

# Computer Vision

CS308

Feng Zheng

SUSTech CS Vision Intelligence and Perception



南方科技大学  
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY



# Content

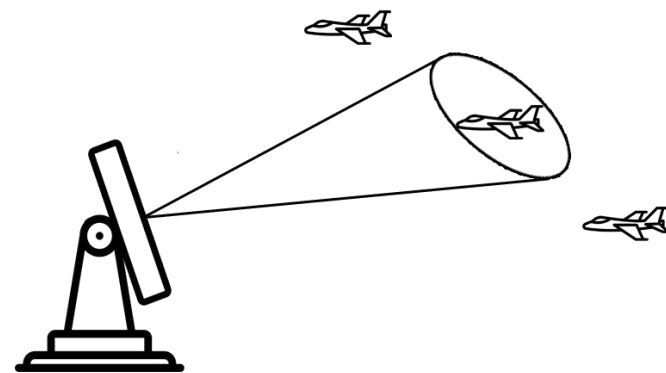
- Object Tracking: *General*
- Siamese-based Tracker
- GOTURN Tracker
- Deep SORT Tracker

# Object Tracking



# What is tracking about?

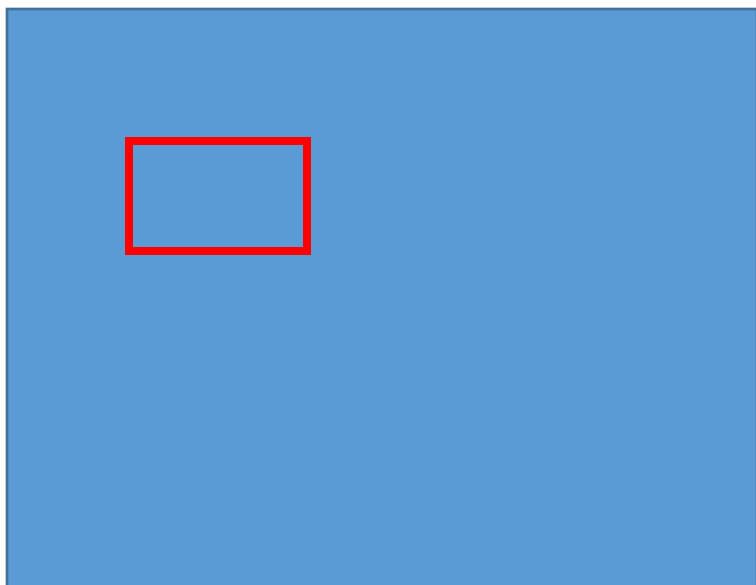
- All started in the early 60s
  - With Kalman filter for military
- Methods
  - Data association
  - Similarity measurement
  - Correlation
  - Matching/Retrieval
- Reasoning with "strong" **priors**: model-based
- Detection with very **similar** examples



CONNECT  
THE DOTS

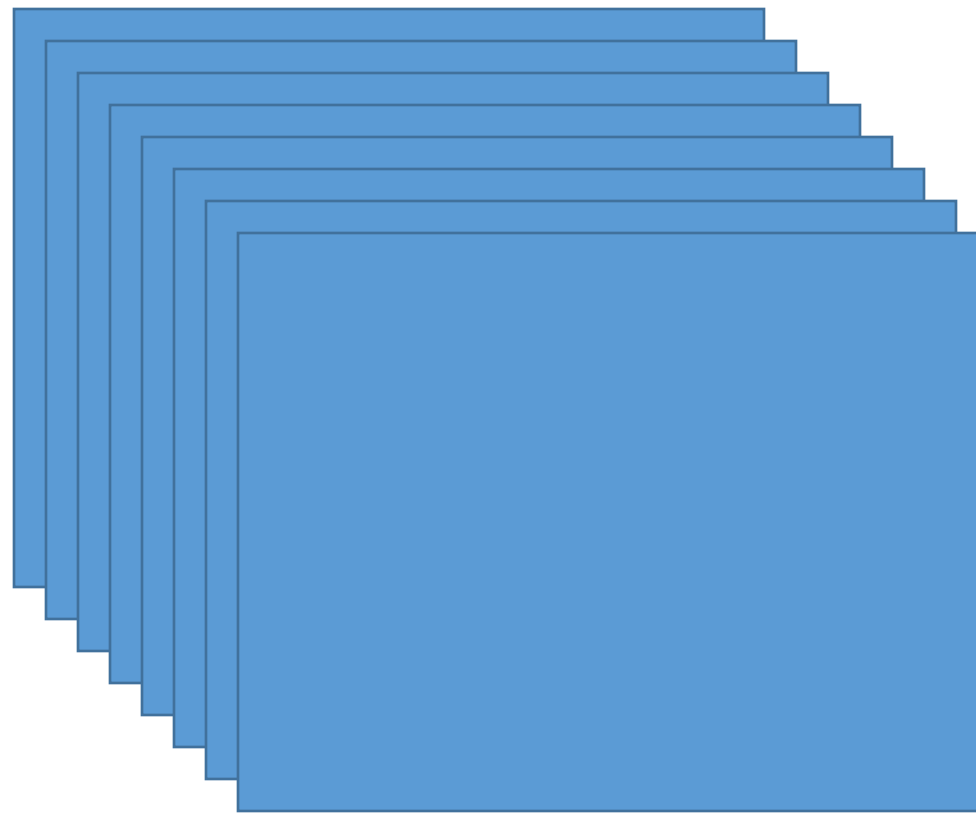


# General Setting



Initialization in the first frame

.....



Localization in the following frames



# 4 stages of the Object Tracking process

- Target initialization
- Appearance modeling
  - **Visual representation:** construct robust features and representation that can describe the object
  - **Statistical modeling:** use statistical learning techniques to build mathematical models for object identification effectively.
- Motion estimation
- Target positioning





# Levels of Object Tracking

- Single Object Tracking(SOT)
- Multiple Object Tracking(MOT): it aims to track objects of multiple classes as we see in self-driving cars



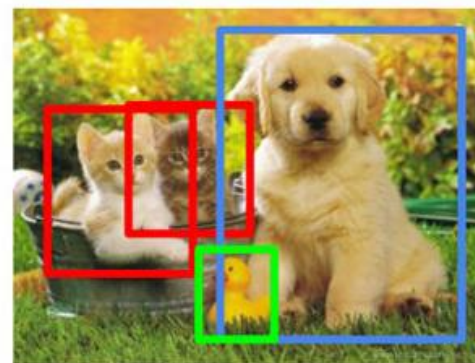


# Object Tracking vs. Object detection

- Object tracking is to estimate or predict the position of a target object in each consecutive frame in a video once the initial position of the target object is defined
- Object detection is the process of detecting a target object in an image or a single frame of the video



**Object Detection**



CAT, DOG, DUCK





# Problem statement

- Input: **target**
- Objective: Estimate target state over **time**
- **State**:
  - Position, Appearance, Shape, Velocity...
  - And affine transformation w.r.t. previous patch
- Choice: (O. S. S.)
  - **O**bject representation
  - **S**imilarity measure
  - **S**earching process



# What are the challenges?

- Illumination Variation - the illumination in the target region is significantly changed
- Scale Variation - the ratio of the bounding boxes of the first frame and the current frame is out of the range
- Occlusion - the target is partially or fully occluded
- Deformation - non-rigid object deformation
- Motion Blur - the target region is blurred due to the motion of target or camera
- Fast Motion - the motion of the ground truth is larger than 20 pixels

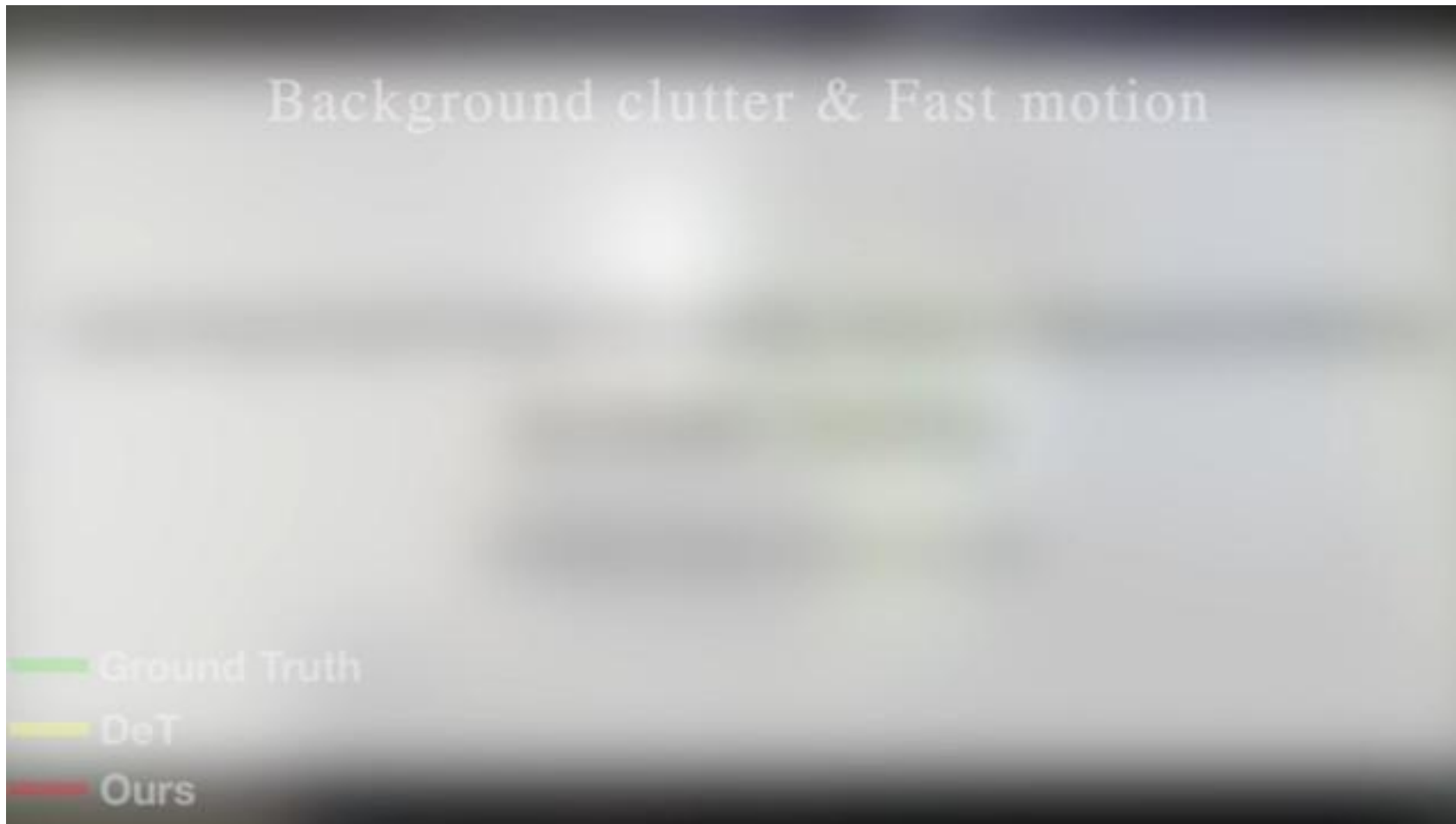


# What are the challenges?

- In-Plane Rotation - the target rotates in the image plane
- Out-of-Plane Rotation - the target rotates out of the image plane
- Out-of-View - some portion of the target leaves the view
- Background Clutters - the background near the target has the similar color or texture as the target
- Low Resolution - the number of pixels inside the ground-truth bounding box is less than 400



