



Robot Perception

Optical Flow & Tracking

Dr. Felix Juefei Xu and Dr. Chen Feng

cfeng@nyu.edu

ROB-GY 6203, Fall 2022



Overview

- Optical flow & Tracking
 - + Basics definitions
 - ++ KLT
 - ++ Mean-shift
 - + Correlation filter
- CNNs for Tracking
 - + Supervised
 - + Self-supervised
- *: know how to code
- ++: know how to derive
- +: know the concept



References

- Sz: Ch 7.1.5, 9.4
- Baker, Simon, and Iain Matthews. "Lucas-kanade 20 years on: A unifying framework." *International journal of computer vision* 56.3 (2004): 221-255.
- Comaniciu, Dorin, and Peter Meer. "Mean shift: A robust approach toward feature space analysis." *IEEE Transactions on Pattern Analysis & Machine Intelligence* 5 (2002): 603-619.
- Feichtenhofer, Christoph, Axel Pinz, and Andrew Zisserman. "Detect to track and track to detect." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3038-3046. 2017.
- Wang, Xiaolong, Allan Jabri, and Alexei A. Efros. "Learning correspondence from the cycle-consistency of time." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2566-2576. 2019.



Tracking an Object



Gif from <https://www.learnopencv.com/object-tracking-using-opencv-cpp-python/>



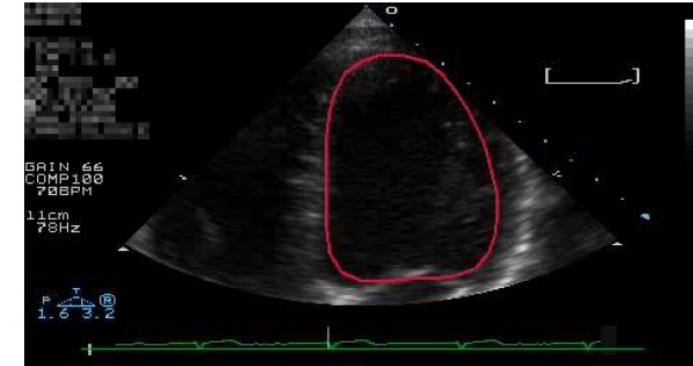
Tracking Applications

- monitoring, assistance, surveillance, control, defense
- robotics, autonomous car driving, rescue
- measurements: medicine, sport, biology, meteorology
- **human computer interaction**
- **augmented reality**
- **film production and postproduction: motion capture, editing**
- **management of video content: indexing, search**
- **action and activity recognition**
- image stabilization
- mobile applications





Tracking Applications





Tracking Definition

- [Forsyth and Ponce, Computer Vision: A modern approach, 2003]
 - “Tracking is the problem of generating an inference about the motion of an object given a sequence of images.”
- Another definition
 - Establishing point-to-point correspondences in consecutive frames of an image sequence
- Yet another definition
 - Given an initial estimate of its position, locate X in a sequence of images
 - X can be a region, an interest point, or an object



Tracking Definition

Given an initial estimate of the pose and state of X :

In all images in a sequence, (in a causal manner)

1. *estimate the pose and state of X*
2. *(optionally) update the model of X*

- Pose: any geometric parameter (position, scale, ...)
- State: appearance, shape/segmentation, visibility, articulations
- Model update: essentially a semi-supervised learning problem
 - a priori information (appearance, shape, dynamics, ...)
 - labeled data (“track this”) + unlabeled data = the sequences
- Causal: for estimation at T, use information from time $t \leq T$



Two Types of Tracking

Short-term Trackers:

- primary objective: “where is X?” = precise estimation of pose
- secondary: be fast; don’t lose track
- evaluation methodology: frame number where failure occurred
- examples: Lucas Kanade tracker, mean-shift tracker

Long-term Tracker-Detectors:

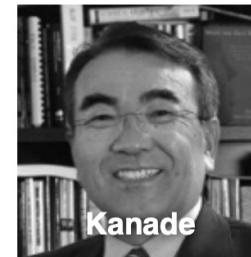
- primary objective: unsupervised learning of a detector, since *every (short-term) tracker fails, sooner or later* (disappearance from the field of view, full occlusion)
- avoid the “*first failure means lost forever*” problem
- close to online-learned detector, but assumes and exploits the fact that a sequence with temporal pose/state dependence is available
- evaluation methodology: precision/recall, false positive/negative rates (i.e. like detectors)
- note: the detector part may help even for short-term tracking problems, provides robustness to fast, unpredictable motions.



