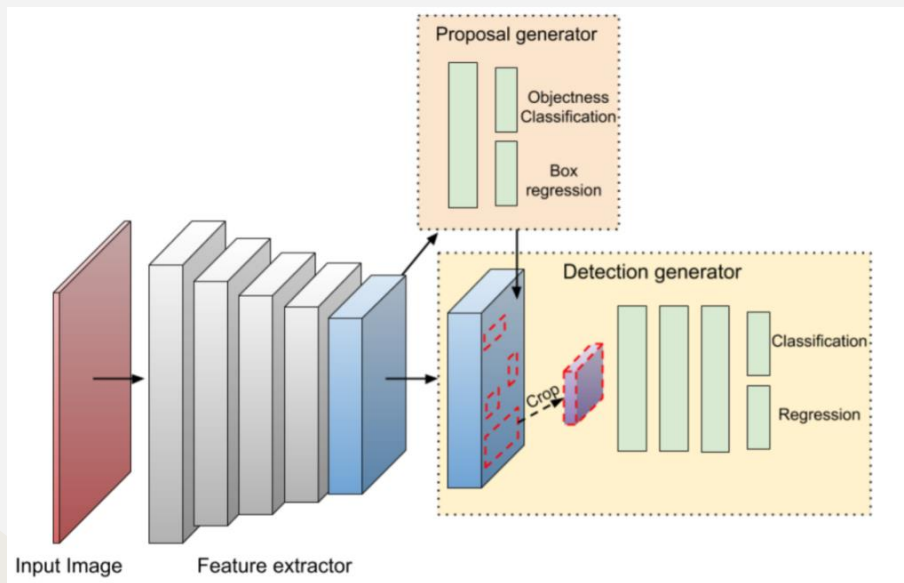
YOLOX



# Anchor-based Two-Stage

# Life of 计算机图像实体检测

## Anchor-based Two-Stage



2014

2015

2016

2017

2018

2019

RCNN  
SPPNet

Fast RCNN

Faster RCNN

FPN

Cascade RCNN

Grid RCNN

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

RCNN: Region CNN

### Region Proposals: Selective Search

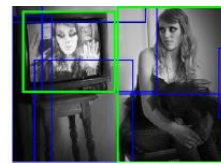
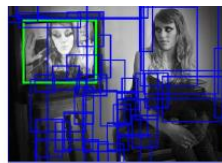
BBox Regressor

AlexNet

NMS

Bottom-up segmentation, merging regions at multiple scales

VOC 2010 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
DPM v5 [17] <sup>†</sup>	49.2	53.8	13.1	15.3	35.5	53.4	49.7	27.0	17.2	28.8	14.7	17.8	46.4	51.2	47.7	10.8	34.2	20.7	43.8	38.3	33.4
UVA [32]	56.2	42.4	15.3	12.6	21.8	49.3	36.8	46.1	12.9	32.1	30.0	36.5	43.5	52.9	32.9	15.3	41.1	31.8	47.0	44.8	35.1
Regionlets [35]	65.0	48.9	25.9	24.6	24.5	56.1	54.5	51.2	17.0	28.9	30.2	35.8	40.2	55.7	43.5	14.3	43.9	32.6	54.0	45.9	39.7
SegDPM [15] <sup>†</sup>	61.4	53.4	25.6	25.2	35.5	51.7	50.6	50.8	19.3	33.8	26.8	40.4	48.3	54.4	47.1	14.8	38.7	35.0	52.8	43.1	40.4
R-CNN	67.1	64.1	46.7	32.0	30.5	56.4	57.2	65.9	27.0	47.3	40.9	66.6	57.8	65.9	53.6	26.7	56.5	38.1	52.8	50.2	50.2
R-CNN BB	<b>71.8</b>	<b>65.8</b>	<b>53.0</b>	<b>36.8</b>	<b>35.9</b>	<b>59.7</b>	<b>60.0</b>	<b>69.9</b>	<b>27.9</b>	<b>50.6</b>	<b>41.4</b>	<b>70.0</b>	<b>62.0</b>	<b>69.0</b>	<b>58.1</b>	<b>29.5</b>	<b>59.4</b>	<b>39.3</b>	<b>61.2</b>	<b>52.4</b>	<b>53.7</b>

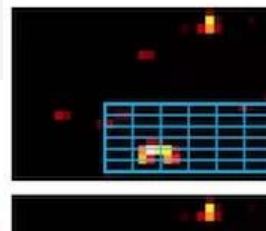


Input Image

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### SPPNet: Spatial Pyramid Pooling



method	mAP	areo	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
SPP-net (1)	59.2	<b>68.6</b>	69.7	57.1	41.2	40.5	66.3	71.3	72.5	34.4	<b>67.3</b>	61.7	63.1	71.0	69.8	57.6	29.7	59.0	50.2	65.2	68.0
SPP-net (2)	59.1	65.7	71.4	57.4	<b>42.4</b>	39.9	67.0	71.4	70.6	32.4	66.7	61.7	64.8	71.7	70.4	56.5	30.8	59.9	53.2	63.9	64.6
combination	<b>60.9</b>	68.5	<b>71.7</b>	<b>58.7</b>	41.9	<b>42.5</b>	<b>67.7</b>	<b>72.1</b>	<b>73.8</b>	<b>34.7</b>	67.0	<b>63.4</b>	<b>66.0</b>	<b>72.5</b>	<b>71.3</b>	<b>58.9</b>	<b>32.8</b>	<b>60.9</b>	<b>56.1</b>	<b>67.9</b>	<b>68.8</b>

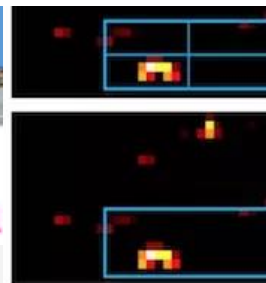


R-CNN

2000 nets on image regions



1



# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### Fast RCNN

method	train set	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	mAP
SPPnet BB [11] <sup>†</sup>	07 \ diff	73.9	72.3	62.5	51.5	44.4	74.4	73.0	74.4	42.3	73.6	57.7	70.3	74.6	74.3	54.2	34.0	56.4	56.4	67.9	73.5	63.1
R-CNN BB [10]	07	73.4	77.0	63.4	45.4	<b>44.6</b>	75.1	78.1	79.8	40.5	73.7	62.2	79.4	78.1	73.1	64.2	<b>35.6</b>	66.8	67.2	70.4	<b>71.1</b>	66.0
FRCN [ours]	07	74.5	78.3	69.2	53.2	36.6	77.3	78.2	82.0	40.7	72.7	67.9	79.6	79.2	73.0	69.0	30.1	65.4	70.2	75.8	65.8	66.9
FRCN [ours]	07 \ diff	74.6	<b>79.0</b>	68.6	57.0	39.3	79.5	<b>78.6</b>	81.9	<b>48.0</b>	74.0	67.4	80.5	80.7	74.1	69.6	31.8	67.1	68.4	75.3	65.5	68.1
FRCN [ours]	07+12	<b>77.0</b>	78.1	<b>69.3</b>	<b>59.4</b>	38.3	<b>81.6</b>	<b>78.6</b>	<b>86.7</b>	42.8	<b>78.8</b>	<b>68.9</b>	<b>84.7</b>	<b>82.0</b>	<b>76.6</b>	<b>69.9</b>	31.8	<b>70.1</b>	<b>74.8</b>	<b>80.4</b>	70.4	<b>70.0</b>

Table 1. VOC 2007 test detection average precision (%). All methods use VGG16. Training set key: **07**: VOC07 trainval, **07 \ diff**: **07** without “difficult” examples, **07+12**: union of **07** and VOC12 trainval. <sup>†</sup>SPPnet results were prepared by the authors of [11].

method	train set	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	mAP
BabyLearning	Prop.	77.7	73.8	62.3	48.8	45.4	67.3	67.0	80.3	41.3	70.8	49.7	79.5	74.7	78.6	64.5	36.0	69.9	55.7	70.4	61.7	63.8
R-CNN BB [10]	12	79.3	72.4	63.1	44.0	44.4	64.6	66.3	84.9	38.8	67.3	48.4	82.3	75.0	76.7	65.7	35.8	66.2	54.8	69.1	58.8	62.9
SegDeepM	12+seg	<b>82.3</b>	75.2	67.1	50.7	<b>49.8</b>	71.1	69.6	88.2	42.5	71.2	50.0	85.7	76.6	81.8	69.3	<b>41.5</b>	<b>71.9</b>	62.2	73.2	<b>64.6</b>	67.2
FRCN [ours]	12	80.1	74.4	67.7	49.4	41.4	74.2	68.8	87.8	41.9	70.1	50.2	86.1	77.3	81.1	70.4	33.3	67.0	63.3	77.2	60.0	66.1
FRCN [ours]	07++12	82.0	<b>77.8</b>	<b>71.6</b>	<b>55.3</b>	42.4	<b>77.3</b>	<b>71.7</b>	<b>89.3</b>	<b>44.5</b>	<b>72.1</b>	<b>53.7</b>	<b>87.7</b>	<b>80.0</b>	<b>82.5</b>	<b>72.7</b>	36.6	68.7	<b>65.4</b>	<b>81.1</b>	62.7	<b>68.8</b>

Table 2. VOC 2010 test detection average precision (%). BabyLearning uses a network based on [17]. All other methods use VGG16. Training set key: **12**: VOC12 trainval, **Prop.**: proprietary dataset, **12+seg**: **12** with segmentation annotations, **07++12**: union of VOC07 trainval, VOC07 test, and VOC12 trainval.

method	train set	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	mAP
BabyLearning	Prop.	78.0	74.2	61.3	45.7	42.7	68.2	66.8	80.2	40.6	70.0	49.8	79.0	74.5	77.9	64.0	35.3	67.9	55.7	68.7	62.6	63.2
NUS_NIN.c2000	Unk.	80.2	73.8	61.9	43.7	<b>43.0</b>	70.3	67.6	80.7	41.9	69.7	51.7	78.2	75.2	76.9	65.1	<b>38.6</b>	<b>68.3</b>	58.0	68.7	63.3	63.8
R-CNN BB [10]	12	79.6	72.7	61.9	41.2	41.9	65.9	66.4	84.6	38.5	67.2	46.7	82.0	74.8	76.0	65.2	35.6	65.4	54.2	67.4	60.3	62.4
FRCN [ours]	12	80.3	74.7	66.9	46.9	37.7	73.9	68.6	87.7	41.7	71.1	51.1	86.0	77.8	79.8	69.8	32.1	65.5	63.8	76.4	61.7	65.7
FRCN [ours]	07++12	<b>82.3</b>	<b>78.4</b>	<b>70.8</b>	<b>52.3</b>	38.7	<b>77.8</b>	<b>71.6</b>	<b>89.3</b>	<b>44.2</b>	<b>73.0</b>	<b>55.0</b>	<b>87.5</b>	<b>80.5</b>	<b>80.8</b>	<b>72.0</b>	35.1	<b>68.3</b>	<b>65.7</b>	<b>80.4</b>	<b>64.2</b>	<b>68.4</b>