# Object Tracking: challenges





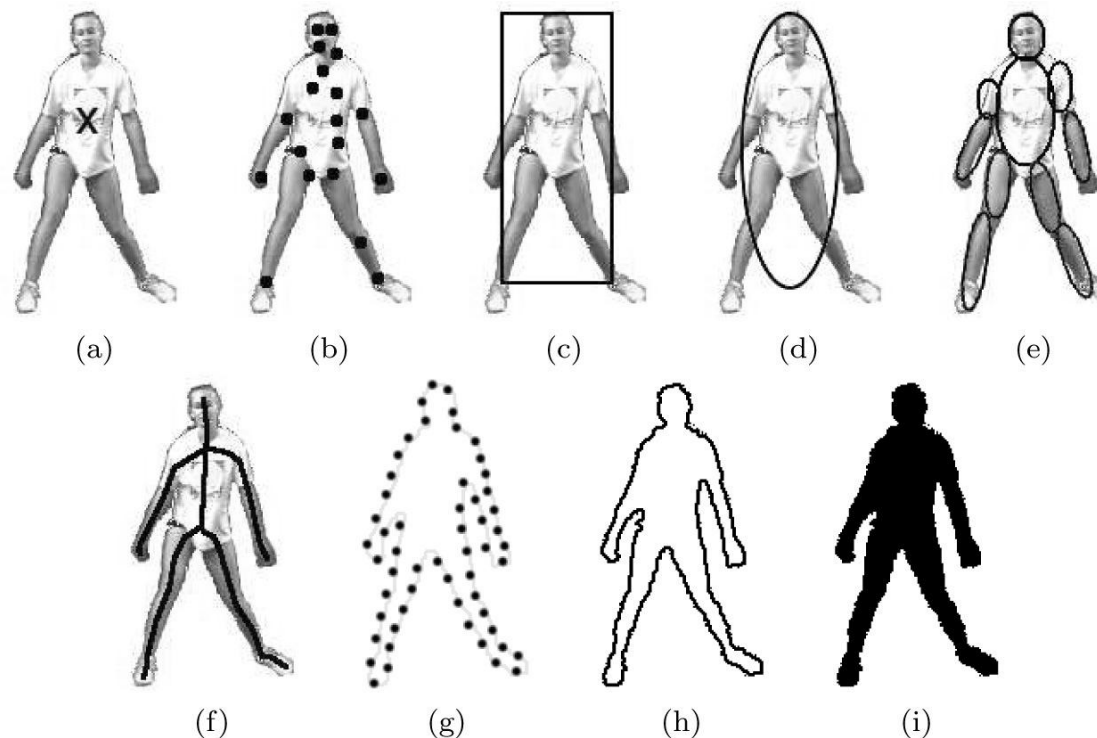
# Object representation

- Taxonomy
  - Low/mid/high level features
  - Grid/Pyramid/Cascade
  - Patch/keypoints
- Goal: we want a representation that is:
  - Descriptive **enough to disambiguate** target VS background
  - Flexible enough to cope with:
    - ✓ **Scale**
    - ✓ **Pose**
    - ✓ **Illumination**
    - ✓ Partial **occlusions**



# Object representation

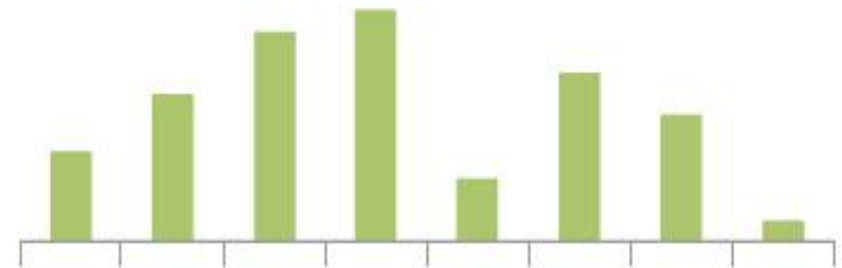
- Object approximation:
  - Segmentation / Polygonal approximation
  - Bounding ellipse/box
  - Position only
- Goal: Measure **affinity**





# Object representation: From light to useful information

- Low/mid/high level features

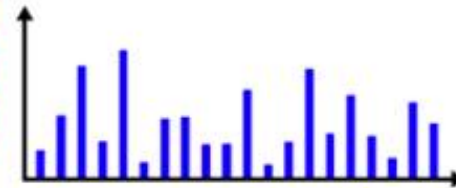
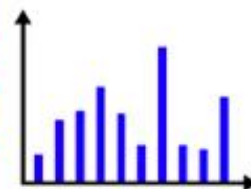
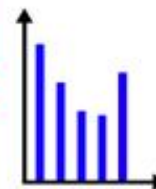


histograms



# Object representation: Sampling strategies

- Grid/pyramid/cascade of coarse-to-fine

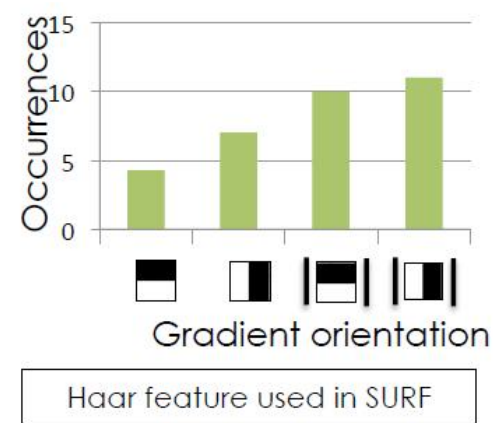
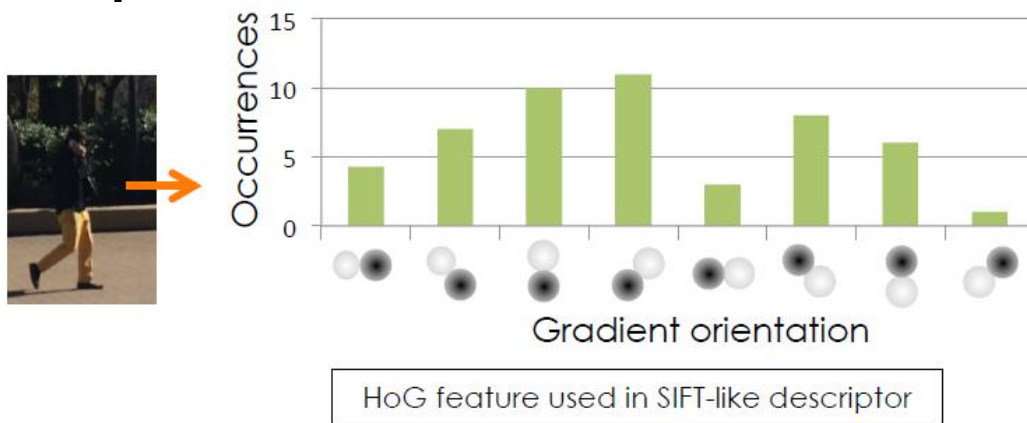




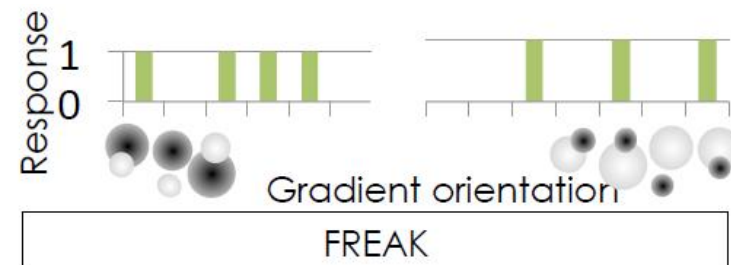
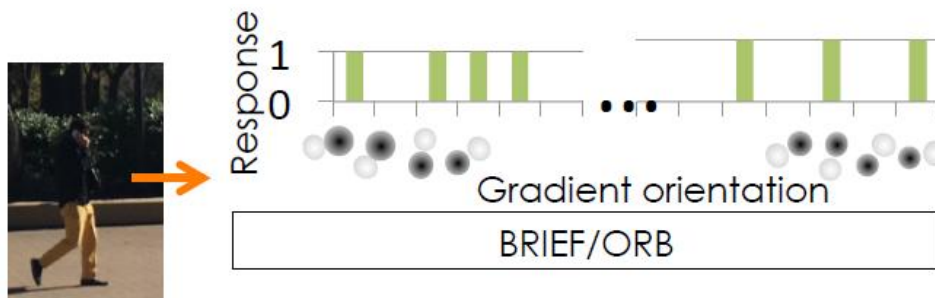


