

# Object Tracking

**Danna Gurari**

University of Colorado Boulder  
Fall 2021



# Review

- Last lecture:
  - Semantic segmentation problem
  - Semantic segmentation applications
  - Semantic segmentation datasets
  - Semantic segmentation evaluation metrics
  - Computer vision models: fully convolutional networks
- Assignments (Canvas)
  - Reading assignment due earlier today
  - Two reading assignments out that are due next Monday and Wednesday
- Questions?

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models

# Object Tracking: Today's Topics

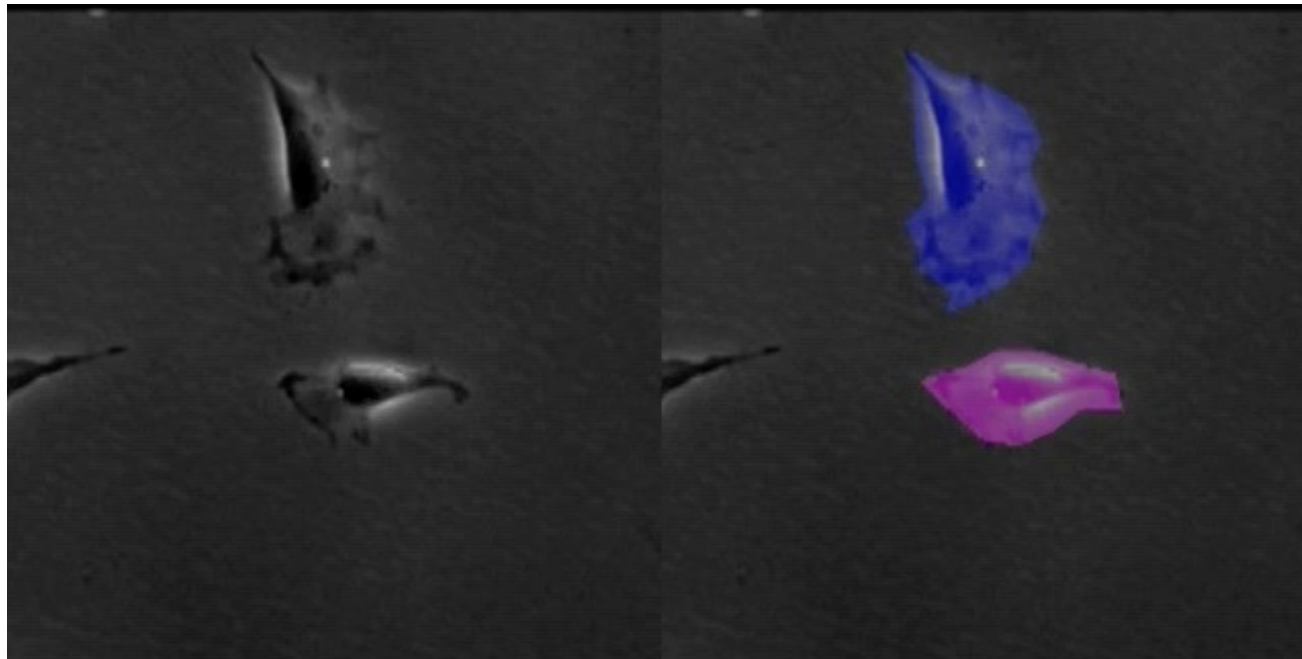
- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models

# Definition

- Identification of the trajectory of an object over time;
  - Single object
  - Multiple objects; e.g.,

Input

Output masks  
overlaid on video



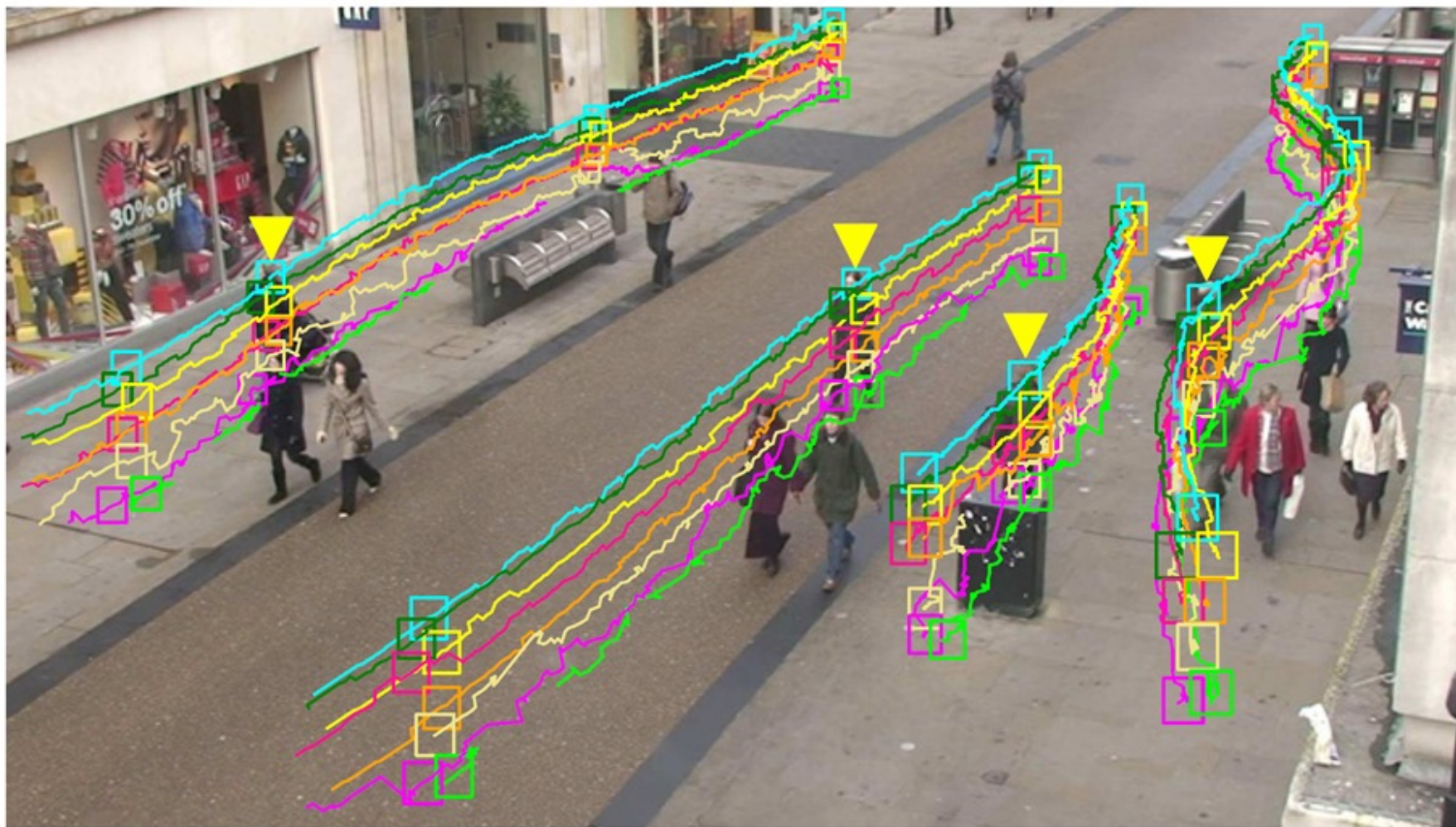
# Definition

- Identification of the trajectory of an object over time
  - Single object
  - Multiple objects
  
- How can the trajectory of an object be represented?
  - Bounding box or ellipse
  - Segmentation or coarse outline
  - Position (e.g., object centroid, corner, salient point)

# Object Tracking: Today's Topics

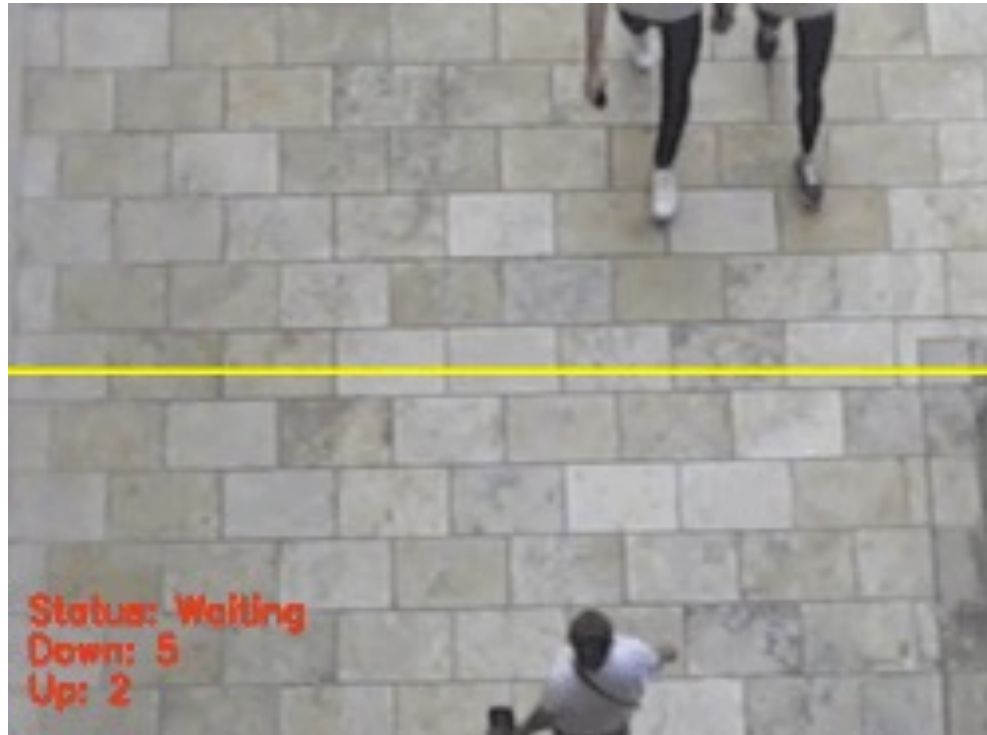
- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models

# Surveillance

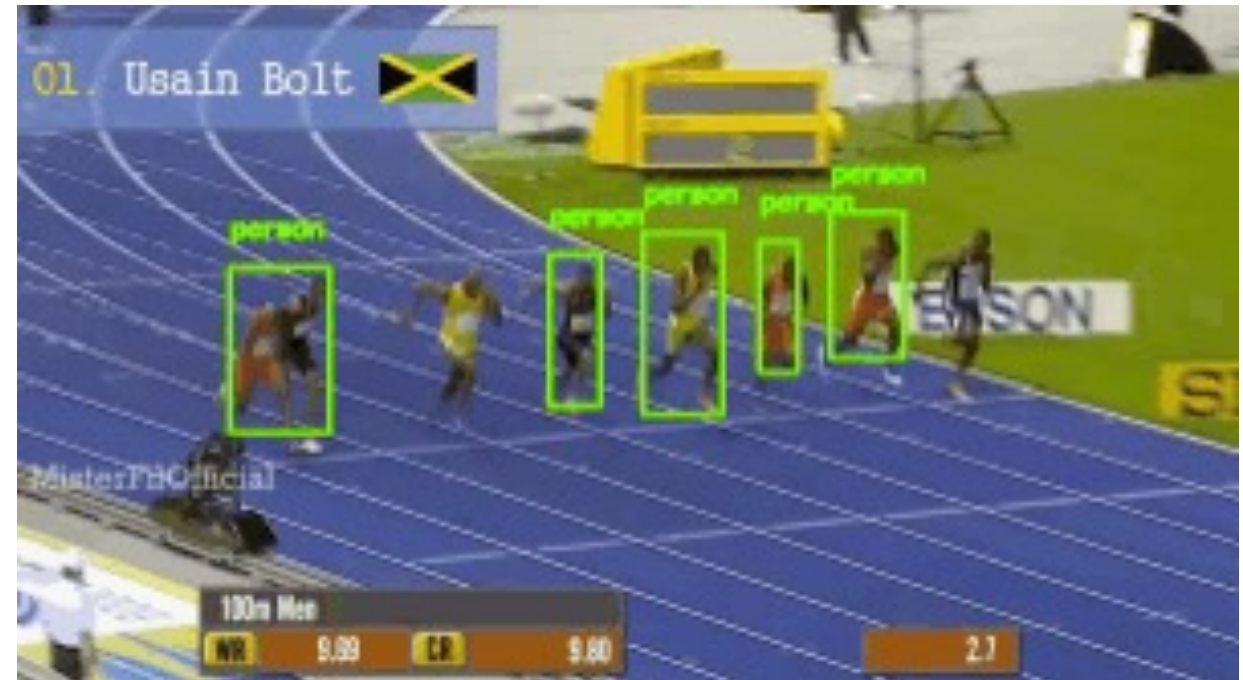




# Business Marketing: People Analytics



# Sports Analysis



<https://www.pyimagesearch.com/2018/10/29/multi-object-tracking-with-dlib/>  
<https://www.pyimagesearch.com/2018/08/06/tracking-multiple-objects-with-opencv/>

# Sports Performance Analytics

Calculate Bat speed from video!