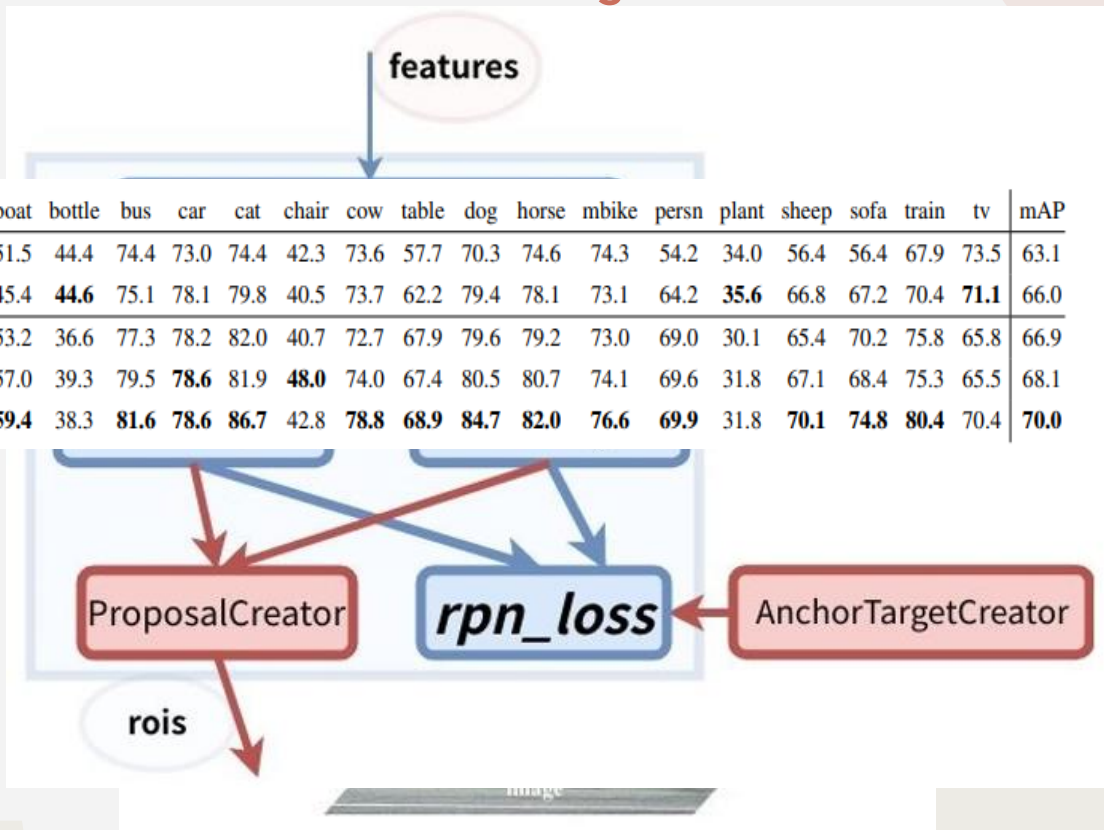
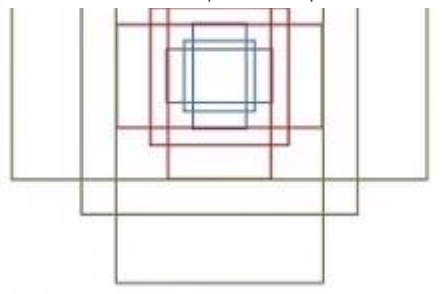
Table 3. VOC 2012 test detection average precision (%). BabyLearning and NUS\_NIN.c2000 use networks based on [17]. All other methods use VGG16. Training set key: see Table 2, **Unk.**: unknown.

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### Faster RCNN

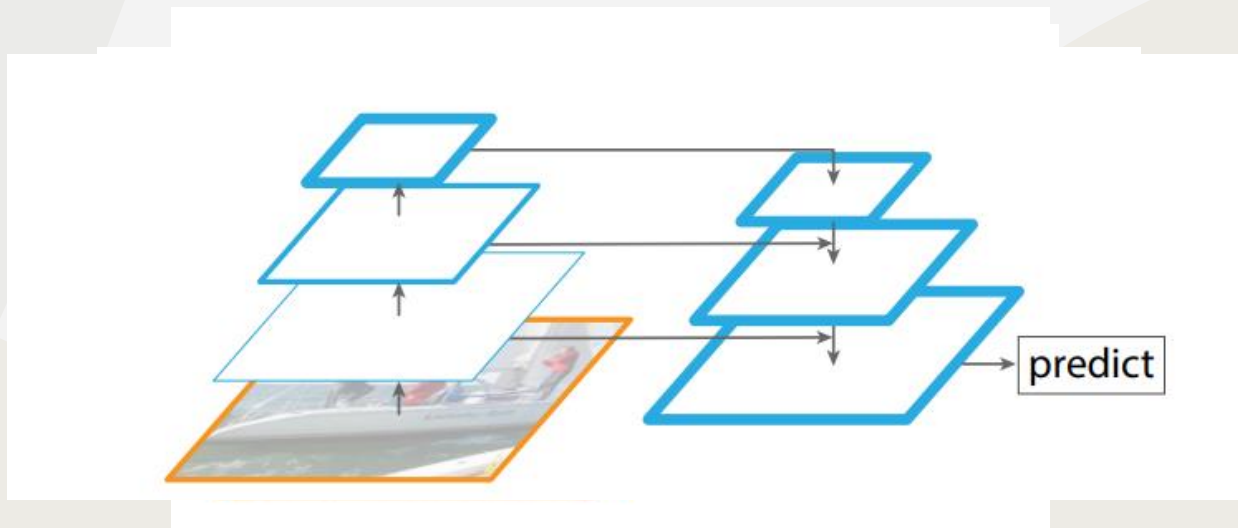
method	train set	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	mAP
SPPnet BB [11] <sup>†</sup>	07 \ diff	73.9	72.3	62.5	51.5	44.4	74.4	73.0	74.4	42.3	73.6	57.7	70.3	74.6	74.3	54.2	34.0	56.4	56.4	67.9	73.5	63.1
R-CNN BB [10]	07	73.4	77.0	63.4	45.4	<b>44.6</b>	75.1	78.1	79.8	40.5	73.7	62.2	79.4	78.1	73.1	64.2	<b>35.6</b>	66.8	67.2	70.4	<b>71.1</b>	66.0
FRCN [ours]	07	74.5	78.3	69.2	53.2	36.6	77.3	78.2	82.0	40.7	72.7	67.9	79.6	79.2	73.0	69.0	30.1	65.4	70.2	75.8	65.8	66.9
FRCN [ours]	07 \ diff	74.6	<b>79.0</b>	68.6	57.0	39.3	79.5	<b>78.6</b>	81.9	<b>48.0</b>	74.0	67.4	80.5	80.7	74.1	69.6	31.8	67.1	68.4	75.3	65.5	68.1
FRCN [ours]	07+12	<b>77.0</b>	78.1	<b>69.3</b>	<b>59.4</b>	38.3	<b>81.6</b>	<b>78.6</b>	<b>86.7</b>	42.8	<b>78.8</b>	<b>68.9</b>	<b>84.7</b>	<b>82.0</b>	<b>76.6</b>	<b>69.9</b>	31.8	<b>70.1</b>	<b>74.8</b>	<b>80.4</b>	70.4	<b>70.0</b>



# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### 基于金字塔的实体检测

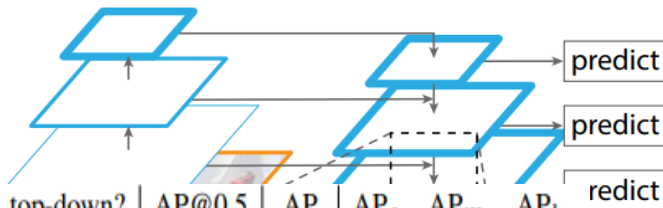


单尺度特征金字塔

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### FPN: Feature Pyramid Networks



Fast R-CNN	proposals	feature	head	lateral?	top-down?	AP@0.5	AP	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	redict
(a) baseline on conv4	RPN, {P <sub>k</sub> }	C <sub>4</sub>	conv5			54.7	31.9	15.7	36.5	45.5	
(b) baseline on conv5	RPN, {P <sub>k</sub> }	C <sub>5</sub>	2fc			52.9	28.8	11.9	32.4	43.4	
(c) FPN	RPN, {P <sub>k</sub> }	{P <sub>k</sub> }	2fc	✓	✓	<b>56.9</b>	<b>33.9</b>	<b>17.8</b>	<b>37.7</b>	<b>45.8</b>	

*Ablation experiments follow:*

(d) bottom-up pyramid	RPN, {P <sub>k</sub> }	{P <sub>k</sub> }	2fc	✓		44.9	24.9	10.9	24.4	38.5	
(e) top-down pyramid, w/o lateral	RPN, {P <sub>k</sub> }	{P <sub>k</sub> }	2fc		✓	54.0	31.3	13.3	35.2	45.3	
(f) only finest level	RPN, {P <sub>k</sub> }	P <sub>2</sub>	2fc	✓	✓	56.3	33.4	17.3	37.3	45.6	



## 多层特征融合金字塔

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### Cascade RCNN

	backbone	cascade	train speed	test speed	param	val (5k)						test-dev (20k)					
						AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Faster R-CNN	VGG	✗	0.12s	0.075s	278M	23.6	43.9	23.0	8.0	26.2	35.5	23.5	43.9	22.6	8.1	25.1	34.7
		✓	0.14s	0.115s	704M	27.0	44.2	27.7	8.6	29.1	42.2	26.9	44.3	27.8	8.3	28.2	41.1
R-FCN	ResNet-50	✗	0.19s	0.07s	133M	27.0	48.7	26.9	9.8	30.9	40.3	27.1	49.0	26.9	10.4	29.7	39.2
		✓	0.24s	0.075s	184M	31.1	49.8	32.8	10.4	34.4	48.5	30.9	49.9	32.6	10.5	33.1	46.9
R-FCN	ResNet-101	✗	0.23s	0.075s	206M	30.3	52.2	30.8	12.0	34.7	44.3	30.5	52.9	31.2	12.0	33.9	43.8
		✓	0.29s	0.083s	256M	33.3	52.0	35.2	11.8	37.2	51.1	33.3	52.6	35.2	12.1	36.2	49.3
FPN+	ResNet-50	✗	0.30s	0.095s	165M	36.5	58.6	39.2	20.8	40.0	47.8	36.5	59.0	39.2	20.3	38.8	46.4
		✓	0.33s	0.115s	272M	40.3	59.4	43.7	22.9	43.7	54.1	40.6	59.9	44.0	22.6	42.7	52.1
FPN+	ResNet-101	✗	0.38s	0.115s	238M	38.5	60.6	41.7	22.1	41.9	51.1	38.8	61.1	41.9	21.3	41.8	49.8
		✓	0.41s	0.14s	345M	42.7	61.6	46.6	23.8	46.2	57.4	42.8	62.1	46.3	23.7	45.5	55.2

# Life of 计算机图像实体检测

## Anchor-based Two-Stage

### Grid RCNN

method	backbone	AP	AP <sub>.5</sub>	AP <sub>.75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Faster R-CNN	ResNet-50	33.8	55.4	35.9	17.4	37.9	45.3
Grid R-CNN	ResNet-50	<b>35.9</b>	54.0	38.0	18.6	40.2	47.8
Faster R-CNN w FPN	ResNet-50	37.4	59.3	40.3	21.8	40.9	47.9
Grid R-CNN w FPN	ResNet-50	<b>39.6</b>	58.3	42.4	22.6	43.8	51.5
Faster R-CNN w FPN	ResNet-101	39.5	61.2	43.1	22.7	43.7	50.8
Grid R-CNN w FPN	ResNet-101	<b>41.3</b>	60.3	44.4	23.4	45.8	54.1