KLT Tracker



Lucas



Kanade

An Iterative Image Registration Technique
with an Application to Stereo Vision.

1981



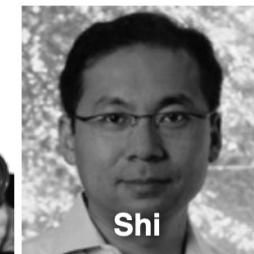
Detection and Tracking of Feature Points.

1991

The original KLT algorithm



Tomasi



Shi

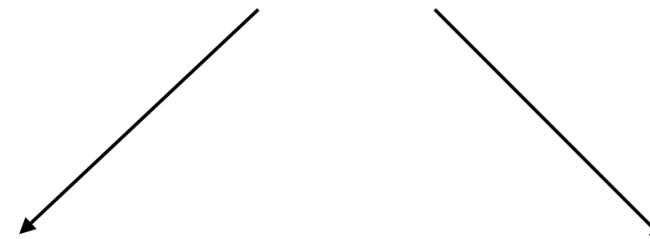
Good Features to Track.

1994



KLT Tracker

Kanade-Lucas-Tomasi



How should we track them from frame to frame?

Lucas-Kanade

Method for aligning (tracking) an image patch

How should we select features?

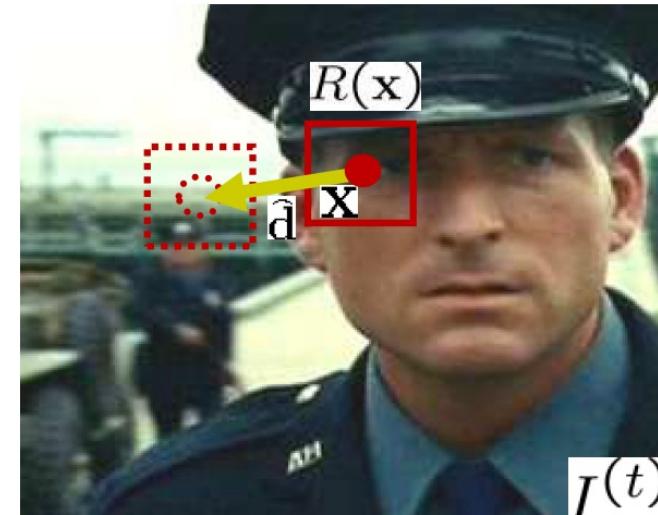
Tomasi-Kanade

Method for choosing the best feature (image patch) for tracking

KLT Tracker

- Problem: tracking “key points” (SIFT, SURF, STAR, RIFF, FAST), or random image patches, as long as possible
 - Input: detected/chosen patches
 - Output: *tracklets* of various life-spans