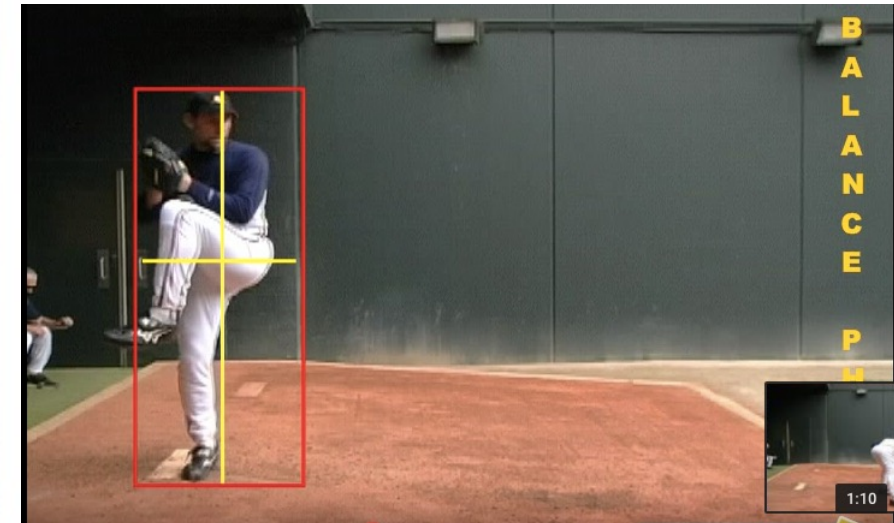
NEW! Track Bowling Ball Path!



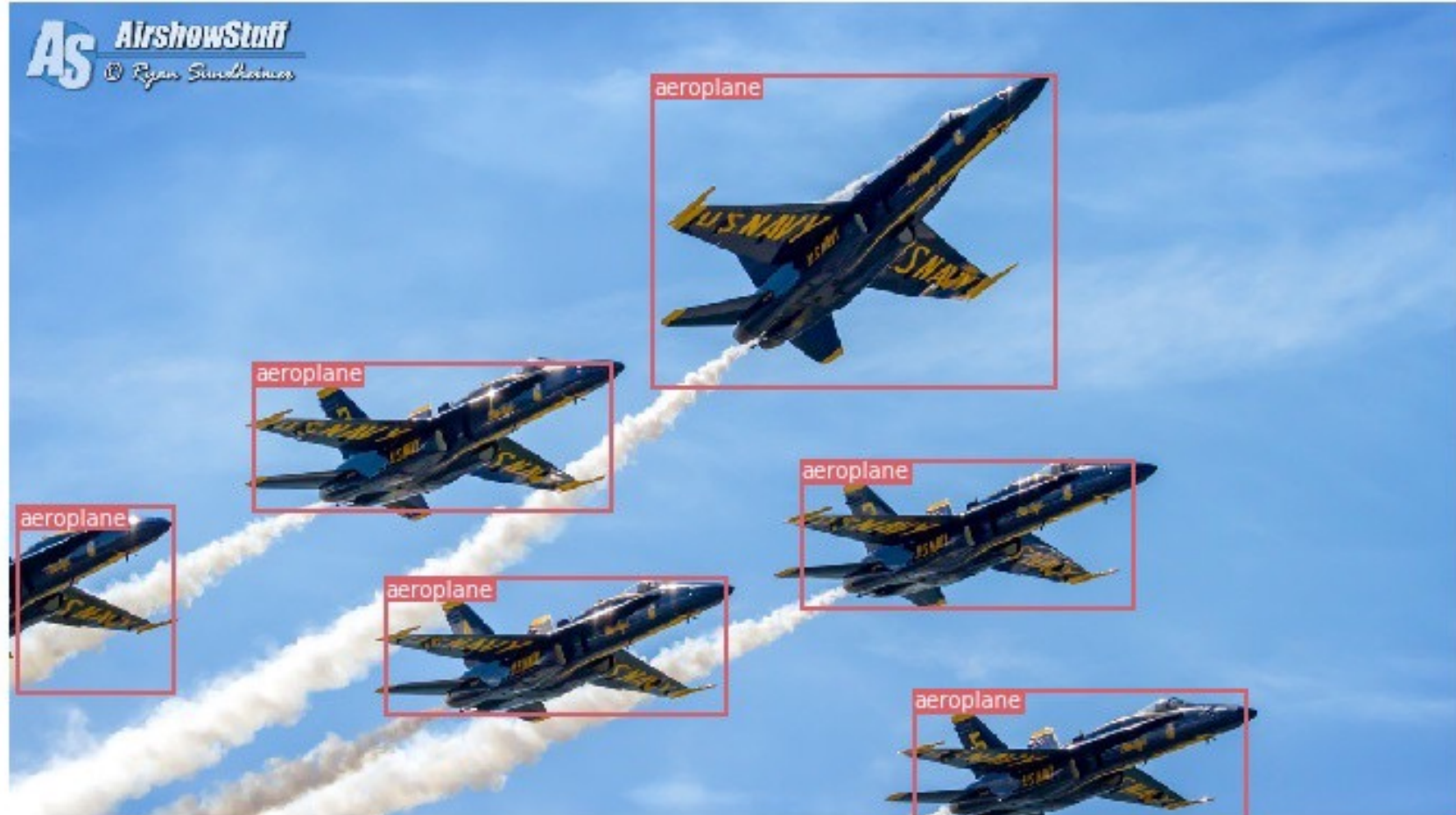
Works great for putting!



<http://www.motionprosoftware.com/>



# Military Defense



# Self-driving Cars

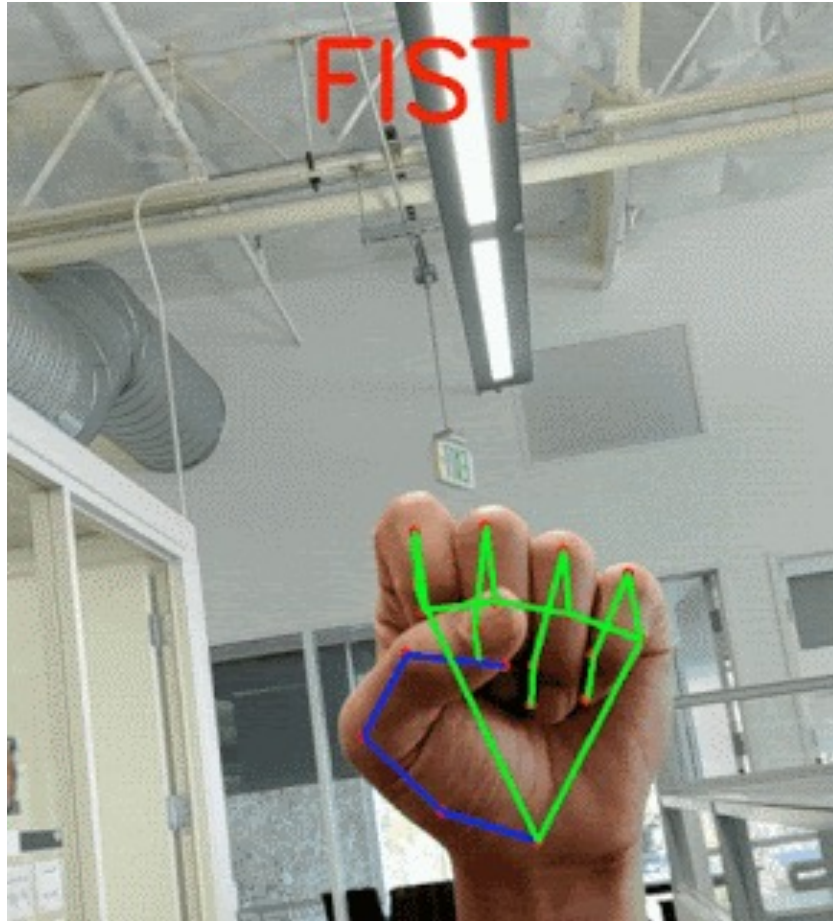


# Human Computer Interaction



Roboceptionist

# Sign Language Recognition



# Biological Monitoring

Counting bats exiting  
a cave in Texas:

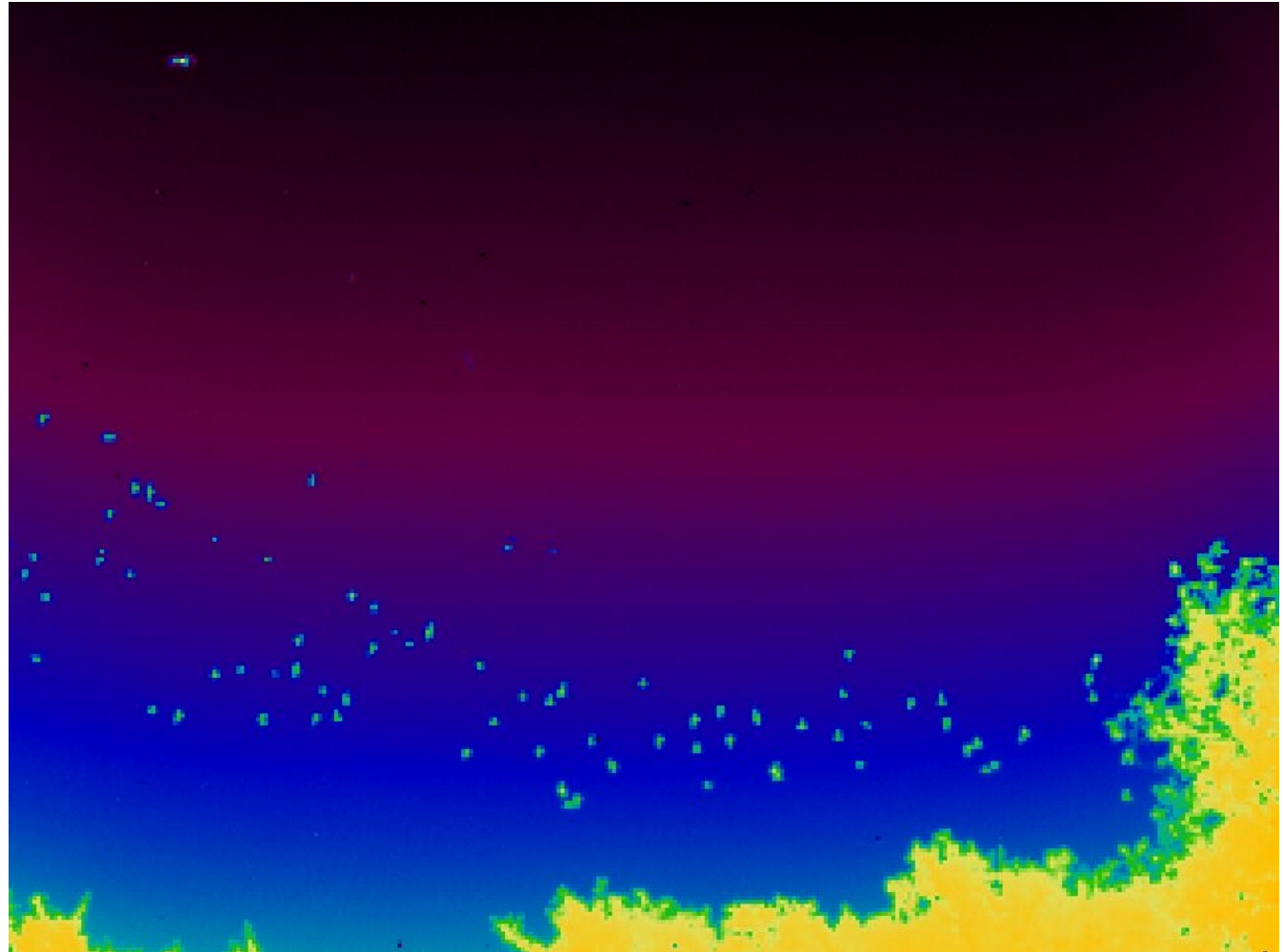


Image source: <https://www.cs.bu.edu/fac/betke/research/bats/images2.html>



# Augmented Reality



<https://virtualrealitypop.com/object-recognition-in-augmented-reality-8f7f17127a7a>

<https://www.geekwire.com/2017/augmented-reality-shopping-phone-patent-hints-amazons-aspirations/>

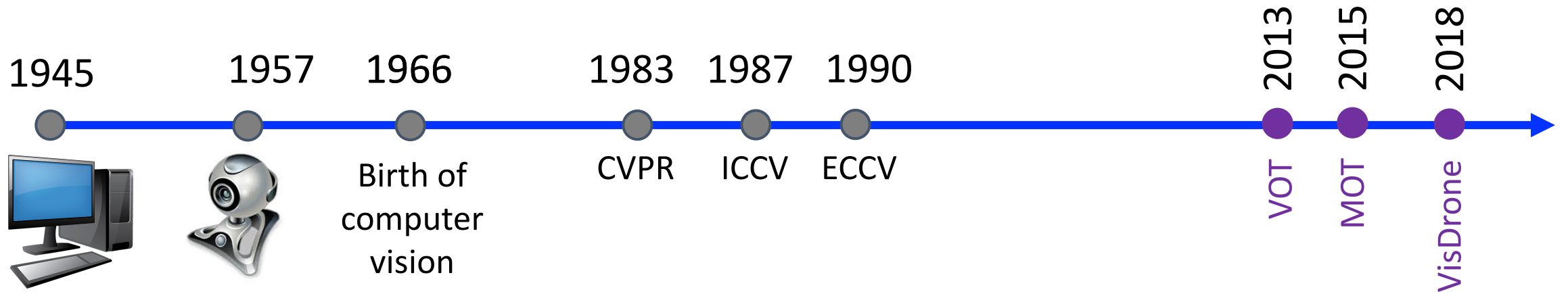
# Applications

What other applications can you think of where object tracking could be useful?