Grid R-CNN w FPN (ours) | ResNeXt-101 | 43.2 | 63.0 | 46.6 | 25.1 | 46.5 | 55.2





# Anchor-based

## One-Stage

# Life of 计算机图像实体检测

## One Stage

2016

YOLOv1  
SSD

2017

YOLOv2

2018

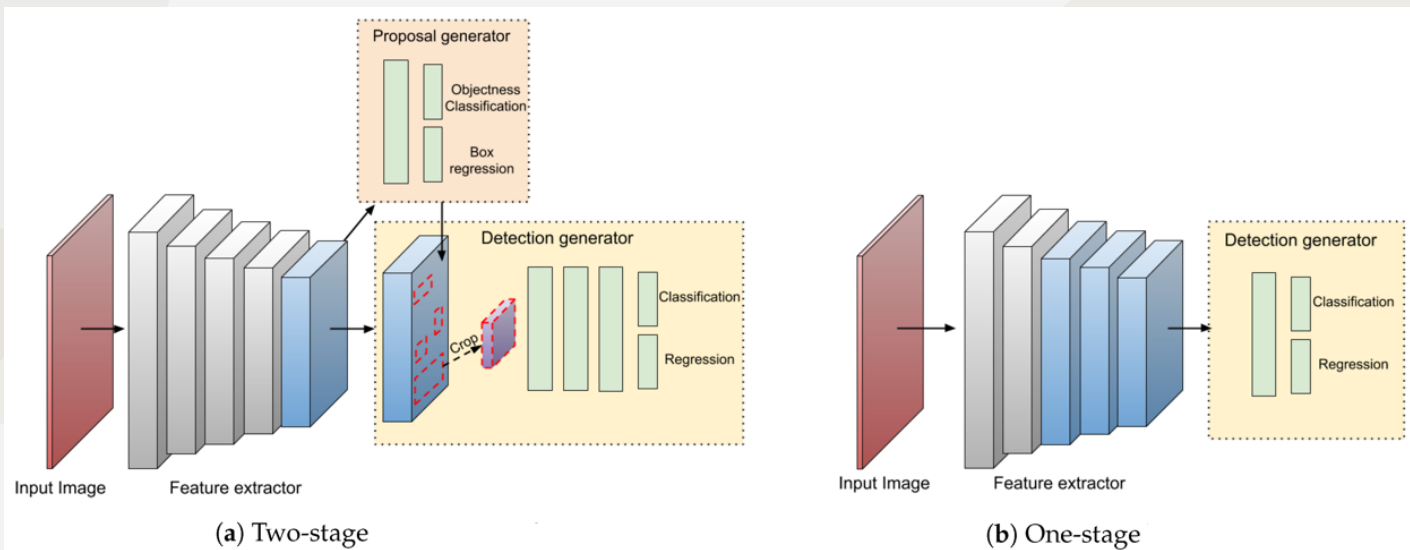
RetinaNet  
YOLOv3

2020

YOLOv4

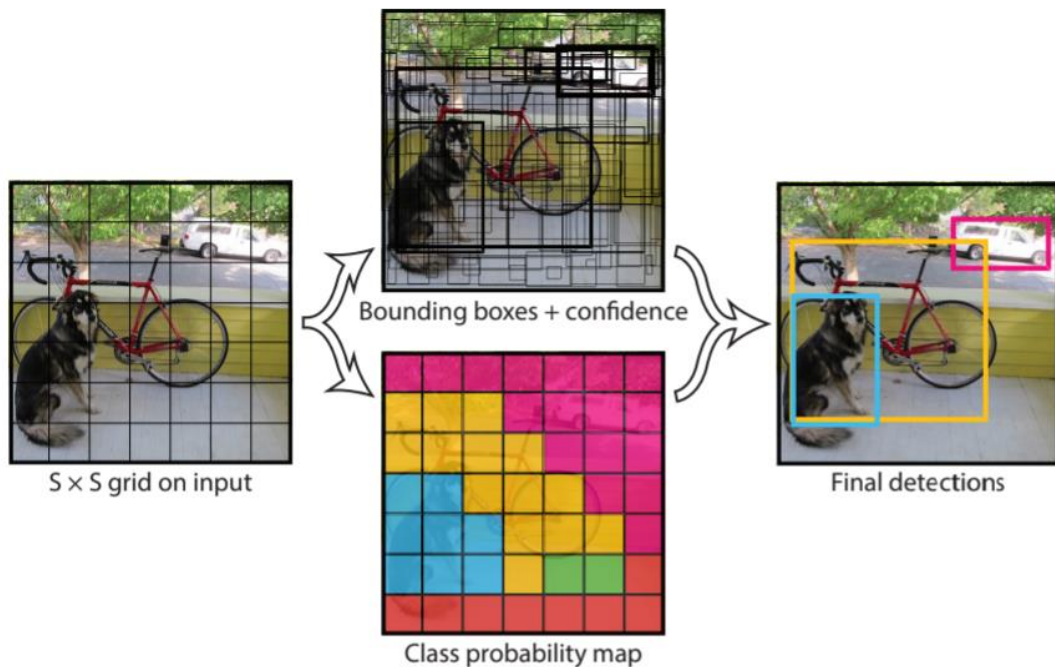
2021

YOLOv5



# Life of 计算机图像实体检测

YOU ONLY LOOK ONCE! 2016

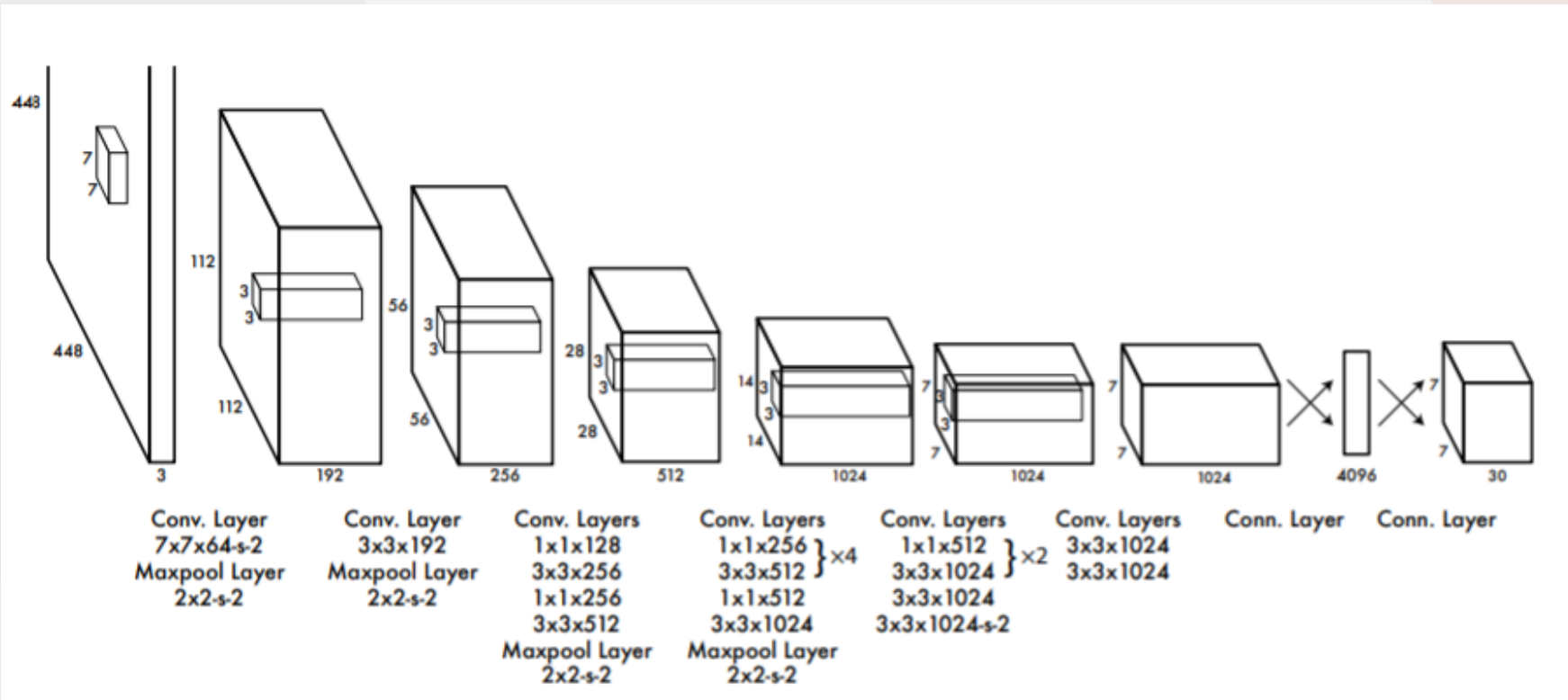


Grid Detection

Anchor-free?

# Life of 计算机图像实体检测

YOU ONLY LOOK ONCE! 2016



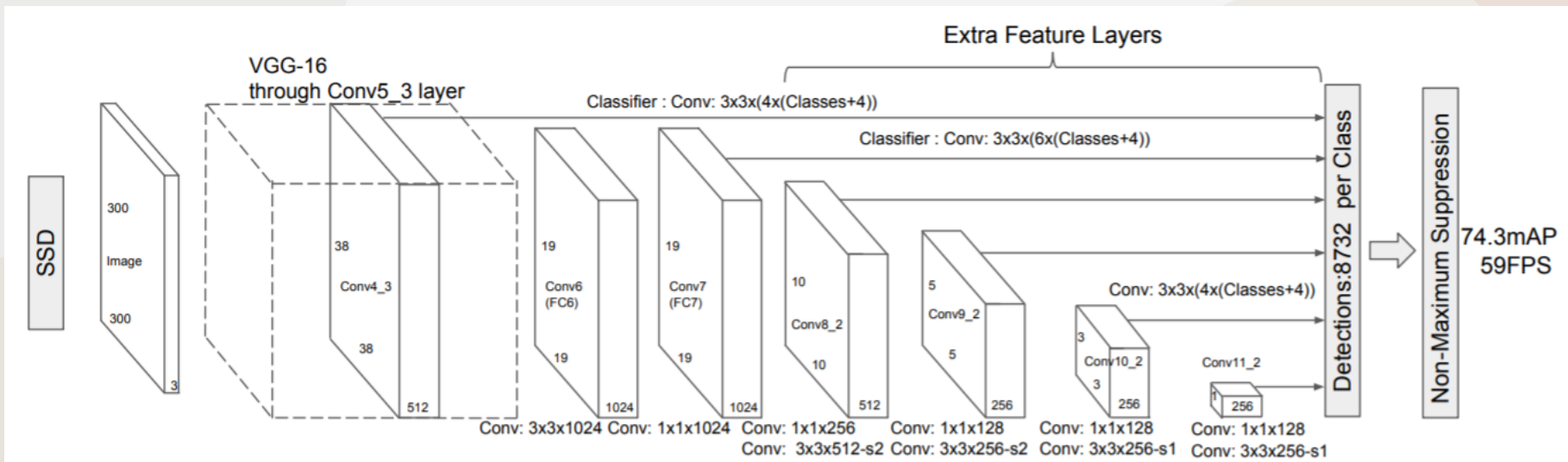
# Life of 计算机图像实体检测

YOU ONLY LOOK ONCE! 2016

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [31]	2007	16.0	100
30Hz DPM [31]	2007	26.1	30
Fast YOLO	2007+2012	52.7	<b>155</b>
YOLO	2007+2012	<b>63.4</b>	45
Less Than Real-Time			
Fastest DPM [38]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[28]	2007+2012	73.2	7
Faster R-CNN ZF [28]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

# Life of 计算机图像实体检测

## SSD: Single Shot MultiBox Detector 2016



# Life of 计算机图像实体检测

## SSD: Single Shot MultiBox Detector 2016

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	74.3	46	1	8732	300 × 300
SSD512	76.8	19	1	24564	512 × 512
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

