

UNIVERSITAT ROVIRA I VIRGILI (URV) Y UNIVERSITAT OBERTA DE CATALUNYA (UOC)

Computational Engineering and Mathematics

Master's Thesis

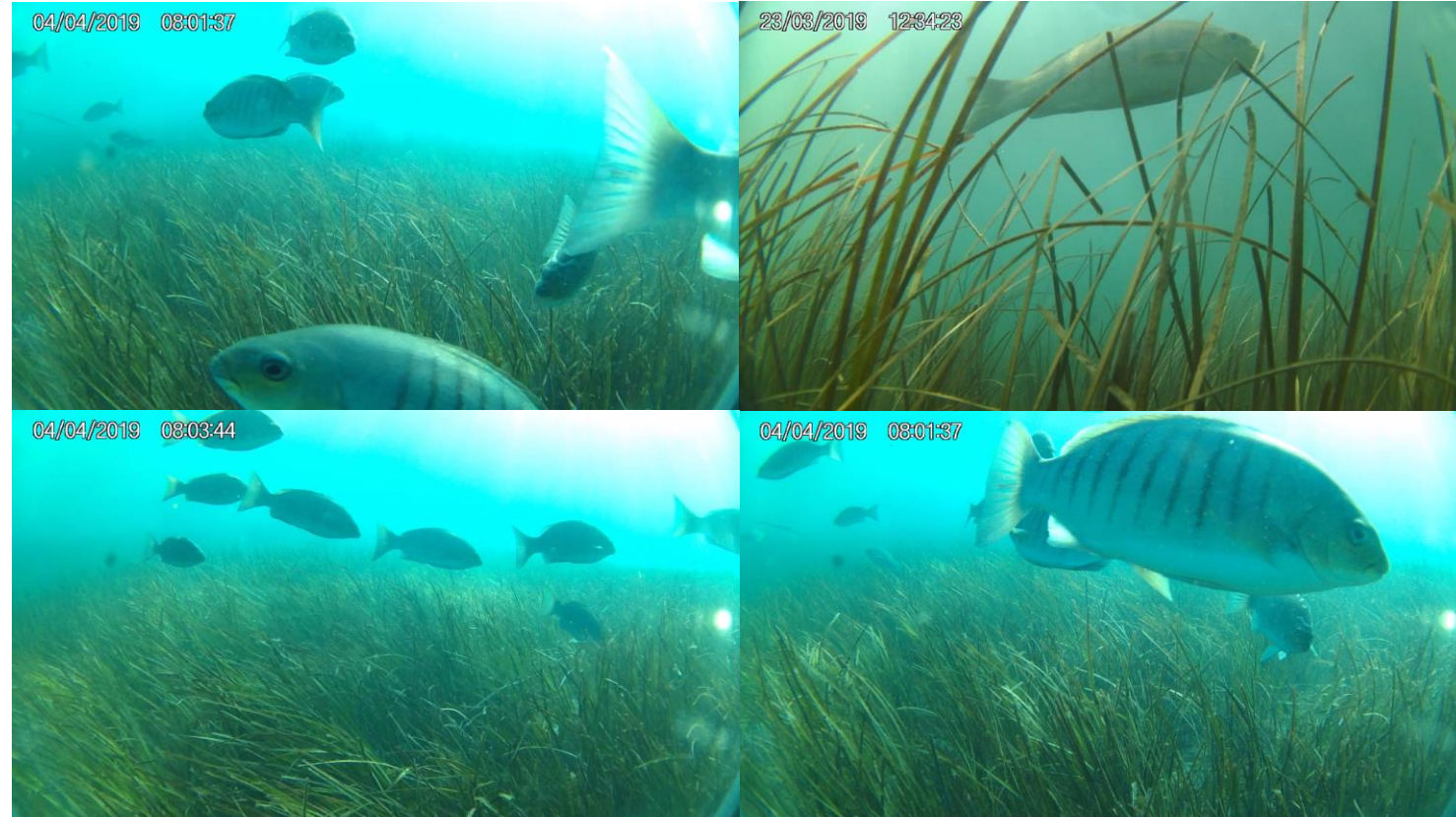
Area: Artificial Intelligence

Underwater Species Detection in Images and Videos

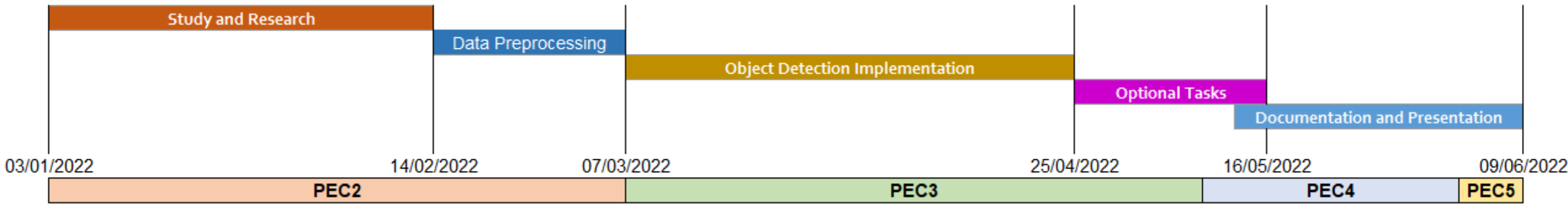
Author:	Ricard Fos Serdà
Consultant teacher:	Antonio Burguera Burguera
Submission date:	01/06/2022

Introduction

- **Deep Learning:**
 - Classify images
 - Detect Objects
 - Natural Language Processing
 - Motion Planning
 - More...
- **Underwater Species:**
 - Identify
 - Count
 - Study
 - Follow



