

# Object Tracking

**Danna Gurari**

University of Colorado Boulder

Fall 2024



# Review

- Last lecture: instance segmentation
  - Motivation
  - Datasets
  - Evaluation metric
  - Mask R-CNN
  - YOLACT
- Assignments (Canvas)
  - Reading assignment was due earlier today
  - Next reading assignments due for next two lectures
  - Project outline due in 1.5 weeks
- Questions?

# Object Tracking: Today's Topics

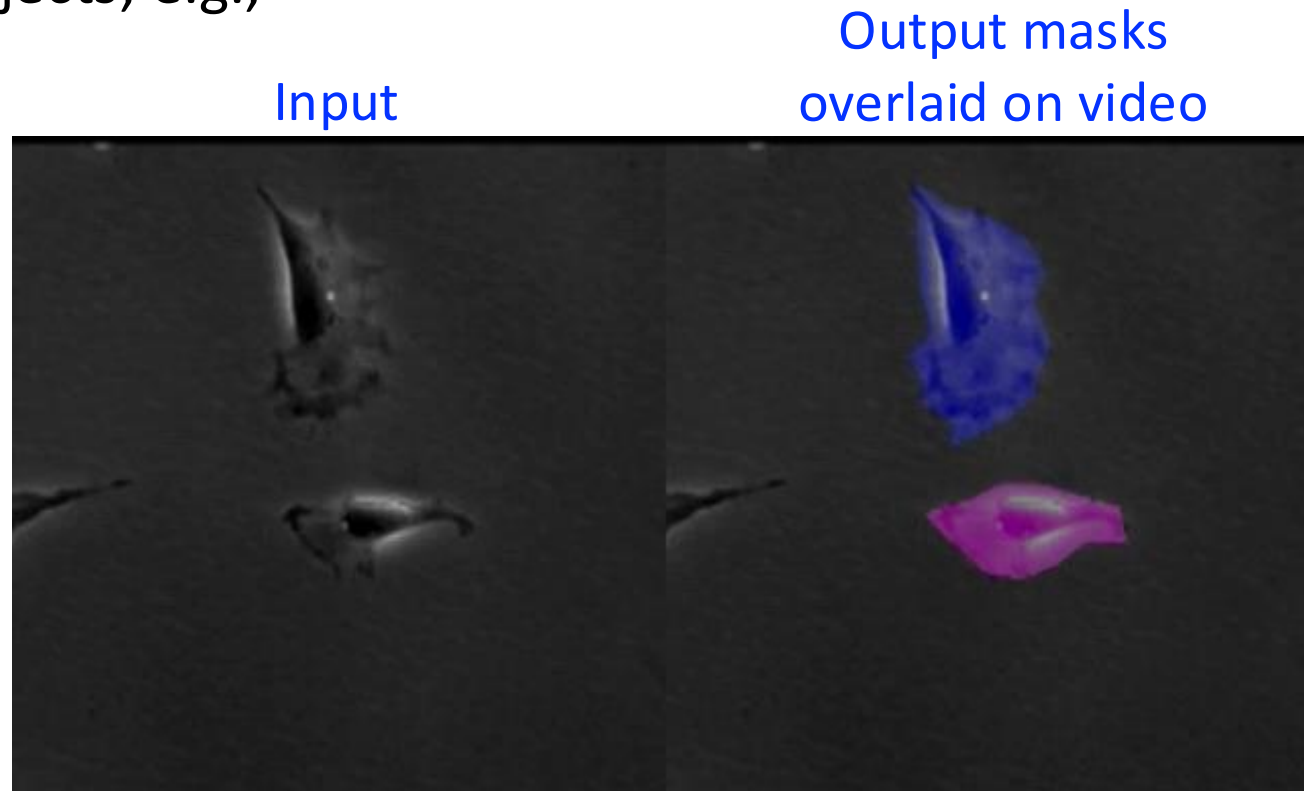
- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Definition

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



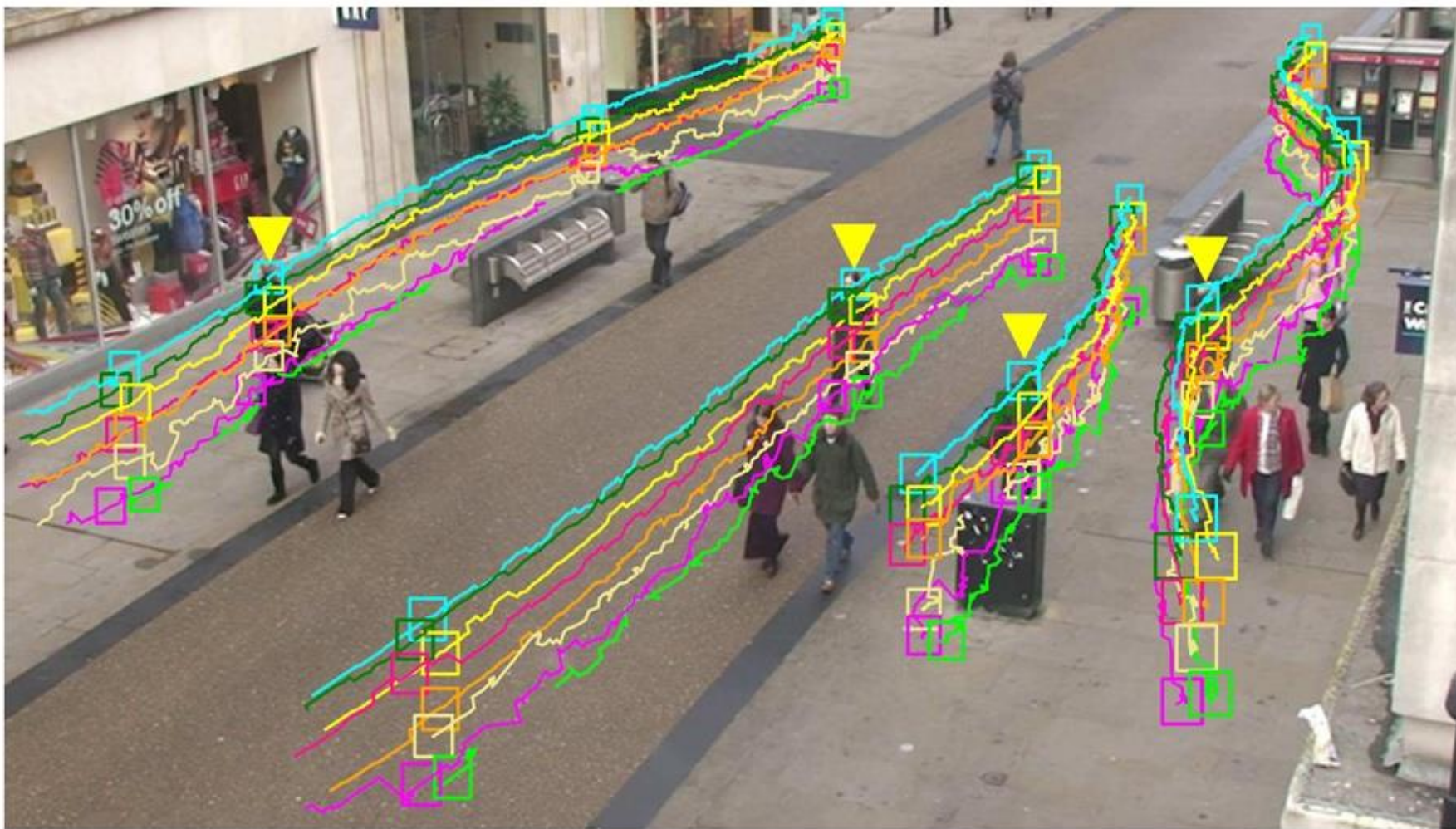
# Definition

- 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

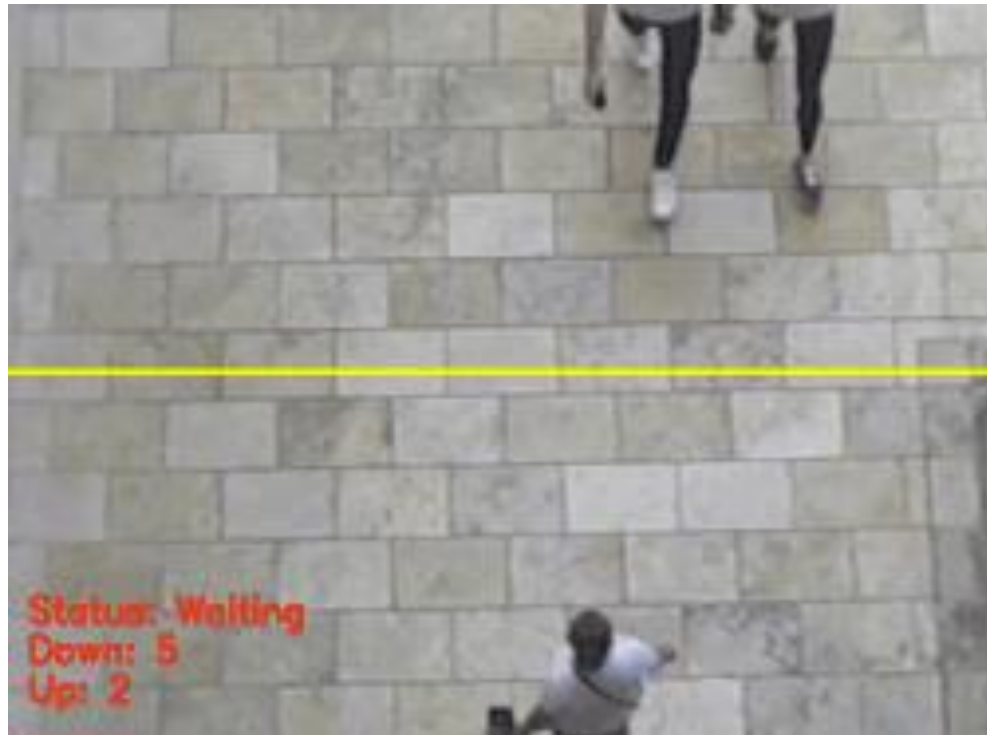
- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Surveillance