# Low-level features

- Integer responses

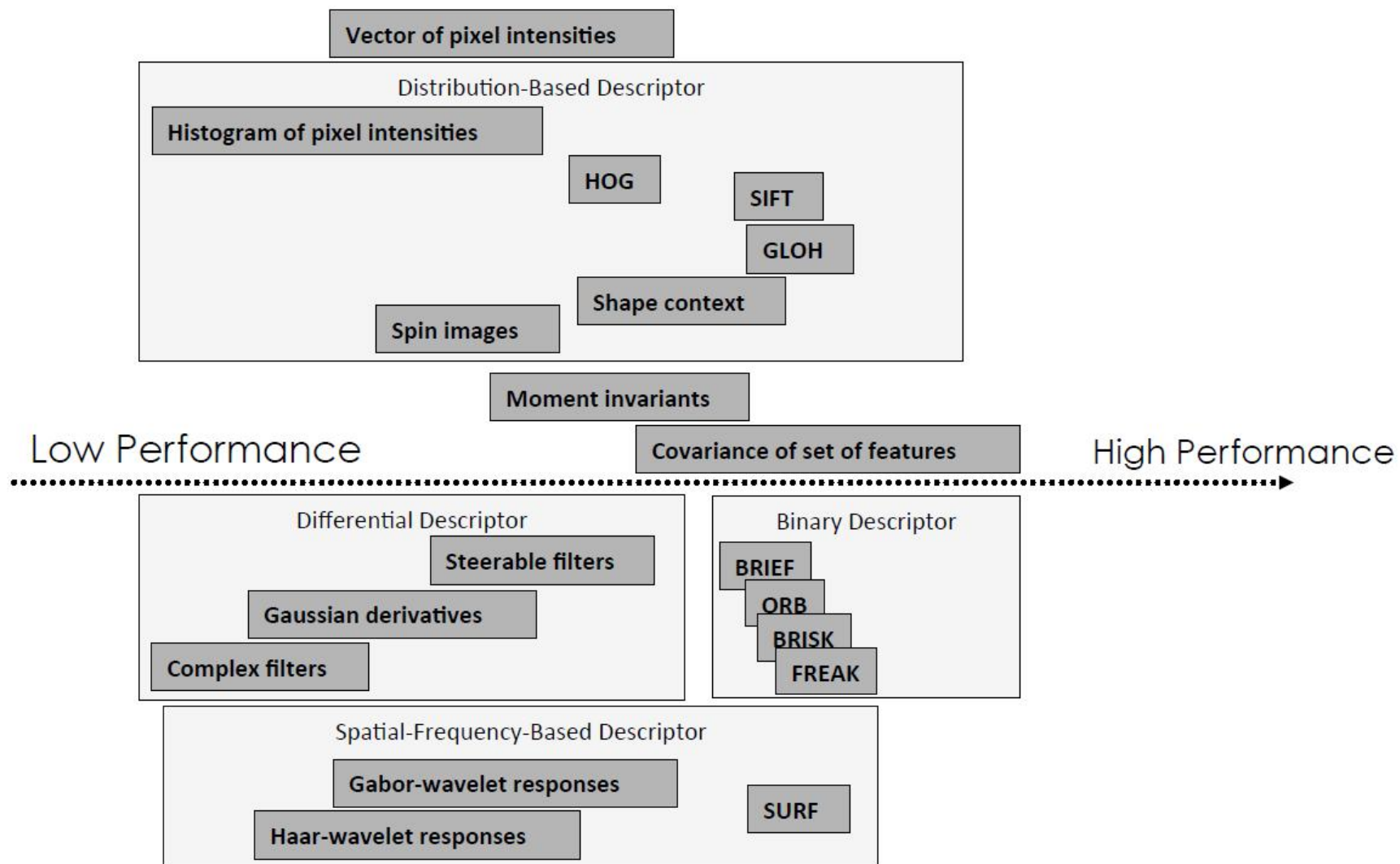


- Binary responses



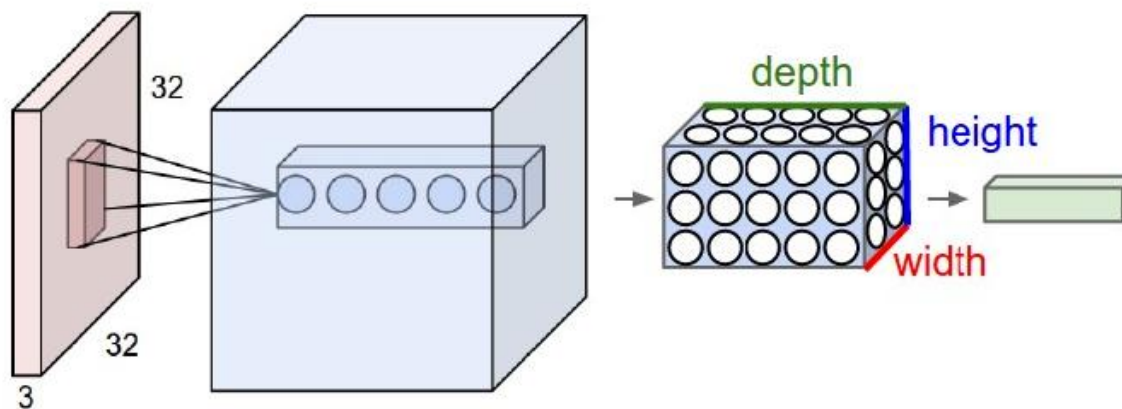


# A bulk of Low-level features





# Recent trend: CNN features





# Measuring Affinity

- In general:  $aff(x, y) = \exp \left( -\frac{1}{2\sigma_d^2} \|f(x) - f(y)\|^2 \right)$
- Examples:
  - Distance:  $f(x) = location(x)$
  - Intensity:  $f(x) = intensity(x)$
  - Color:  $f(x) = color(x)$
  - Texture:  $f(x) = filterbank(x)$

**Pixels => Regions**

- Note: Can also modify distance metric

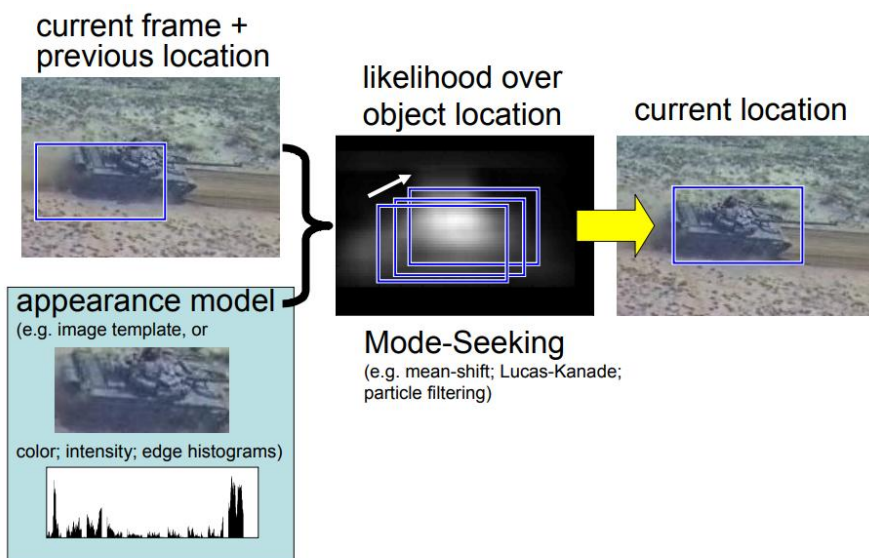


# Classical Methods

- Optimal flow
- Kalman filter
- Particle filter
- Mean-shift



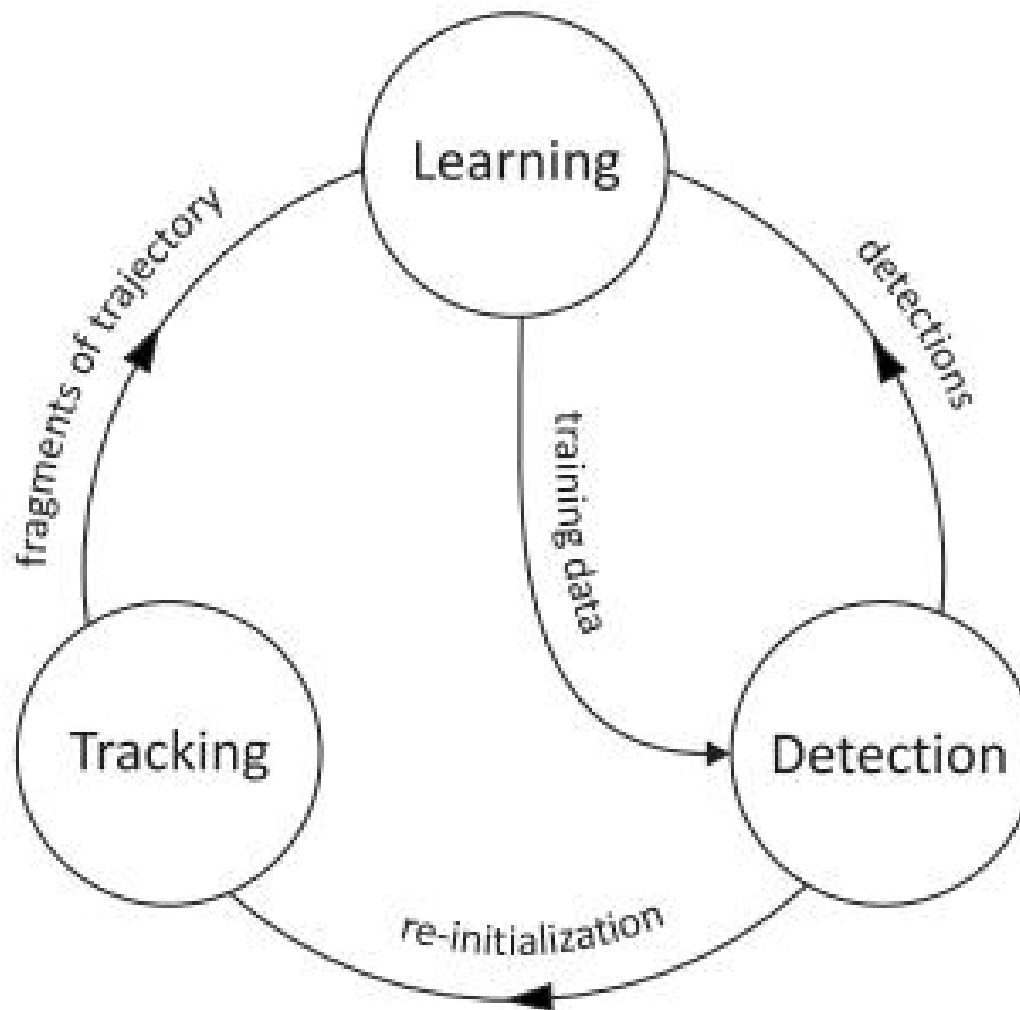
<https://www.kalmanfilter.net/default.aspx>