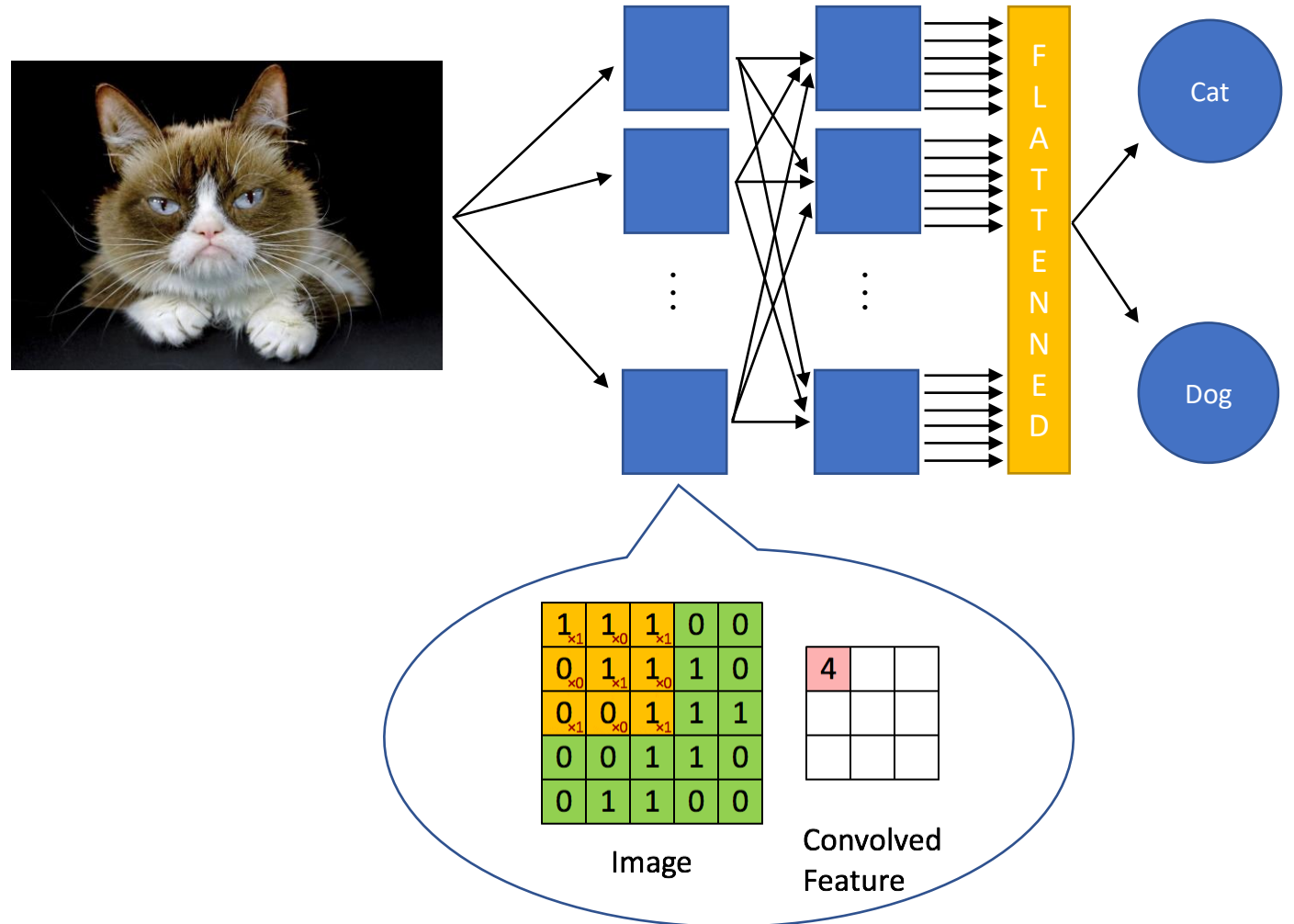
Work Planning



Convolutional Neural Networks (CNN)

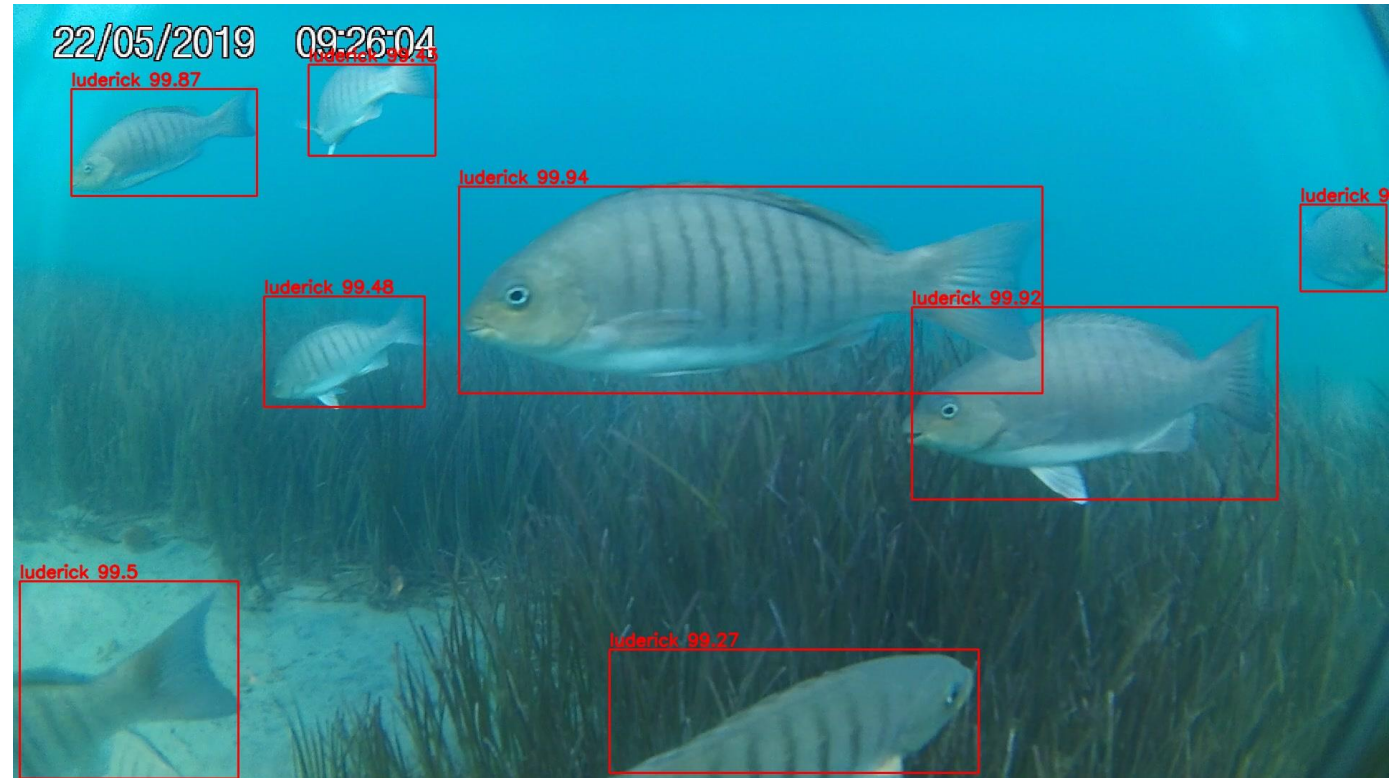
- **Image Classification:**

1. Convolutional Layers
2. Pooling Layers
3. Flatten
4. Classify



Two-Stage Object Detection

- Bounding Boxes
 - [xmin, ymin, xmax, ymax]
- Class label
- Confidence



Region-Based CNN (R-CNN)