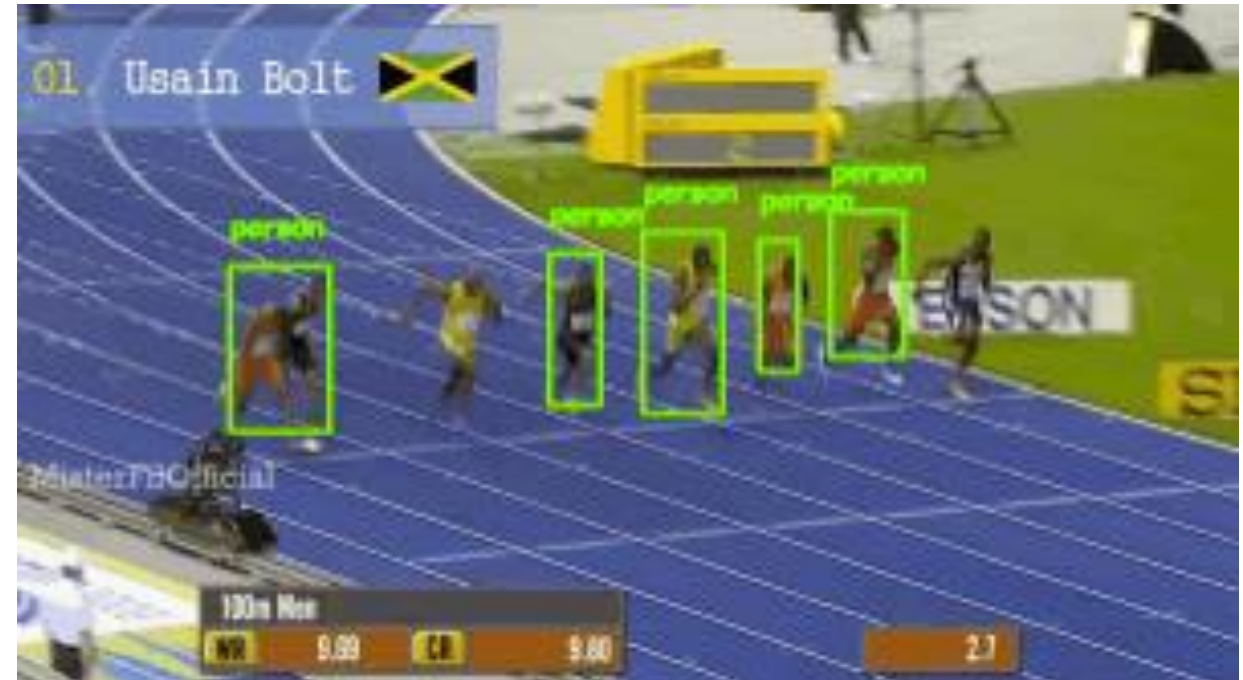




# Business Marketing: People Analytics



# 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!



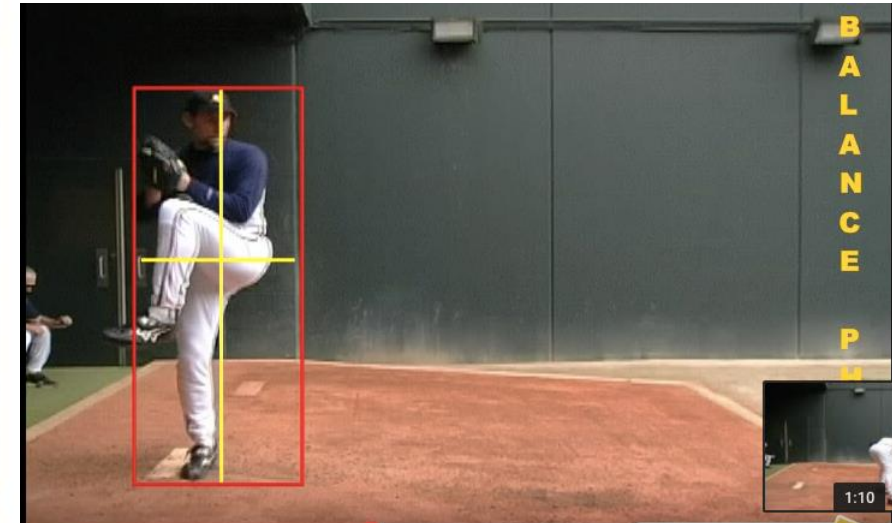
NEW! Track Bowling Ball Path!



Works great for putting!



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



# Military Defense



# Self-driving Cars

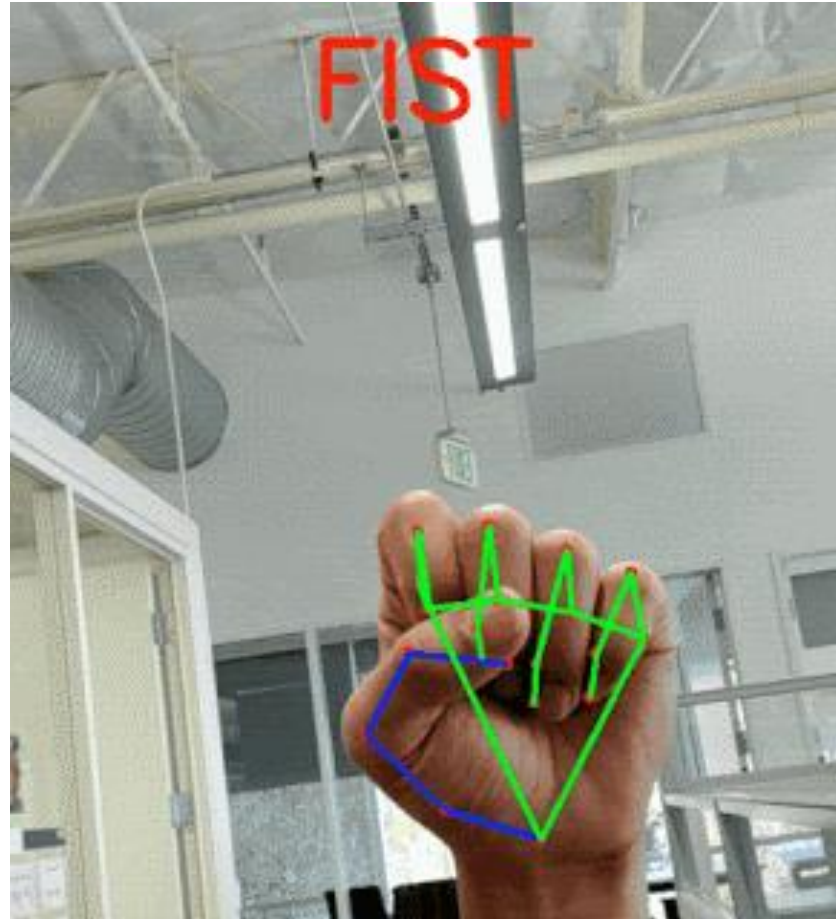


# Human Computer Interaction



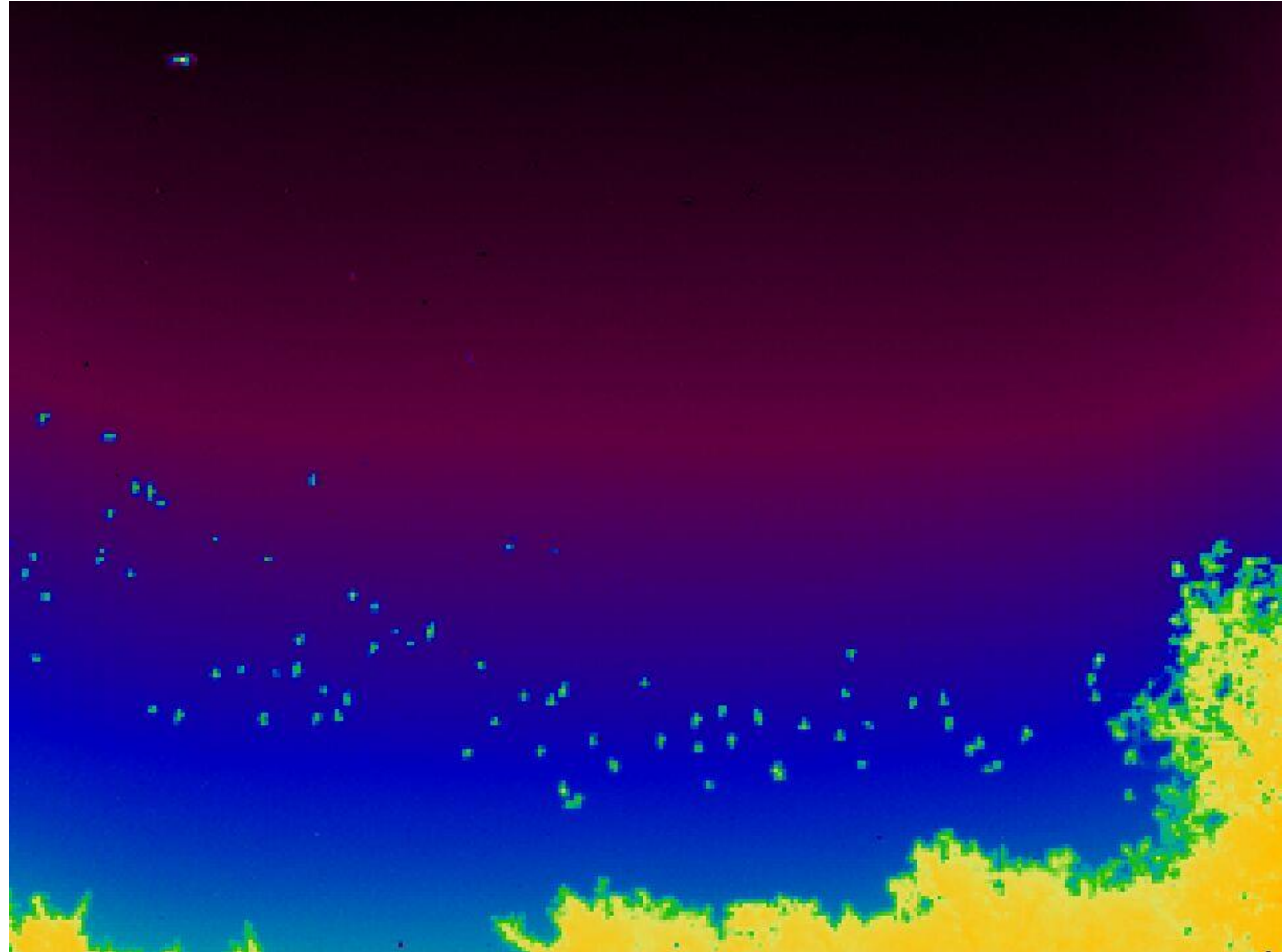
Roboceptionist

# Sign Language Recognition



# Biological Monitoring

Counting bats exiting  
a cave in Texas:





# Augmented Reality



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

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

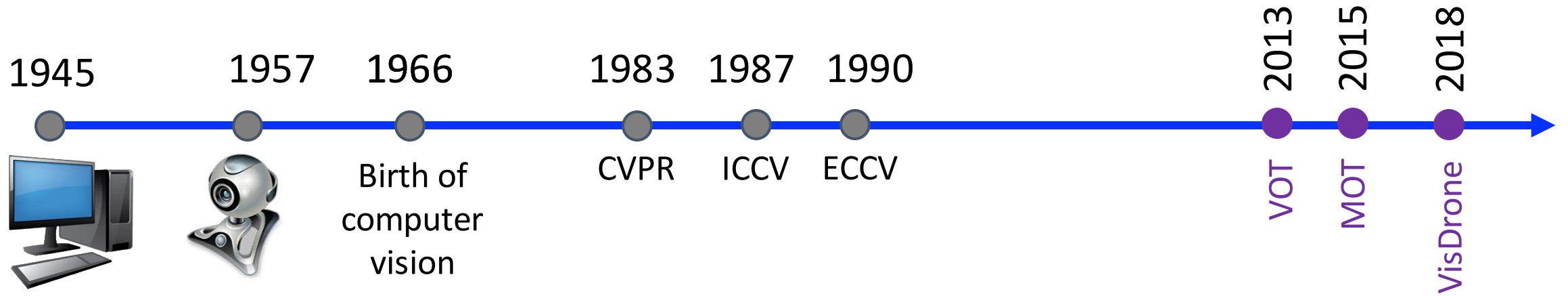
# Applications

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

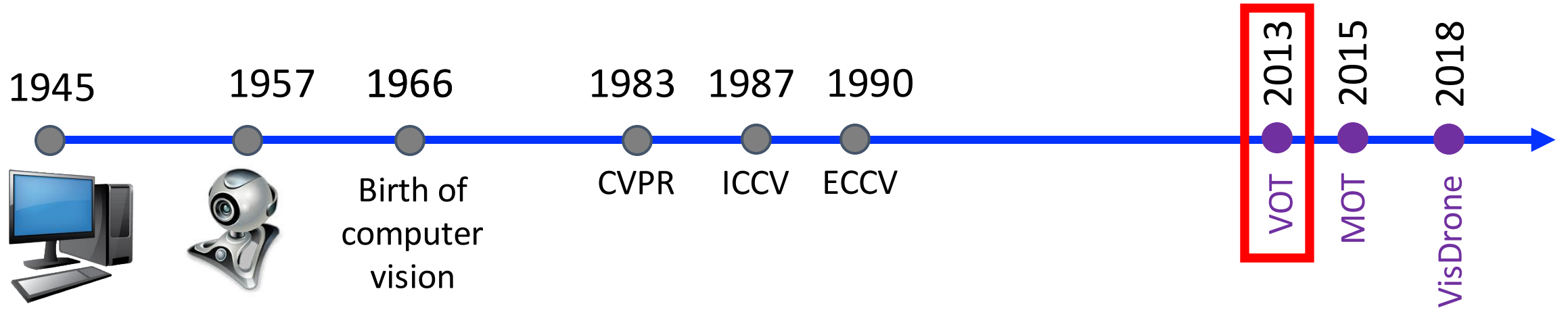
# Object Tracking: Today's Topics

- Problem
- Applications
- **Datasets**
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Object Tracking Datasets