# Life of 计算机图像实体检测

## YOLO9000: Better, Faster, Stronger 2017

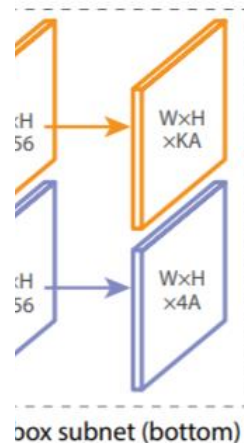
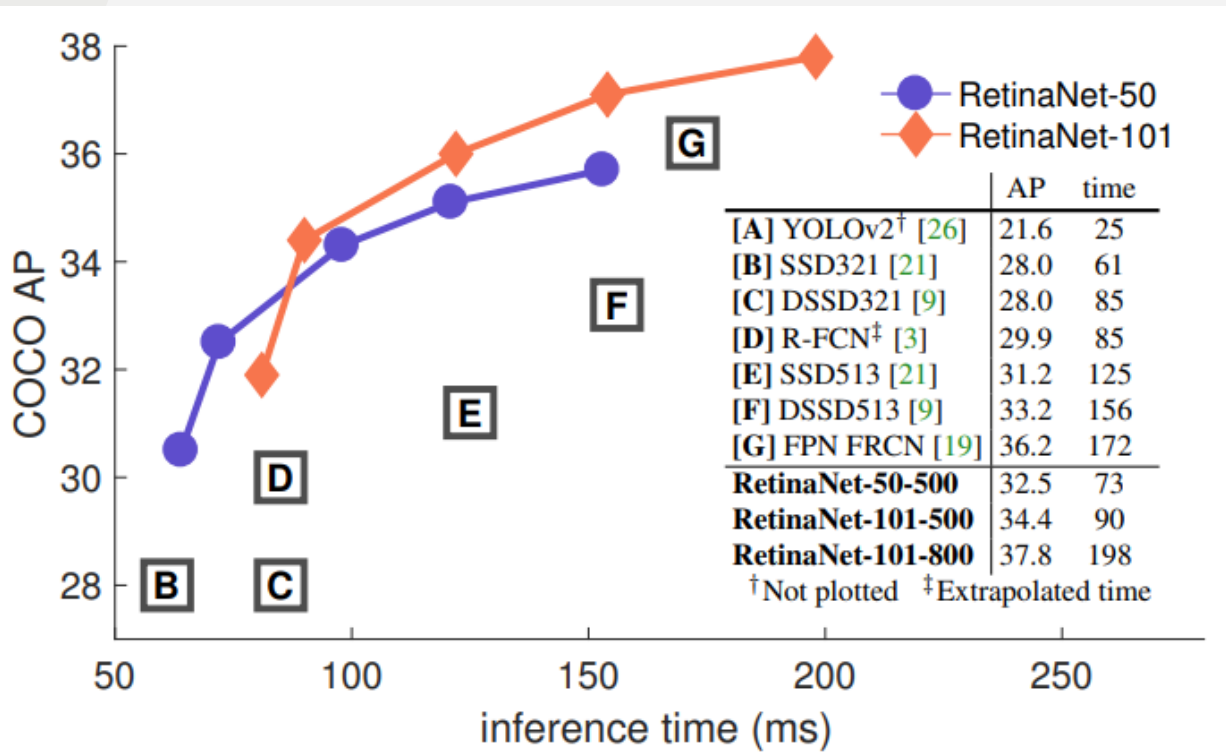
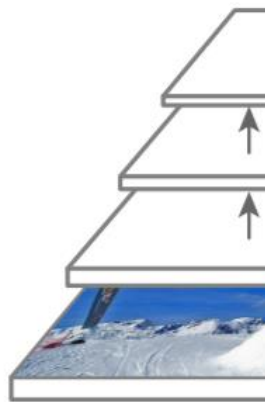
	Detection Frameworks	Train	mAP	FPS
batch r	Fast R-CNN [5]	2007+2012	70.0	0.5
hi-res class	Faster R-CNN VGG-16[15]	2007+2012	73.2	7
convoluti	Faster R-CNN ResNet[6]	2007+2012	76.4	5
anchor b	YOLO [14]	2007+2012	63.4	45
new netw	SSD300 [11]	2007+2012	74.3	46
dimension p	SSD500 [11]	2007+2012	76.8	19
location predic	YOLOv2 288 × 288	2007+2012	69.0	91
passthro	YOLOv2 352 × 352	2007+2012	73.7	81
multi-s	YOLOv2 416 × 416	2007+2012	76.8	67
hi-res deta	YOLOv2 480 × 480	2007+2012	77.8	59
VOC2007	YOLOv2 544 × 544	2007+2012	<b>78.6</b>	40

Hi-resclassifier...

Darknet19

# Life of 计算机图像实体检测

## Focal Loss for Dense Object Detection 2018



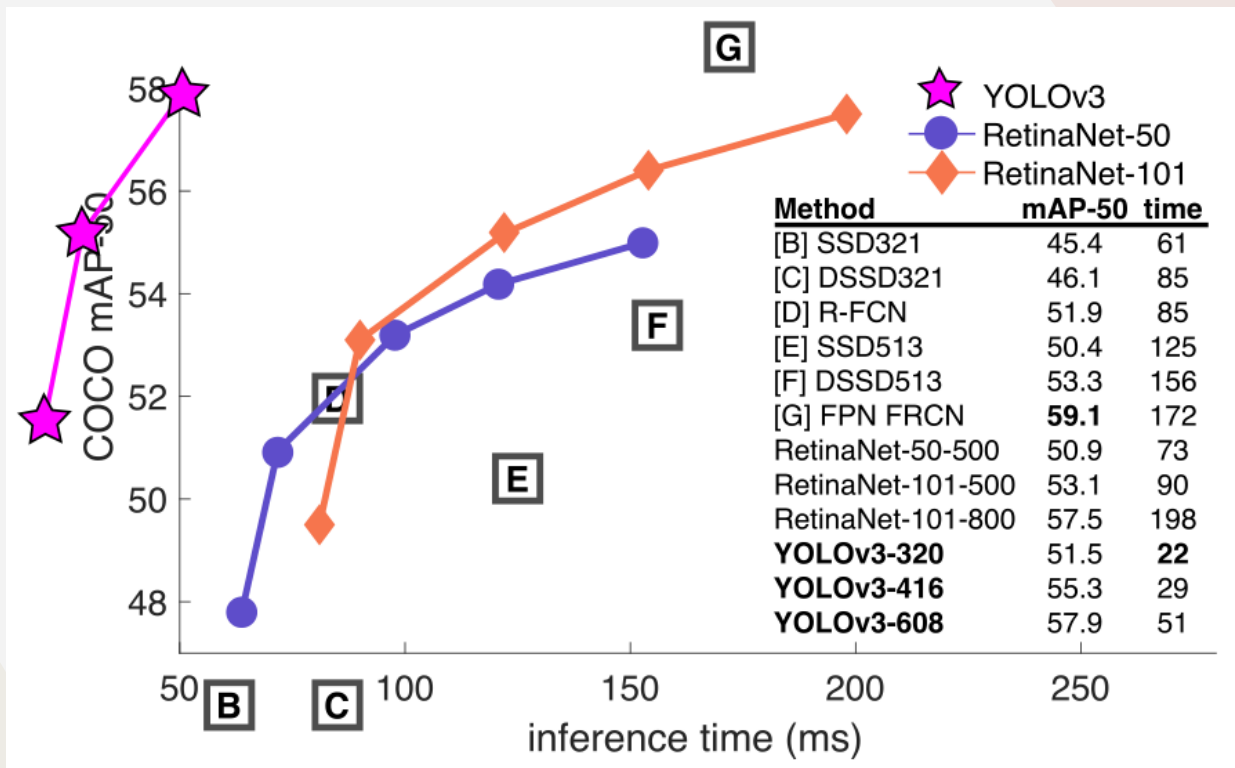
# Life of 计算机图像实体检测

## YOLOv3: An Incremental Improvement 2018

Darknet53

多尺度

忽略样本



# Life of 计算机图像实体检测

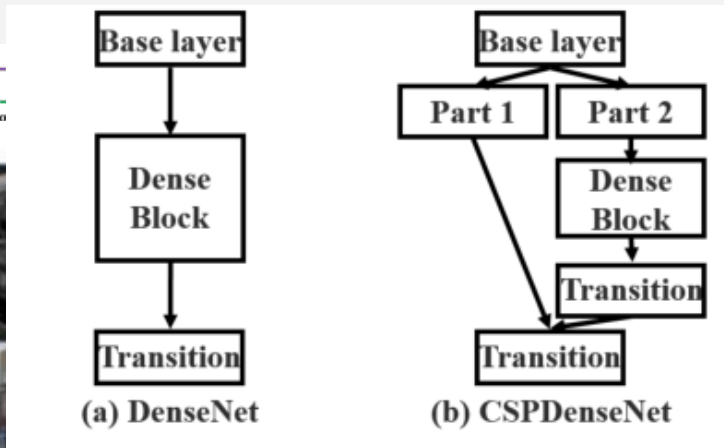
## YOLOv4: Optimal Speed and Accuracy of Object Detection 2020

Mosaic数据增强

CSPDarknet53

Bag-of-Freebies

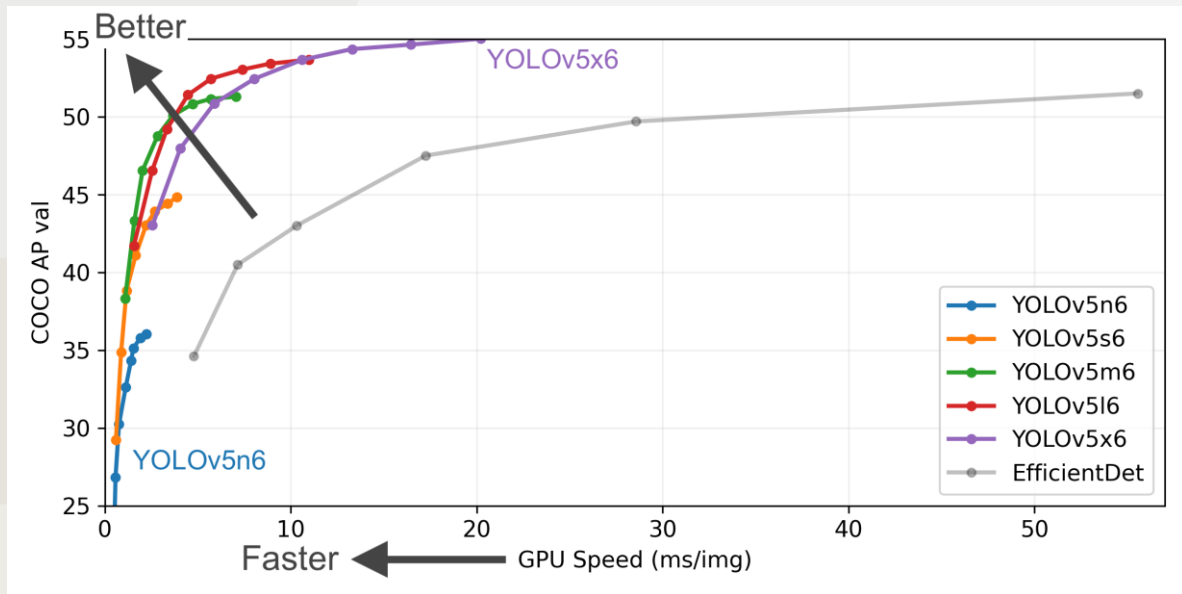
Bag-of-Special



Backbone	Top-1	Top-5	Bn Ops	BFLOP/s	FPS
Darknet-19 [15]	74.1	91.8	7.29	1246	<b>171</b>
ResNet-101[5]	77.1	93.7	19.7	1039	53
ResNet-152 [5]	<b>77.6</b>	<b>93.8</b>	29.4	1090	37
Darknet-53	77.2	<b>93.8</b>	18.7	<b>1457</b>	78

# Life of 计算机图像实体检测

## YOLOv5



Mosaic数据增强

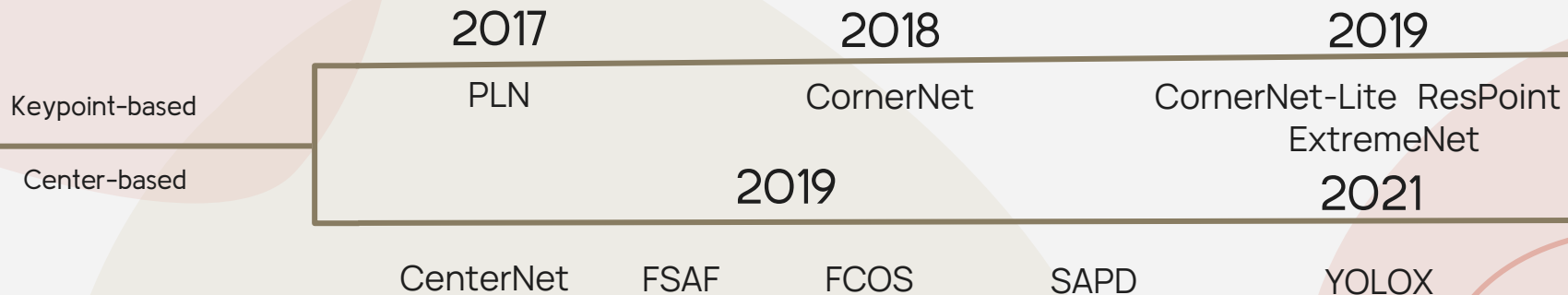
自适应锚框计算

自适应图片缩放

Focus结构

CSP结构.....