1. Region proposals

- Selective Search

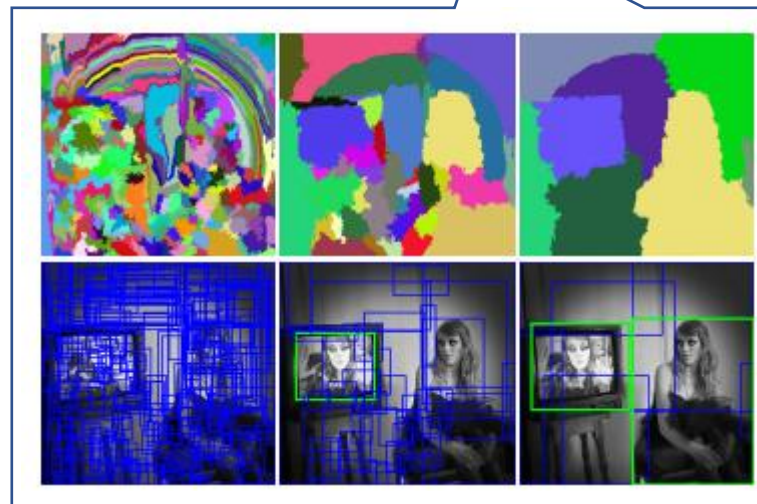
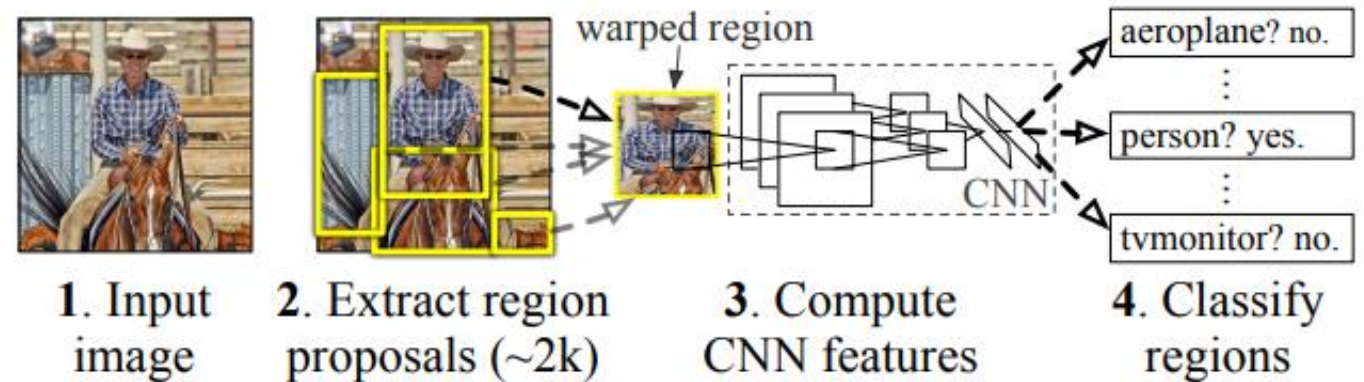
2. CNN Backbone