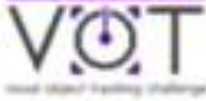





# Object Tracking Datasets



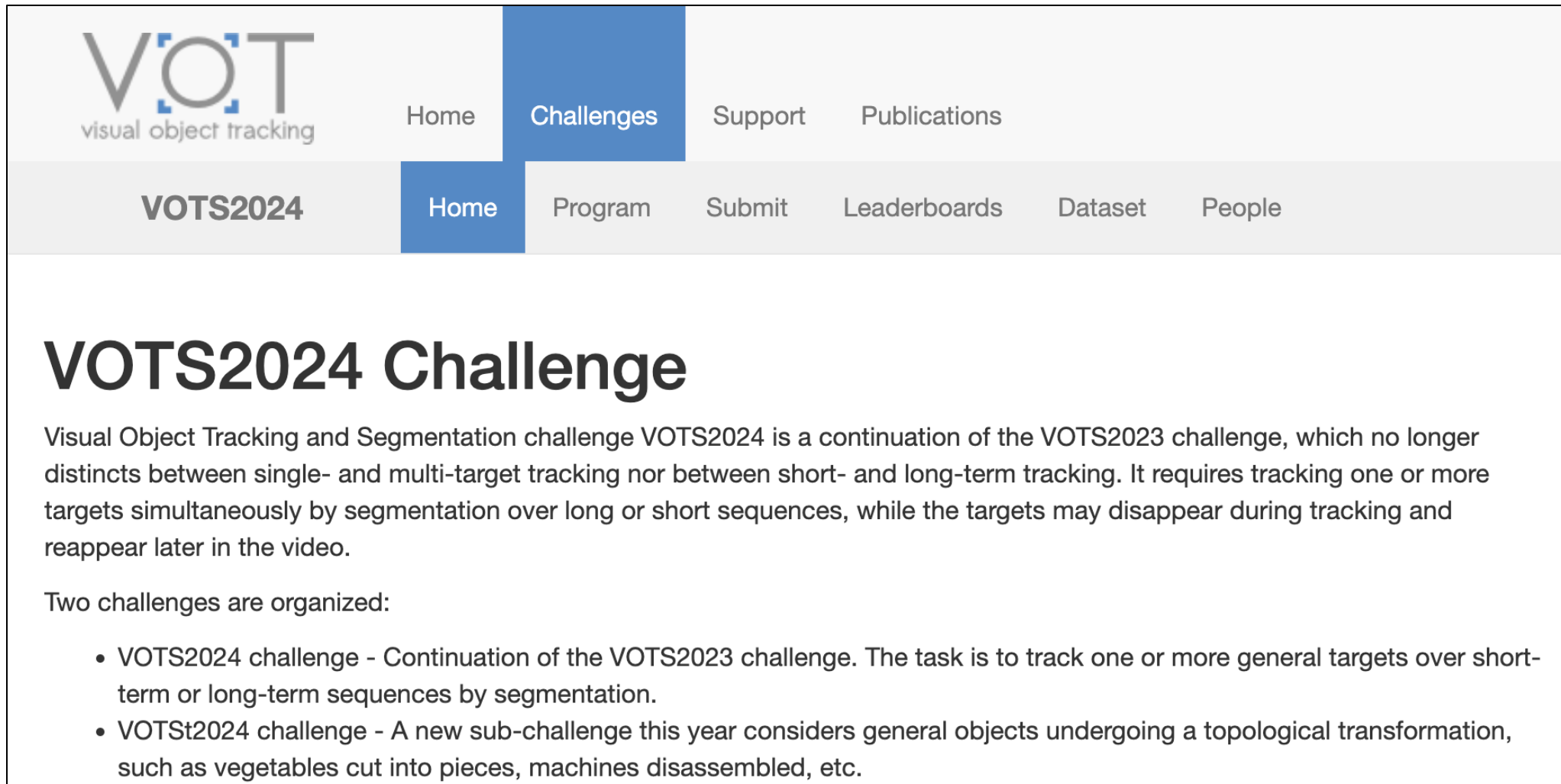
# Single Object Tracking Dataset: VOT

- 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 believed had 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 (12<sup>th</sup> year in 2024)



The image shows a screenshot of the VOTS2024 Challenge website. The top navigation bar includes links for Home, Challenges (highlighted in blue), Support, and Publications. Below this, a secondary navigation bar features VOTS2024, Home (highlighted in blue), Program, Submit, Leaderboards, Dataset, and People. The main content area displays the title 'VOTS2024 Challenge' and a descriptive paragraph about the challenge. It also lists two organized challenges: VOTS2024 and VOTSt2024.

**VOTS2024**

## VOTS2024 Challenge

Visual Object Tracking and Segmentation challenge VOTS2024 is a continuation of the VOTS2023 challenge, which no longer distinguishes between single- and multi-target tracking nor between short- and long-term tracking. It requires tracking one or more targets simultaneously by segmentation over long or short sequences, while the targets may disappear during tracking and reappear later in the video.

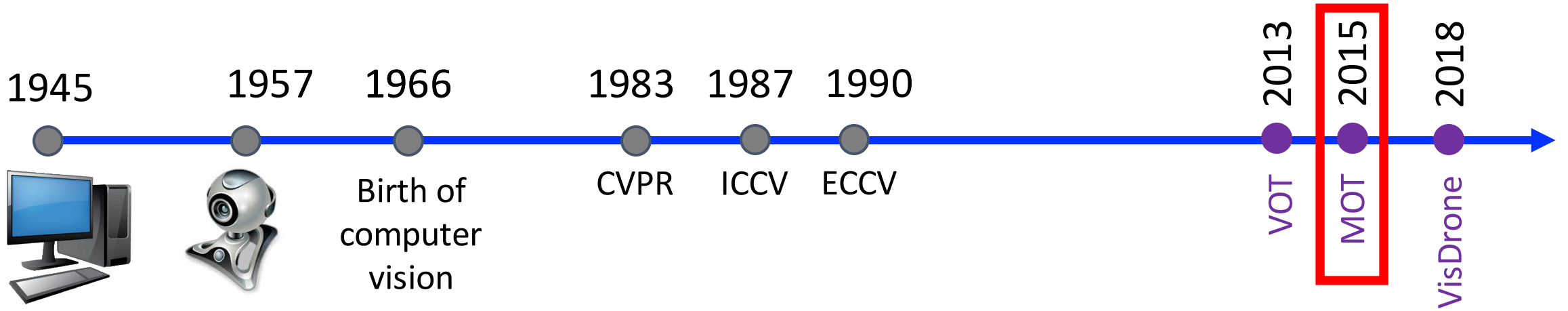
Two challenges are organized:

- VOTS2024 challenge - Continuation of the VOTS2023 challenge. The task is to track one or more general targets over short-term or long-term sequences by segmentation.
- VOTSt2024 challenge - A new sub-challenge this year considers general objects undergoing a topological transformation, such as vegetables cut into pieces, machines disassembled, etc.

<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 vs on ground looking up)
  - **Multiple weather conditions**; e.g., sunny vs cloudy vs night time
- 16 videos from existing datasets and other 6 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 vs 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

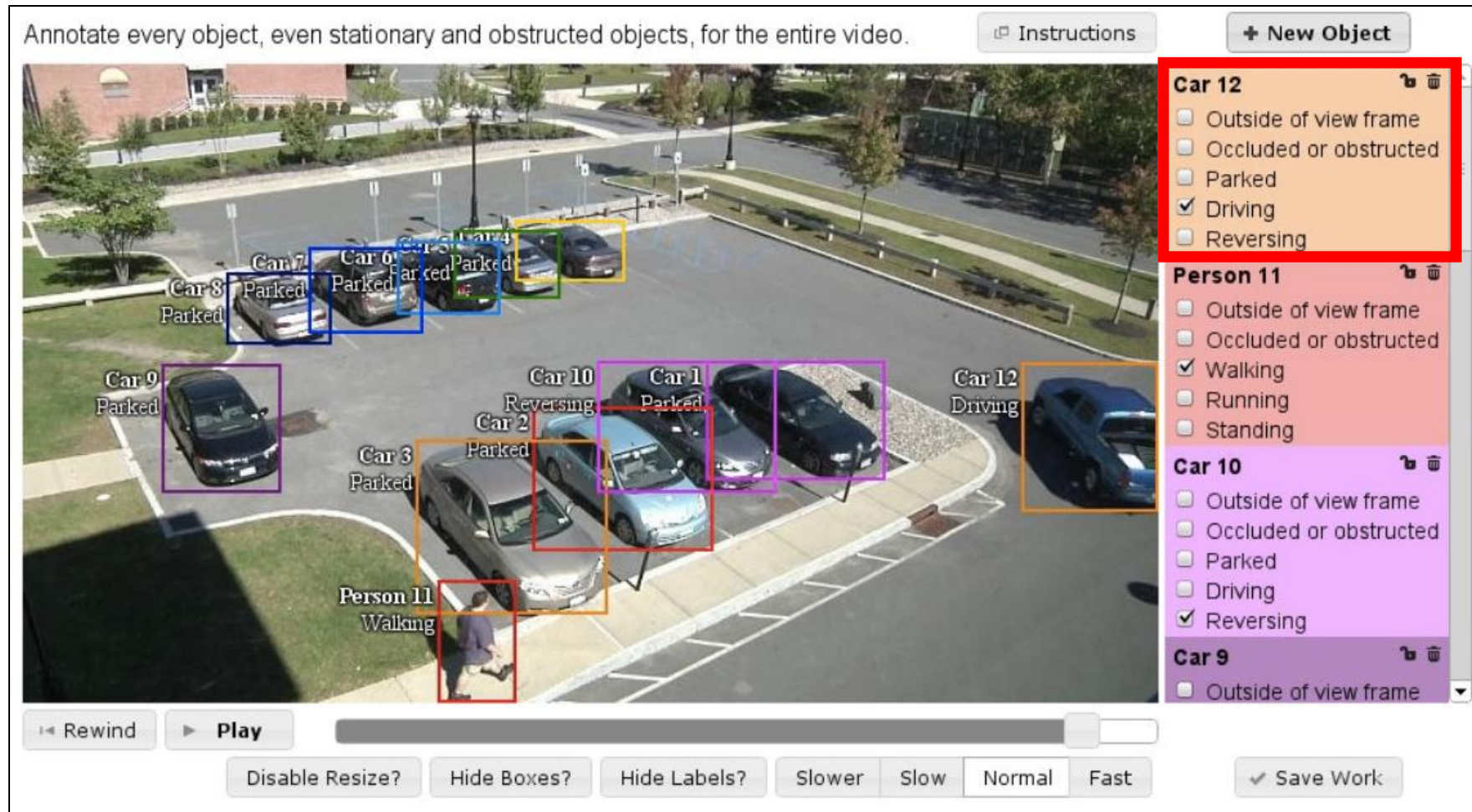
# Multiple Object Tracking Annotation: VATIC



Demo: <https://www.youtube.com/watch?v=ljl5pAowACc>

# Multiple Object Tracking Annotation: VATIC

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



The screenshot displays the VATIC annotation interface. The main window shows a video frame of a parking lot with several cars and a person. Each object is enclosed in a colored bounding box and labeled with its ID and activity. The sidebar on the right provides a detailed view of the metadata for selected objects. The 'Car 12' entry is highlighted with a red border.