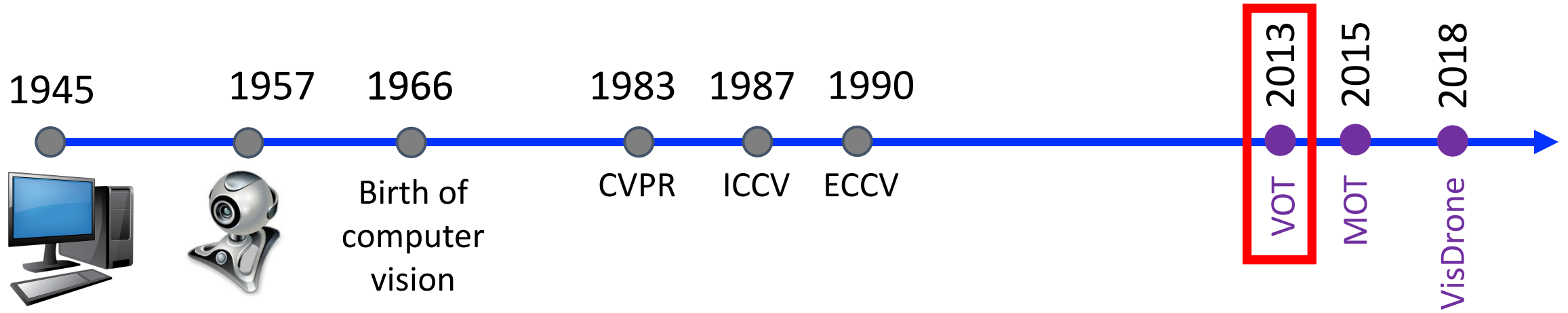
# Object Tracking: Today's Topics

- Problem
- Applications
- **Datasets**
- Evaluation metrics
- Computer vision models

# Object Tracking Datasets









# Object Tracking Datasets



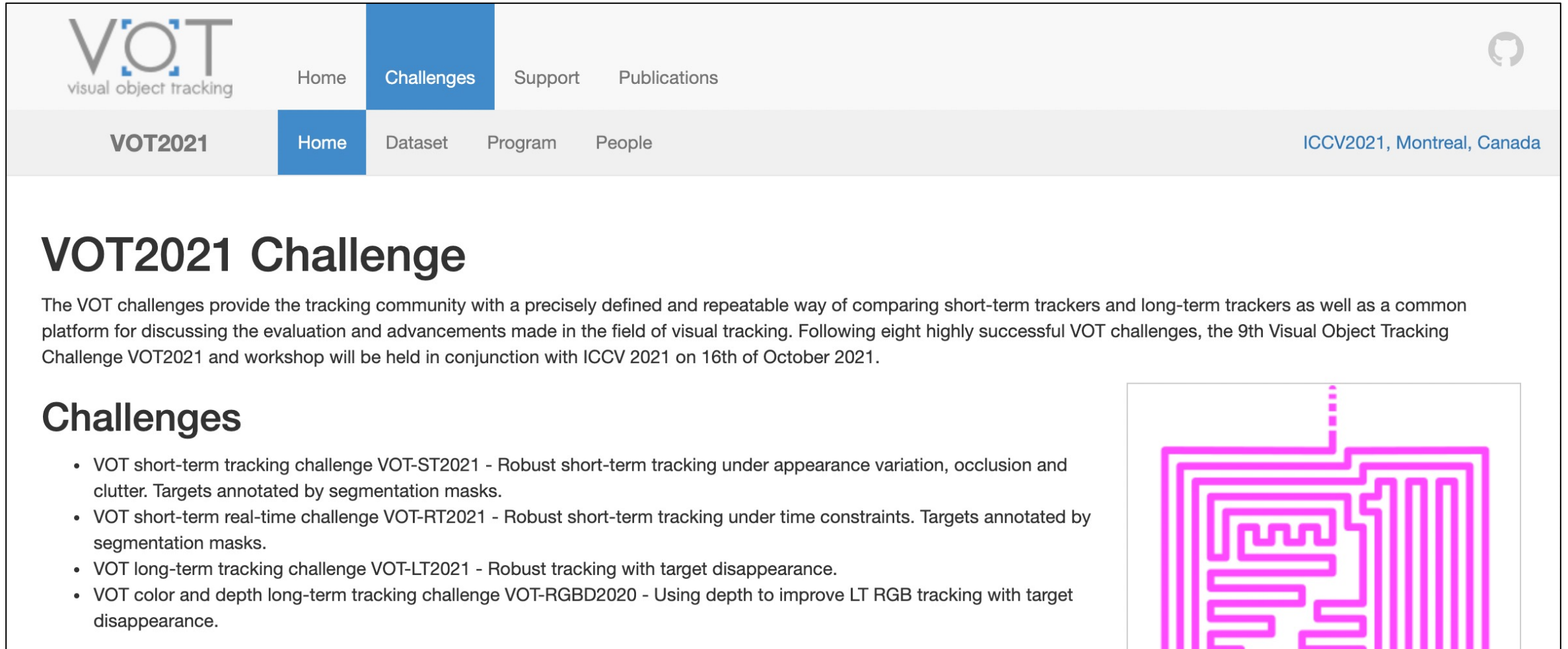
# Single Object Tracking Dataset: VOT

- Aggregated 16 videos from existing datasets that used bounding boxes to track a single object in each video
  - Limitation: inconsistent annotation methodologies across videos (e.g., different bounding box criteria)
- Authors re-annotated object tracking for videos they deemed to have unsuitable annotations

# Single Object Tracking Dataset: VOT's Evolution

 <p>visual object tracking challenge</p>	 <p>visual object tracking challenge</p>
 <p><b>VOT2016 benchmark</b></p> <p>The fourth challenge updated the dataset of 60 sequences with new annotations. The results were published in a joint paper presented at a workshop at ECCV2016.</p>	 <p><b>VOT2015 benchmark</b></p> <p>The third challenge introduced a dataset of 60 challenging sequences, a formalized sequence selection methodology and improvements to evaluation methodology. The results were published in a joint paper presented at an ICCV2015 workshop.</p>
 <p><b>VOT2014 benchmark</b></p> <p>The second challenge introduced several improvements in annotations and testing of statistical significance, new set of 25 sequences and an improved evaluation kit. The results were published in a joint paper presented at an ECCV2014 workshop.</p>	 <p><b>VOT2013 benchmark</b></p> <p>The first challenge introduced a new evaluation kit plus 16 well-known short videos. 27 single-target trackers submitted by 51 participants participated at the challenge. The results were published in a joint paper presented at an ICCV2013 workshop which was attended by over 70 researchers.</p>

# Single Object Tracking Annual Challenge (9<sup>th</sup> year now)



**VOT**  
visual object tracking

Home **Challenges** Support Publications

**VOT2021** Home Dataset Program People


ICCV2021, Montreal, Canada

## VOT2021 Challenge

The VOT challenges provide the tracking community with a precisely defined and repeatable way of comparing short-term trackers and long-term trackers as well as a common platform for discussing the evaluation and advancements made in the field of visual tracking. Following eight highly successful VOT challenges, the 9th Visual Object Tracking Challenge VOT2021 and workshop will be held in conjunction with ICCV 2021 on 16th of October 2021.

### Challenges

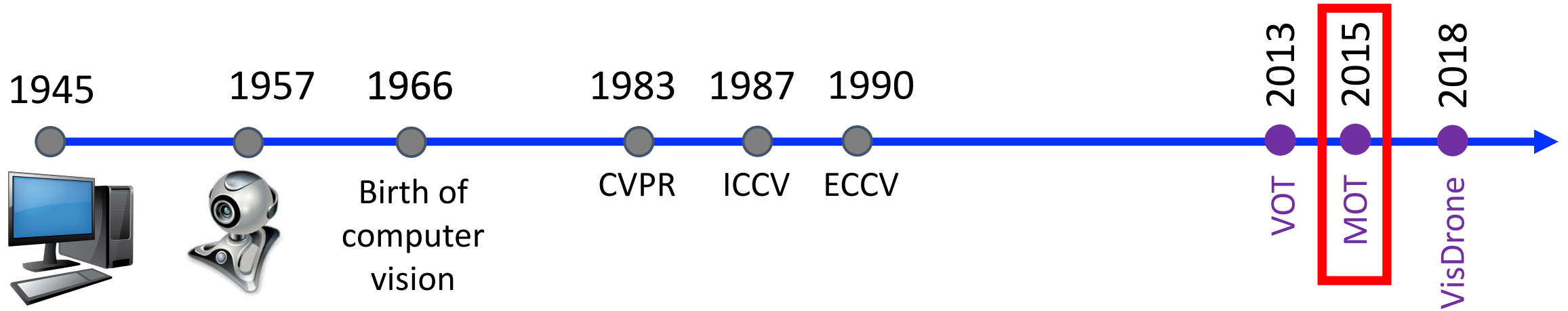
- VOT short-term tracking challenge VOT-ST2021 - Robust short-term tracking under appearance variation, occlusion and clutter. Targets annotated by segmentation masks.
- VOT short-term real-time challenge VOT-RT2021 - Robust short-term tracking under time constraints. Targets annotated by segmentation masks.
- VOT long-term tracking challenge VOT-LT2021 - Robust tracking with target disappearance.
- VOT color and depth long-term tracking challenge VOT-RGBD2020 - Using depth to improve LT RGB tracking with target disappearance.



<https://www.votchallenge.net/vot2021/>



# Object Tracking Datasets



# Multiple Object Tracking Dataset: MOT

- Authors aggregated 22 videos that contain a total of 11,286 frames associated with 61,440 annotated bounding boxes
  - **Static and moving camera;** e.g., held by a person, stroller, and car
  - **Multiple viewpoints;** e.g., cameras positioned at a high, medium, and low position (e.g., person's height versus on the ground looking up)
  - **Multiple weather conditions;** e.g., sunny versus cloudy versus night time
- 16 of the videos came from existing datasets while the other 6 were generated by the authors; tracked objects were people and vehicles



# Multiple Object Tracking Dataset: MOT

- Authors aggregated 22 videos that contain a total of 11,286 frames associated with 61,440 annotated bounding boxes
  - **Static and moving camera**; e.g., held by a person, stroller, and car
  - **Multiple viewpoints**; e.g., cameras positioned at a high, medium, and low position (e.g., person's height versus on the ground looking up)
  - **Multiple weather conditions**; e.g., sunny versus cloudy versus night time
- Annotations:
  - Automatically-generated detections for the dataset provided
  - For existing videos, there GT was used
  - For new videos, the VATIC annotation tool was used to generate tracks

# Multiple Object Tracking Annotation: VATIC

Annotate every object, even stationary and obstructed objects, for the entire video. Instructions + New Object

The screenshot displays a video player interface for the VATIC annotation tool. The main video frame shows a parking lot with several cars and a person. The cars are labeled with IDs and their status: Car 7 (Parked), Car 8 (Parked), Car 9 (Parked), Car 10 (Reversing), Car 11 (Parked), Car 12 (Driving), and Person 11 (Walking). The right sidebar shows a list of annotated objects with their status and control options. The video player controls include a Play button, a progress bar, and buttons for Slower, Slow, Normal, and Fast. There are also buttons for Disable Resize?, Hide Boxes?, Hide Labels?, and Save Work.

**Car 12**

- Outside of view frame
- Occluded or obstructed
- Parked
- Driving
- Reversing

**Person 11**

- Outside of view frame
- Occluded or obstructed
- Walking
- Running
- Standing

**Car 10**

- Outside of view frame
- Occluded or obstructed
- Parked
- Driving
- Reversing

**Car 9**

- Outside of view frame

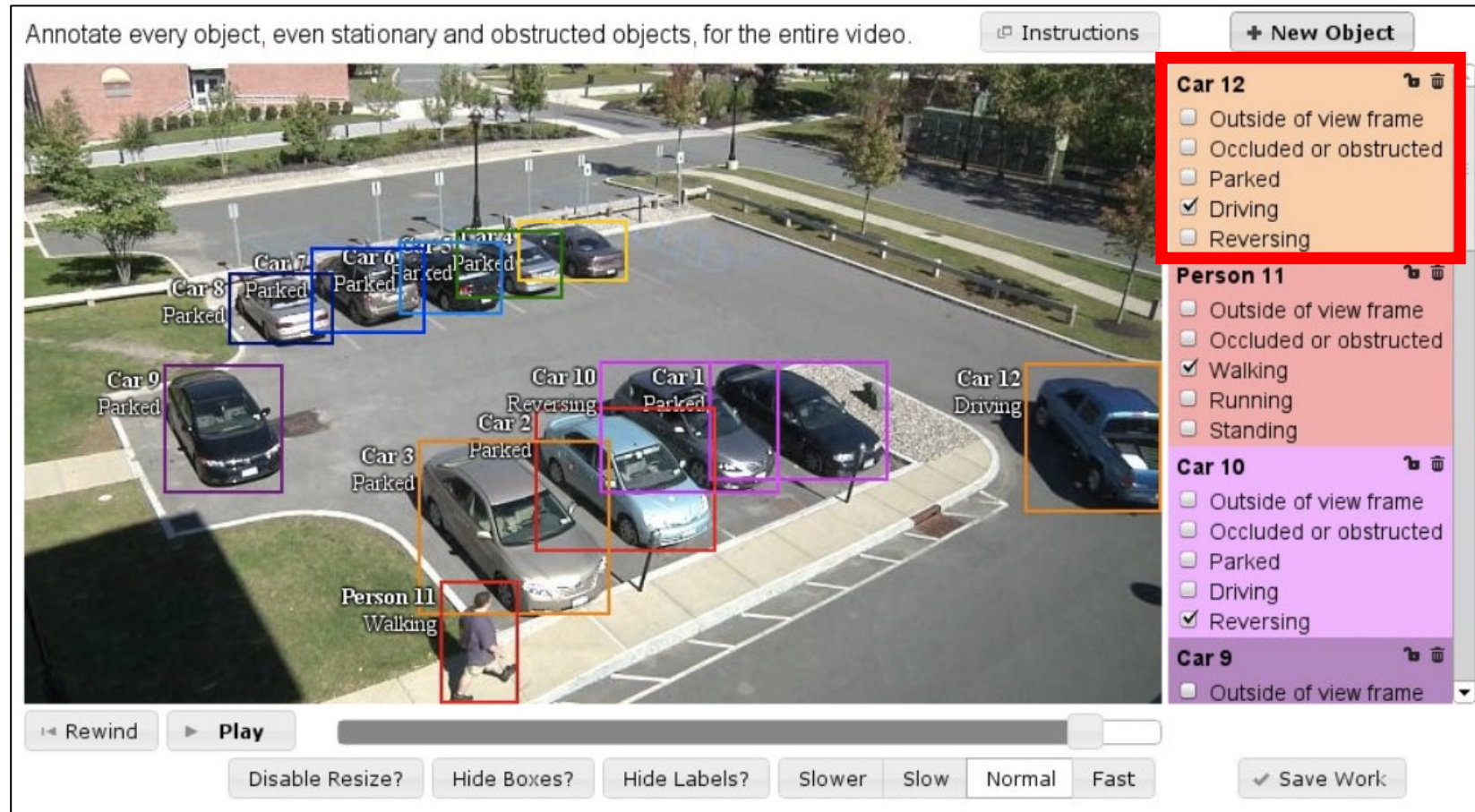
◀ Rewind ▶ **Play** [Progress Bar]

Disable Resize? Hide Boxes? Hide Labels? Slower Slow Normal Fast ✓ Save Work

Demo: <https://www.youtube.com/watch?v=ljI5pAowACc>

Carl Vondrick, Donald Patterson, and Deva Ramanan. Efficiently Scaling Up Crowdsourced Video Annotation: A Set of Best Practices for High Quality, Economical Video Labeling. IJCV 2012.