3. Classify

4. Box regression

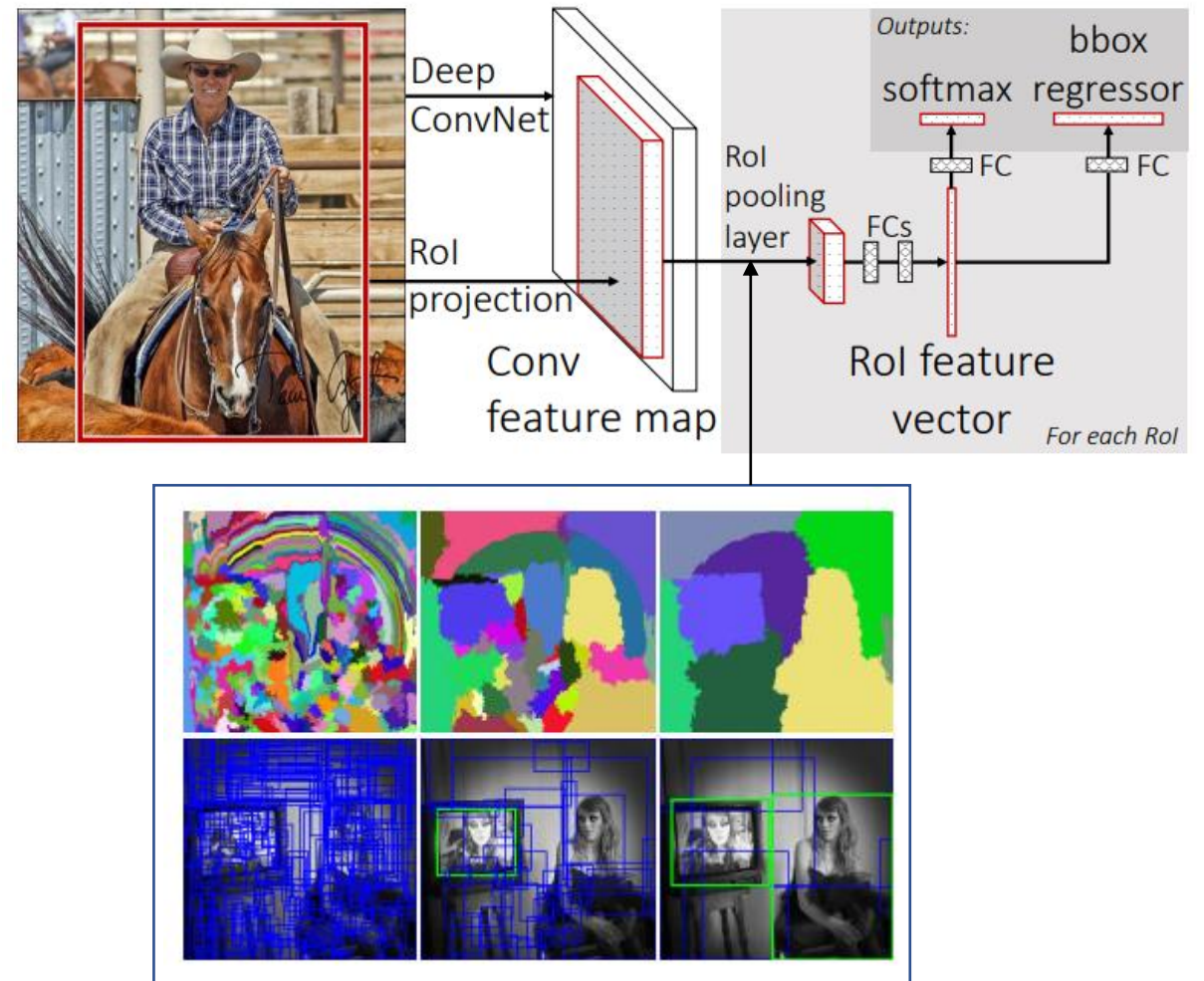
5. Non-Max Suppresion

R-CNN: *Regions with CNN features*



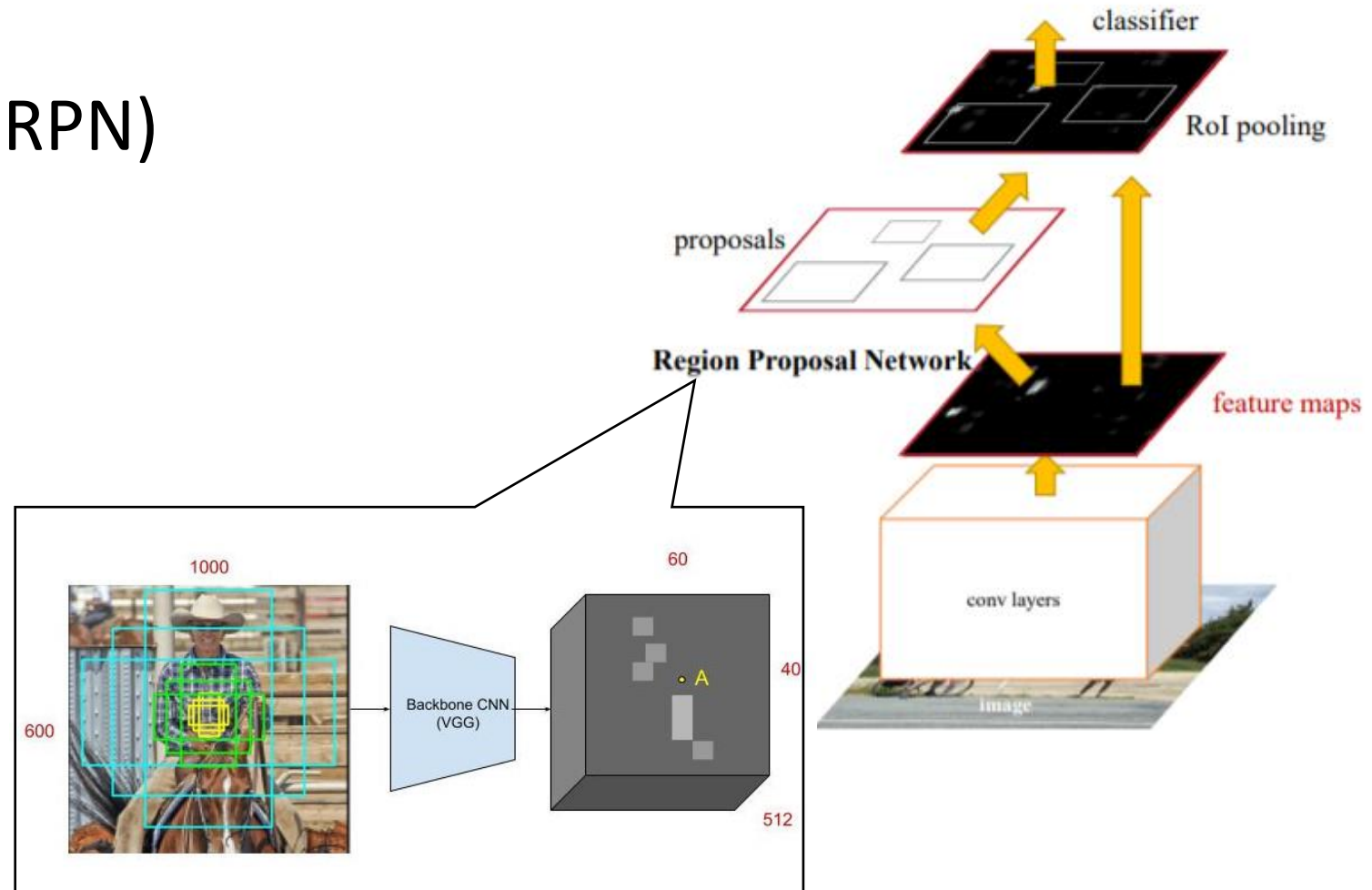
Fast R-CNN

- CNN backbone used once
- Region proposals
 - Selective Search
 - Feature Map
- Region of Interest Pooling
 - Fixed-Size ROIs




Faster R-CNN

- Region Proposal Network (RPN)
 - Anchor Boxes
 - Objectness score
 - Generate Region Proposals



Evaluation Metrics for Object Detection

- Intersection over Union (IoU)

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


- Precision $Precision(c) = \frac{TPc}{TPc + FPc}$

- Recall $Recall(c) = \frac{TPc}{TPc + FNC}$

- Mean Average Precision (mAP)

- Area under the Precision-Recall Curve

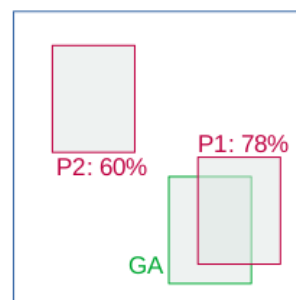


Image 1

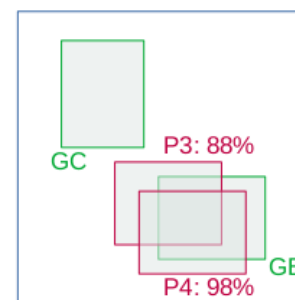


Image 2

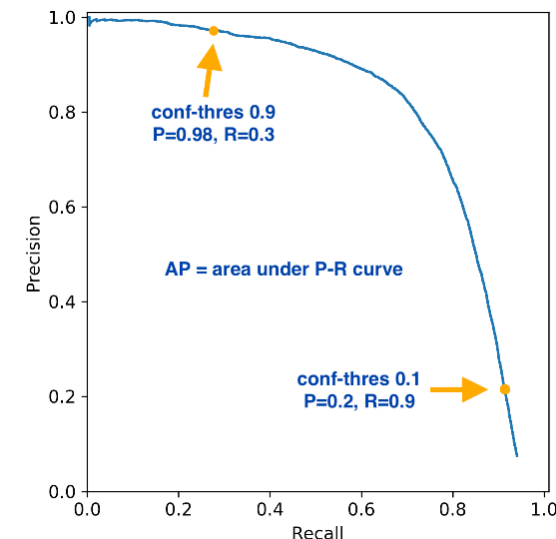
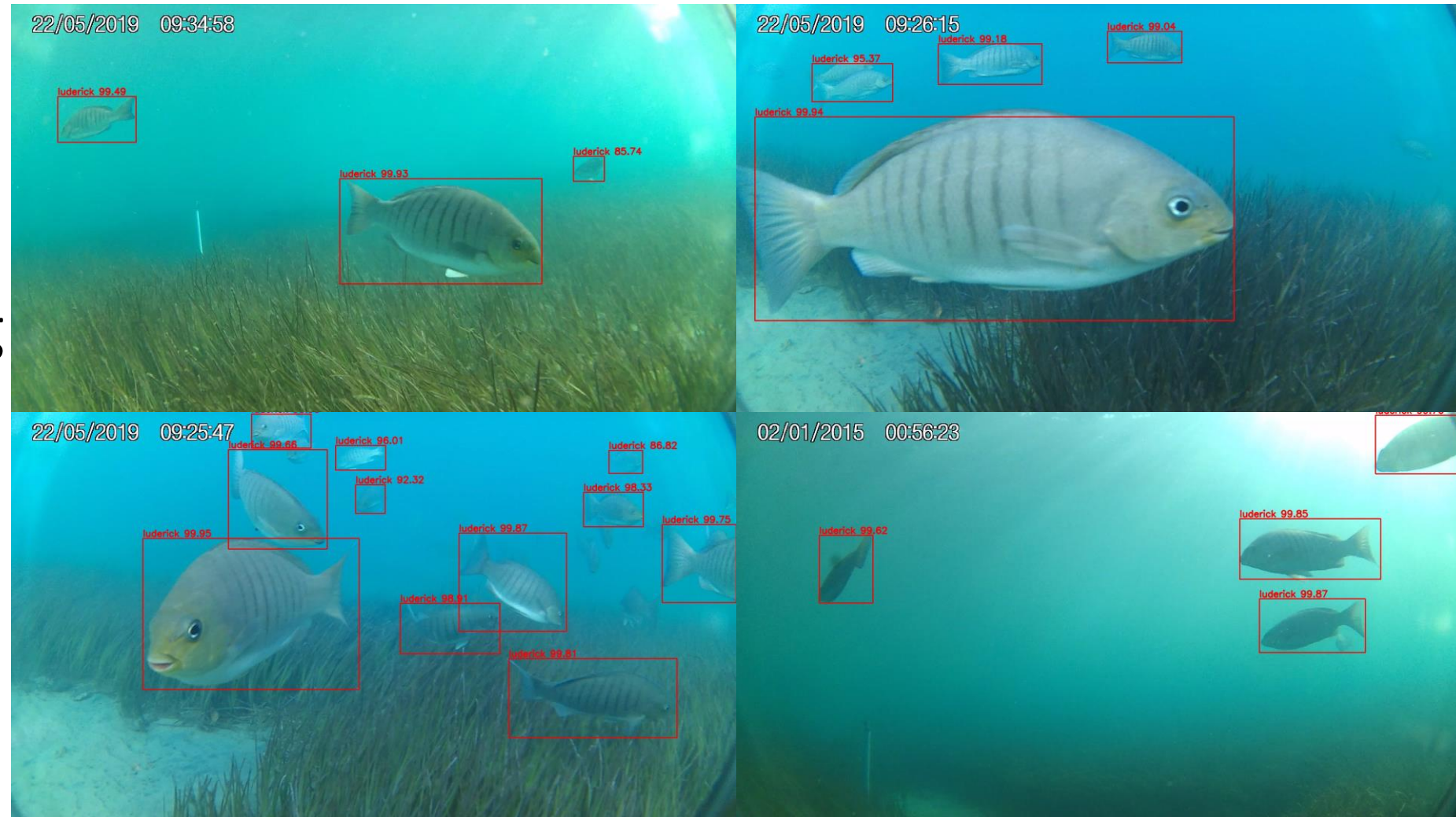