Object ID	Activity	Attributes
Car 12	Driving	Outside of view frame, Occluded or obstructed, Parked, Reversing
Person 11	Walking	Outside of view frame, Occluded or obstructed, Running, Standing
Car 10	Reversing	Outside of view frame, Occluded or obstructed, Parked, Driving
Car 9	Outside of view frame	

Car 12 metadata (highlighted in red):

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

Person 11 metadata:

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

Car 10 metadata:

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

Car 9 metadata:

- Outside of view frame

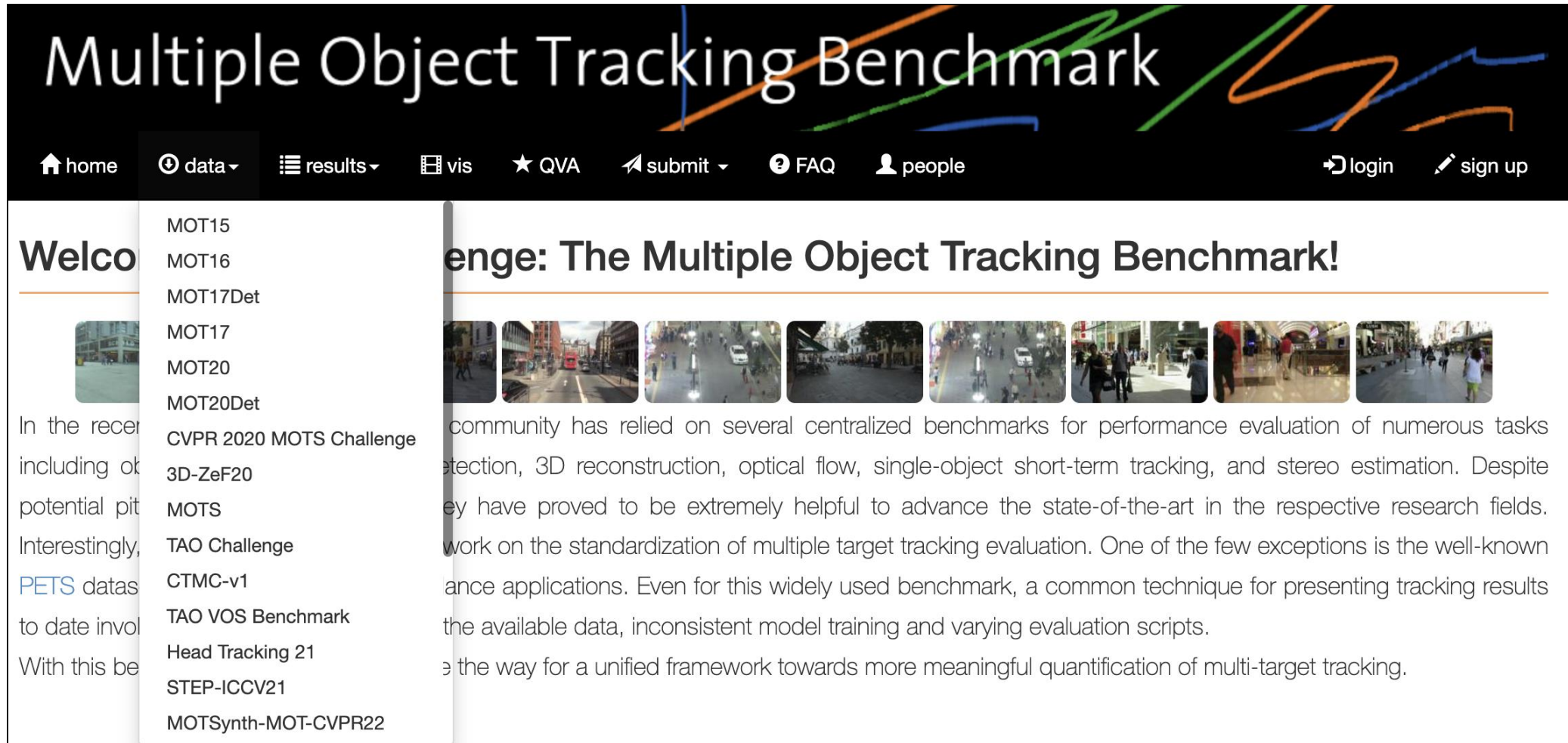
Video player controls: Rewind, Play, Slower, Slow, Normal, Fast, Save Work

Metadata about each object includes activity and attributes

# Multiple Object Tracking Annotation: VATIC

- How to handle occlusions?
  - 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"*
  - "**visibility**" flag: 0-1 with 1 fully visible and less than 1 indicating occlusions
  - "**confidence**" flag: 1 when box should be considered for evaluation and 0 otherwise (e.g., a pedestrian is too small)
  - Non-tracked categories: "**class**" value is occluder and ignored during evaluation

# Multiple Object Tracking Annual Challenge (10<sup>th</sup> year in 2024)



The screenshot shows the homepage of the Multiple Object Tracking Benchmark website. The header features the title "Multiple Object Tracking Benchmark" in white text on a black background, with colorful abstract lines. Below the header is a navigation bar with links for home, data, results, vis, QVA, submit, FAQ, people, login, and sign up. The main content area is divided into three columns. The left column has a "Welcome" section with a small image and a list of challenge datasets. The middle column has a "Challenge: The Multiple Object Tracking Benchmark!" section with a row of eight small images showing various tracking scenarios. The right column contains a paragraph of text describing the challenge and its goals.

## Multiple Object Tracking Benchmark

home data results vis QVA submit FAQ people login sign up

### Welcome

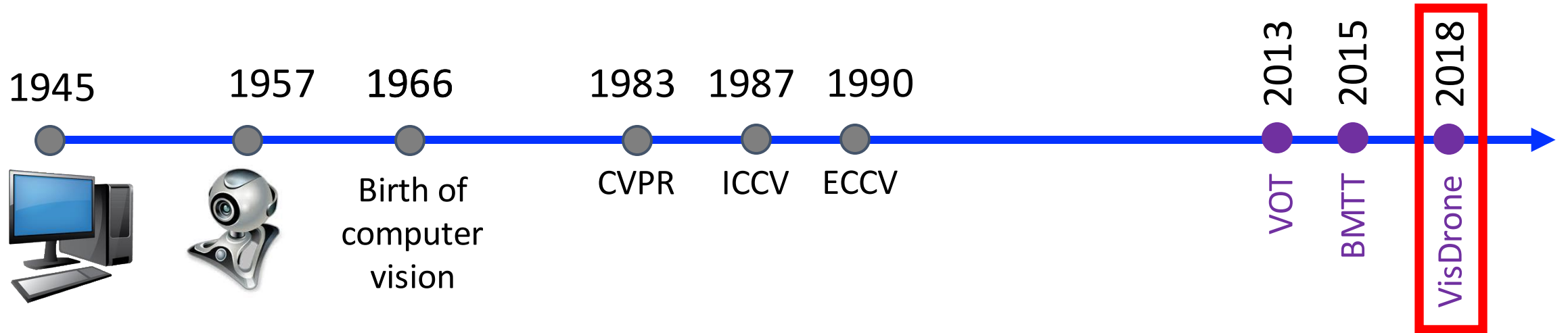
- MOT15
- MOT16
- MOT17Det
- MOT17
- MOT20
- MOT20Det
- CVPR 2020 MOTS Challenge
- 3D-ZeF20
- MOTS
- TAO Challenge
- CTMC-v1
- TAO VOS Benchmark
- Head Tracking 21
- STEP-ICCV21
- MOTSynth-MOT-CVPR22

### Challenge: The Multiple Object Tracking Benchmark!

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. We are now working on the standardization of multiple target tracking evaluation. One of the few exceptions is the well-known MOT17 challenge applications. Even for this widely used benchmark, a common technique for presenting tracking results is to use the available data, inconsistent model training and varying evaluation scripts. We aim to provide 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, without description of how 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

#00073 #00098 #00143  
#00250 #00286 #00342

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

# Discussion

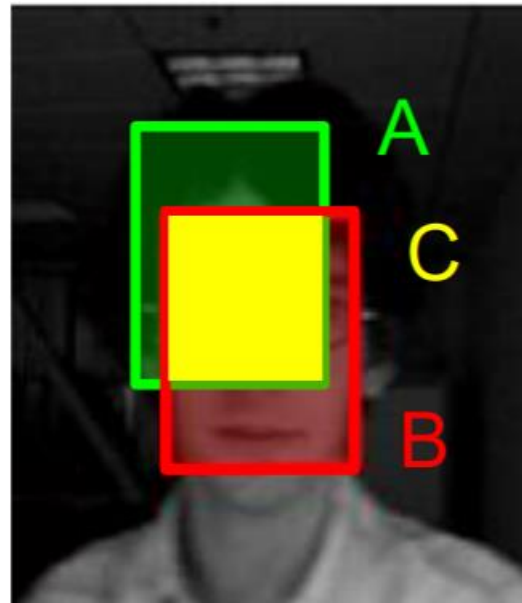
- When designing an annotation protocol, how should these scenarios be handled:
  - Partially visible object
  - Occluded object
  - Object is reflected in reflective surfaces such as mirrors or windows

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Accuracy

Average IoU from a tracker across all video frames



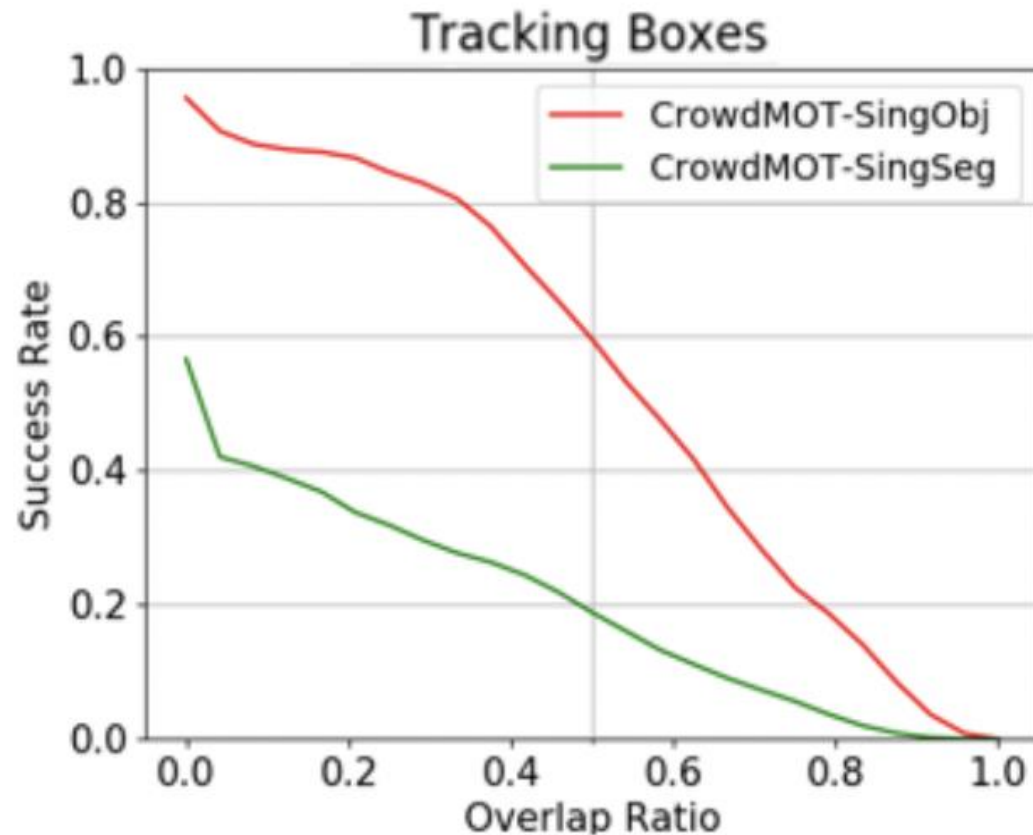
A = Ground Truth  
B = Predicted Track  
C = Intersection