# Online object tracking : Tracking-Learning- Detection



# Tracking-Learning-Detection



Built in the first  
or last frame

Focus more on:  
**Similarity**

Built offline or in  
past frames

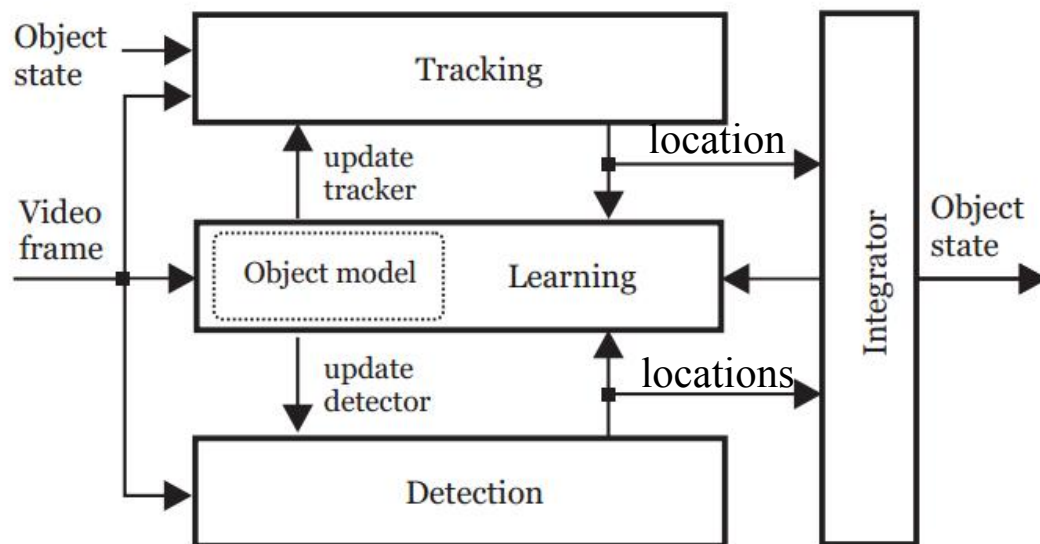
Focus more on:  
**categories**

<https://github.com/gnebehay/TLD>

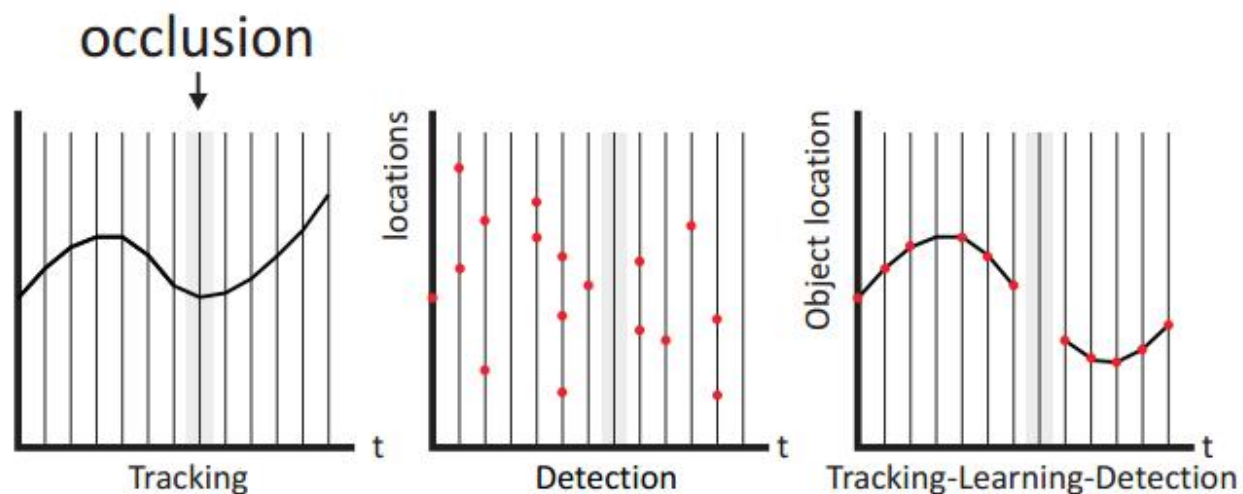




# Tracking-Learning-Detection



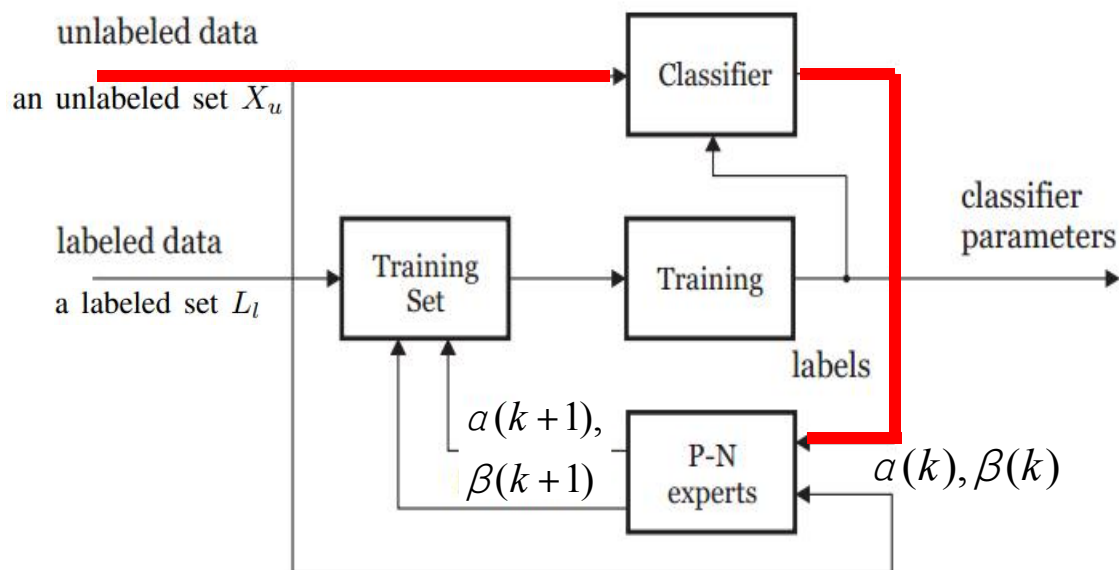
Multiple models







# PN-learning



- (i) A classifier to be learned
- (ii) Training set - a collection of labeled training examples
- (iii) Supervised training- a method that trains a classifier from training set
- (iv) P-N experts

GOTURN Tracker



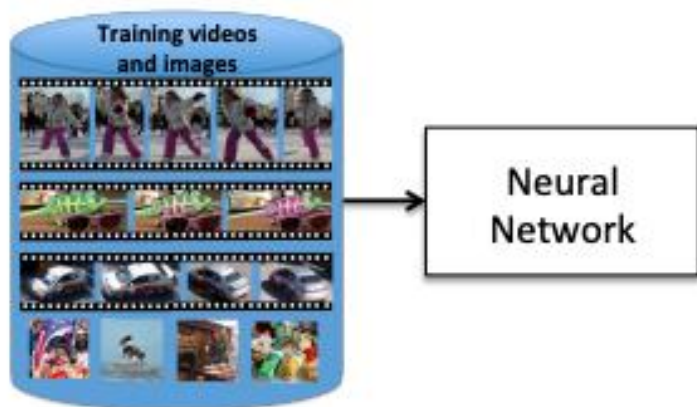
# GOTURN Tracker

- GOTURN: Generic Object Tracking Using Regression Network
- Motivation
  - Past trackers are still trained from scratch **online** and do not benefit from the large number of videos
  - No real-time speed (<20FPS)
- Solution
  - **Offline** training
  - Simple feed-forward network (no online training required, **100FPS**)



# GOTURN Tracker

Training:

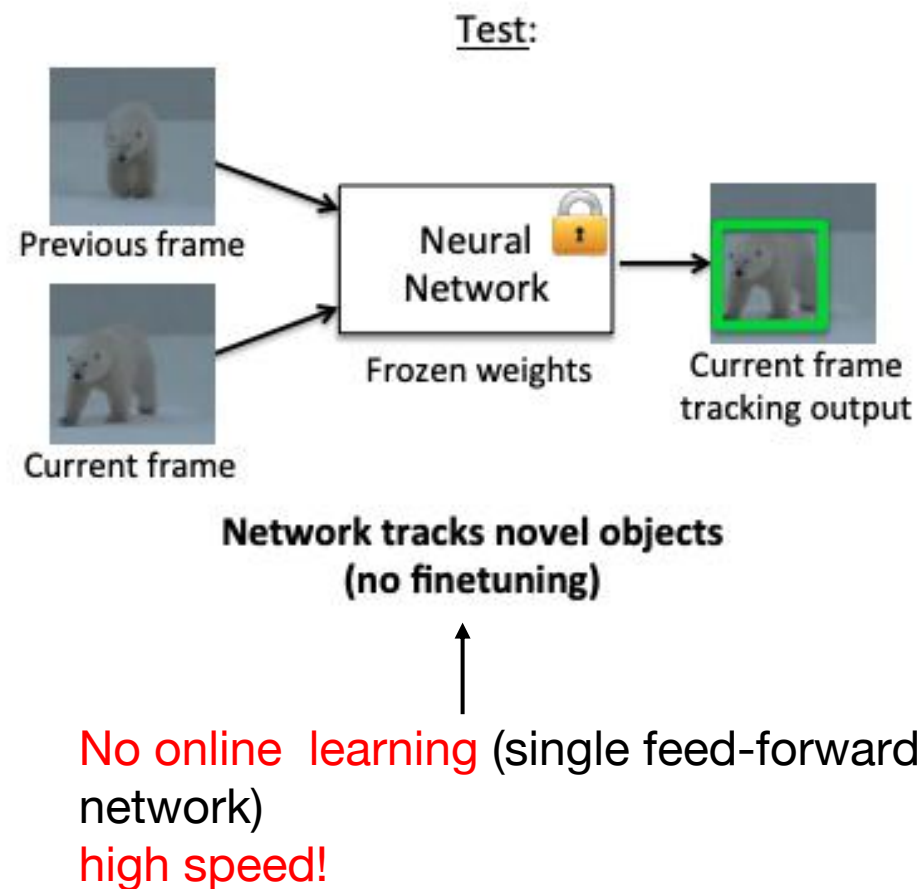
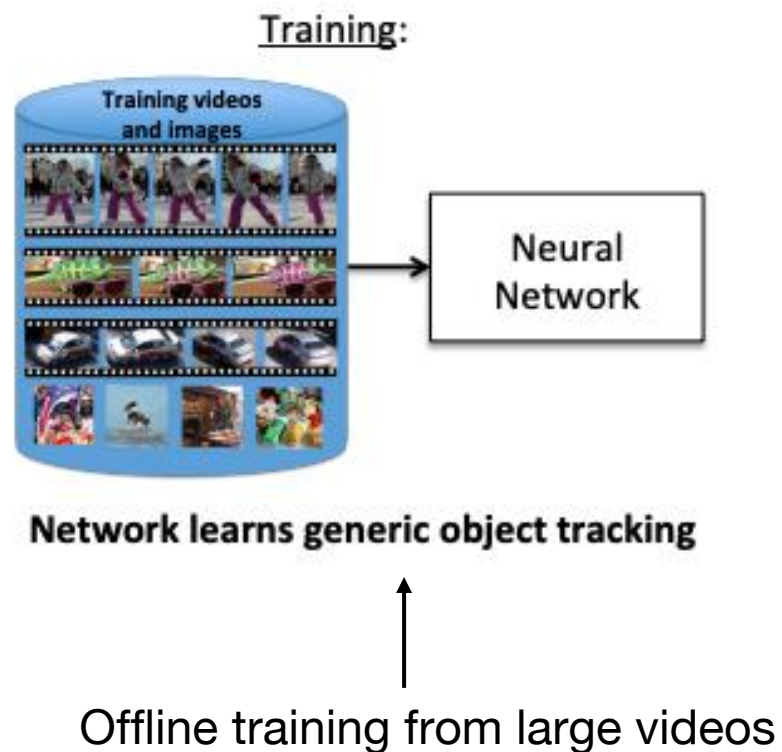


**Network learns generic object tracking**

Offline training from large videos



# GOTURN Tracker





# GOTURN Tracker

Current frame



Search Region



Crop



Previous frame

Crop



What to track



What to track?

- **Crop** and **scale** the previous frame to be centered on the target object
- **Padding**-receive extra contextual information about the surroundings



# GOTURN Tracker

Current frame



Search Region



Crop



Previous frame

Crop



What to track

What to track?

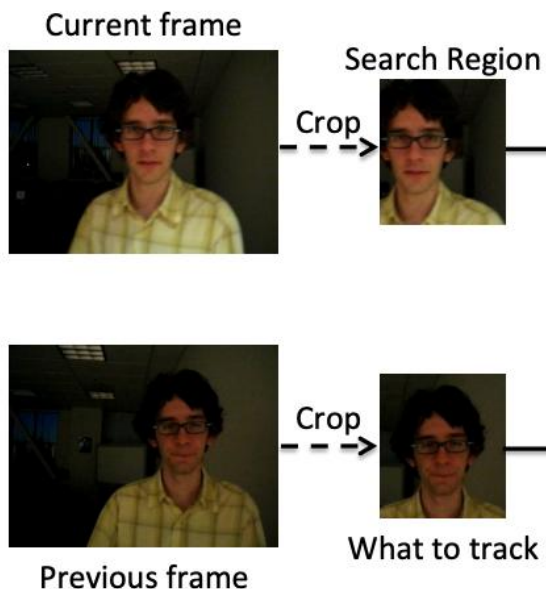
- **Crop** and **scale** the previous frame to be centered on the target object
- **Padding**-receive extra contextual information about the surroundings

Where to look?

- Hypothesis: objects tend to **move smoothly** through space
- Choose a search region in current frame based on the object's **previous** location —> **crop**



# GOTURN Tracker



Where to look?