# Multiple Object Tracking Annotation: VATIC



Metadata about each object: e.g., activity, attributes, etc.

Carl Vondrick, Donald Patterson, and Deva Ramanan. Efficiently Scaling Up Crowdsourced Video Annotation: A Set of Best Practices for High Quality, Economical Video Labeling. IJCV 2012.

# Multiple Object Tracking Annotation: VATIC

- How to handle occlusions?
  - Annotation instructions: *"Always annotate during occlusions if the position can be determined unambiguously. If the occlusion is very long and it is not possible to determine the path of the object using simple reasoning (e.g. constant velocity assumption), the object will be assigned a new ID once it reappears"*
  - Annotation "**visibility**" flag: ranges between 0-1 (1 when it's fully visible, and less than 1 when it's occluded)
  - Annotation "**confidence**" flag: set to 1 when box should be considered for evaluation and 0 otherwise (for example, when a pedestrian is too small)
  - For non-tracked categories: annotate object with "**class**" value as occluder and ignore during evaluation

# Single Object Tracking Dataset: MOT's Evolution

**Multiple Object Tracking Benchmark**

home data results vis QVA submit FAQ people login sign up

## Welcome

- MOT15
- MOT16
- MOT17Det
- MOT17
- MOT20
- MOT20Det
- 3D-ZeF20
- MOTS
- TAO Challenge
- CTMC-v1
- TAO VOS Benchmark
- Head Tracking 21 **NEW**
- STEP-ICCV21 **NEW**
- DevKit

## Challenge: The Multiple Object Tracking Benchmark!

In the recent years, the computer vision community has relied on several centralized benchmarks for performance evaluation of numerous tasks including object detection, 3D reconstruction, optical flow, single-object short-term tracking, and stereo estimation. Despite their utility, they have proved to be extremely helpful to advance the state-of-the-art in the respective research fields. Interestingly, there has been limited work on the standardization of multiple target tracking evaluation. One of the few exceptions is the well-known PETS dataset used in surveillance applications. Even for this widely used benchmark, a common technique for presenting tracking results to date involves showing only one of the available data, inconsistent model training and varying evaluation scripts. With this benchmark, we aim to pave the way for a unified framework towards more meaningful quantification of multi-target tracking.

<https://motchallenge.net/>

# Single Object Tracking Annual Challenge (7<sup>th</sup> year now)

## Multiple Object Tracking Benchmark

home data results vis QVA submit FAQ people login sign up

### Welcome

- MOT15
- MOT16
- MOT17Det
- MOT17
- MOT20
- MOT20Det
- 3D-ZeF20
- MOTS
- TAO Challenge
- CTMC-v1
- TAO VOS Benchmark
- Head Tracking 21 **NEW**
- STEP-ICCV21 **NEW**
- DevKit

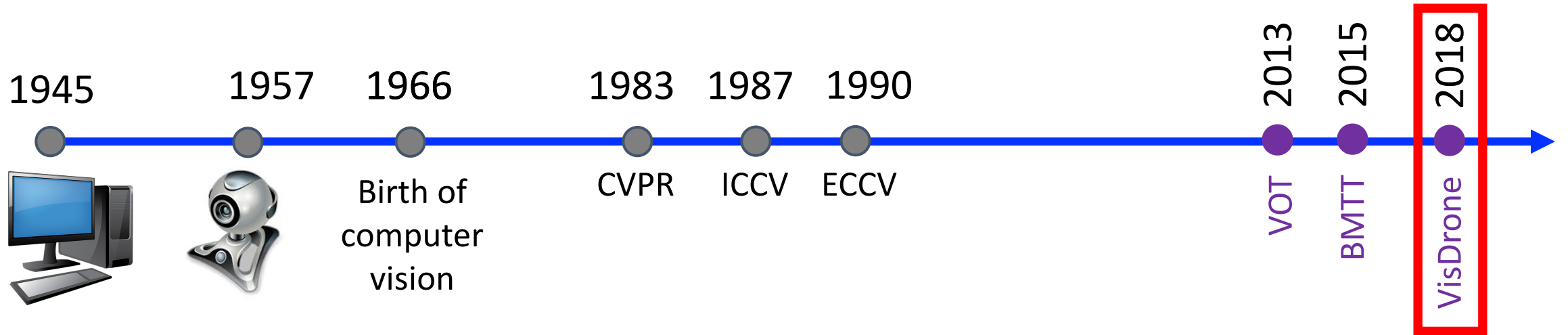
### Challenge: The Multiple Object Tracking Benchmark!

In the recent years, the computer vision community has relied on several centralized benchmarks for performance evaluation of numerous tasks including object detection, 3D reconstruction, optical flow, single-object short-term tracking, and stereo estimation. Despite their success, they have proved to be extremely helpful to advance the state-of-the-art in the respective research fields. However, there has been limited work on the standardization of multiple target tracking evaluation. One of the few exceptions is the well-known PETS dataset used in surveillance applications. Even for this widely used benchmark, a common technique for presenting tracking results to date involves showing only one of the available data, inconsistent model training and varying evaluation scripts. With this benchmark, we aim to pave the way for a unified framework towards more meaningful quantification of multi-target tracking.

<https://motchallenge.net/>



# Object Tracking Datasets



# VisDrone

- Authors collected 263 video clips (179,264 frames) from drones in Asia



- Annotations created for over 2.5 million object instances, however it is unspecified how these annotations were collected

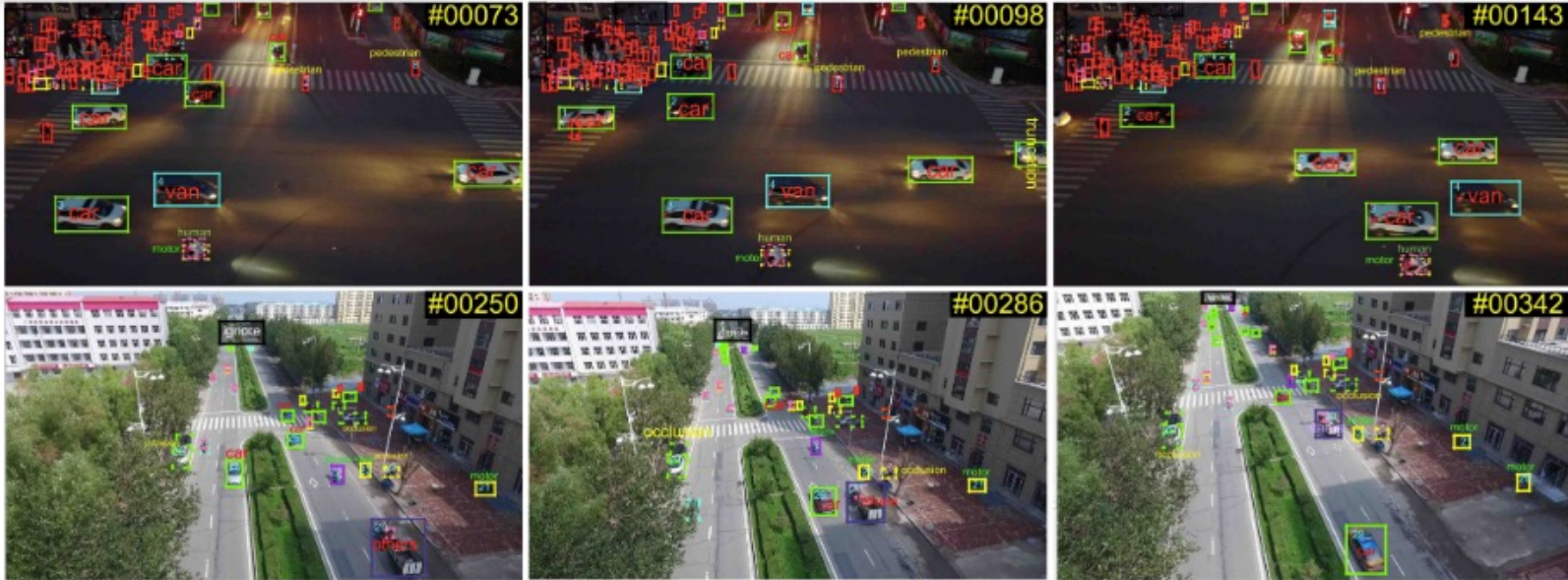
# VisDrone Challenge

Home Challenge ▾ Evaluate ▾ Download Submit FAQ ICCV2019 People Sign in Sign up

Object Detection in Images  
Object Detection in Videos  
**Single-Object Tracking**  
Multi-Object Tracking

## Multi-Object Tracking

...r results here!! Note that the evaluation server on the test-dev set will be open for



<http://www.aiskyeye.com/views/index>

# Discussion

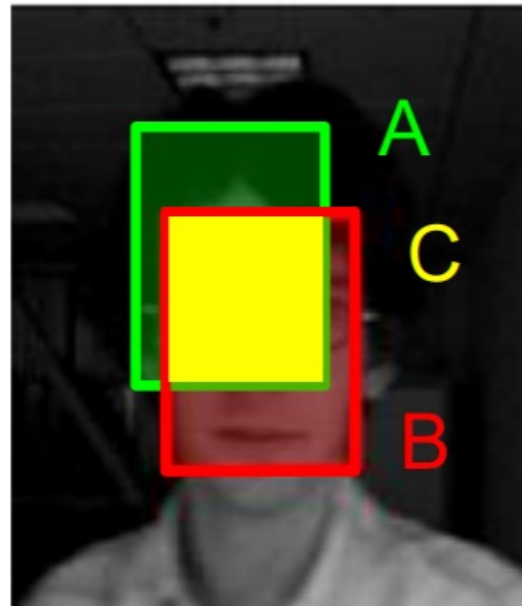
- When designing an annotation protocol to collect **high quality** object tracking annotations, how should these scenarios be handled:
  - Partially visible object
  - Occluded object
  - Object is reflected in reflective surfaces such as mirrors or windows
- What will be the total crowdsourcing task cost to annotate 1,000 1-minute videos where you need to track 5 humans/video (assume 30 frames/second)?

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- **Evaluation metrics**
- Computer vision models

# Accuracy

Average IoU from a tracker across all video frames



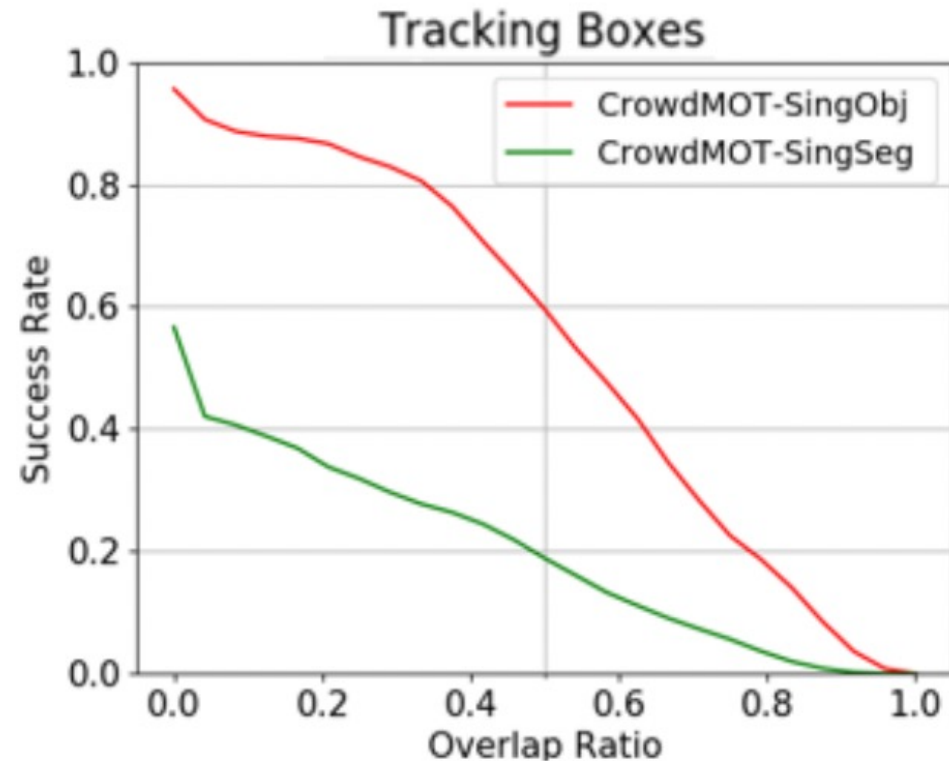
A = Ground Truth  
B = Predicted Track  
C = Intersection