# Anchor-free





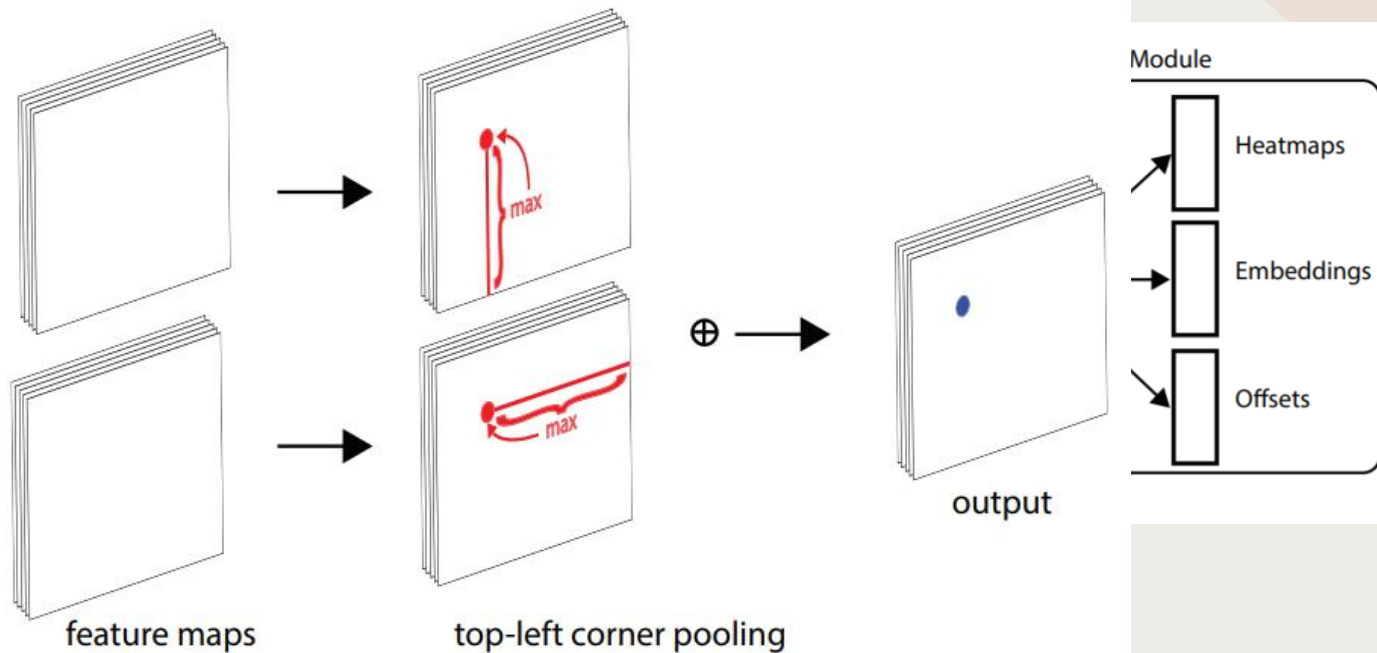
# Anchor-free

## Keypoint-based

# Life of 计算机图像实体检测

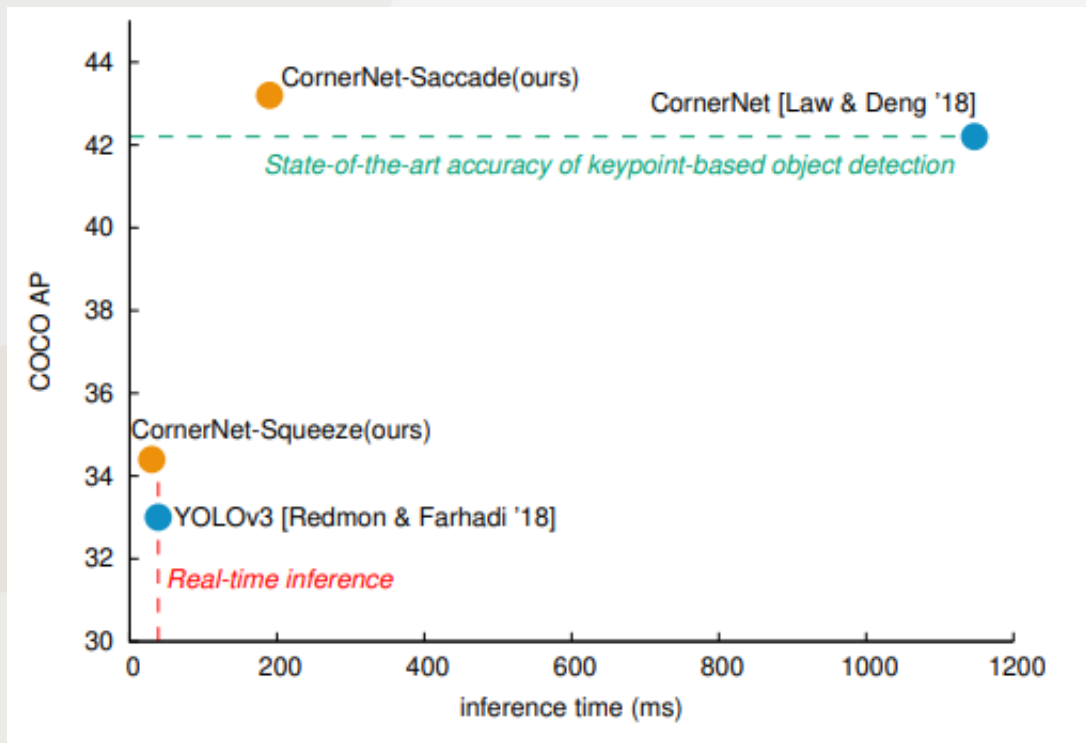
Anchor-free Keypoint-based

CornerNet: Detecting Objects as Paired Keypoints 2018



# Life of 计算机图像实体检测

## CornerNet-Lite: Efficient Keypoint Based Object Detection 2019

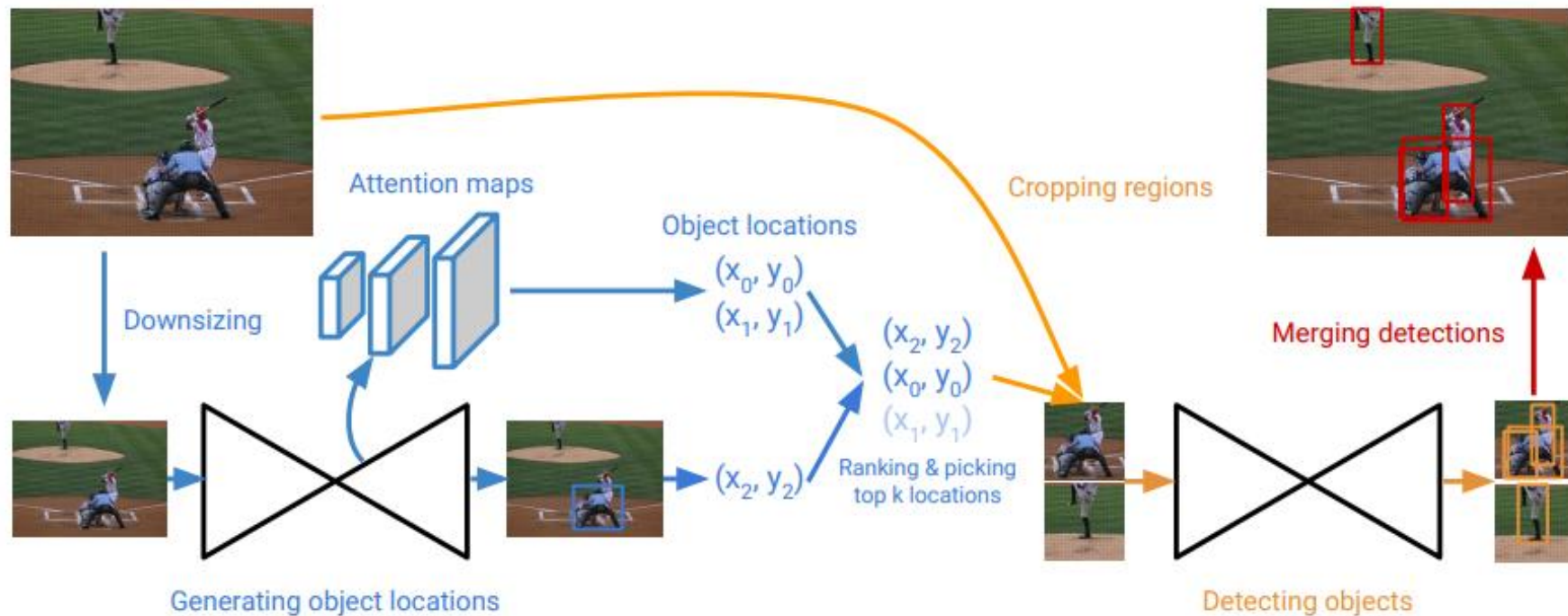


CornerNet-Saccade

CornerNet-Squeeze

# Life of 计算机图像实体检测

## CornerNet-Saccade



# Life of 计算机图像实体检测

## CornerNet-Squeeze

SqueezeNet

MobileNet

Input	Operator	Output
Residual block in CornerNet		
$h \times w \times k$	$3 \times 3$ Conv, ReLU	$h \times w \times k'$
$h \times w \times k'$	$3 \times 3$ Conv, ReLU	$h \times w \times k'$
Fire module in CornerNet-Squeeze		
$h \times w \times k$	$1 \times 1$ Conv	$h \times w \times \frac{k'}{2}$
$h \times w \times \frac{k'}{2}$	$1 \times 1$ Conv + $3 \times 3$ Dwise, ReLU	$h \times w \times k'$

Comparison between the residual block and the new fire module.