- Hypothesis: objects tend to **move smoothly** through space
- Choose a search region in current frame based on the object's **previous** location  $\rightarrow$  **crop**

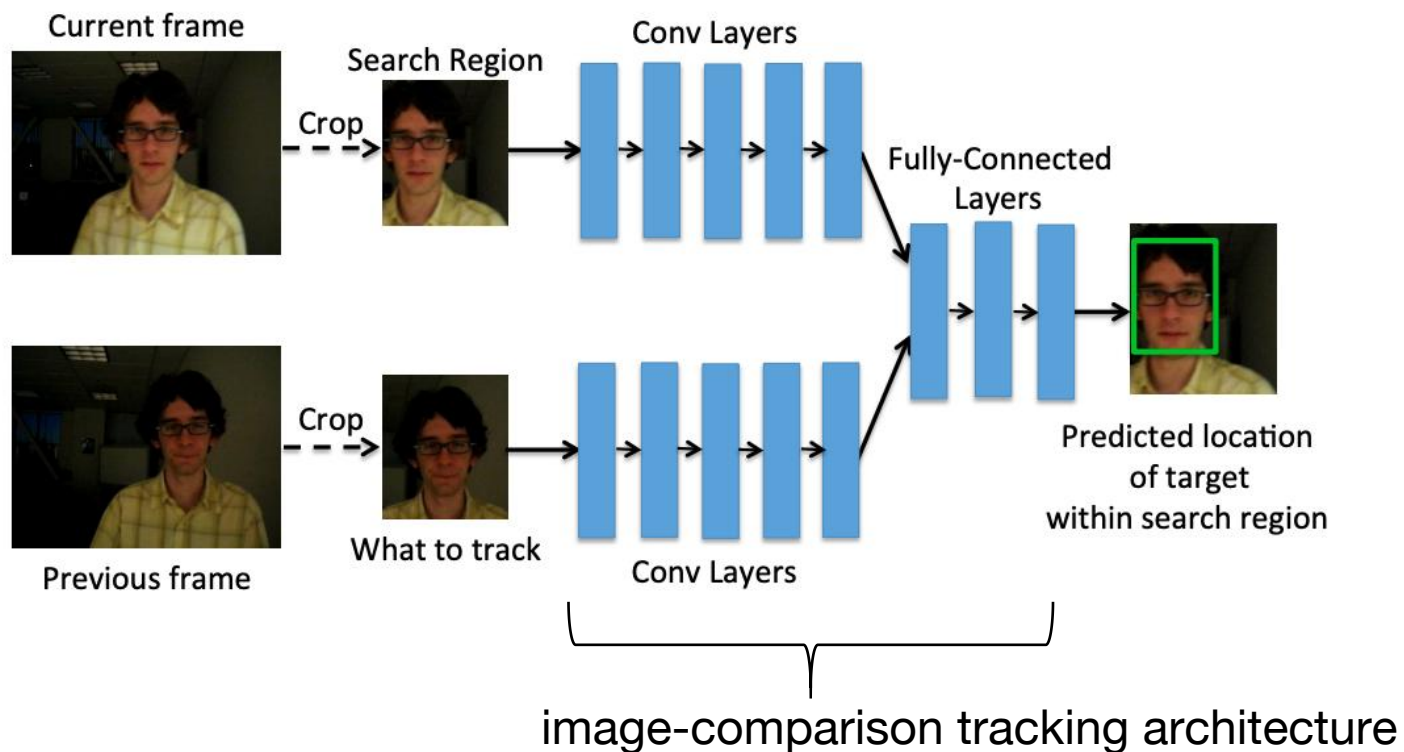
$$\begin{aligned}c'_x &= c_x + w \cdot \Delta x & w' &= w \cdot \gamma_w \\c'_y &= c_y + h \cdot \Delta y & h' &= h \cdot \gamma_h\end{aligned}$$

- $(c'_x, c'_y)$ : the center of the bounding box in the current frame
- $(c_x, c_y)$ : the center of the bounding box in the previous frame
- $w, h$ : the width and height
- $\Delta x, \Delta y$ : position change of the bounding box relative to its size (a **Laplace distribution** with a mean of 0, in practice)
- $\gamma_w, \gamma_h$ : the **size change** of the bounding box (modeled by a Laplace distribution with a mean of 1)





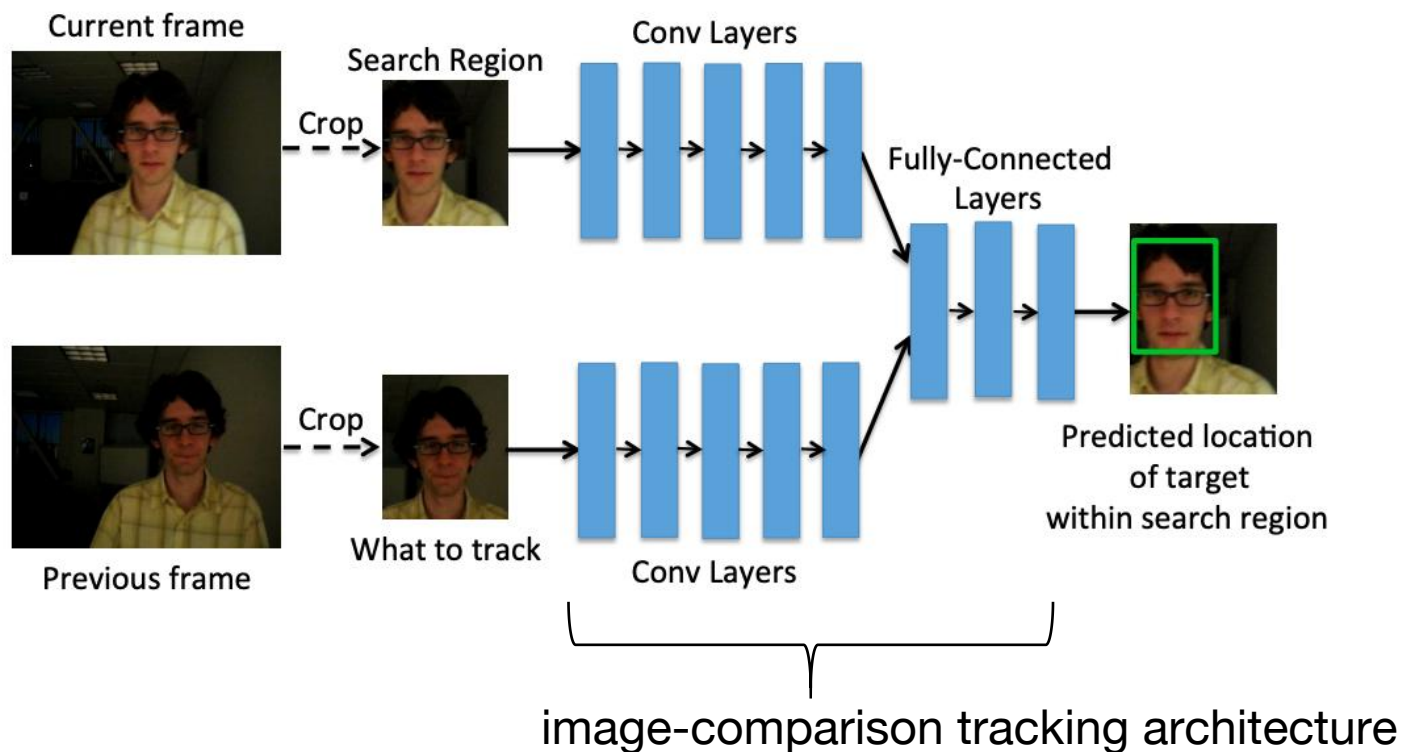
# GOTURN Tracker



- Conv Layers: extract image features
- Fully connected layers: **compare** the features from the target object to the features in the current frame to find where the target object has moved and **output bounding box**.



# GOTURN Tracker

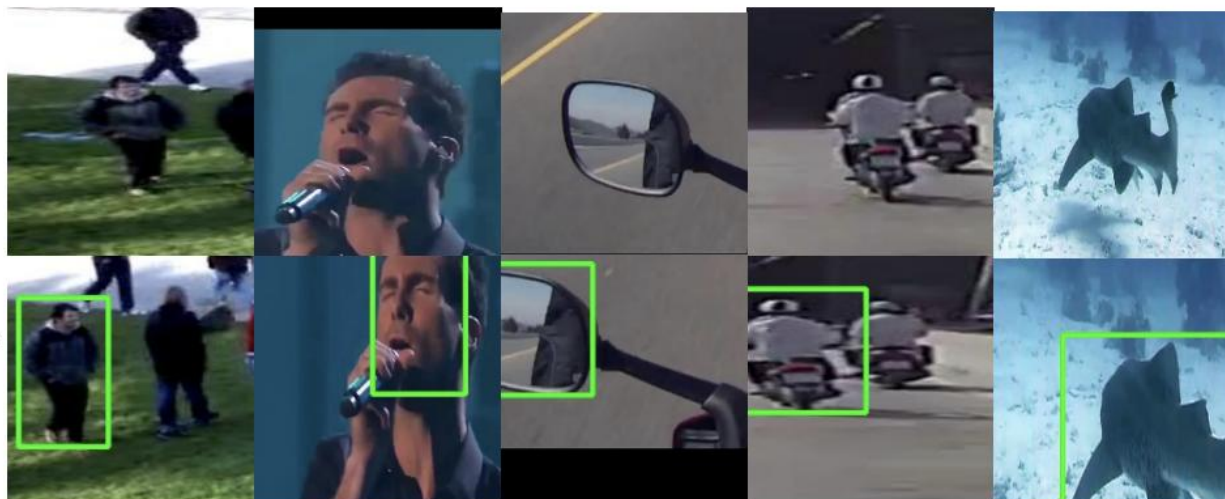


- Loss: L1 loss between predicted bounding box and ground truth.