slide credit:
Patrick Perez

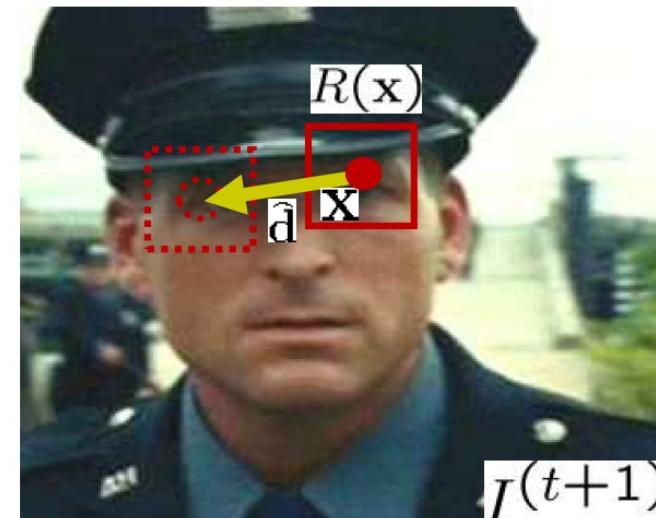


$$\hat{d} = \arg \min_d \underbrace{\sum_{p \in R(x)} |I^{(t+1)}(p + d) - I^{(t)}(p)|^2}_{\text{SSD}}$$

KLT Tracker

- Problem: tracking “key points” (SIFT, SURF, STAR, RIFF, FAST), or random image patches, as long as possible
 - Input: detected/chosen patches
 - Output: *tracklets* of various life-spans

slide credit:
Patrick Perez



$$\hat{d} = \arg \min_d \underbrace{\sum_{p \in R(x)} |I^{(t+1)}(p + d) - I^{(t)}(p)|^2}_{\text{SSD}}$$



KLT Tracker

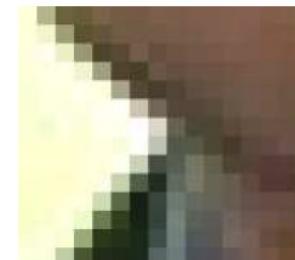
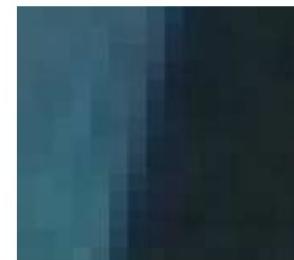
$$f(x) = f(a) + f'(a)(x - a) + R_2$$

– First assuming small displacement: 1st-order Taylor expansion inside SSD

$$\hat{d} = \operatorname{argmin}_d \sum_{p \in R(x)} |I^{(t+1)}(p) + \nabla I^{(t+1)}(p)^T d - I^{(t)}(p)|^2$$
$$\hat{d} = - \left(\sum_{p \in R(x)} \nabla I^{(t+1)}(p) \nabla I^{(t+1)}(p)^T \right)^{-1} \sum_{p \in R(x)} \nabla I^{(t+1)}(p) [I^{(t+1)}(p) - I^{(t)}(p)]$$

For good conditioning, patch must be textured/structured enough:

- Uniform patch: no information
- Contour element: aperture problem (one dimensional information)
- Corners, blobs and texture: best estimate



[Lucas and Kanda 1981][Tomasi and Shi, CVPR'94]

$$\underset{\beta}{\operatorname{minimize}} \boxed{\|Y - X\beta\|^2}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

KLT Tracker

- Translation is usually sufficient for small fragments, but:
 - Perspective transforms and occlusions cause drift and loss
- Two complementary options
 - Kill tracklets when minimum SSD too large
 - Compare as well with *initial patch under affine transform (warp) assumption*

slide credit:
Patrick Perez

$$\hat{\mathbf{d}} = \arg \min_{\mathbf{d}} \sum_{\mathbf{p} \in R_t} |I^{(t+1)}(\mathbf{p} + \mathbf{d}) - I^{(t)}(\mathbf{p})|^2$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{\mathbf{p} \in R_0} |I^{(t+1)}(\mathbf{w}[\mathbf{p}]) - I^{(0)}(\mathbf{p})|^2$$