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

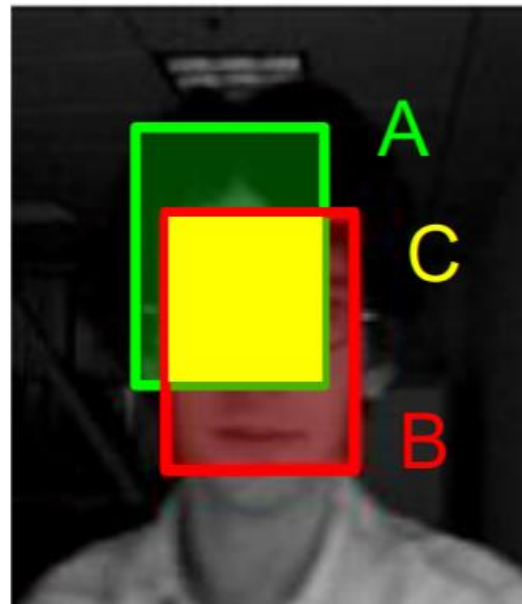
# Success Plot

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



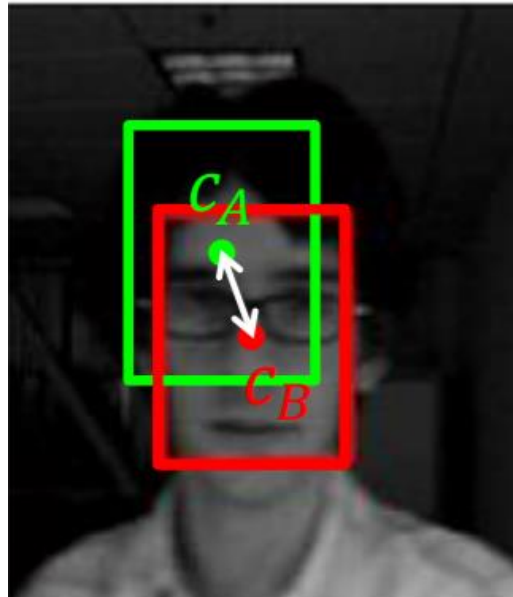
A = Ground Truth  
B = Predicted Track  
C = Intersection

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



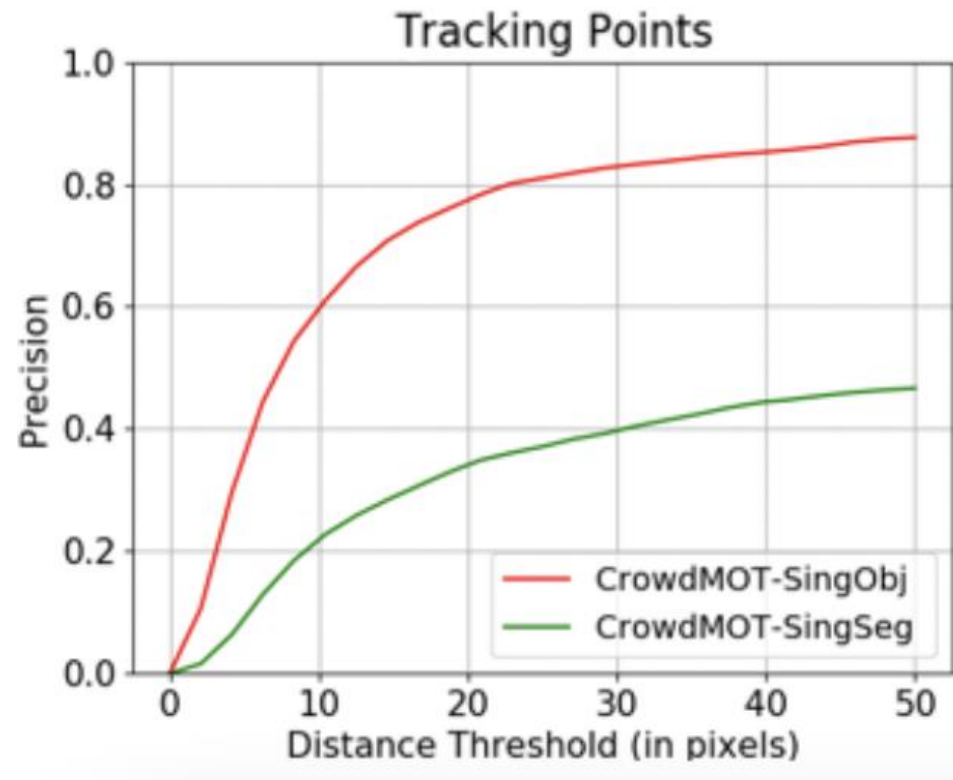
A = Ground Truth  
B = Predicted Track

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

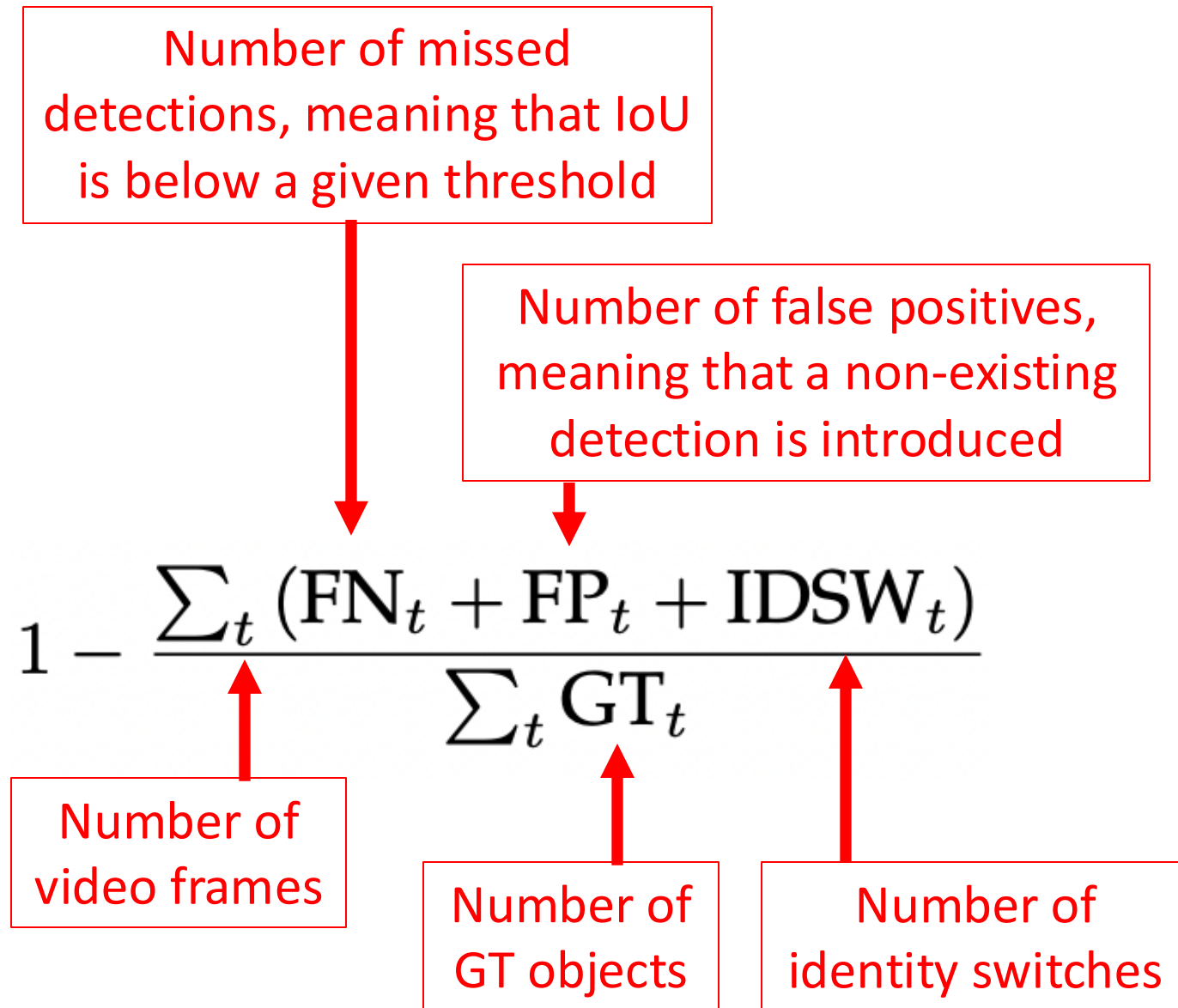


# Precision Plot

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



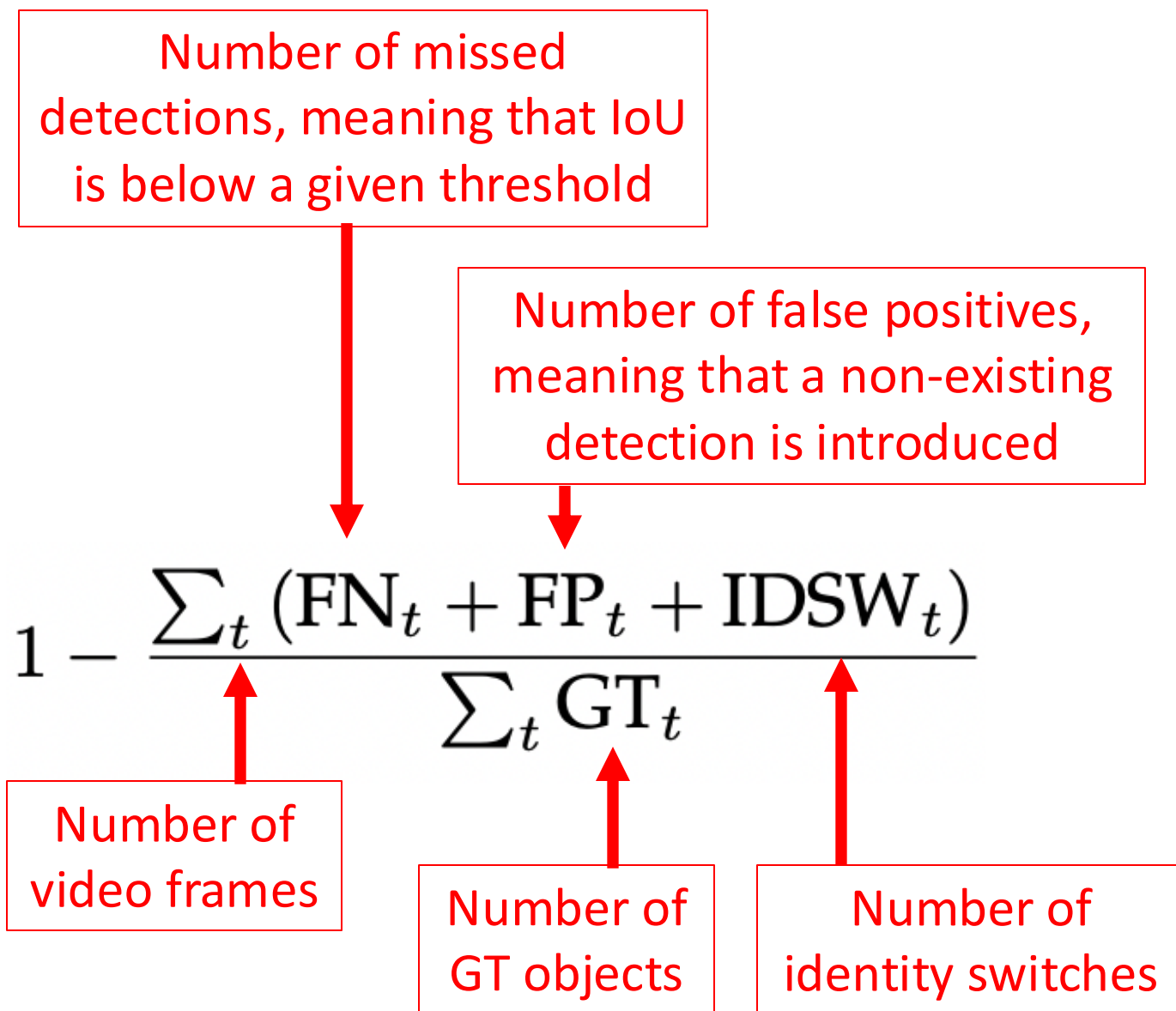
# MOTA

What is the range of possible values?

- $(-\infty, 100]$  (original value usually multiplied by 100)

When is MOTA negative?