# Life of 计算机图像实体检测

## Anchor-free Keypoint-based

Point

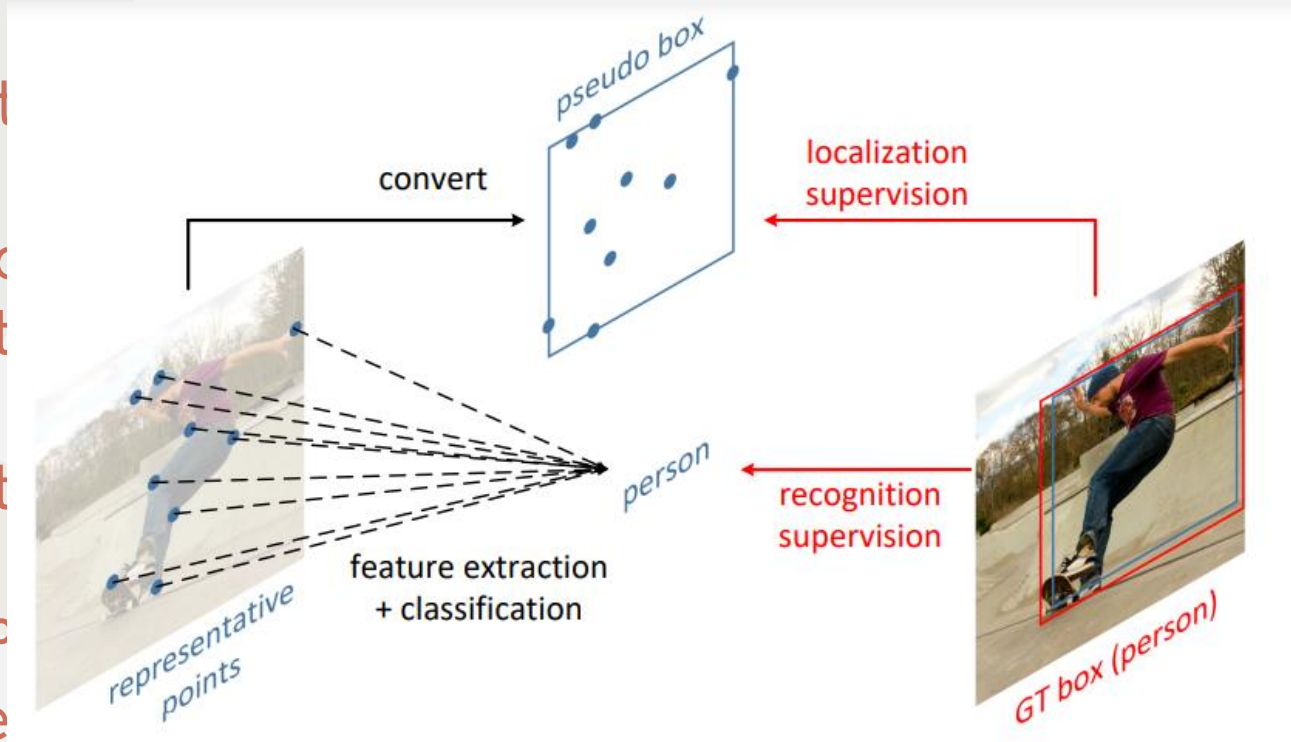
Bottom

Center

Center

RepF

Dete



e and

2019



# Anchor-free

## Center-based

# Life of 计算机图像实体检测

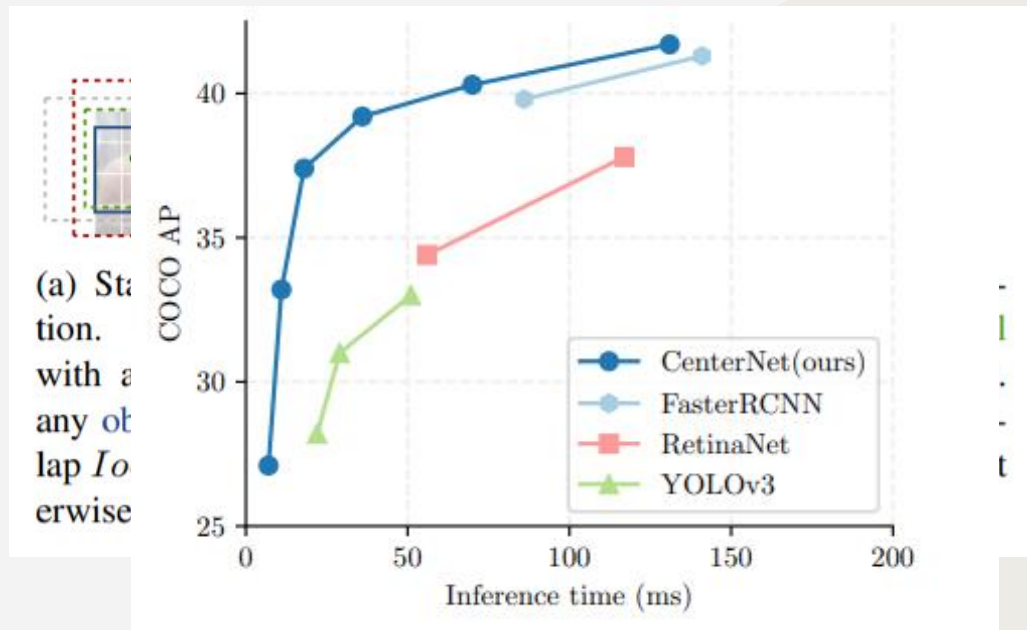
## Anchor-free Center-based

## CenterNet: Objects as Points 2019

前后景分类 ×

NMS ×

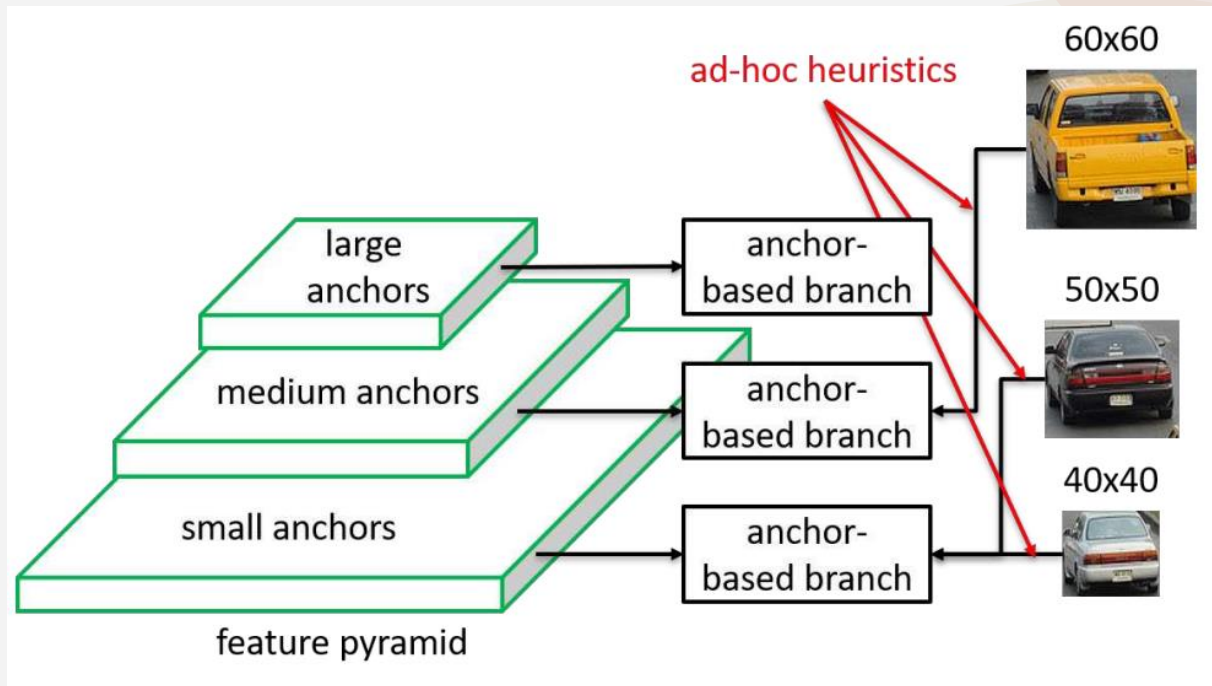
FPN ×



# Life of 计算机图像实体检测

## Feature Selective Anchor-Free Module for Single-Shot Object Detection 2019

人为设定的特征水平？



# Life of 计算机图像实体检测

## Feature Selective Anchor-Free Module for Single-Shot Object Detection 2019

	anchor-						
CornerNet511 [17] (single-scale)	Hourglass-104	40.5	56.5	43.1	19.4	42.7	53.9
CornerNet [17] (multi-scale)		42.1	57.8	45.3	20.8	44.8	56.7
GHM800 [18]	ResNeXt-101	41.6	62.8	44.2	22.3	45.1	<b>55.3</b>
<b>Ours800</b> (single-scale)		42.9	63.8	46.3	26.6	46.2	52.7
<b>Ours</b> (multi-scale)		<b>44.6</b>	<b>65.2</b>	<b>48.6</b>	<b>29.7</b>	<b>47.1</b>	54.6