# GOTURN Tracker

Previous  
video frame  
centered on  
object



Current video frame,  
shifted, with  
ground-truth  
bounding box

## Two types of training sets

Image  
centered on  
object



Shifted image  
with ground-truth  
bounding box

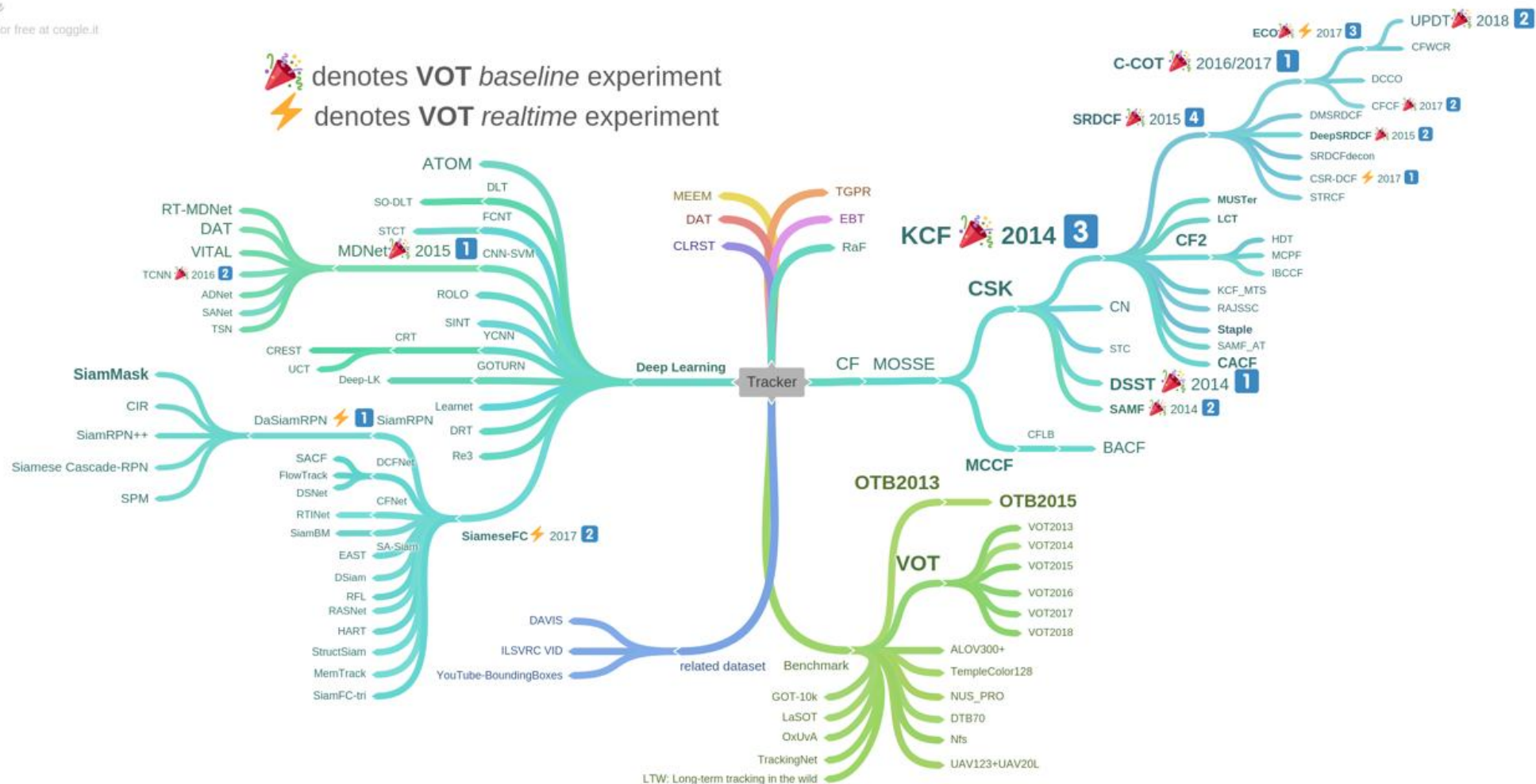
# Siamese-based Tracker



coggle

made for free at coggle.it

# Siamese-based Tracker





# Siamese-based Tracker

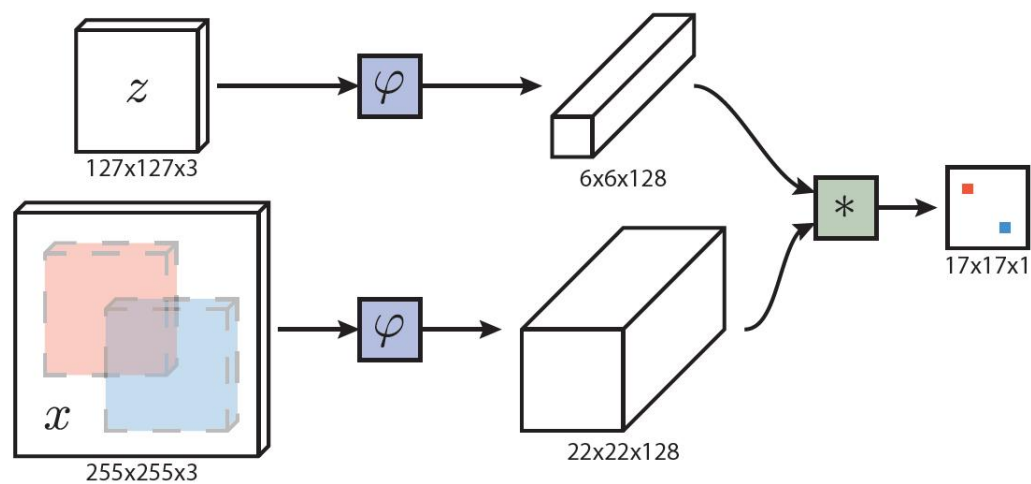
- SiamFC
- SiamRPN
- SiamRPN++





# SiamFC

## • Architecture

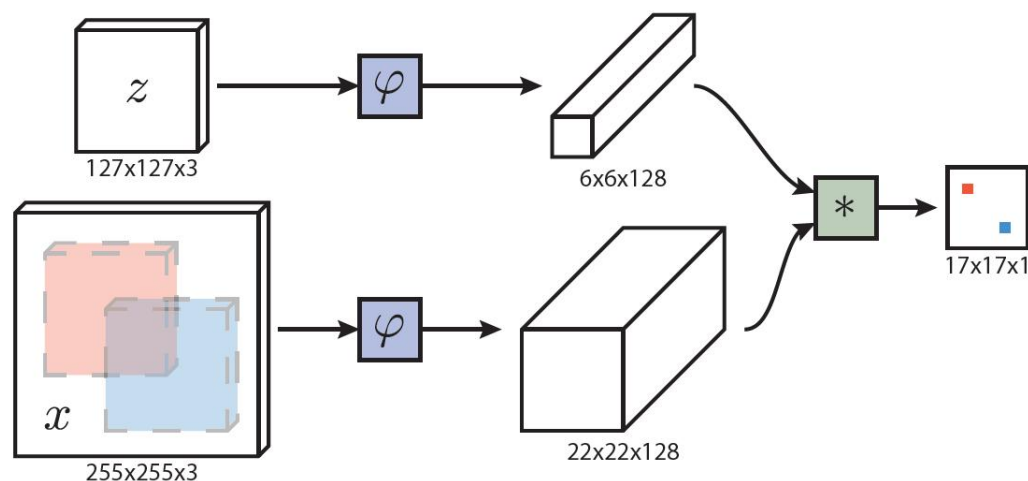


Layer	Support	Chan. map	Stride	Activation size		
				for exemplar	for search	chans.
				127 × 127	255 × 255	×3
conv1	11 × 11	96 × 3	2	59 × 59	123 × 123	×96
pool1	3 × 3		2	29 × 29	61 × 61	×96
conv2	5 × 5	256 × 48	1	25 × 25	57 × 57	×256
pool2	3 × 3		2	12 × 12	28 × 28	×256
conv3	3 × 3	384 × 256	1	10 × 10	26 × 26	×192
conv4	3 × 3	384 × 192	1	8 × 8	24 × 24	×192
conv5	3 × 3	256 × 192	1	6 × 6	22 × 22	×128



# SiamFC

## • Architecture



Images are **scaled** such that the bounding box, plus an added margin for context, has a fixed area.

### \* Input image

$$p = (w + h)/4$$

$$s(w + 2p) \times s(h + 2p) = A$$

scale factor  $s$

$$A = 127^2$$

bounding box has size  $(w, h)$

context margin is  $p$

### \* Network architecture

### \* Similarity measurement

$$f(z, x) = \varphi(z) * \varphi(x) + b \mathbb{1}$$

### \* Object location





# SiamFC

- **\* Loss function**

$$\ell(y, v) = \log(1 + \exp(-yv))$$

$$L(y, v) = \frac{1}{|\mathcal{D}|} \sum_{u \in \mathcal{D}} \ell(y[u], v[u])$$

$$y[u] = \begin{cases} +1 & \text{if } k\|u - c\| \leq R \\ -1 & \text{otherwise} \end{cases}.$$

the centre

- $v$  real-valued score of a single exemplar-candidate pair
- $y$  ground-truth label
- $D$  score map
- $R$  radius



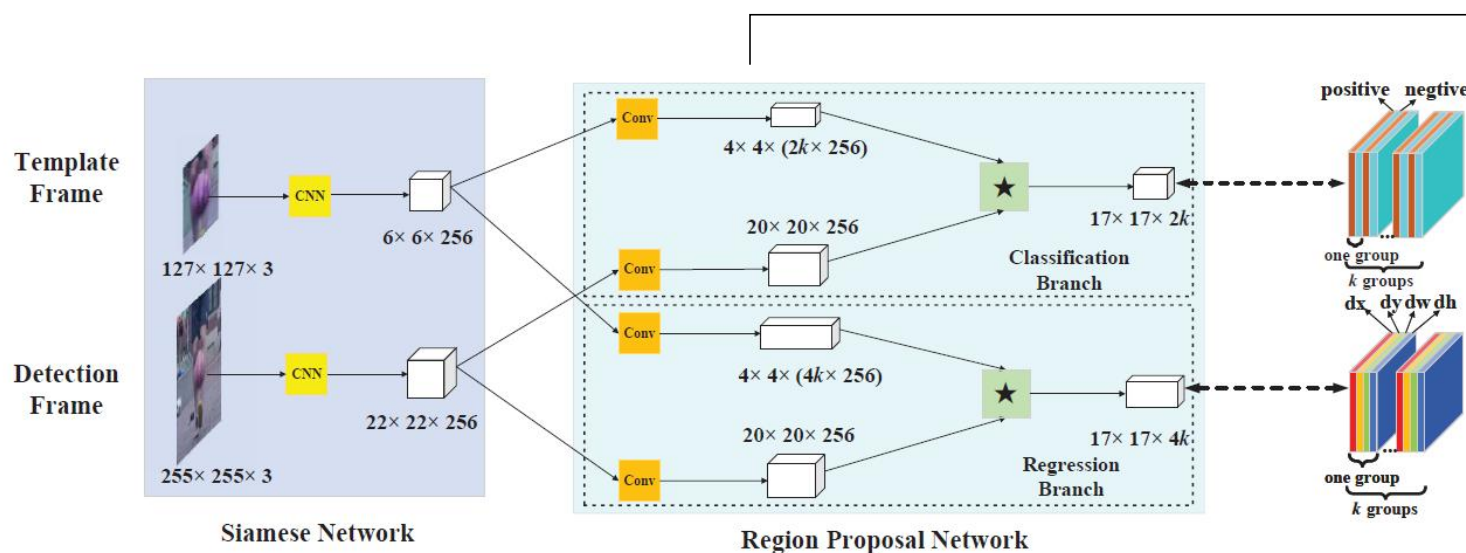
# Siamese-based Tracker

- SiamFC
  - Lack of bounding box regression
  - Need to do multi-scale test
- SiamRPN
- SiamRPN++



# SiamRPN

## • Architecture



$2k$  represents for negative and positive activation of each anchor at corresponding location on original map

If there are  $k$  anchors, network needs to output  $2k$  channels for classification and  $4k$  channels for regression.

$4k$  represents for  $dx, dy, dw, dh$  measuring the distance between anchor and corresponding groundtruth

$$A_{w \times h \times 2k}^{cls} = [\varphi(x)]_{cls} \star [\varphi(z)]_{cls}$$

$$A_{w \times h \times 4k}^{reg} = [\varphi(x)]_{reg} \star [\varphi(z)]_{reg}$$

Can we have a deeper backbone? AlexNet→Resnet50/101



# SiamRPN

- Proposal generation

- Collect the top  **$K$  points** in all  $A_{w \times h \times 2k}^{cls}$
- Get the corresponding refinement coordinates

$$A_{w \times h \times 2k}^{cls} = \{(x_i^{cls}, y_j^{cls}, c_l^{cls})\} \rightarrow ANC^* = \{(x_i^{an}, y_j^{an}, w_l^{an}, h_l^{an})\} \rightarrow \begin{aligned} x_i^{pro} &= x_i^{an} + dx_l^{reg} * w_l^{an} \\ y_j^{pro} &= y_j^{an} + dy_l^{reg} * h_l^{an} \\ w_l^{pro} &= w_l^{an} * e^{dw_l} \\ h_l^{pro} &= h_l^{an} * e^{dh_l} \end{aligned}$$

- Proposal selection

- Discard the bounding boxes generated by the anchors too far away from the center
- Use cosine window and scale change penalty to re-rank the proposals' score to get the best one

$$penalty = e^{k * \max(\frac{r}{r'}, \frac{r'}{r}) * \max(\frac{s}{s'}, \frac{s'}{s})} \quad \text{Ratio and overall scale between proposal and last frame}$$



# Siamese-based Tracker

- SiamFC
- SiamRPN
  - Too shallow to meet the strict translation invariance (no padding)
  - Imbalance of parameter distribution (i.e. the RPN module contains 20M parameters while the feature extractor only contains 4M parameters)
- SiamRPN++