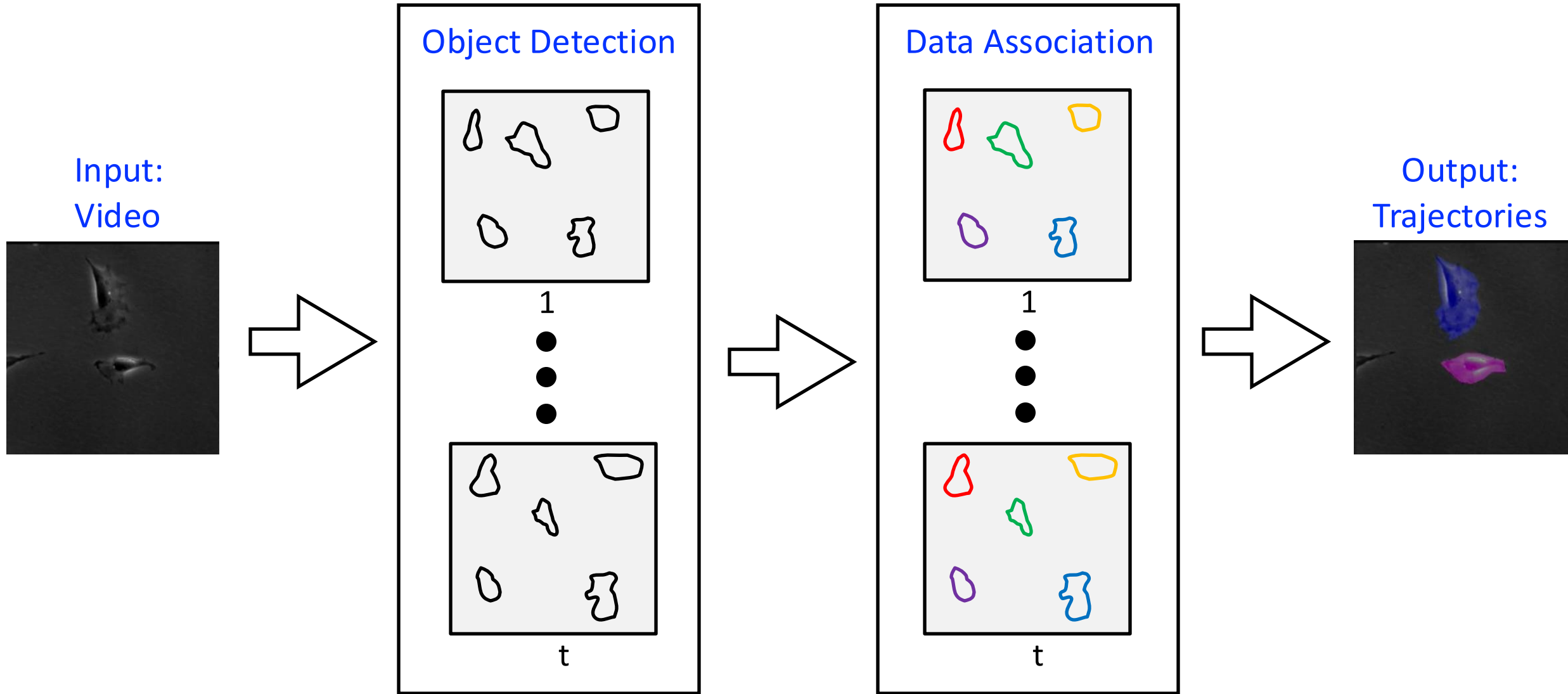
- When the number of errors exceed the number of objects in the frames



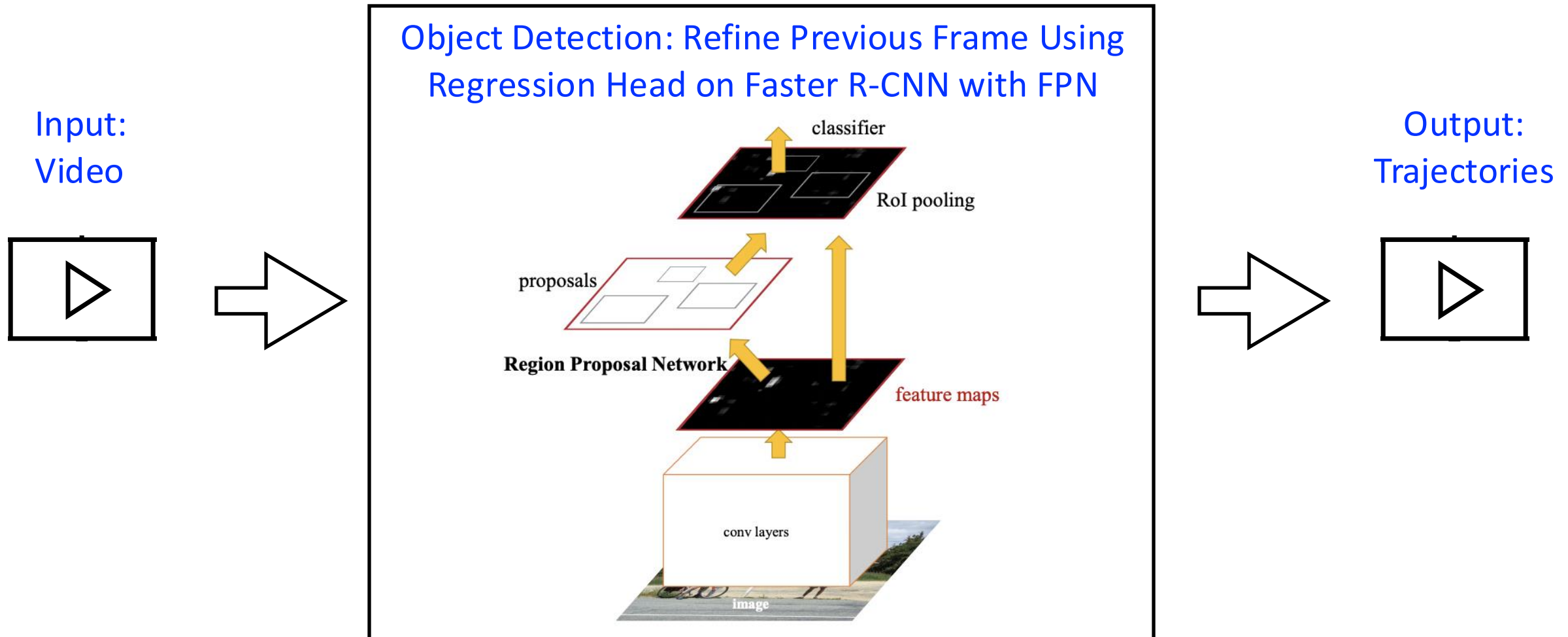
# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- **Computer vision models**
- Discussion (chosen by YOU 😊)

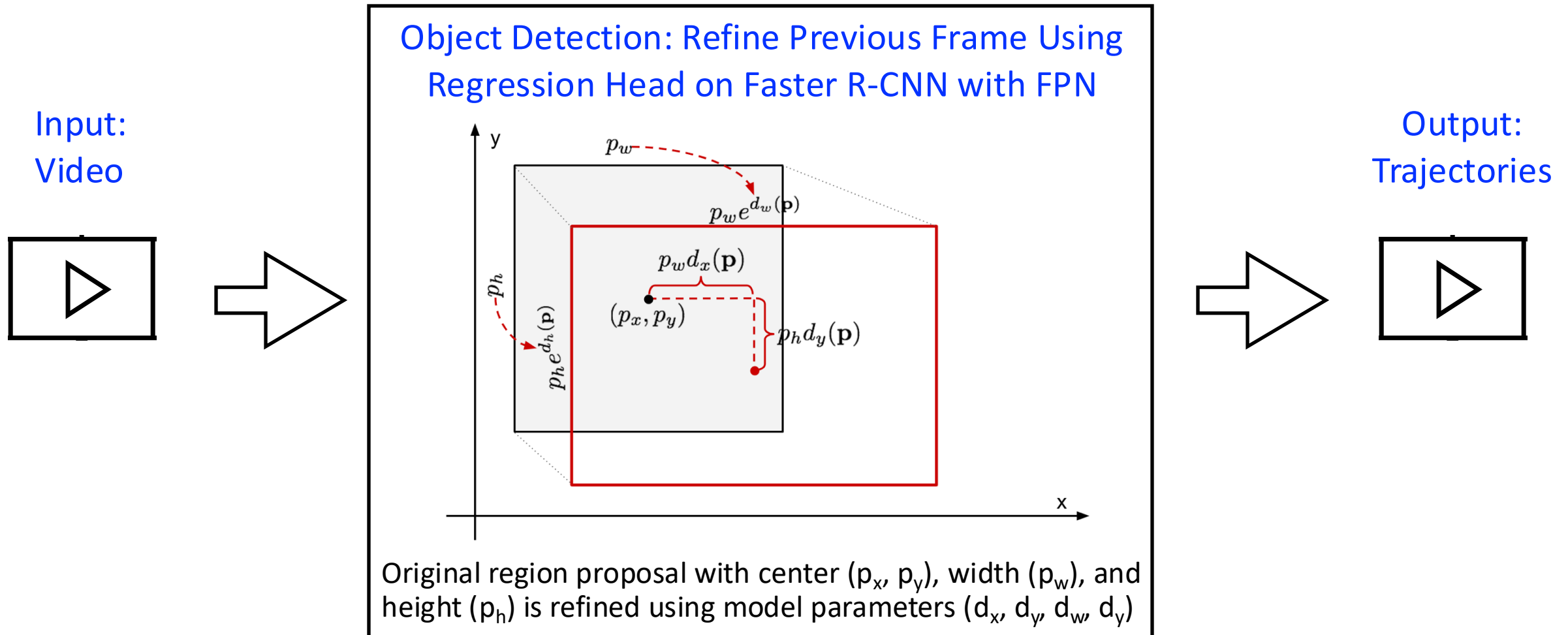
# A Common Approach: Tracking-by-Detection



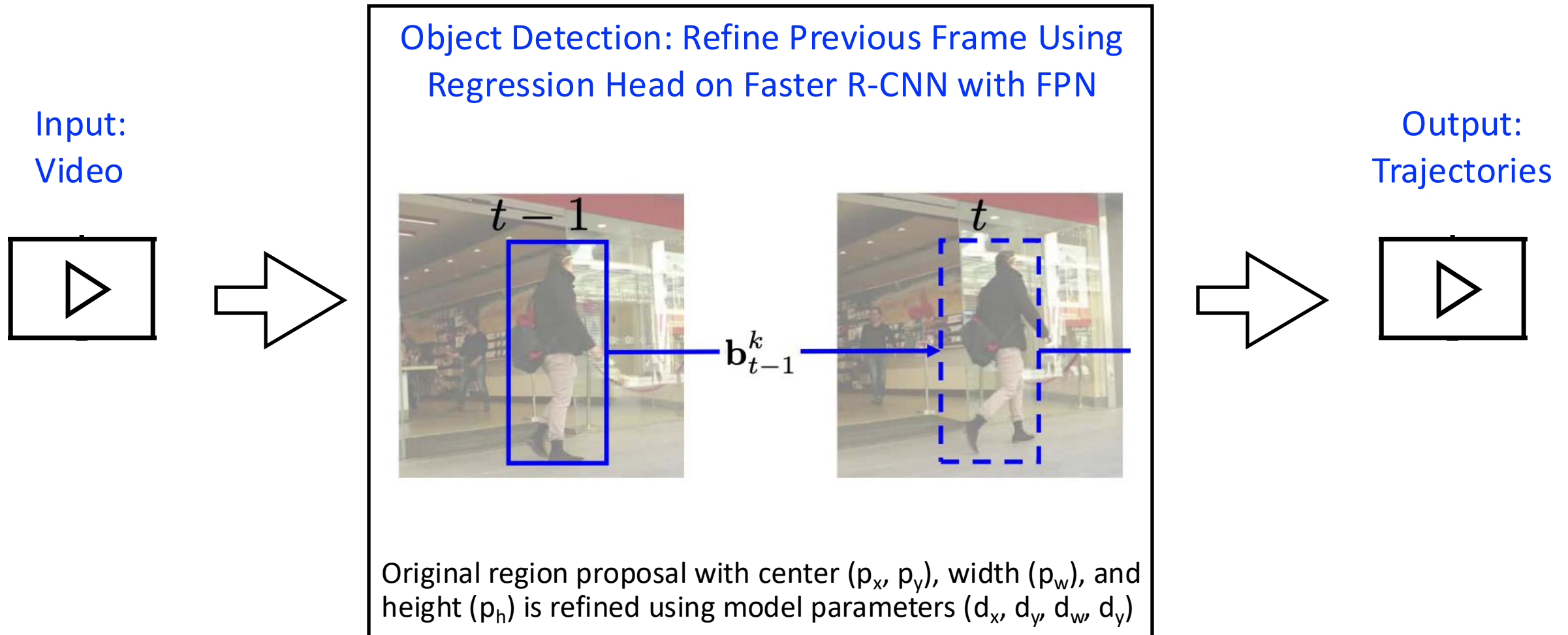
# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)