Image	Detection	Confidence	IoU	Ground Truth	TP/FP	Acc TP	Acc FP	Precision	Recall
Image 2	P4	98%	> 0.5	GB	TP	1	0	1	0.33
Image 2	P3	88%	> 0.5	GB	FP	1	1	0.5	0.33
Image 1	P1	78%	> 0.5	GA	TP	2	1	0.67	0.67
Image 1	P2	60%	< 0.5	-	FP	2	2	0.5	0.67

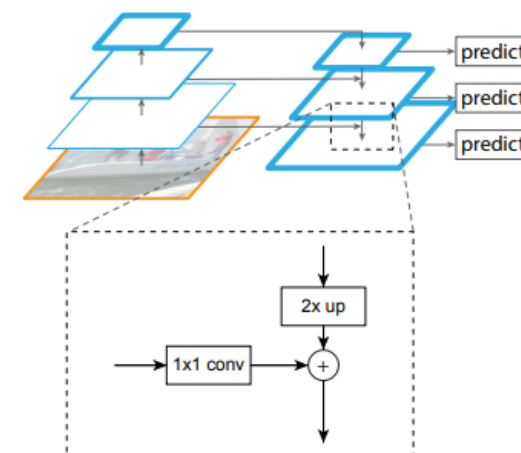
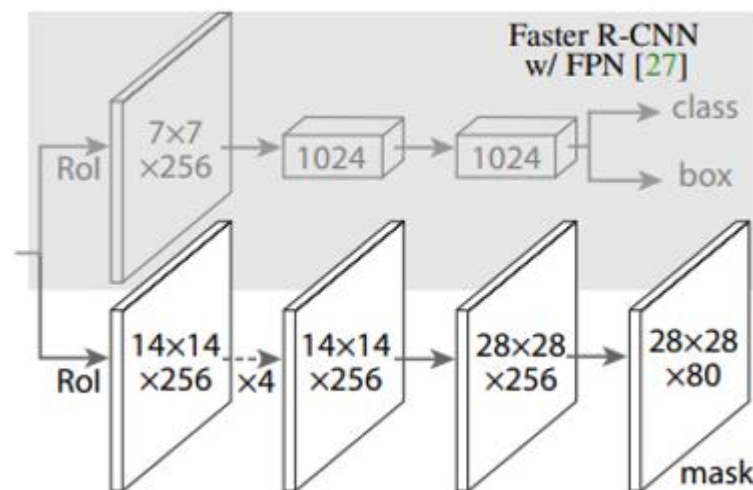
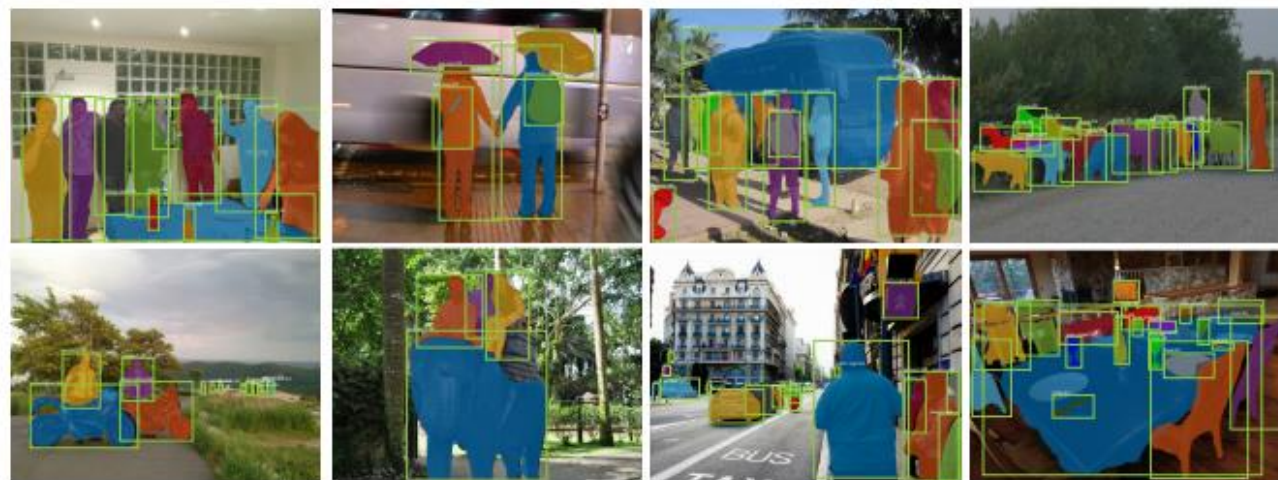
Faster R-CNN Results

- Pytorch Library
- Transfer Learning
- **7.5 average FPS**
- **61.6 mAP**



Mask R-CNN

- Instance Segmentation
- Segmentation CNN
- Feature Pyramid Networks
- ROI Align



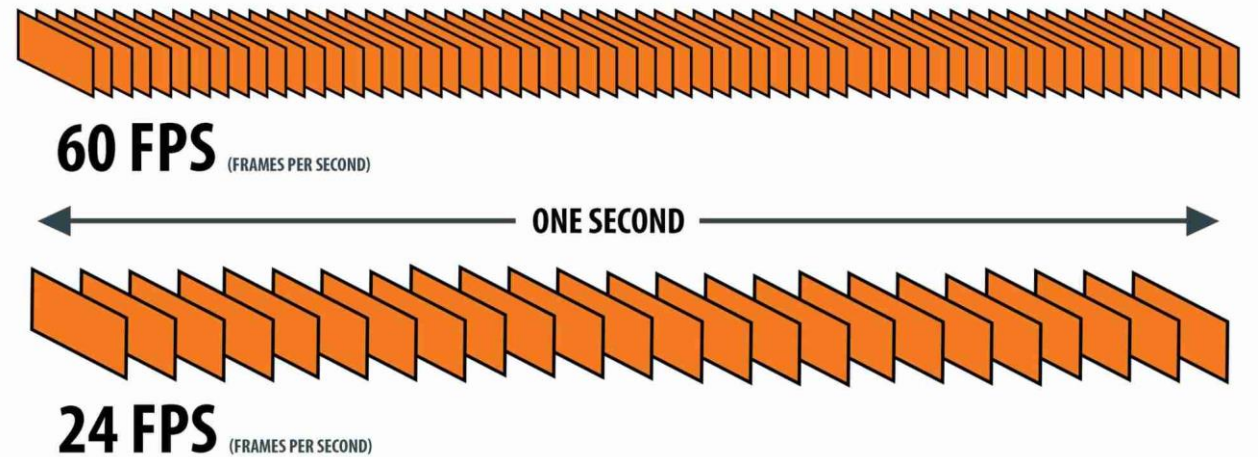
Mask R-CNN Results

- Mask data
- Pytorch Library
- Transfer Learning
- **7 average FPS**
- **61.3 mAP**