# SiamRPN++

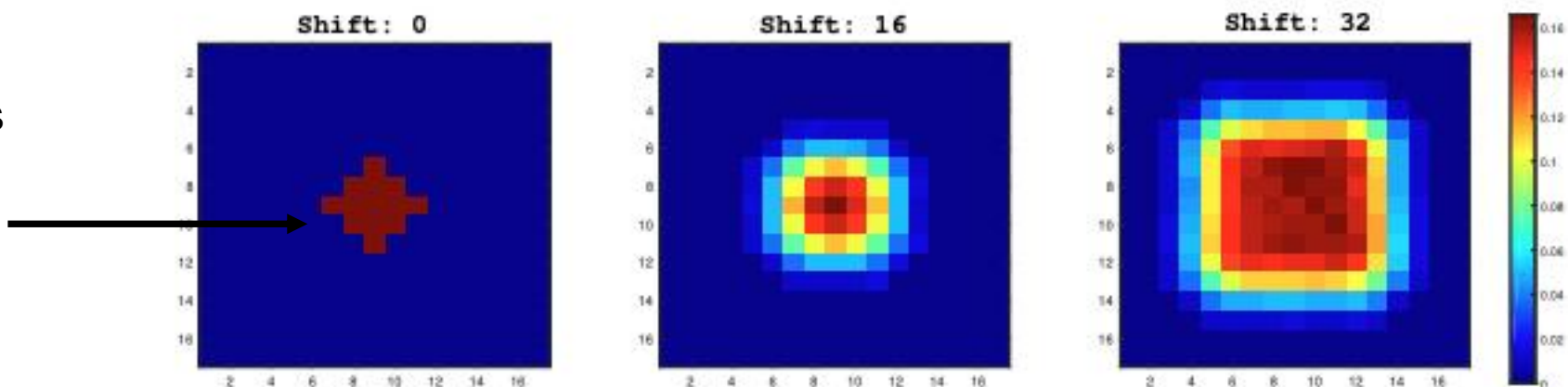
- Motivation: Can we have a **deeper** backbone?
  - The deeper, the better?
  - But deeper networks need padding
- Two restrictions in correlation :  $f(z, x) = \Phi(z) * \Phi(x) + b$ 
  - Strict translation invariance. **Padding X**
  - Correlation symmetry. RPN requires **asymmetrical** for Cls and Reg.
- Solution
  - Spatial aware sampling strategy
  - Depthwise Cross Correlation



# Spatial aware **sampling** strategy

- Positive samples are evenly distributed within a certain range, **not always** at the **center**
- The range is the distance from the center point, which means **shift**

The disadvantage of a **deep** network that breaks strict translational invariance is that it learns **positional bias**.

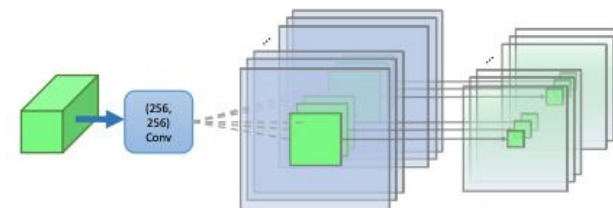
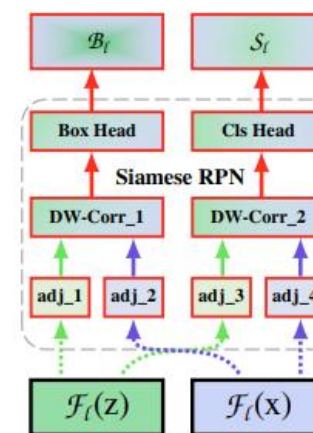
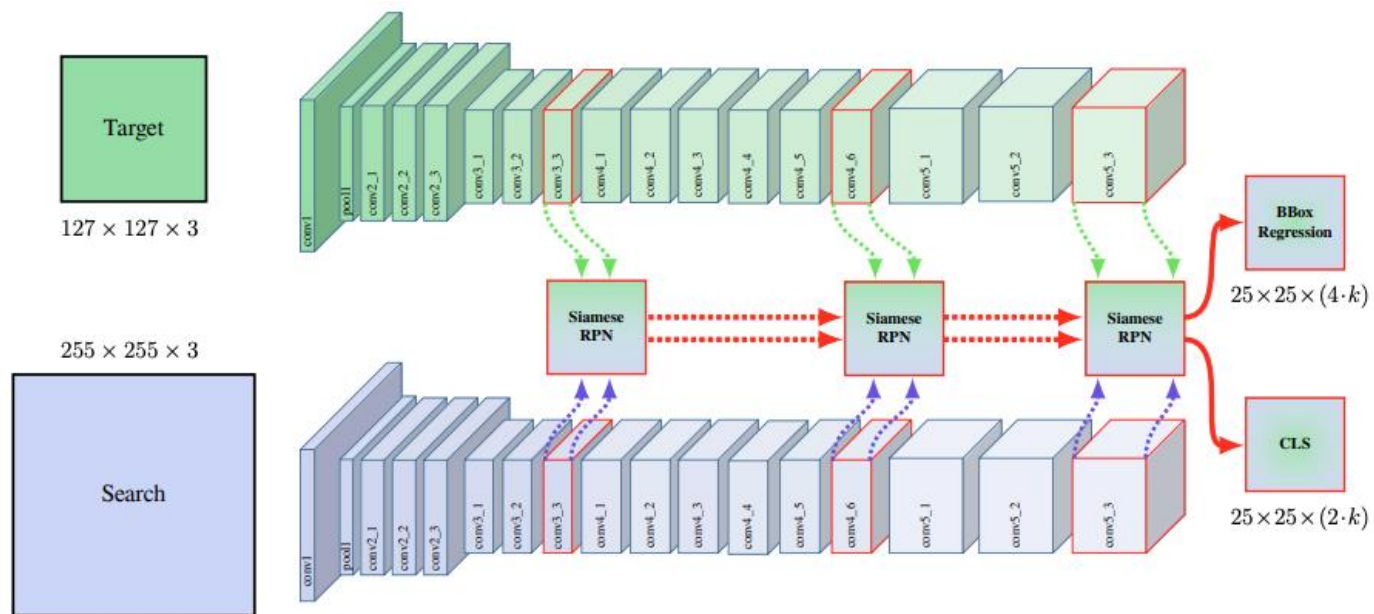


Target may appear at any position in the search region and learned feature representation should stay spatial invariant



# Layer-wise Aggregation

- Architecture



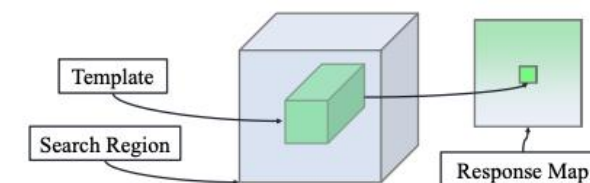
(c) Depth-wise Cross Correlation Layer



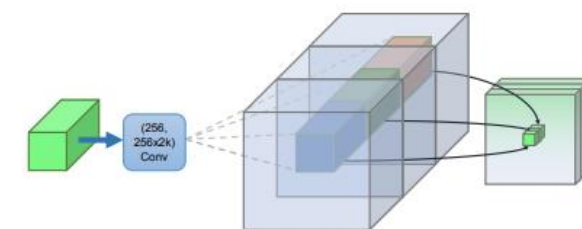


# Depthwise Cross Correlation

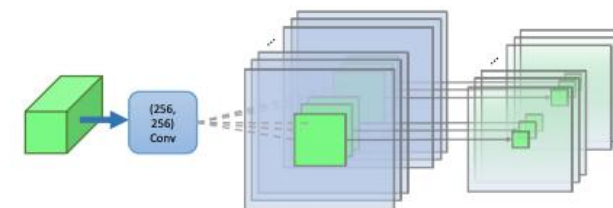
- (a) SiamFC: Cross Correlation (XCorr) layer predicts a **single** channel similarity map between target template and search patches
- (b) SiamRPN: Up-Channel Cross Correlation (UP-XCorr) layer outputs a **multi-channel correlation** features by cascading a **heavy** convolutional layer with several independent XCorr layers
- (c) SiamRPN++: Depth-wise Cross Correlation (DW-XCorr) layer predicts multi-channel correlation features between a template and search patches



(a) Cross Correlation Layer



(b) Up-Channel Cross Correlation Layer



(c) Depth-wise Cross Correlation Layer

The same category have high response on same channels



# Results

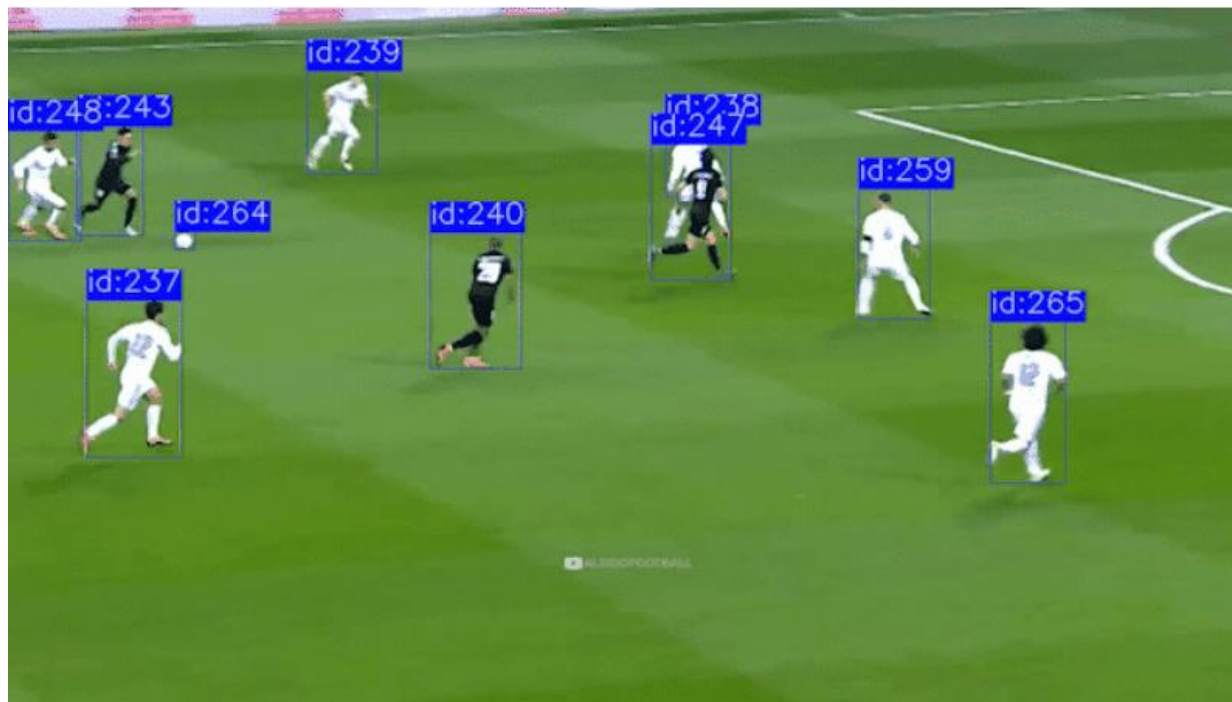
	CSRDCF	ECO	SiamFC	CFNet	MDNet	DaSiamRPN	Ours
	[28]	[5]	[1]	[41]	[32]	[52]	
AUC (%)	53.4	55.4	57.1	57.8	60.6	63.8	<b>73.3</b>
P (%)	48.0	49.2	53.3	53.3	56.5	59.1	<b>69.4</b>
P <sub>norm</sub> (%)	62.2	61.8	66.3	65.4	70.5	73.3	<b>80.0</b>

Deep SORT Tracker



# What is multiple object tracking(MOT)?

- What is MOT?
  - Locate **multiple objects** of interest in a given video simultaneously, **maintain their IDs** and **record their trajectories**