Figure credit: [https://ags.cs.uni-kl.de/fileadmin/inf\\_ags/opt-ss15/OPT\\_SS2015\\_lec11.pdf](https://ags.cs.uni-kl.de/fileadmin/inf_ags/opt-ss15/OPT_SS2015_lec11.pdf)

Matej Kristan et al. "A Novel Performance Evaluation Methodology for Single-Target Trackers." PAMI 2016

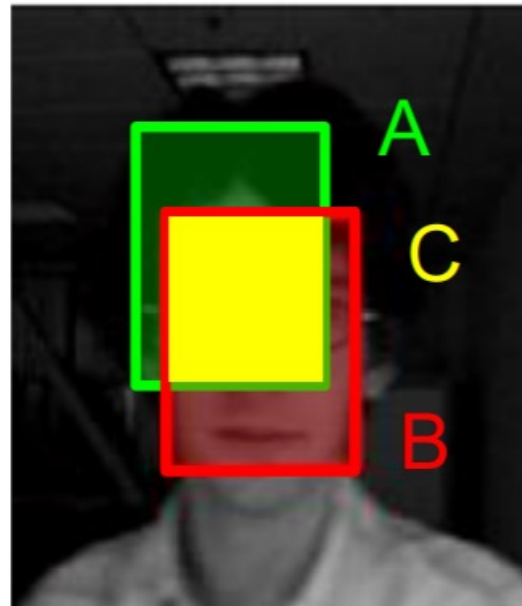
# Success Plot

Percentage of frames where the IoU is larger than a given threshold (e.g., 0.5); can create a plot by varying the threshold amount



# Robustness

Average number of times a tracker drifts to an IoU value of 0 and so needs to be re-initialized to the ground truth bounding box per video



A = Ground Truth  
B = Predicted Track  
C = Intersection

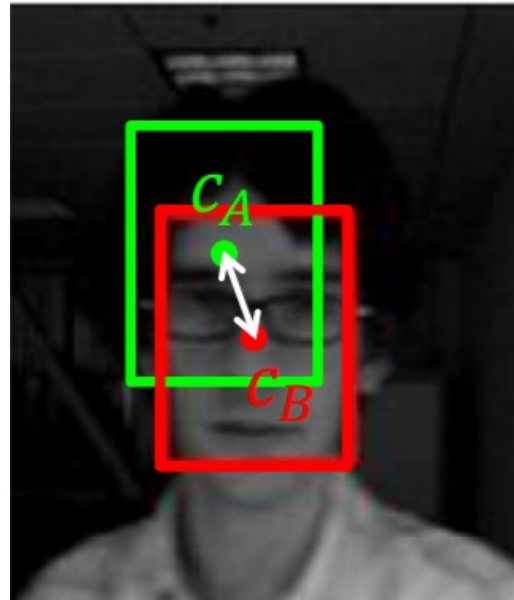
Figure credit: [https://ags.cs.uni-kl.de/fileadmin/inf\\_ags/opt-ss15/OPT\\_SS2015\\_lec11.pdf](https://ags.cs.uni-kl.de/fileadmin/inf_ags/opt-ss15/OPT_SS2015_lec11.pdf)

Matej Kristan et al. "A Novel Performance Evaluation Methodology for Single-Target Trackers." PAMI 2016



# Precision

Distance between the centers of bounding boxes for each frame

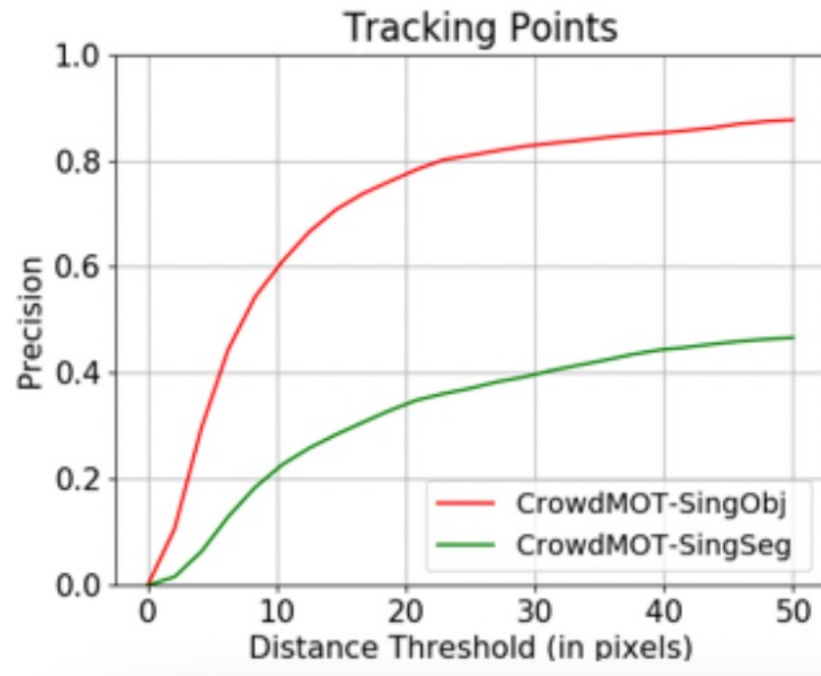


A = Ground Truth  
B = Predicted Track

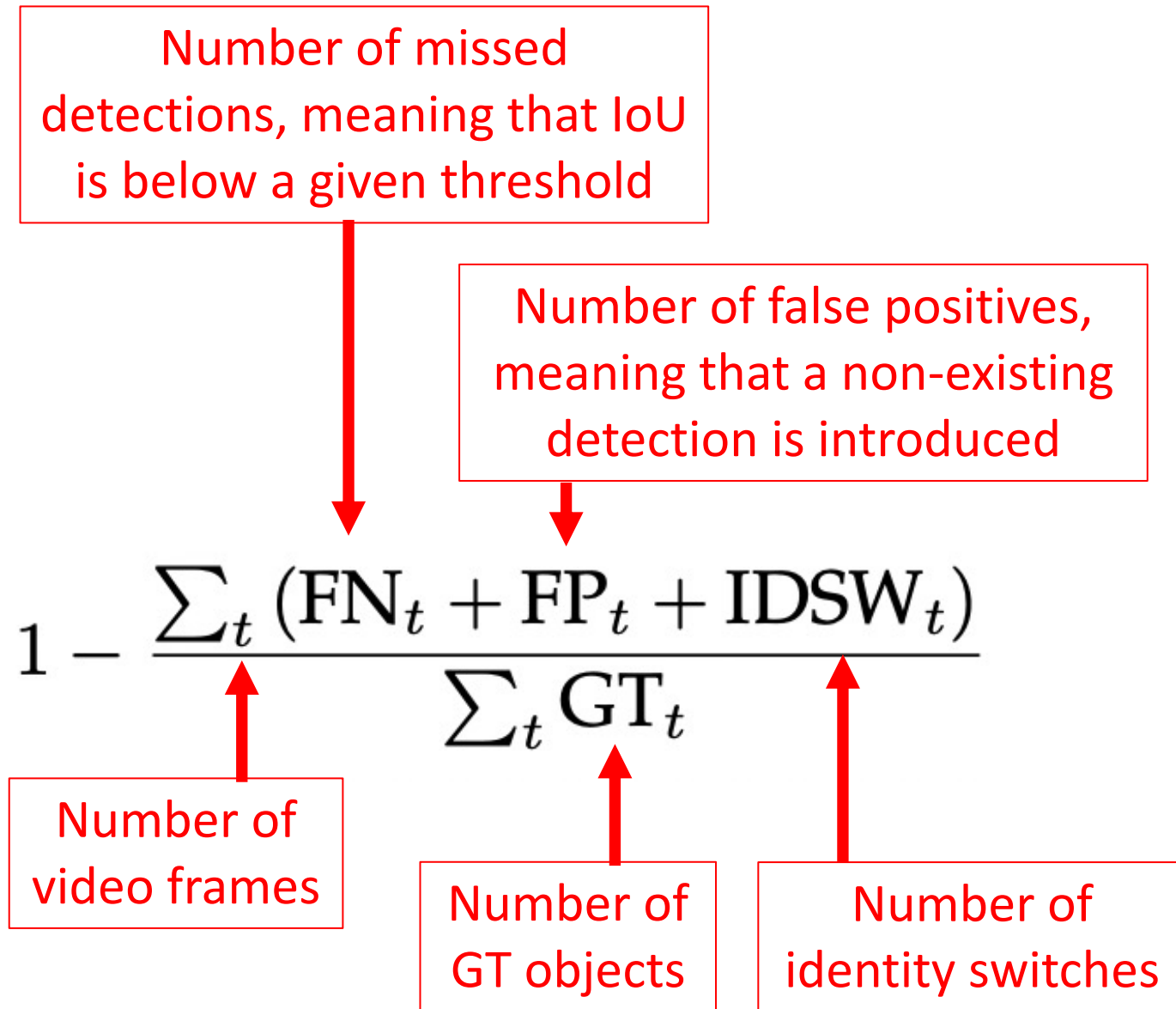
$$p = \|c_A - c_B\|$$

# Precision Plot

Percentage of frames with predicted location within a given threshold distance of ground truth (e.g., 20 pixels); can create a plot by varying the threshold amount



# MOTA



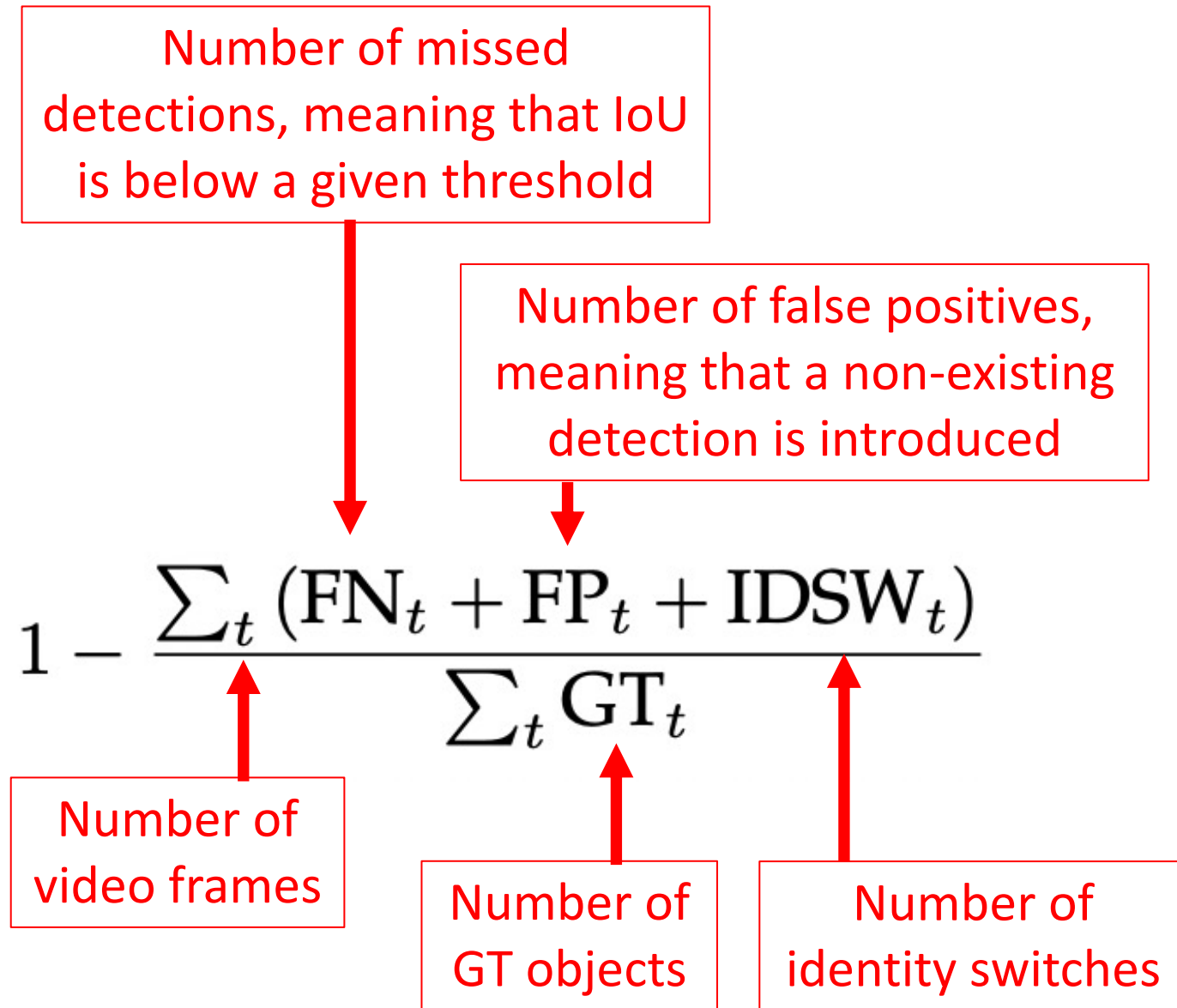
# MOTA

What is the range of possible values?

- (- infinite, 100] (original value is multiplied by 100)

When is MOTA negative?