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

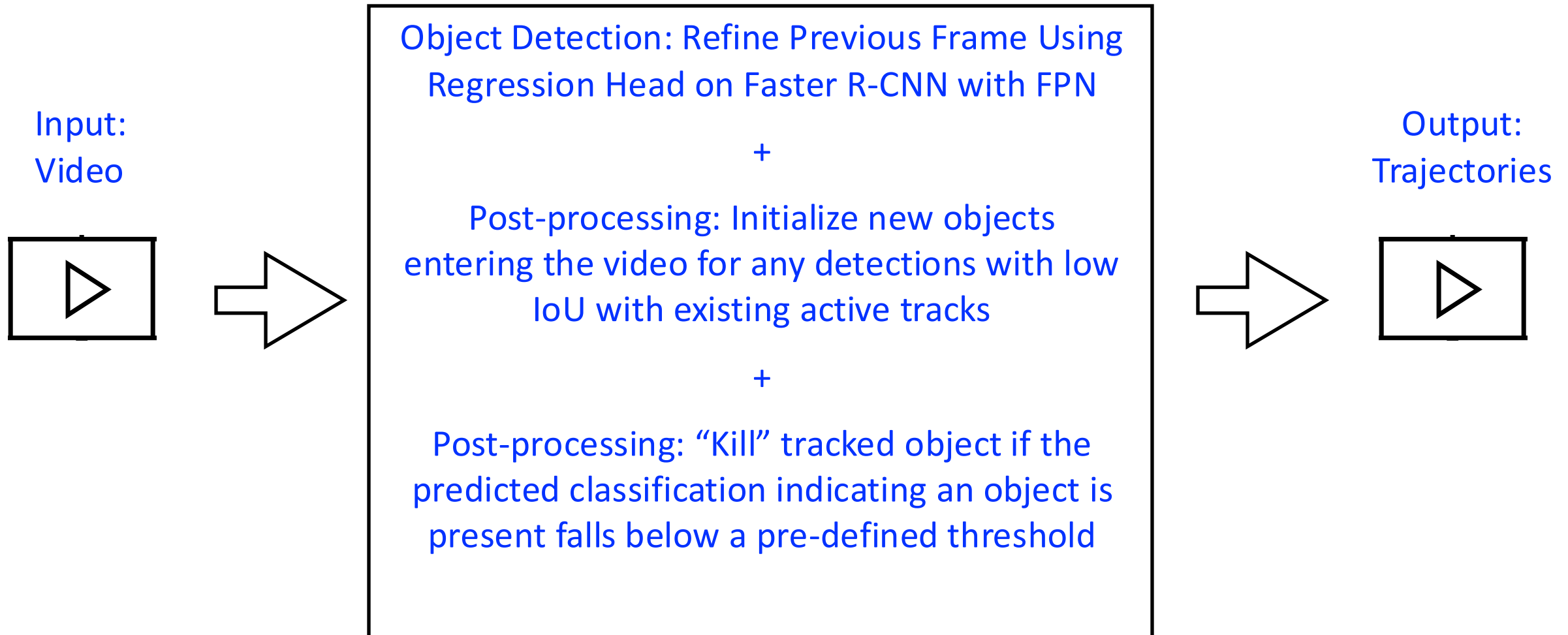


# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

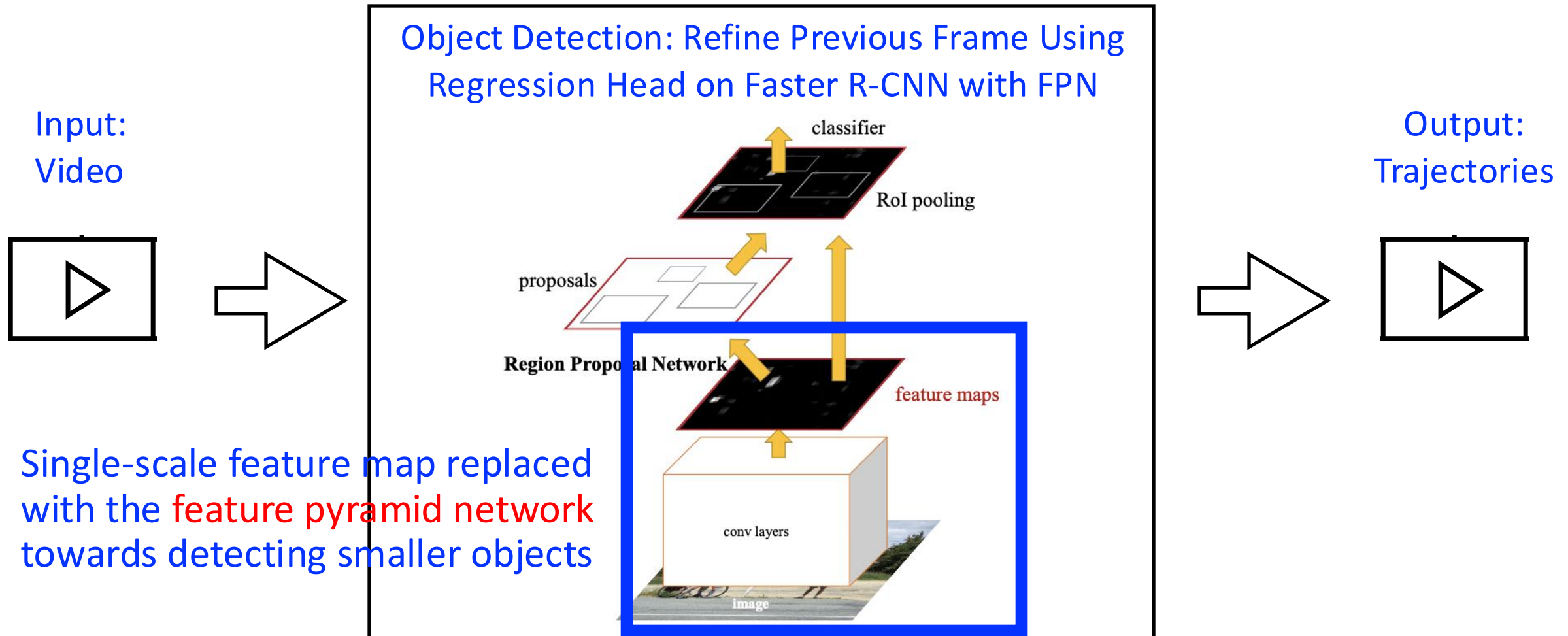




# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)

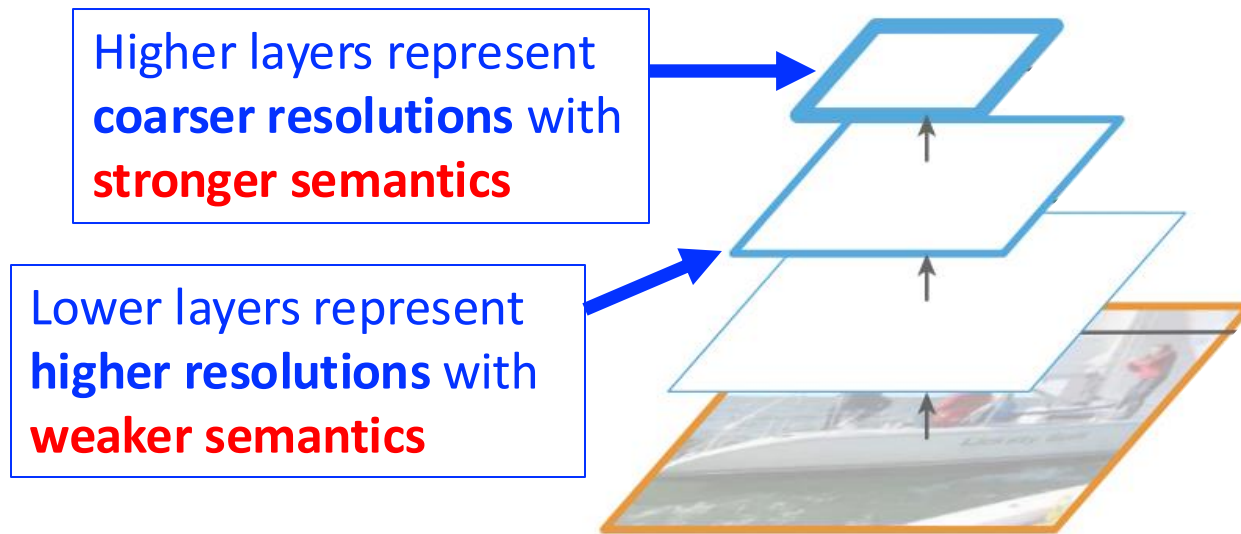


# Tracktor – Base Architecture (FPN Variant of Faster R-CNN)



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

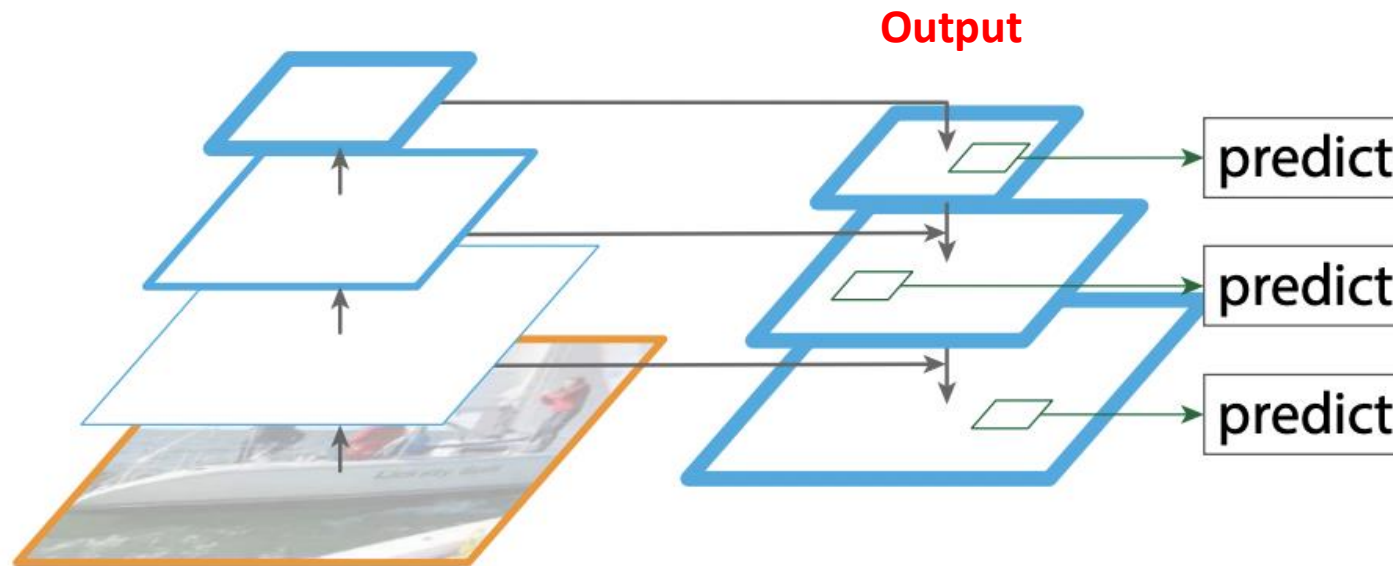
**Step 1.** Compute hierarchy of feature maps at several scales with your favorite backbone architecture (e.g., ResNet)  
(Feature Pyramid Network)