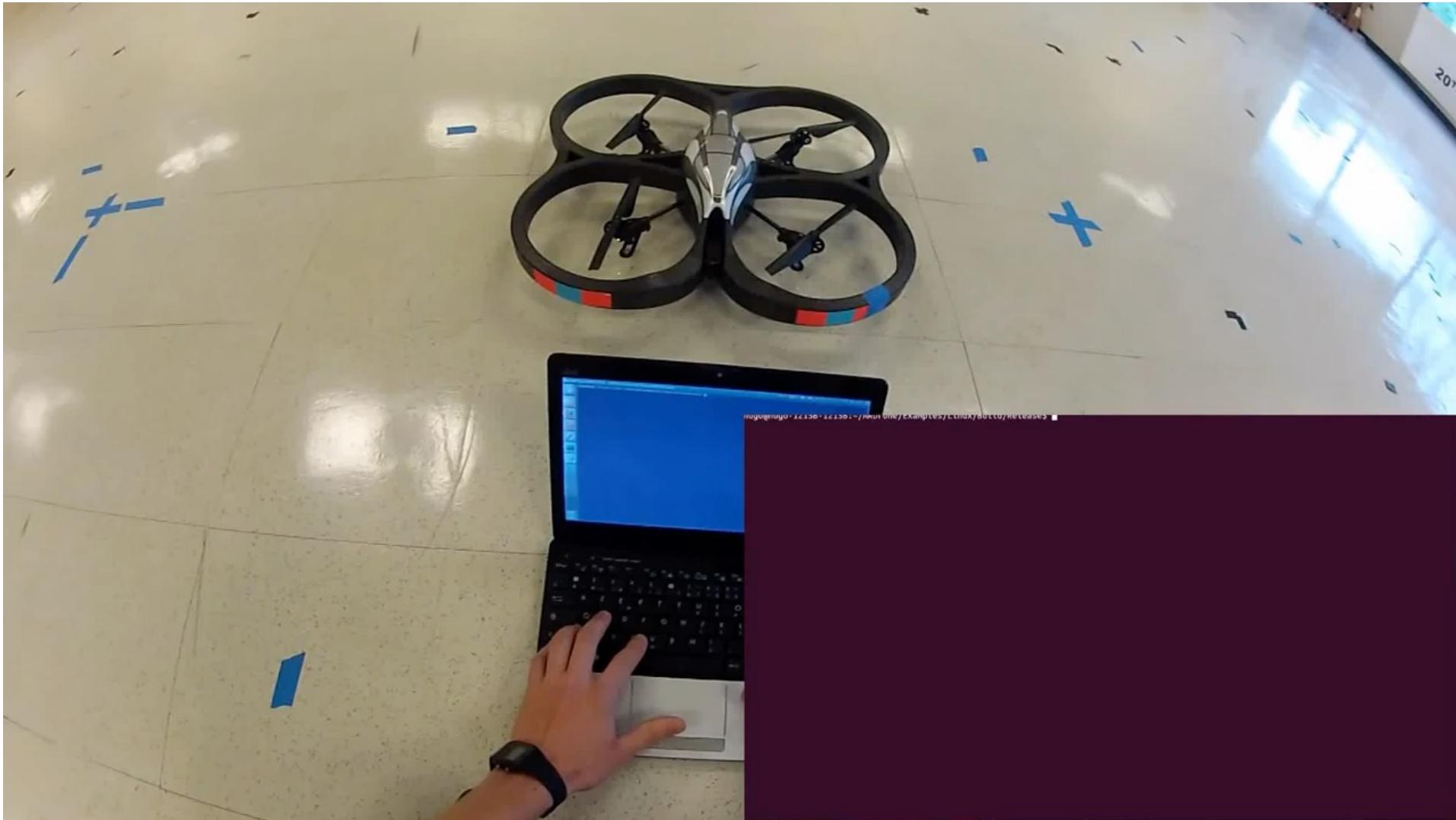


KLT Tracker

- cost function: sum of squared intensity differences between template and window
 - optimization technique: gradient descent
 - model learning: no update
-
- attractive properties:
 - fast
 - easily extended to image-to-image transformations with multiple parameters