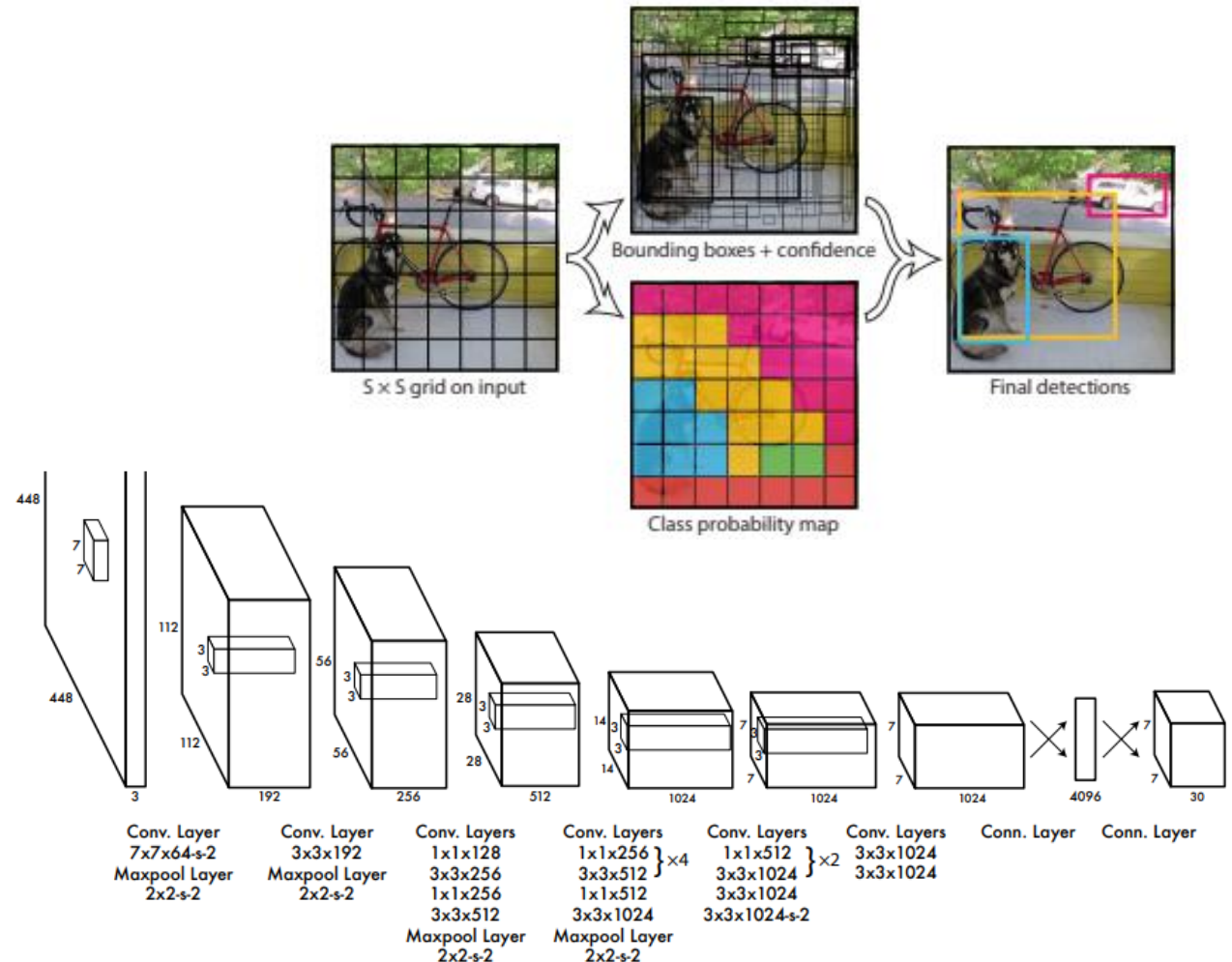
One stage detection

- Two-Stage detection is slow
 - 7 fps achieved
 - Want 30 fps minimum
- Avoid Region Proposal Stage



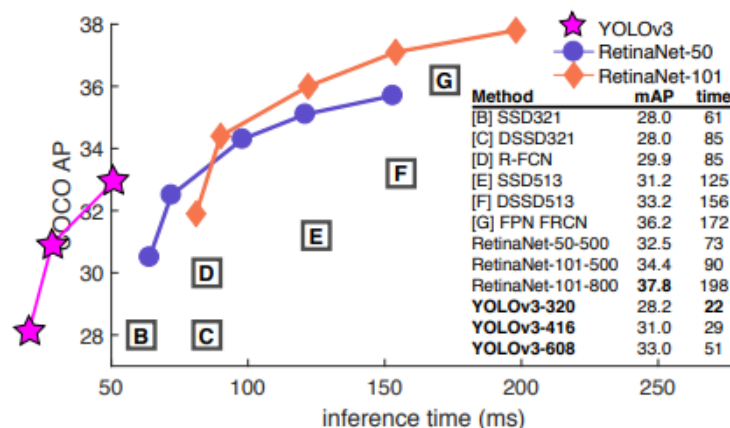
You Only Look Once (YOLO)

- Divide image in 7x7 grid cells
 - 1 Class predicted per cell
 - 2 boxes predicted per cell
- No anchor boxes
- No Region Proposals
- 45 FPS



YOLOv2 and YOLOv3

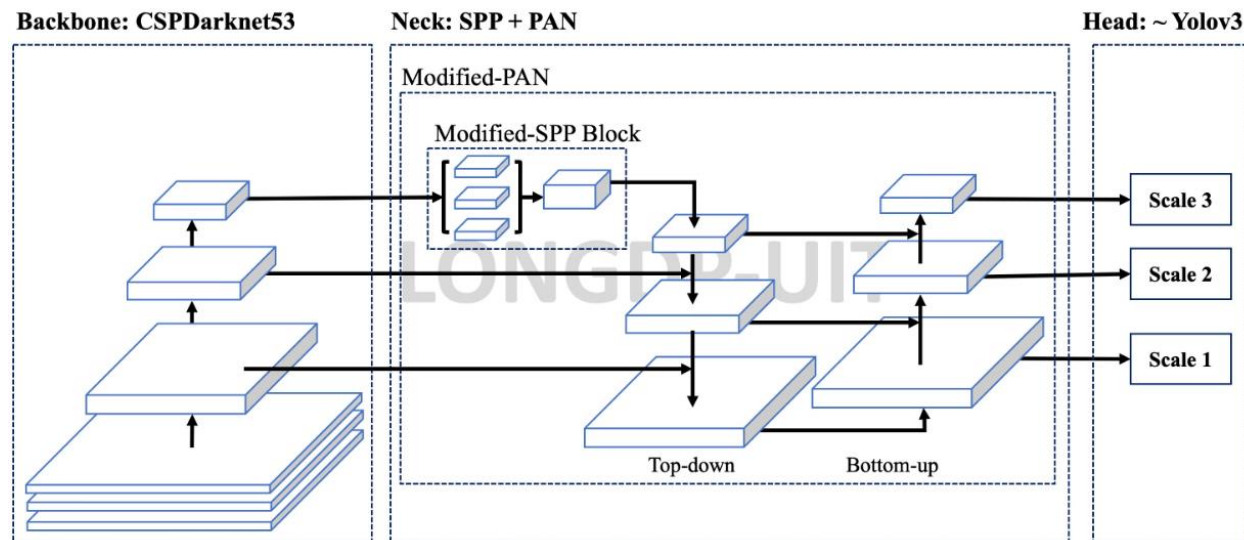
- Small improvements
- Anchor Boxes
- Darknet-53 backbone
- Feature Pyramid Networks (FPN)