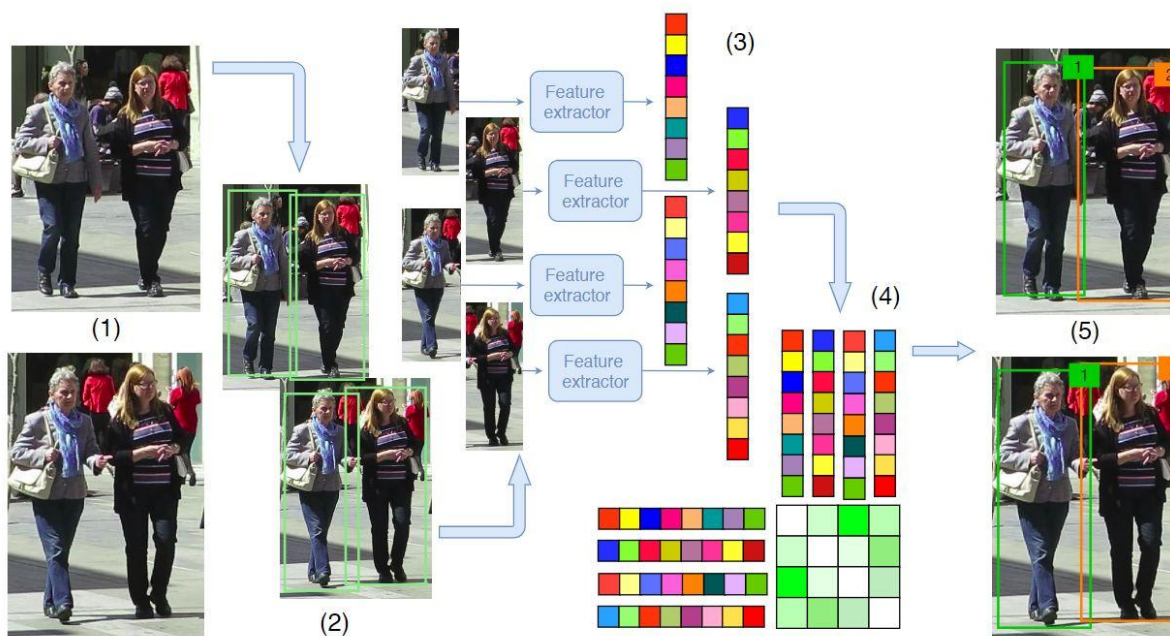




# What is multiple object tracking(MOT)?

- Popular solution:

- First, all **detections** are extracted via YOLO, Faster RCNN, and so on.
- Then, an **association** algorithm is performed to link these detections to different tracks. Usually, the association algorithm considers motion (direction, speed, ...) information and appearance information.





# What is multiple object tracking(MOT)?

- Compared to single object tracking (SOT), MOT requires two additional tasks to be solved:
  - (1) determining the **number** of objects
  - (2) maintaining their **identities**, since in each moment, the algorithm should decide whether a new track appears and an existing track vanishes.
- Single object tracking
  - The basic assumption is that there is a **single object** in the scene to be tracked



# Motivation

- Deep SORT: Simple online and realtime tracking with a **deep association** metric
- The motivation of deep sort:
  - Integrate appearance information and motion information (Mahalanobis distance for motion information, cosine distance for appearance information)





# Motion information

- The state of each target at some point is modelled as:

$$(u, v, \gamma, h, \dot{x}, \dot{y}, \dot{\gamma}, \dot{h})$$

- $(u, v)$  is bounding box center position,
  - $\gamma$  is aspect ratio,
  - $h$  is height
  - Overdot means their respective **velocities** in image coordinates.
- A standard Kalman filter with **constant velocity** motion and linear observation model, is used to update the above target state.





# Motion information

- A standard Kalman filter with constant velocity motion and linear observation model, is used to update the above target state
- What is Kalman filter?
  - Kalman filtering is an algorithm that uses a series of measurements observed over time, and produces estimates of unknown variables

$$x = x_0 + v_0 \Delta t + \frac{1}{2} a \Delta t^2$$

Where:

$x$  is the target position

$x_0$  is the initial target position

$v_0$  is the initial target velocity

$a$  is the target acceleration

$\Delta t$  is the time interval (5 seconds in our example)

**Measurement Noise**  
**Process Noise**



# Motion information

- Link detections to existing tracks:

$$d^{(1)}(i, j) = (\mathbf{d}_j - \mathbf{y}_i)^T \mathbf{S}_i^{-1} (\mathbf{d}_j - \mathbf{y}_i) \quad (\text{cost})$$

- $\mathbf{d}_j$  is the  $j$ -th bounding box detection
  - $\mathbf{S}_i$  is the covariance matrix of the Kalman filter prediction ( $i$ -th tracker)
  - $\mathbf{y}_i$  is the Kalman filter prediction bounding box
  - The equation calculates the Mahalanobis distance between groundtruth detection and the Kalman filter prediction
- IOU is another kind of cost on motion information [1]
    - IOU simply calculates the **maximum overlap ratio** of any bounding box in the track and the new detection bounding box



# Motion information

- Link detections to existing tracks:

$$d^{(1)}(i, j) = (\mathbf{d}_j - \mathbf{y}_i)^T \mathbf{S}_i^{-1} (\mathbf{d}_j - \mathbf{y}_i)$$

- Then apply a threshold  $t^{(1)}$  to  $d^{(1)}(i, j)$ , to check whether it's feasible to accept this link:

$$b_{i,j}^{(1)} = \mathbb{1}[d^{(1)}(i, j) \leq t^{(1)}] \quad (\text{gate})$$

- For each bounding box detection  $\mathbf{d}_j$ 
  - Compute an appearance descriptor  $\mathbf{r}_j$  with  $\|\mathbf{r}_j\| = 1$  by L2 normalization, where  $\mathbf{r}_j$  comes from a **convolutional neural network** (wide residual).



# Appearance information

- Keep a gallery  $\mathcal{R}_i = \{\mathbf{r}_k^{(i)}\}_{k=1}^{100}$  of the last 100 associated appearance descriptors for each track  $i$ , i.e, only keep the last 100 descriptors.
- Distance between the  $i$ -th track and  $j$ -th detection in appearance space is the smallest cosine distance the  $i$ -th track and  $j$ -th detection that:

$$d^{(2)}(i, j) = \min\{1 - \mathbf{r}_j^T \mathbf{r}_k^{(i)} \mid \mathbf{r}_k^{(i)} \in \mathcal{R}_i\}$$

- Then, apply a threshold  $t^{(2)}$  to  $d^{(2)}(i, j)$ , check whether it's feasible to accept this link:

$$b_{i,j}^{(2)} = \mathbb{1}[d^{(2)}(i, j) \leq t^{(2)}]$$



# Combine motion & appearance information

- **Combine** both metrics using a weighted sum:

$$c_{i,j} = \lambda d^{(1)}(i,j) + (1 - \lambda) d^{(2)}(i,j) \quad (5)$$

- This term is interpreted as the cost of associating the i-th track and j-th detection.

- And check whether motion and appearance are both less than the threshold:

$$b_{i,j} = \prod_{m=1}^2 b_{i,j}^{(m)}. \quad (6)$$

- This term is interpreted as the **gate** of associating the i-th track and j-th detection.



# Symbols in pipeline

---

## Listing 1 Matching Cascade

---

**Input:** Track indices  $\mathcal{T} = \{1, \dots, N\}$ , Detection indices  $\mathcal{D} = \{1, \dots, M\}$ , Maximum age  $A_{\max}$

```
1: Compute cost matrix  $\mathbf{C} = [c_{i,j}]$  using Eq. 5
2: Compute gate matrix  $\mathbf{B} = [b_{i,j}]$  using Eq. 6
3: Initialize set of matches  $\mathcal{M} \leftarrow \emptyset$ 
4: Initialize set of unmatched detections  $\mathcal{U} \leftarrow \mathcal{D}$ 
5: for  $n \in \{1, \dots, A_{\max}\}$  do
6:   Select tracks by age  $\mathcal{T}_n \leftarrow \{i \in \mathcal{T} \mid a_i = n\}$ 
7:    $[x_{i,j}] \leftarrow \text{min\_cost\_matching}(\mathbf{C}, \mathcal{T}_n, \mathcal{U})$ 
8:    $\mathcal{M} \leftarrow \mathcal{M} \cup \{(i, j) \mid b_{i,j} \cdot x_{i,j} > 0\}$ 
9:    $\mathcal{U} \leftarrow \mathcal{U} \setminus \{j \mid \sum_i b_{i,j} \cdot x_{i,j} > 0\}$ 
10: end for
11: return  $\mathcal{M}, \mathcal{U}$ 
```

---

- $\mathcal{T} = \{1, \dots, N\}$ : every track contains all the past detections in that track.
- $\mathcal{D} = \{1, \dots, M\}$ : detections in all frames.
- $A_{\max}$ : tracks has no new added frame in the past  $A_{\max}$  frames are thought dead.
- $\mathbf{C} = [c_{i,j}]$ :  $c_{i,j}$  is the cost of associating the  $i$ -th track and  $j$ -th detection.
- $\mathbf{B} = [b_{i,j}]$ :  $b_{i,j}$  is the gate of associating the  $i$ -th track and  $j$ -th detection.





# Pipeline

---

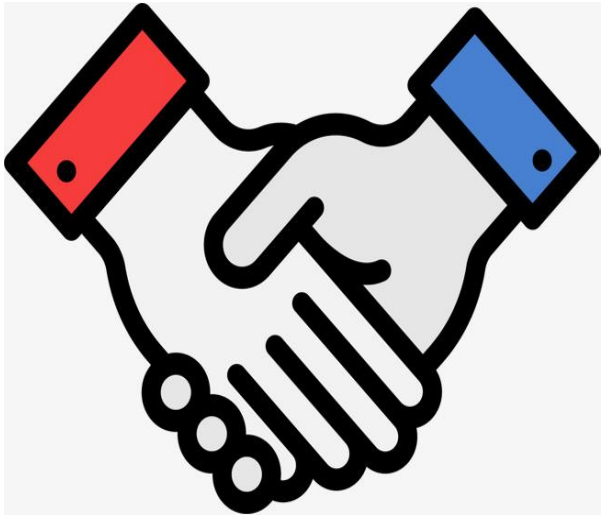
## Listing 1 Matching Cascade

---

**Input:** Track indices  $\mathcal{T} = \{1, \dots, N\}$ , Detection indices  $\mathcal{D} = \{1, \dots, M\}$ , Maximum age  $A_{\max}$  Get detections from another algorithm(YOLO, faster RCNN, ...)

- 1: Compute cost matrix  $\mathbf{C} = [c_{i,j}]$  using Eq. 5 Compute the association cost matrix and
- 2: Compute gate matrix  $\mathbf{B} = [b_{i,j}]$  using Eq. 6 the matrix of admissible associations.
- 3: Initialize set of matches  $\mathcal{M} \leftarrow \emptyset$
- 4: Initialize set of unmatched detections  $\mathcal{U} \leftarrow \mathcal{D}$
- 5: **for**  $n \in \{1, \dots, A_{\max}\}$  **do**  $n$  is in how many frames a tracker has not been updated.
- 6:   Select tracks by age  $\mathcal{T}_n \leftarrow \{i \in \mathcal{T} \mid a_i = n\}$  select the subset of tracks  $\mathcal{T}_n$  that have not been associated with a detection in the last  $n$  frames
- 7:    $[x_{i,j}] \leftarrow \text{min\_cost\_matching}(\mathbf{C}, \mathcal{T}_n, \mathcal{U})$
- 8:    $\mathcal{M} \leftarrow \mathcal{M} \cup \{(i, j) \mid b_{i,j} \cdot x_{i,j} > 0\}$  Solve the linear assignment between tracks in  $\mathcal{T}_n$  and unmatched detections  $\mathcal{U}$
- 9:    $\mathcal{U} \leftarrow \mathcal{U} \setminus \{j \mid \sum_i b_{i,j} \cdot x_{i,j} > 0\}$
- 10: **end for** Update the set of matches and unmatched detections
- 11: **return**  $\mathcal{M}, \mathcal{U}$

---



Thanks



[zhengf@sustc.edu.cn](mailto:zhengf@sustc.edu.cn)