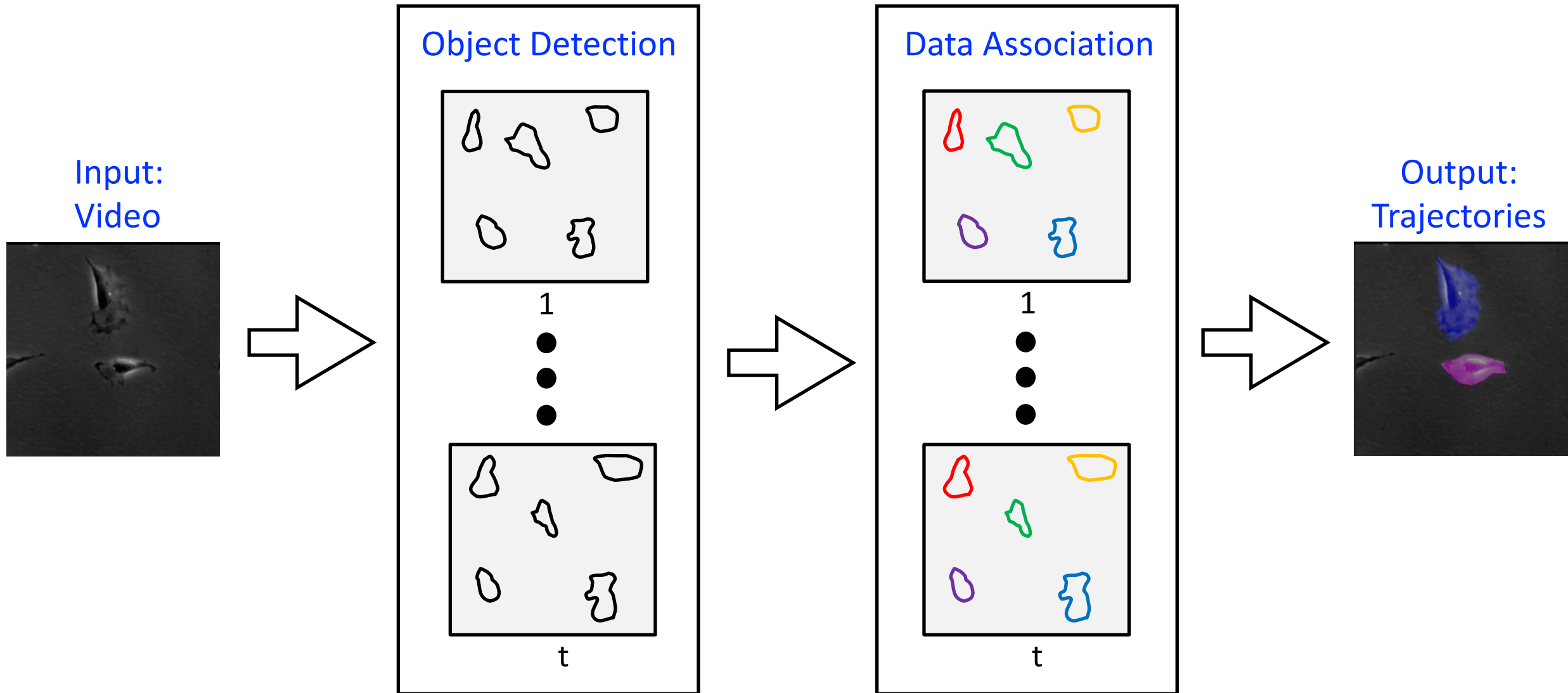
- When the number of errors exceed the number of objects in the frames



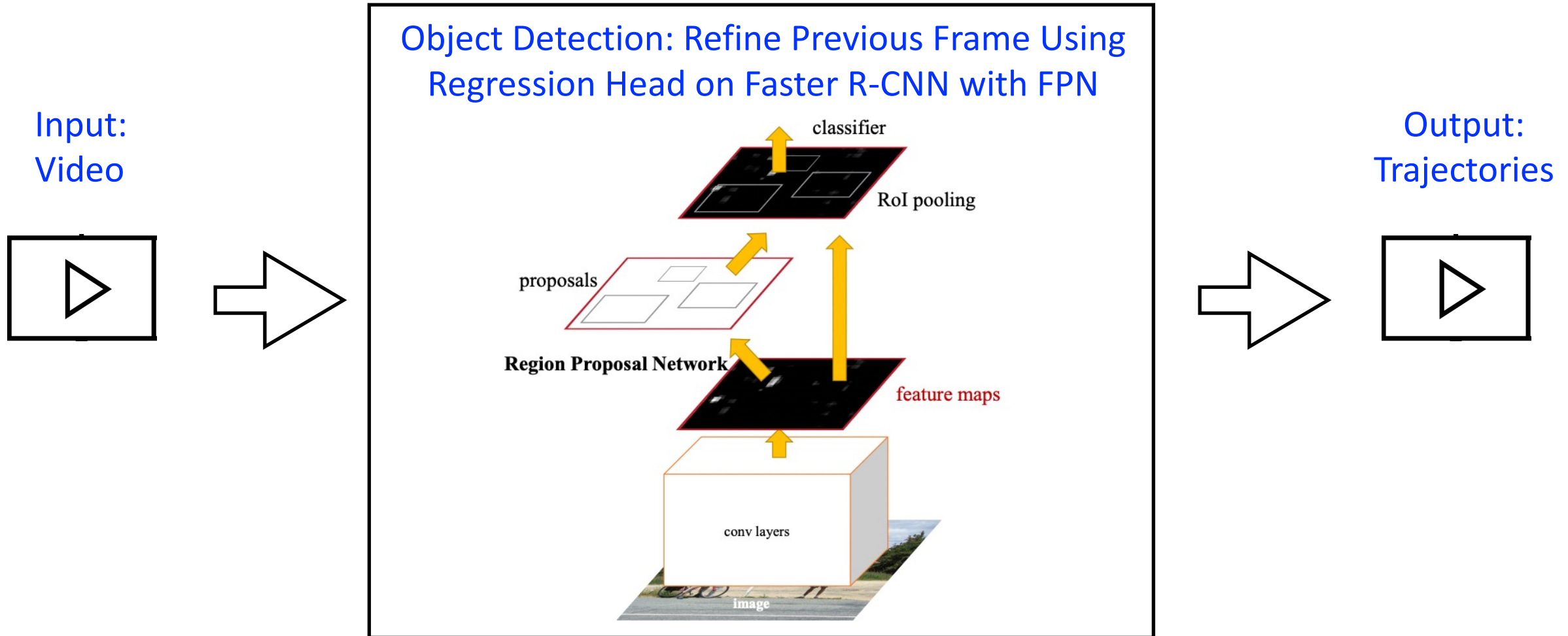
# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models

# Common Approach: Tracking-by-Detection



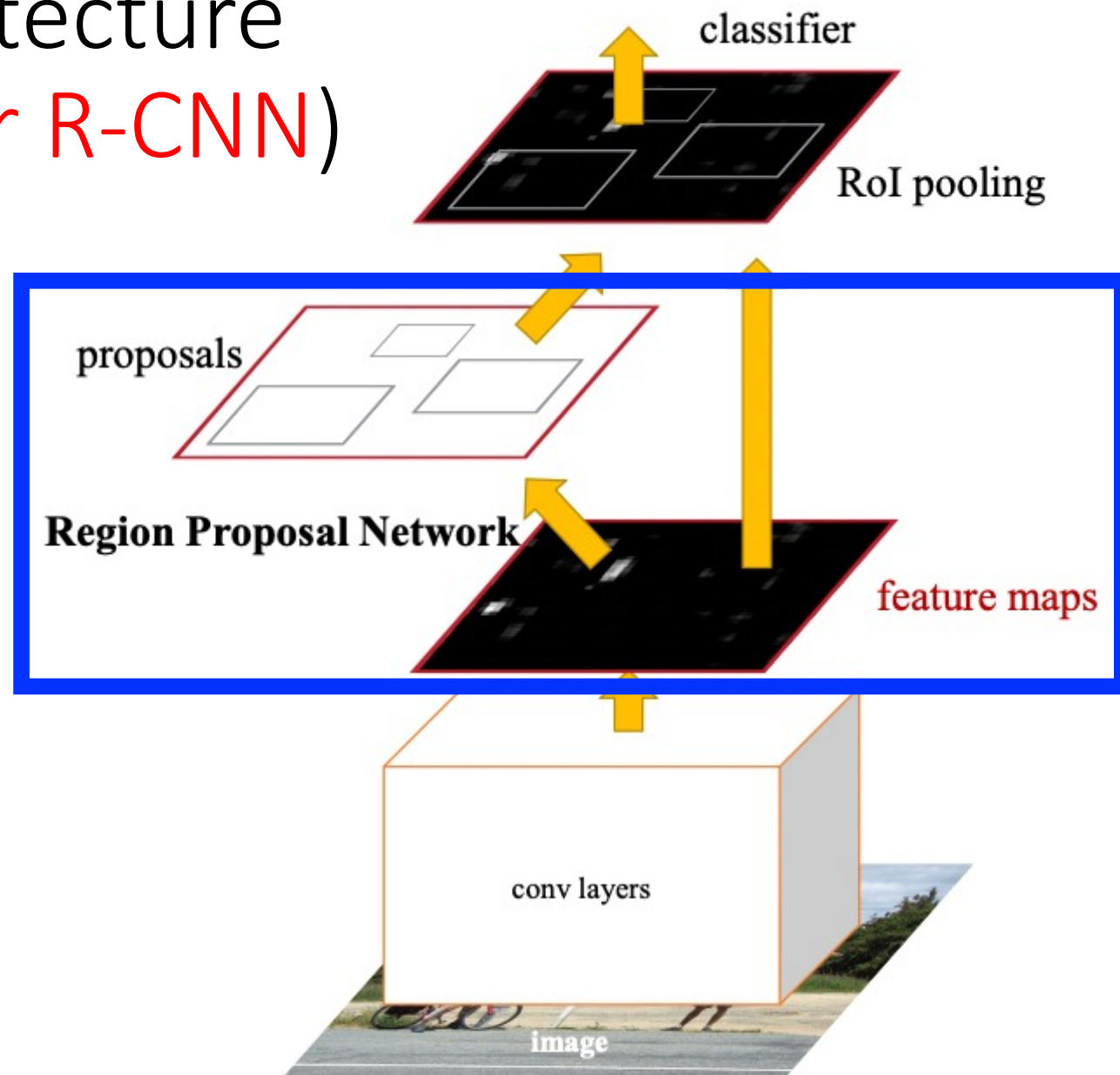
# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)



# Tracktor – Base Architecture (FPN Variant of **Faster R-CNN**)

- Recall:

Embeds a region proposal network in Fast R-CNN using a sliding window approach





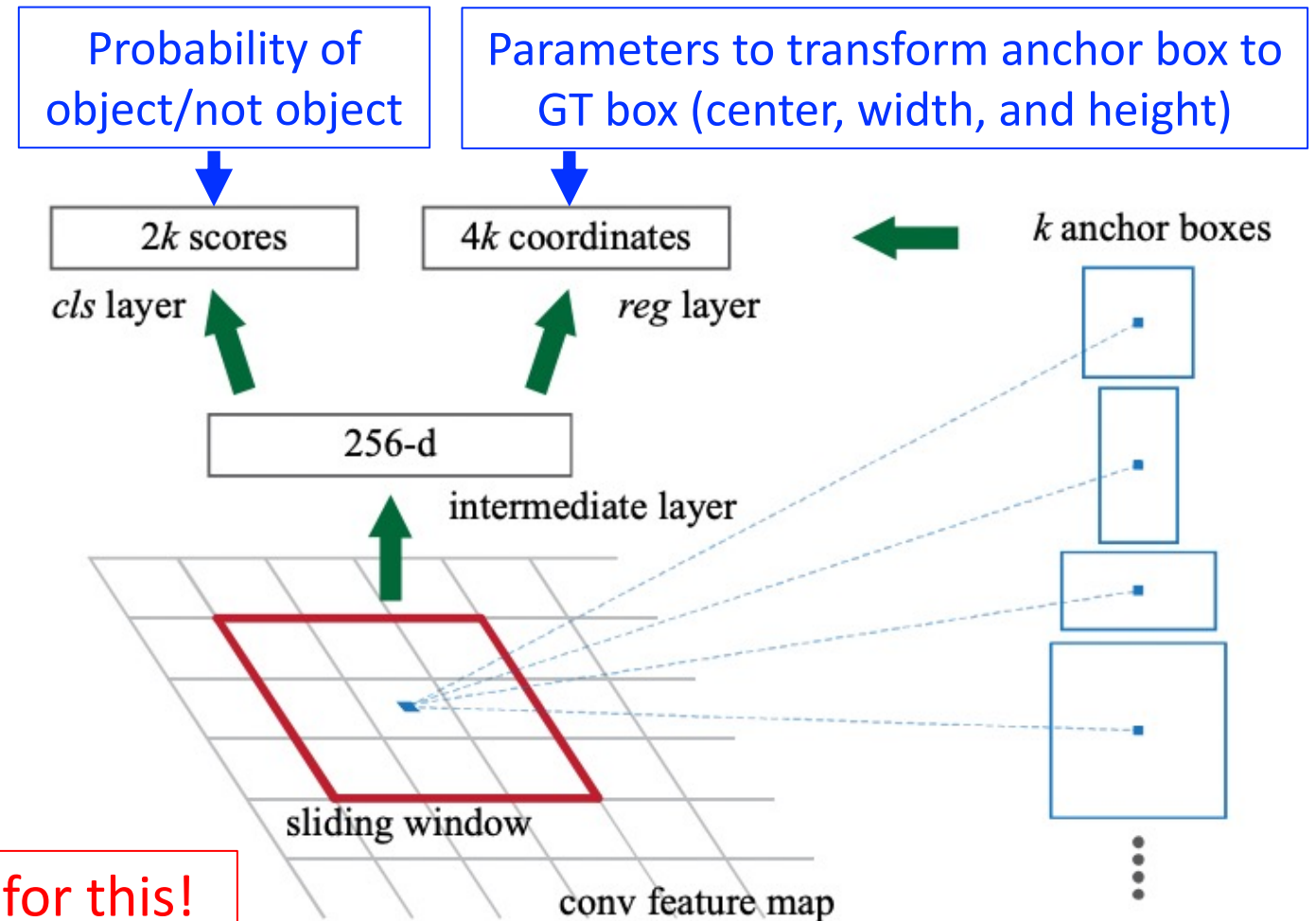
# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