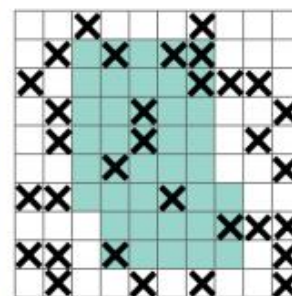
	Type	Filters	Size	Output
1x	Convolutional	32	3 × 3	256 × 256
	Convolutional	64	3 × 3 / 2	128 × 128
	Convolutional	32	1 × 1	
	Convolutional	64	3 × 3	
	Residual			128 × 128
	Convolutional	128	3 × 3 / 2	64 × 64
	Convolutional	64	1 × 1	
	Convolutional	128	3 × 3	
	Residual			64 × 64
	Convolutional	256	3 × 3 / 2	32 × 32
8x	Convolutional	128	1 × 1	
	Convolutional	256	3 × 3	
	Residual			32 × 32
	Convolutional	512	3 × 3 / 2	16 × 16
	Convolutional	256	1 × 1	
	Convolutional	512	3 × 3	
	Residual			16 × 16
	Convolutional	1024	3 × 3 / 2	8 × 8
	Convolutional	512	1 × 1	
	Convolutional	1024	3 × 3	
4x	Residual			8 × 8
	Avgpool		Global	
	Connected		1000	
	Softmax			

YOLOv4

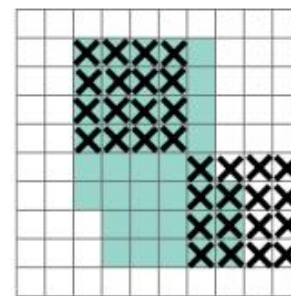
- Complex network
- Regularization
 - Data augmentation
 - DropBlock
- Modified FPN



(a)



(b)



(c)



YOLOv5

- YOLOv4 with improvements
- Managed by a company
- Continuous development
- Multiple model size options
- User-Friendly

YOLOv5 by 

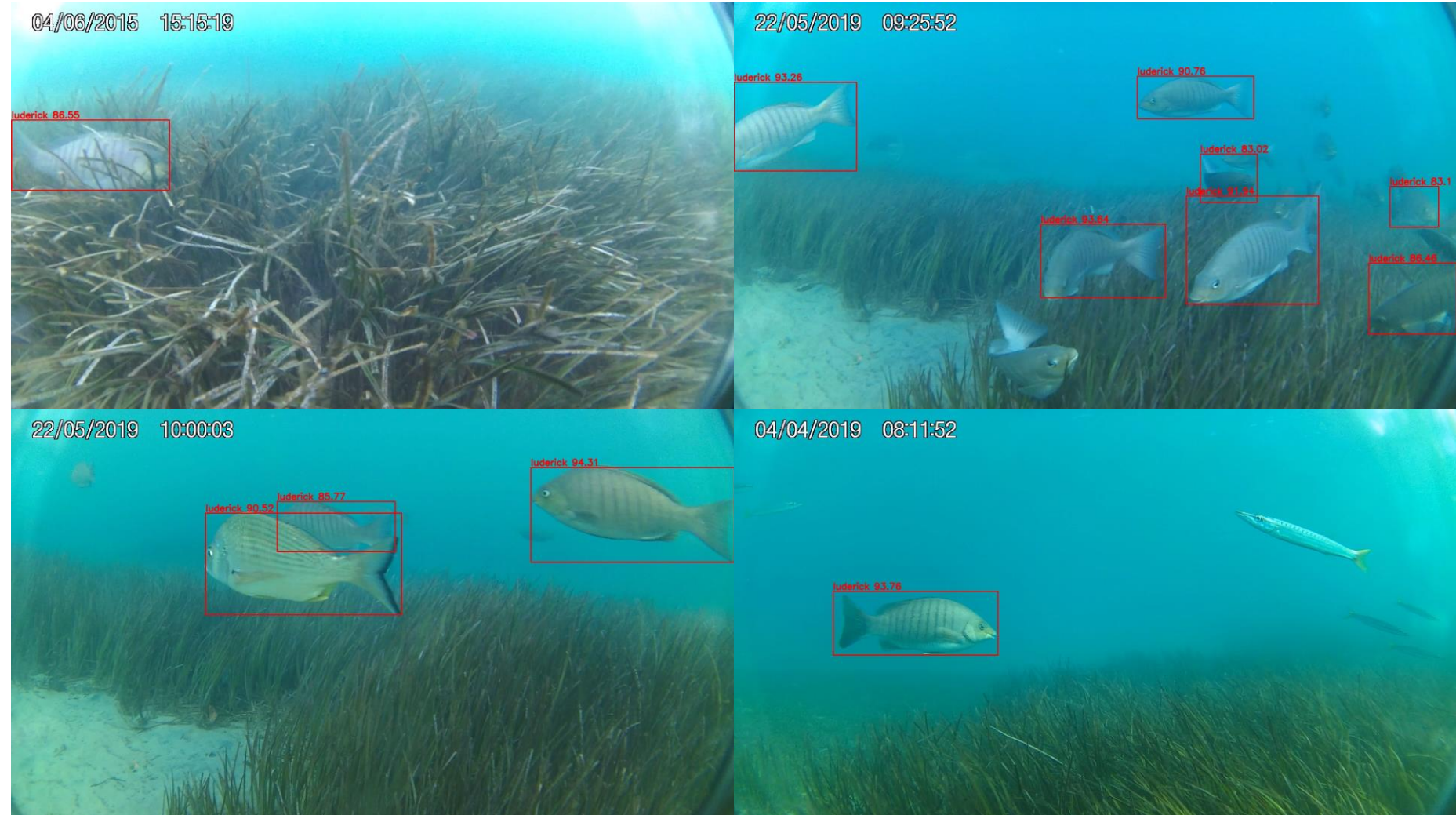


 Download on the App Store  Coming Soon on Google Play

				
Nano	Small	Medium	Large	XLarge
YOLOv5n	YOLOv5s	YOLOv5m	YOLOv5l	YOLOv5x
4 MB _{FP16}	14 MB _{FP16}	41 MB _{FP16}	89 MB _{FP16}	166 MB _{FP16}
6.3 ms _{V100}	6.4 ms _{V100}	8.2 ms _{V100}	10.1 ms _{V100}	12.1 ms _{V100}
28.4 mAP _{COCO}	37.2 mAP _{COCO}	45.2 mAP _{COCO}	48.8 mAP _{COCO}	50.7 mAP _{COCO}