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

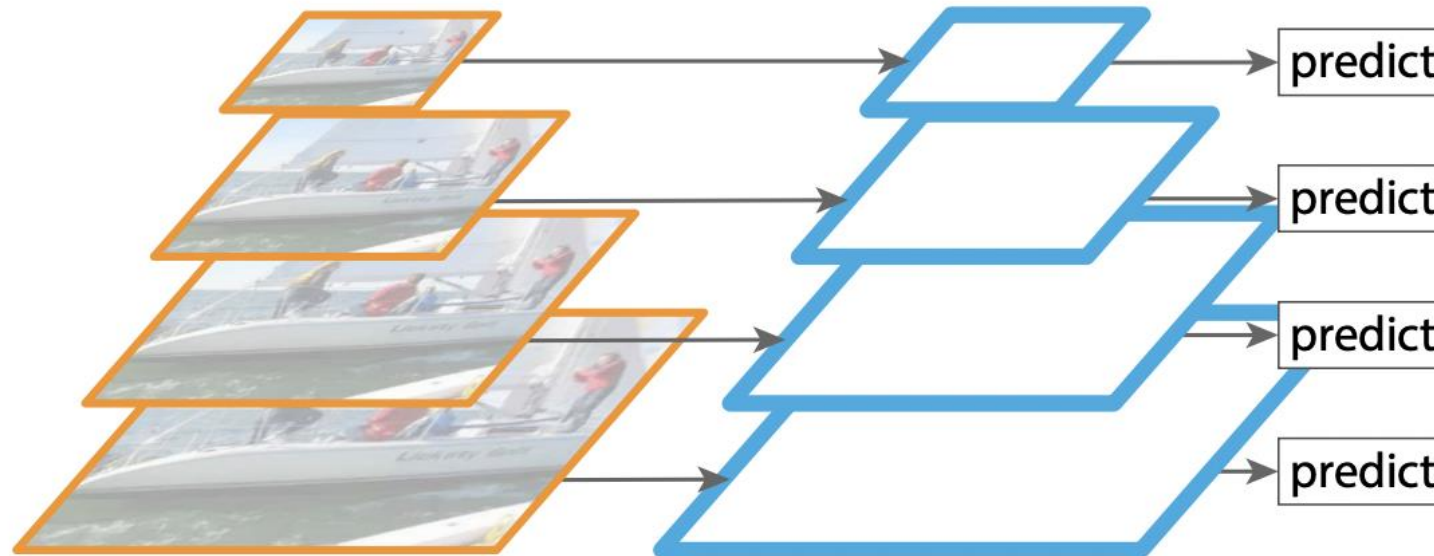
**Step 1.** Compute hierarchy of feature maps at several scales with your favorite backbone architecture (e.g., ResNet)  
(Feature Pyramid Network)

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



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

- Why not use **image pyramids**? (i.e., convert an image into multiple scales and then extract a semantically strong feature for each scale)
  - Relatively slow at test time (must test an image at every scale)



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

1. **Motion model:** for objects' considerably changing positions between frames

- Low frame rate: assume constant velocity for all objects
- Moving camera: apply image registration

2. **Reidentification:** accounts for linking an object that disappears for a short time to itself when it re-appears

- Compare appearance similarity of killed objects to newly tracked objects

# 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





# Tracktor++ Weaknesses

- 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)

# SAM-2: Semantics-Agnostic, Semi-Automated Tracking

Achieves state-of-the-art performance for video object segmentation, when specifying at the first frame what to track (e.g., click, box, mask)

Demo: <https://sam2.metademolab.com/>

Today's focus for the programming tutorial!

# SAM-2: Semantics-Agnostic, Semi-Automated Tracking

Perspective:

What would it cost to annotate 12,000 1-minute videos (i.e., 200 hours),  
with 6 frames sampled per second and 30 seconds to annotate each frame?

## SA-V Dataset

- 642.6 K masklets
- 35.5 M masks
- 50.9 K videos
- 196.0 hours



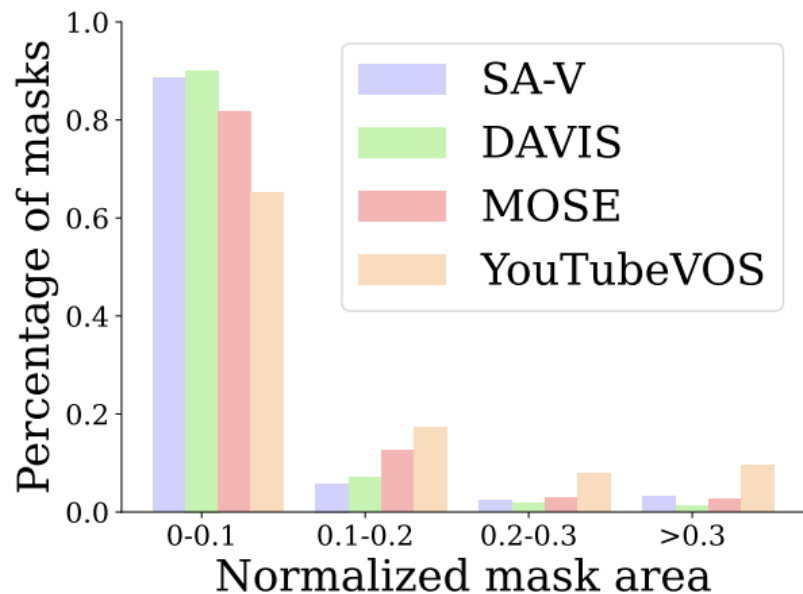
Key idea: huge training dataset (all videos and annotations from crowdworkers)!

# SAM-2: Semantics-Agnostic, Semi-Automated Tracking

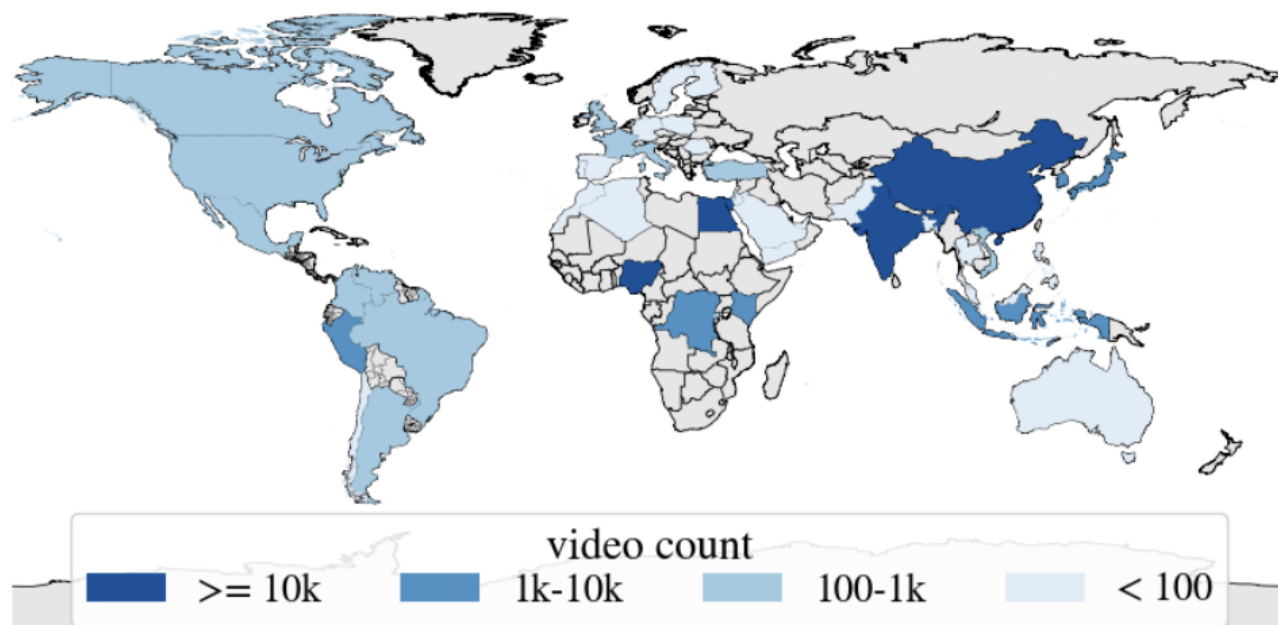
	#Videos	Duration	#Masklets	#Masks	#Frames
DAVIS 2017 ( <a href="#">Pont-Tuset et al., 2017</a> )	0.2K	0.1 hr	0.4K	27.1K	10.7K
YouTube-VOS ( <a href="#">Xu et al., 2018b</a> )	4.5K	5.6 hr	8.6K	197.3K	123.3K
UVO-dense ( <a href="#">Wang et al., 2021b</a> )	1.0K	0.9 hr	10.2K	667.1K	68.3K
VOST ( <a href="#">Tokmakov et al., 2022</a> )	0.7K	4.2 hr	1.5K	175.0K	75.5K
BURST ( <a href="#">Athar et al., 2022</a> )	2.9K	28.9 hr	16.1K	600.2K	195.7K
MOSE ( <a href="#">Ding et al., 2023</a> )	2.1K	7.4 hr	5.2K	431.7K	638.8K
Internal	62.9K	281.8 hr	69.6K	5.4M	6.0M
SA-V Manual	50.9K	196.0 hr	190.9K	10.0M	4.2M
SA-V Manual+Auto	50.9K	196.0 hr	642.6K	35.5M	4.2M

Key idea: huge training dataset (all videos and annotations from crowdworkers)!

# SAM-2: Semantics-Agnostic, Semi-Automated Tracking



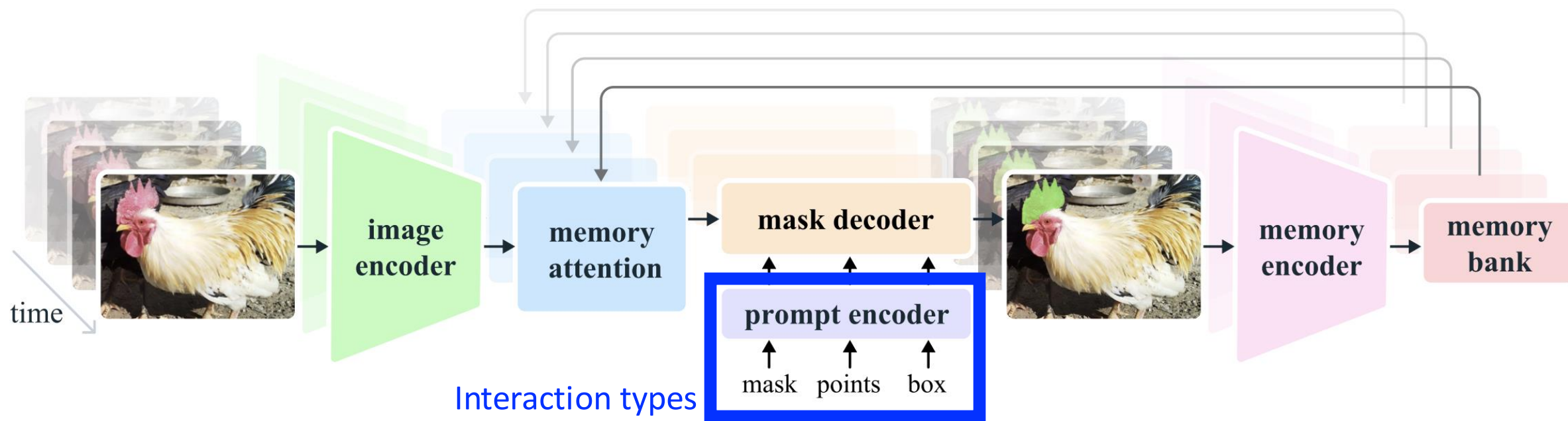
**(a)** Size



**(b)** Geography

Masks tend to occupy 10% or less of frames for videos from around the world

# SAM-2: Semantics-Agnostic, Semi-Automated Tracking



Architecture extends SAM model with memory to retain tracking information from previous frames

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)

# Object Tracking: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metrics
- Computer vision models
- Discussion (chosen by YOU 😊)





*The End*