feature pyramid

anchor-based branch

anchor-free branch

classification

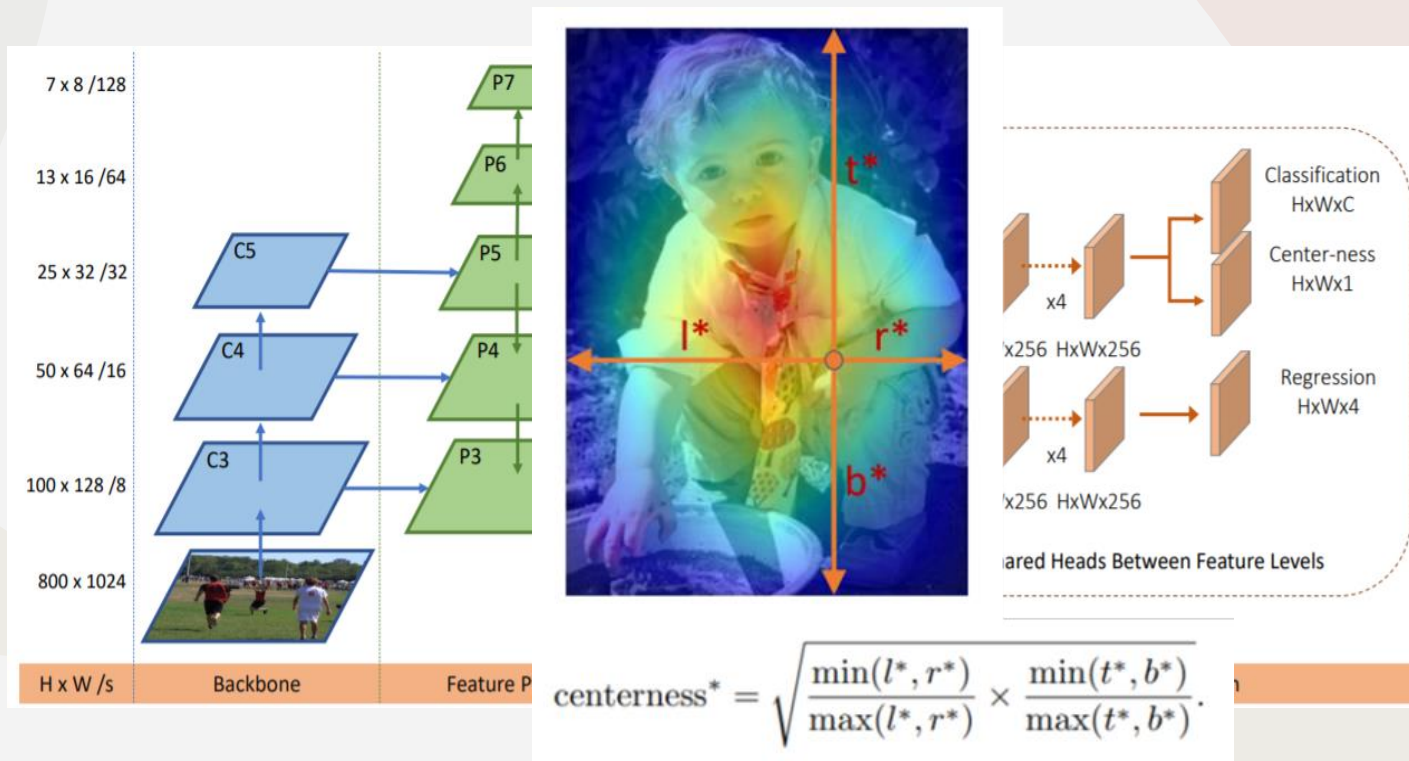
FSAF module

# Life of 计算机图像实体检测

## FCOS: Fully Convolutional One-Stage Object Detection 2019

逐像素回归预测

center-ness



# Life of 计算机图像实体检测

## FCOS: Fully Convolutional One-Stage Object Detection 2019

Method	Backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Two-stage methods:							
Faster R-CNN w/ FPN [14]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [11]	Inception-ResNet-v2 [27]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w/ TDM [25]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	52.1
One-stage methods:							
YOLOv2 [22]	DarkNet-19 [22]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [18]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [5]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
RetinaNet [15]	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
CornerNet [13]	Hourglass-104	40.5	56.5	43.1	19.4	42.7	53.9
FSAF [34]	ResNeXt-64x4d-101-FPN	42.9	63.8	46.3	26.6	46.2	52.7
FCOS	ResNet-101-FPN	41.5	60.7	45.0	24.4	44.8	51.6
FCOS	HRNet-W32-51 [26]	42.0	60.4	45.3	25.4	45.0	51.0
FCOS	ResNeXt-32x8d-101-FPN	42.7	62.2	46.1	26.0	45.6	52.6
FCOS	ResNeXt-64x4d-101-FPN	43.2	62.8	46.6	26.5	46.2	53.3
FCOS w/ improvements	ResNeXt-64x4d-101-FPN	<b>44.7</b>	<b>64.1</b>	<b>48.4</b>	<b>27.6</b>	<b>47.5</b>	<b>55.6</b>

# Life of 计算机图像实体检测

## YOLOX: Exceeding YOLO Series in 2021

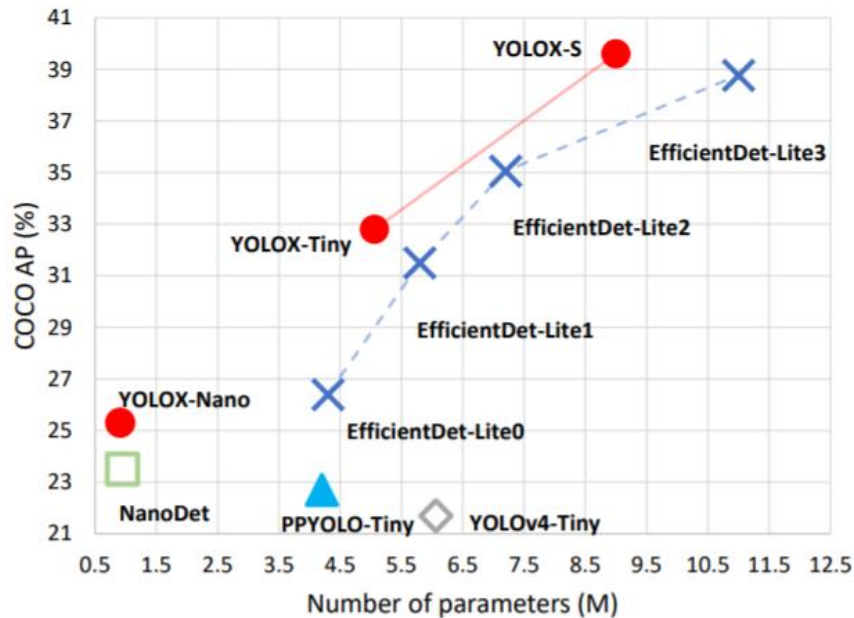
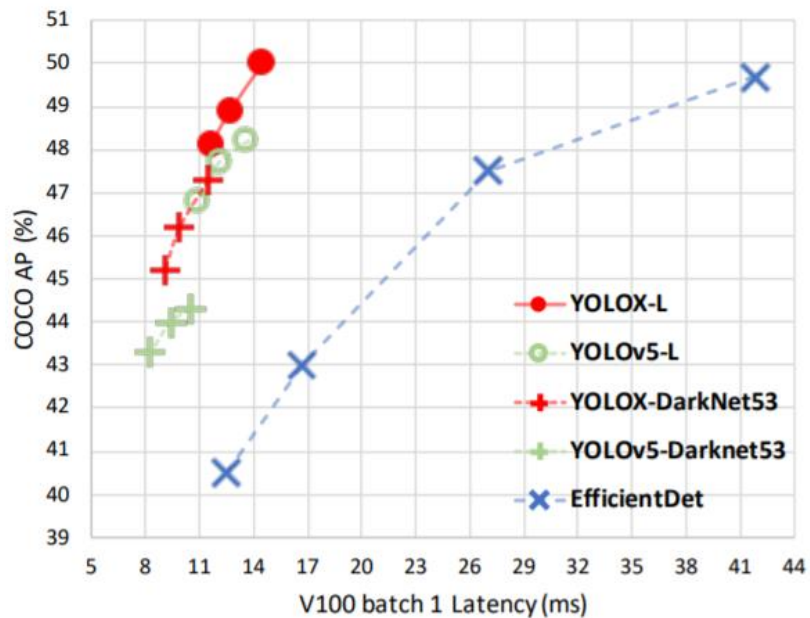


Figure 1: Speed-accuracy trade-off of accurate models (top) and Size-accuracy curve of lite models on mobile devices (bottom) for YOLOX and other state-of-the-art object detectors.

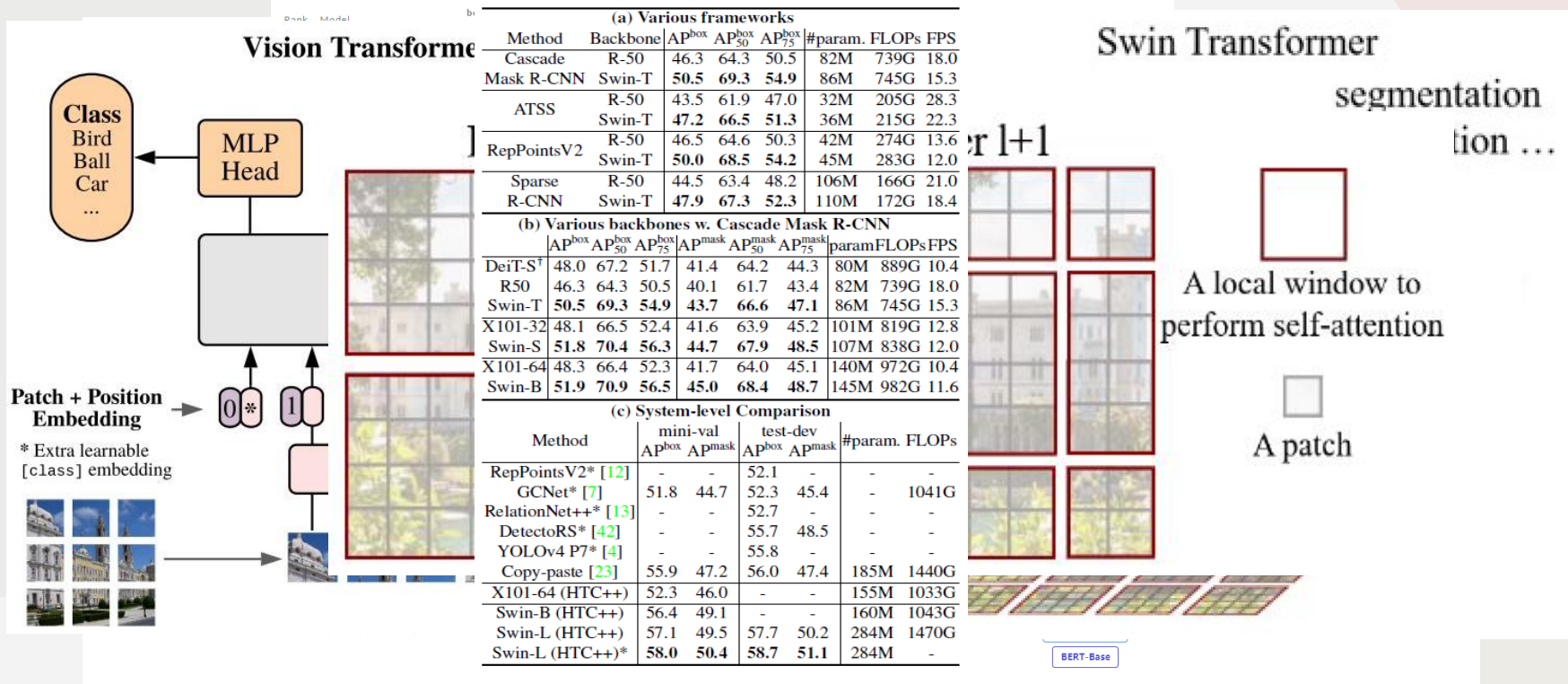


新进展

Transformer

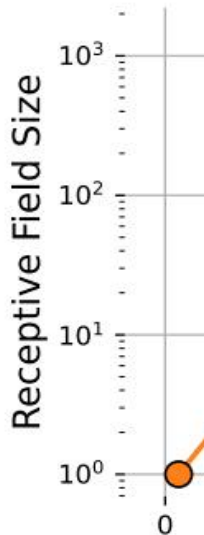
# Life of 计算机图像实体检测

## 新进展——Vision Transformer、Swin Transformer

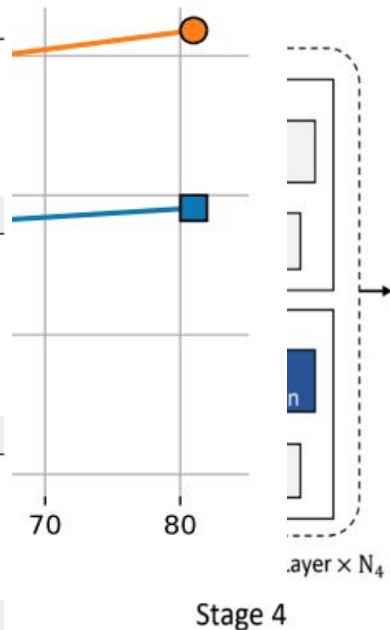


# Life of 计算机图像实体检测

## 新进展——Focal Transformer



Backbone	RetinaNet	Mask R-CNN	
	$AP^b$	$AP^b$	$AP^m$
ResNet-50 [34]	36.3	38.0	34.4
PVT-Small	40.4	40.4	37.8
ViL-Small [80]	41.6	41.8	38.5
Swin-Tiny [44]	42.0	43.7	39.8
<b>Focal-Tiny (Ours)</b>	<b>43.7 (+1.7)</b>	<b>44.8 (+1.1)</b>	<b>41.0 (+1.3)</b>
ResNet-101 [34]	38.5	40.4	36.4
ResNeXt101-32x4d [70]	39.9	41.9	37.5
PVT-Medium [63]	41.9	42.0	39.0
ViL-Medium [80]	42.9	43.4	39.7
Swin-Small [44]	45.0	46.5	42.1
<b>Focal-Small (Ours)</b>	<b>45.6 (+0.6)</b>	<b>47.4 (+0.9)</b>	<b>42.8 (+0.7)</b>
ResNeXt101-64x4d [70]	41.0	42.8	38.4
PVT-Large [63]	42.6	42.9	39.5
ViL-Base [80]	44.3	45.1	41.0
Swin-Base [44]	45.0	46.9	42.3
<b>Focal-Base (Ours)</b>	<b>46.3 (+1.3)</b>	<b>47.8 (+0.9)</b>	<b>43.2 (+0.9)</b>