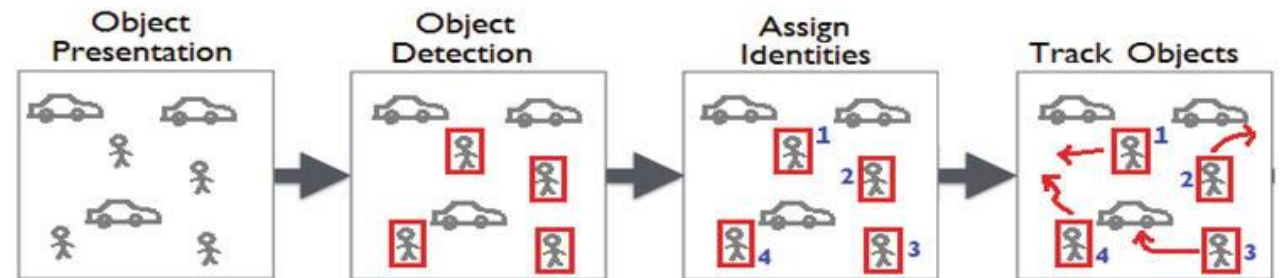
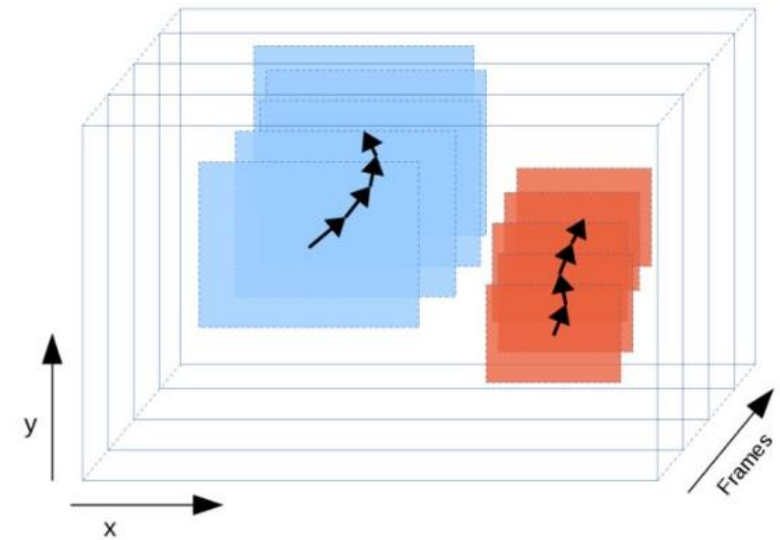
YOLOv5 Results

- Text input format
- YOLOv5 Github
- Small size option
- **47.3 average FPS**
- **63.4 mAP**



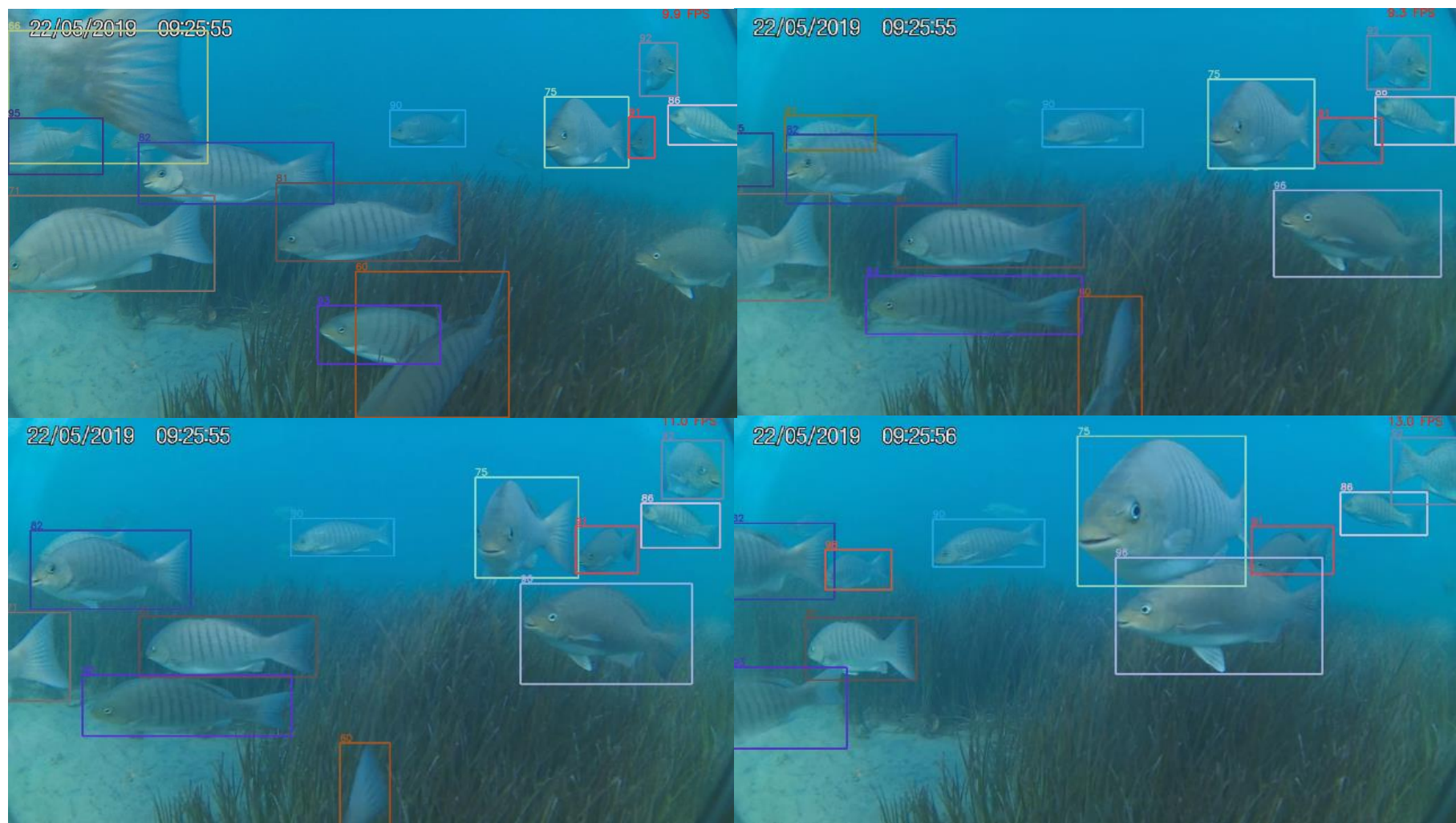
Object Tracking

- Used on top of Detection models
- Track through frames
- IOU tracking
- Hungarian Algorithm
- Kalman Filters
- SORT



Tracking Results

- YOLOv5 detection
- Simple SORT
- No Kalman Filters
- **23.6 average FPS**



Conclusions

- YOLOv5 Achieves precise real-time detection
- Detection can be used for tracking and other purposes
- Deep Learning is replacing traditional computer vision
- New solutions will probably appear each year