- Recall:

Uses sliding window, since based on convolution

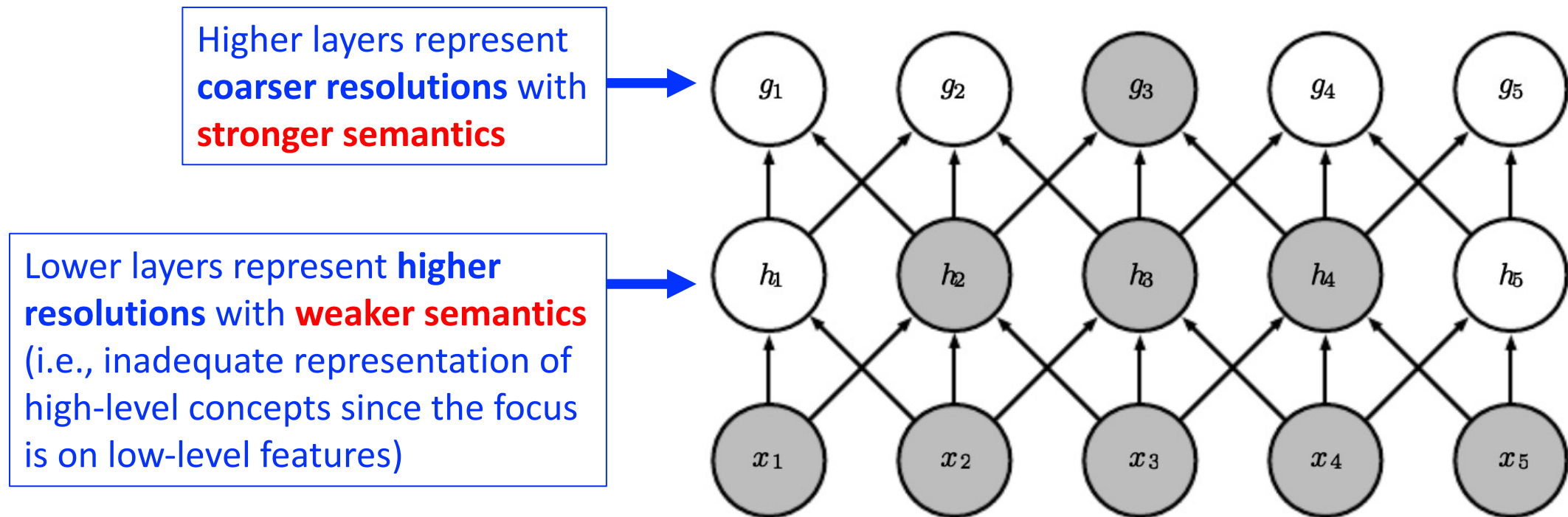
- At each sliding window position, region proposals are predicted with respect to an anchor point (i.e., center of sliding window position)
- At each anchor point,  $k = 9$  anchors are used to represent 3 aspect ratios and 3 scales

Variant addresses and so removes need for this!



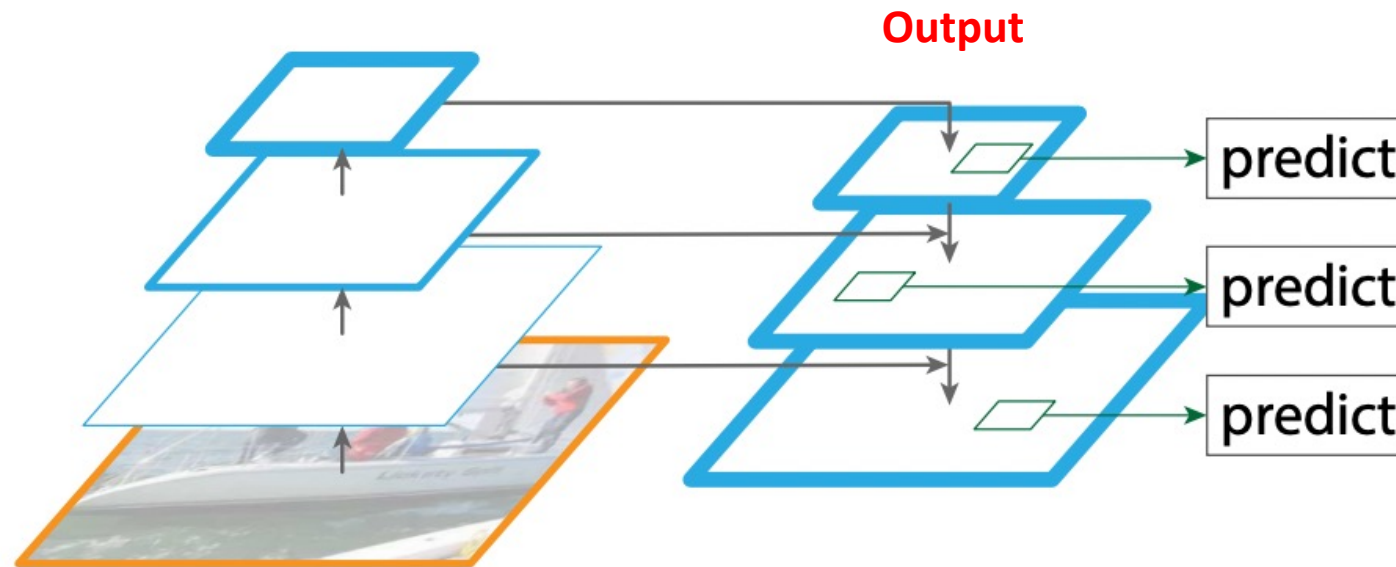
# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

- FPN – feature pyramid network
  - Transforms a convolutional network's pyramidal feature hierarchy to have **stronger semantics** at all scales, so that different object sizes are supported



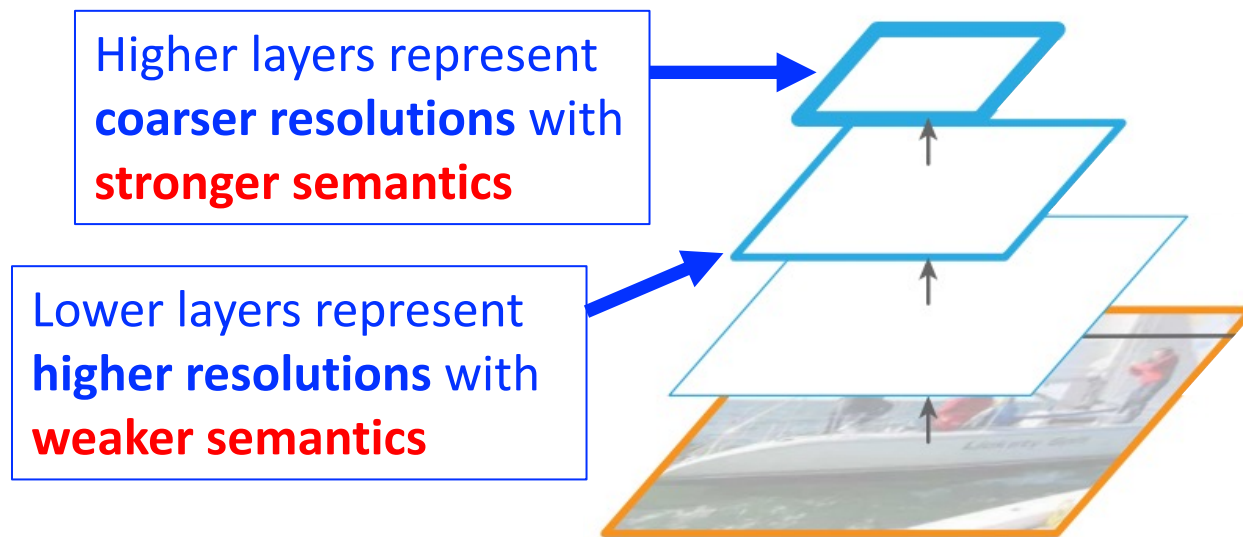
# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

- FPN – feature pyramid network
  - Transforms a convolutional network's pyramidal feature hierarchy to have **stronger semantics** at all scales, so that different object sizes are supported
  - Given a single image scale, its fully convolutional approach generates feature maps with **strong semantics at multiple levels** for use in a downstream task



# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

**Step 1.** Compute a feature hierarchy consisting of feature maps at several scales using your favorite backbone architecture (e.g., ResNet)



# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

**Step 1.** Compute a feature hierarchy consisting of feature maps at several scales using your favorite backbone architecture (e.g., ResNet)

**Step 2.** Fuse semantically stronger, coarser resolution feature maps with higher resolution, semantically weak features maps by upsampling the coarser resolution feature maps

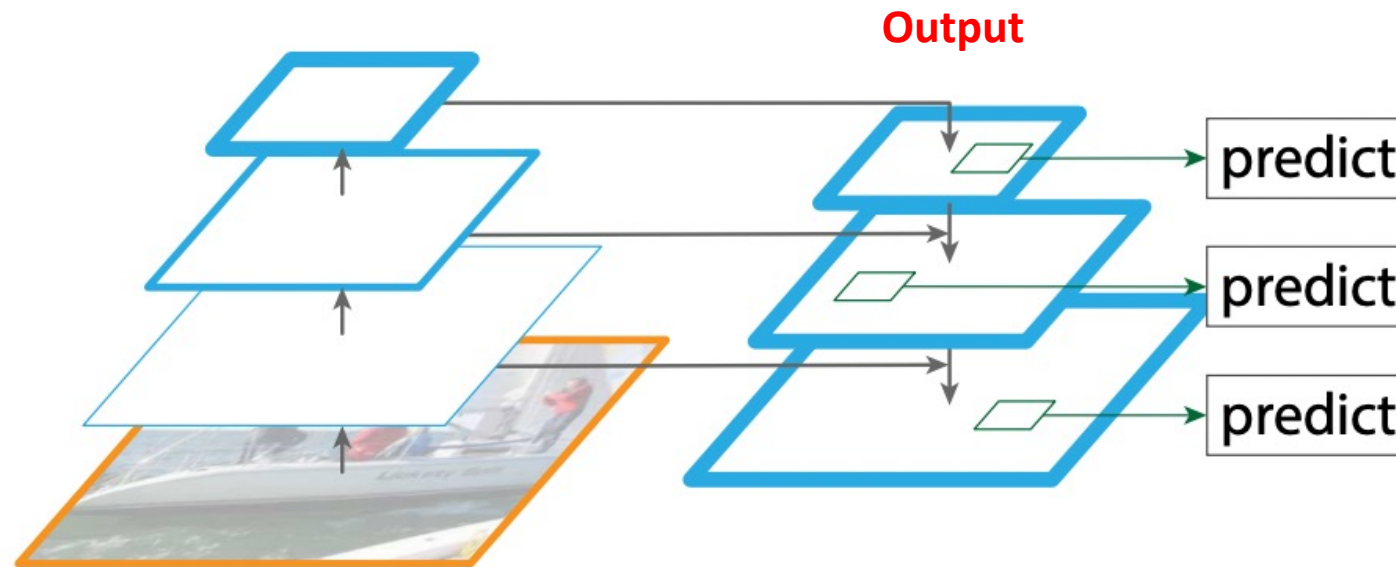
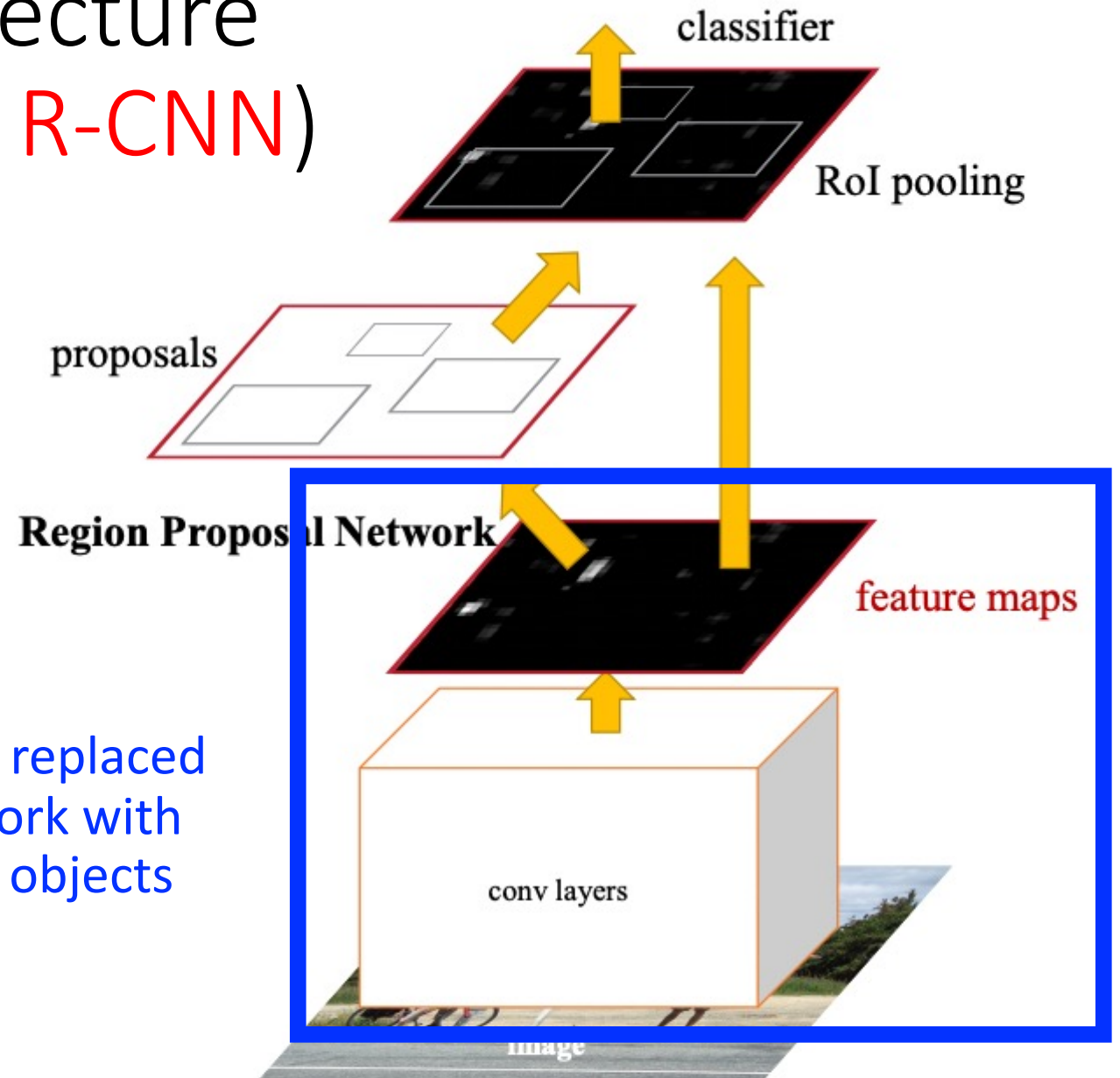