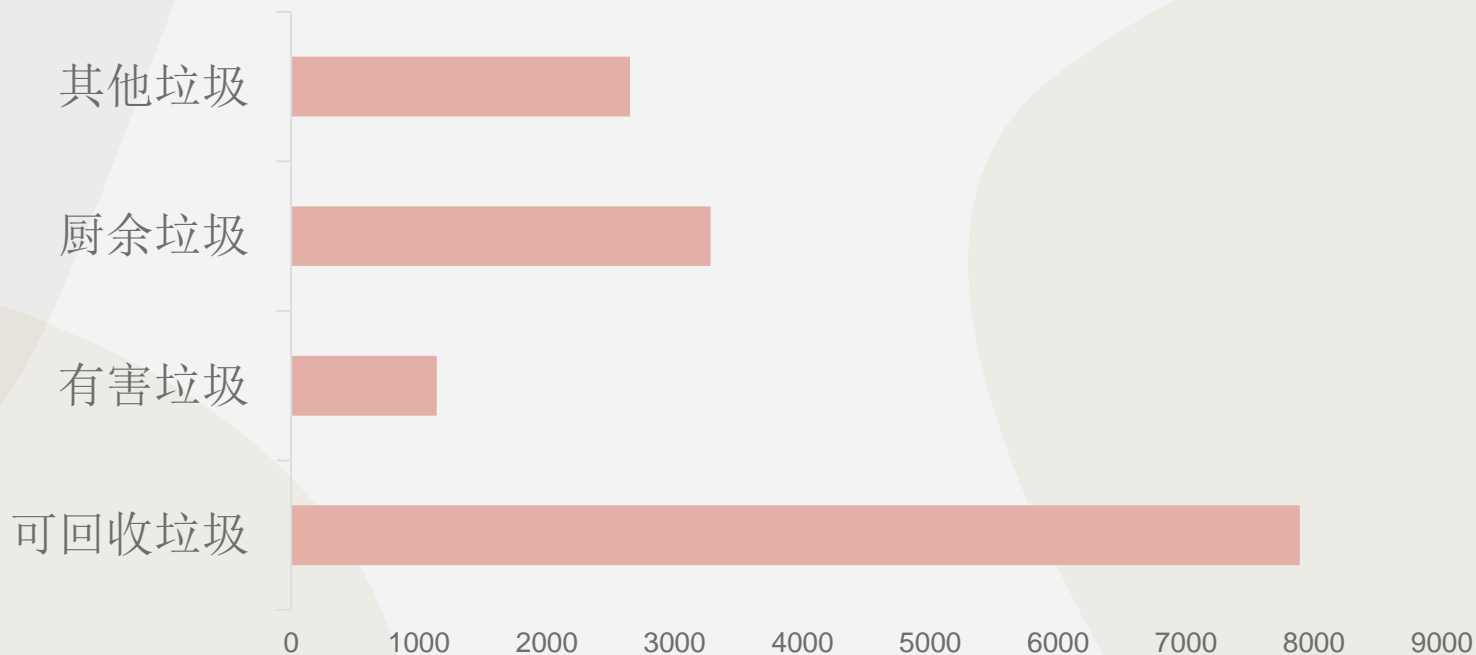
# 项目展示

## 垃圾分类检测Demo

# Demo of 计算机图像实体检测

The Trash Dataset of Our Demo

## 四类垃圾分布



# Demo of 计算机图像实体检测

## The Trash Dataset of Our Demo

可回收垃圾



# Demo of 计算机图像实体检测

## The Trash Dataset of Our Demo



# Demo of 计算机图像实体检测

## The Models of Our Demo

Val Datasets

AP(%)

One-Stage YOLOv5-s

36.7

Two-Stage Faster R-CNN + FPN

COCO

35.9

Anchor-free YOLOX-s

39.6

	Val Datasets	AP(%)
One-Stage YOLOv5-s	COCO	36.7
Two-Stage Faster R-CNN + FPN		35.9
Anchor-free YOLOX-s		39.6

# Demo of 计算机图像实体检测

```
文件(E) 编辑(E) 视图(V) 导航(N) 代码(C) 重构(R) 运行(U) 工具(I) Git(G) 窗口(W) 帮助(H) predict.py - predict.py
2dashujijushu > deep-learning-for-image-processing-master > deep-learning-for-image-processing-master > pytorch_object_detection > faster_rcnn > predict.py
项目
predict.py x demo.py x detect.py x
56 assert os.path.exists(train_weights), "{} file dose not exist.".format(train_weights)
57 model.load_state_dict(torch.load(train_weights, map_location=device)["model"])
58 model.to(device)
59
60 # read class_indict
61 label_json_path = './pascal_voc_classes.json'
62 assert os.path.exists(label_json_path), "json file {} dose not exist.".format(label_json_path)
63 json_file = open(label_json_path, 'r', encoding='UTF-8')
64 class_dict = json.load(json_file)
65 json_file.close()
66 category_index = {v: k for k, v in class_dict.items()}
67
68 # load image
69 original_img = Image.open("./test1.jpg")
70
71 # from pil image to tensor, do not normalize image
72 data_transform = transforms.Compose([transforms.ToTensor()])
73 img = data_transform(original_img)
74 # expand batch dimension
75 img = torch.unsqueeze(img, dim=0)
76
77 model.eval() # 进入验证模式
78 with torch.no_grad():
79     # init
80     img_height, img_width = img.shape[-2:]
81     init_img = torch.zeros((1, 3, img_height, img_width), device=device)
82     model(init_img)
83
84     t_start = time_synchronized()
85     predictions = model(img.to(device))[0]
86     t_end = time_synchronized()
main()
```

4 1 4

Fast Request

结构

Git TODO 问题 终端 Python Packages Python 控制台

72:65 LF UTF-8 4个空格 Python 3.8 (pytorch) (2) master

## Faster R-CNN with FPN

# Demo of 计算机图像实体检测

## The Models of Our Demo

YOLOX

Faster R-CNN

YOLOv5



可回收垃圾63%



可回收垃圾86.2%



可回收垃圾73%

# Demo of 计算机图像实体检测

```
video.py - video.py
yolov5-5.0-trash | video.py
项目
提交
Pull Request
结构
收藏类
1 import time
2 import cv2
3 import numpy as np
4 import torch
5 from models.experimental import attempt_load
6 from utils.datasets import letterbox
7 from utils.general import check_img_size, non_max_suppression, scale_coords, xyxy2xywh, set_logging, check_requirements
8 from utils.plots import colors, plot_one_box
9 from utils.torch_utils import select_device, time_synchronized
10
11 trash_classification = ["可回收垃圾", "有害垃圾", "厨余垃圾", "其他垃圾"]
12 trash_names=[3, 0, 1, 2, 0, 0, 0, 0, 0, 2,
13              1, 0, 0, 0, 0, 0, 2, 0, 3, 3,
14              3, 3, 3, 0, 0, 0, 0, 3, 2, 2,
15              2, 0, 3, 1, 0, 0, 0, 0, 0, 3,
16              0, 0, 0, 2]
17
18 @torch.no_grad()
19 def detect(
20     weights='runs/train/base/weights/best.pt', # 训练好的模型路径 (必改)
21     imgsz=640, # 训练模型设置的尺寸 (必改)
22     cap=0, # 摄像头
23     conf_thres=0.25, # 置信度
24     iou_thres=0.45, # NMS IOU 阈值
25     max_det=1000, # 最大检测的目标数
26     device='', # 设备
27     crop=True, # 显示预测框
28     classes=None, # 种类
29     agnostic_nms=False, # class-agnostic NMS
30     augment=False, # 是否扩充推理
31     half=False, # 使用CPU半精度推理

```

1:1 CRLF UTF-8 4个空格 Python 3.8 (pytorch) (2) master

# Demo of 计算机图像实体检测

## The Models of Our Demo

	mAP	Model	Inference Time
One-Stage YOLOv5-s	80.3%	13.9MB	27ms
Two-Stage Faster R-CNN + FPN	53.9%	316.7MB	108ms
Anchor-free YOLOX-s	71.4%	68.6MB	205ms

Tested on RTX1080

# 04 未来展望

Future

# Challenge of 计算机图像实体检测



· 实体多样

光照、视角、类间 ...

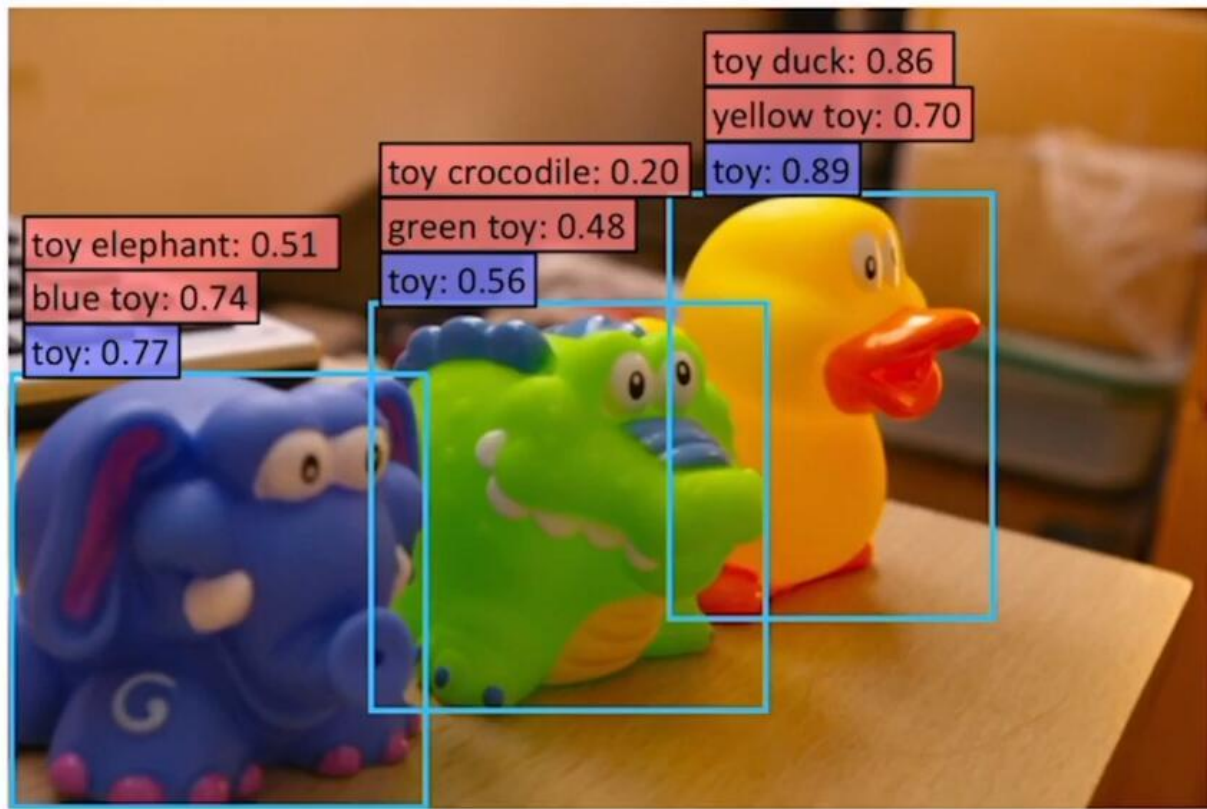


实体

湖图像

# Future of 计算机图像实体检测

- 精
- 自
- 弱



**THANKS!**