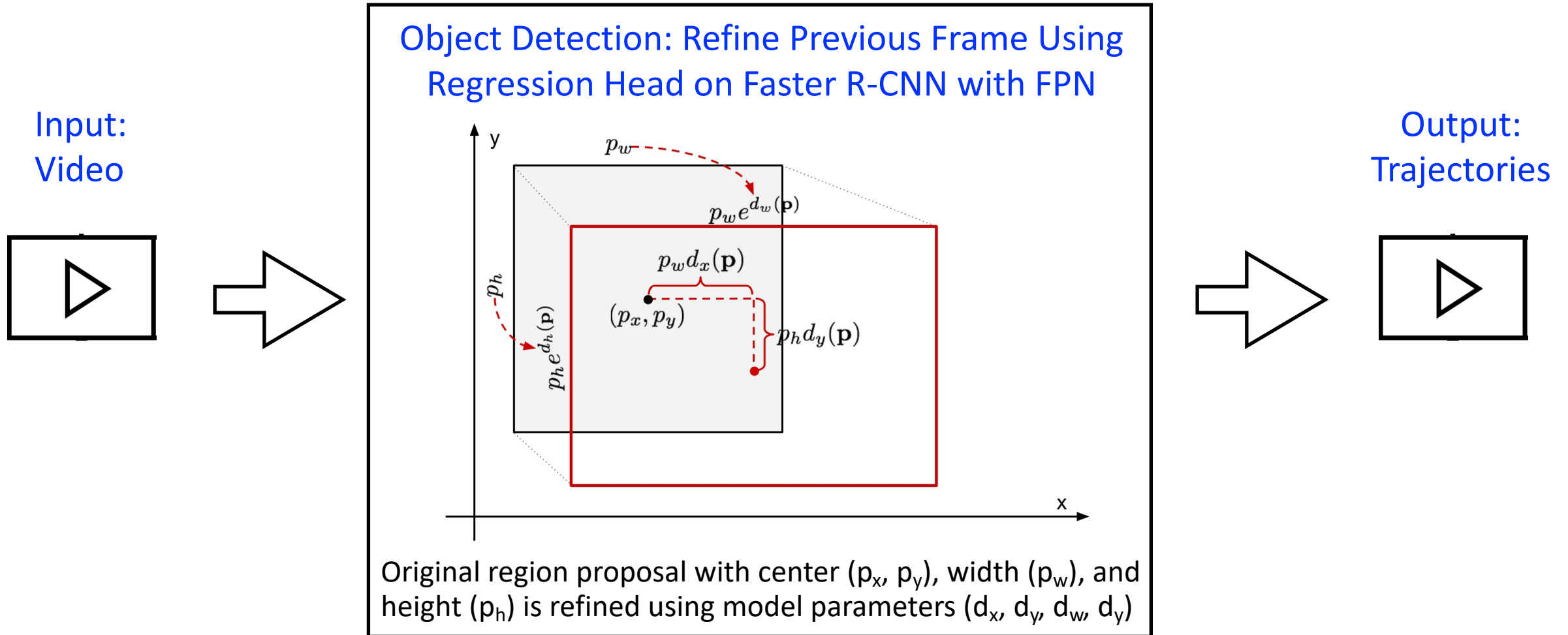
Figure source: Tsung-Yi Lin et al. "Feature Pyramid Networks for Object Detection." CVPR 2017.

# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

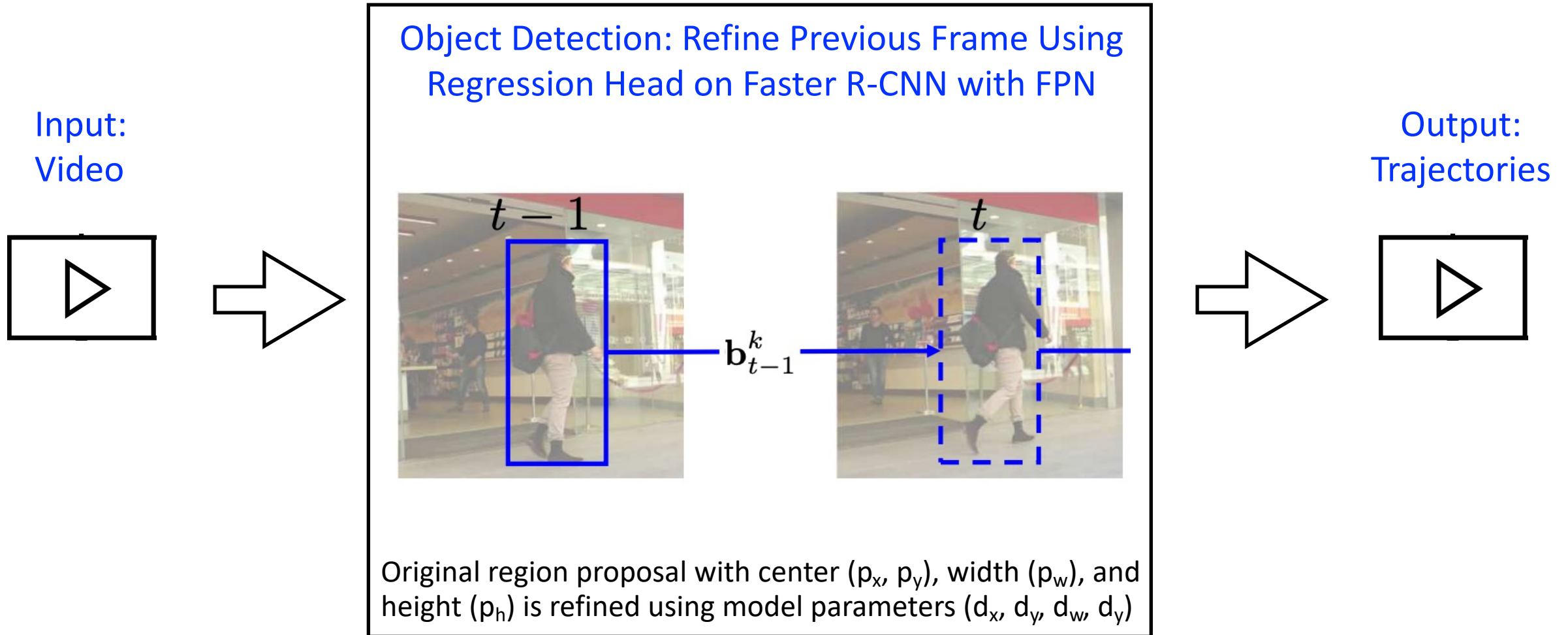


The single-scale feature map is replaced with the feature pyramid network with the aim to detect smaller sized objects

# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

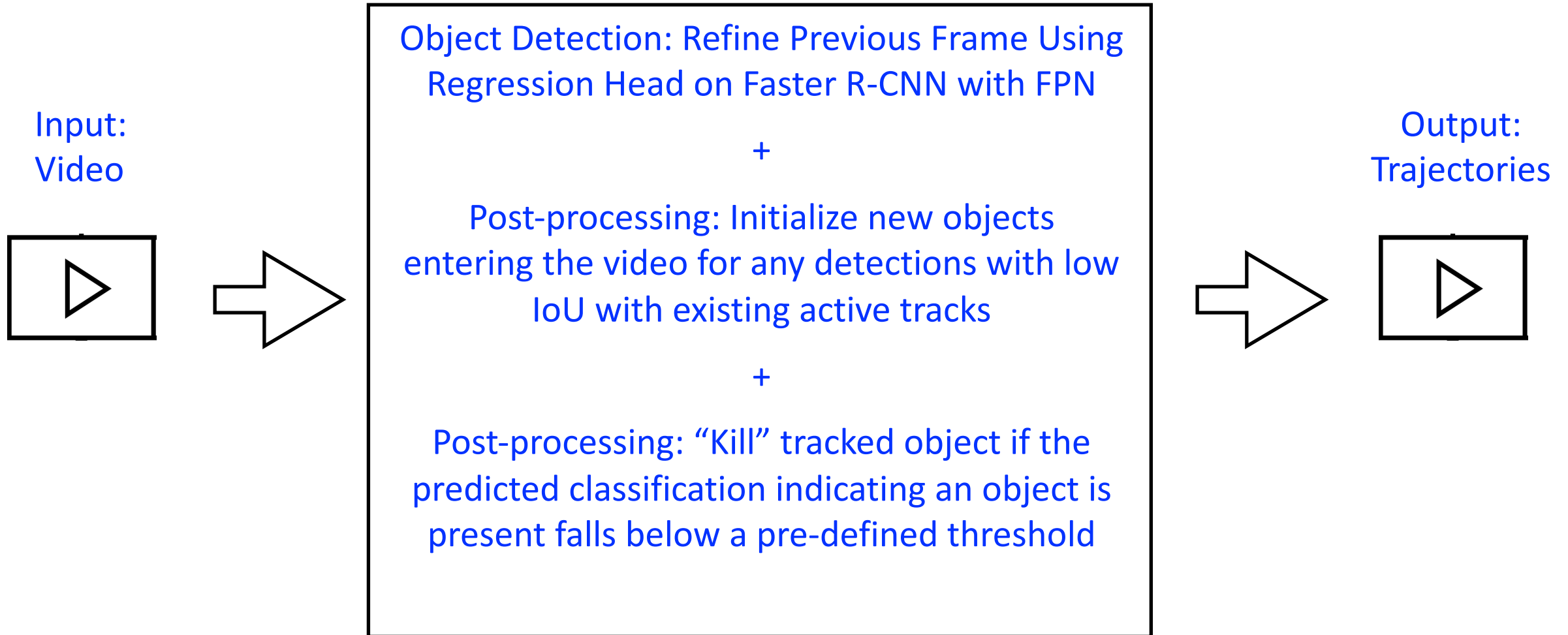


# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)





# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)



# Tracktor++ (i.e., with *More* Post-Processing)

1. Motion model to address when an object's position changes a lot between consecutive frames
  - For moving camera, apply image registration
  - For low frame rate, assume constant velocity for all objects
2. Reidentification to enable an object that disappears for a short time to be linked to itself when it re-appears
  - Compare detected object appearance of a “killed” object to newly tracked objects in future video frames using a Siamese network
  - Given two bounding boxes of a “killed” object and candidate detection with high IoU, the Siamese network computes their similarity to determine whether the two tracks should be linked (i.e., if similarity exceeds a pre-defined threshold)

# Tracktor++ Performance

State-of-art performance on three datasets with respect to MOTA!

	Method	MOTA ↑
MOT17	Tracktor++	<b>53.5</b>
	eHAF [58]	51.8
	FWT [23]	51.3
	jCC [30]	51.2
	MOTDT17 [9]	50.9
	MHT_DAM [32]	50.7
MOT16	Tracktor++	<b>54.4</b>
	HCC [44]	49.3
	LMP [59]	48.8
	GCRA [43]	48.2
	FWT [23]	47.8
	MOTDT [9]	47.6
2D MOT 2015	Tracktor++	<b>44.1</b>
	AP_HWDPL_p [8]	38.5
	AMIR15 [56]	37.6
	JointMC [30]	35.6
	RAR15pub [17]	35.1

# Ablation Study of Tracktor++


- Test set: MOT17 which consists of 7 sequences

Method	MOTA $\uparrow$
D&T [18]	50.1
Tracktor-no-FPN	57.4
Tracktor	61.5
Tracktor+reID	61.5
Tracktor+CMC	<b>61.9</b>
Tracktor++ (reID + CMC)	<b>61.9</b>

Greatest boost in performance comes from using a feature pyramid network



Remainder of performance boost stems from the motion model



# What Makes MOT Difficult?

- When targets have diminished visibility (i.e., from occlusion)
- When objects are small
- When there is a large gap for a tracked object (i.e., missed detections)

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models

A dark gray background with a central circular glow. The glow is a gradient from light gray in the center to dark gray at the edges. The text "The End" is centered within this glow. The entire scene is framed by a white film strip border with sprocket holes on the left and right sides.

*The End*