Optical Flow in Robotics



<https://youtu.be/C95bngCOv9Q>



Mean Shift Theory

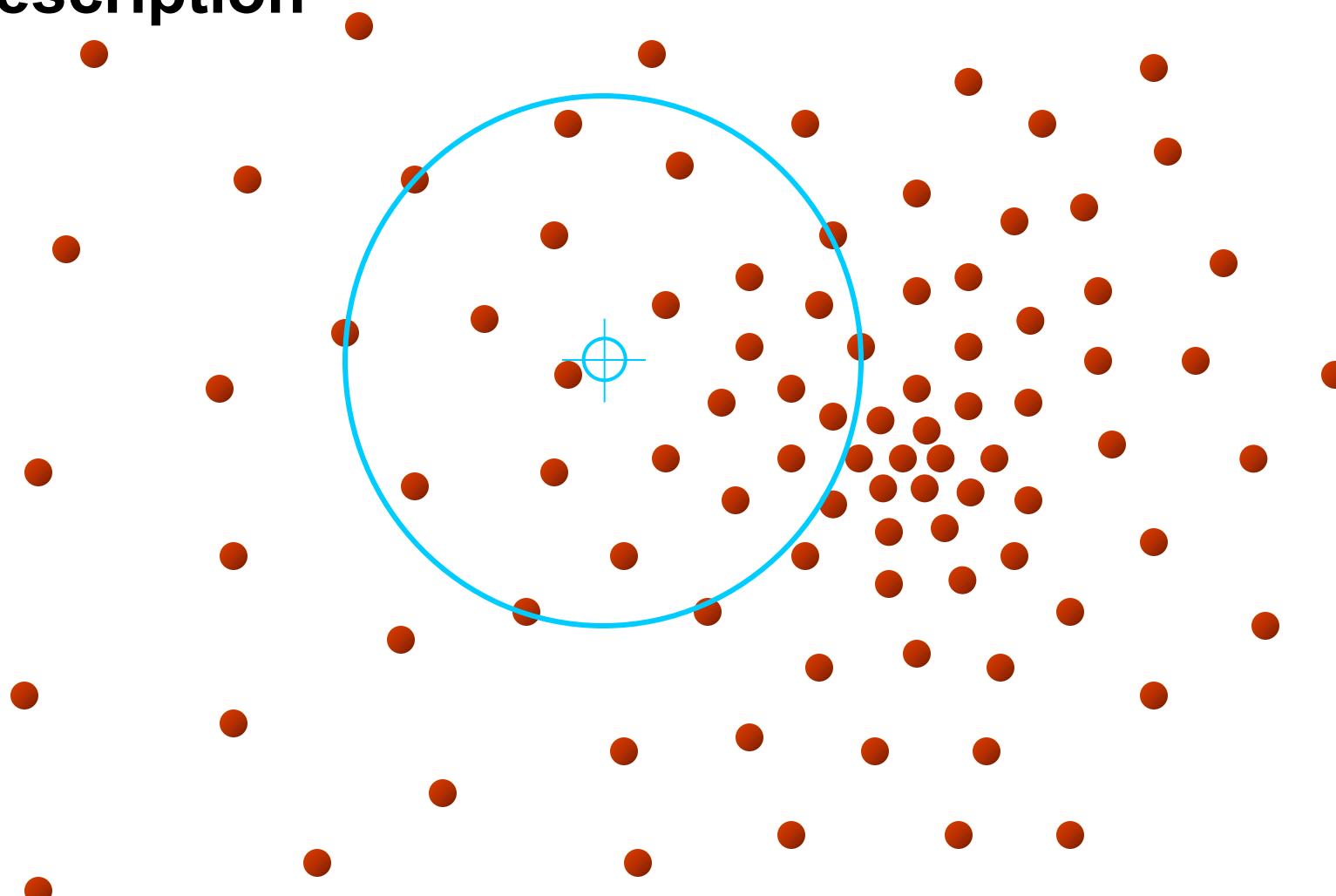
Comaniciu, Dorin, Visvanathan Ramesh, and Peter Meer.
"Real-Time Tracking of Non-Rigid Objects Using Mean Shift."

Proceedings IEEE Conference on Computer Vision and Pattern Recognition.
CVPR 2000. Vol. 2. IEEE, 2000.

Won the Best Paper Award CVPR 2000



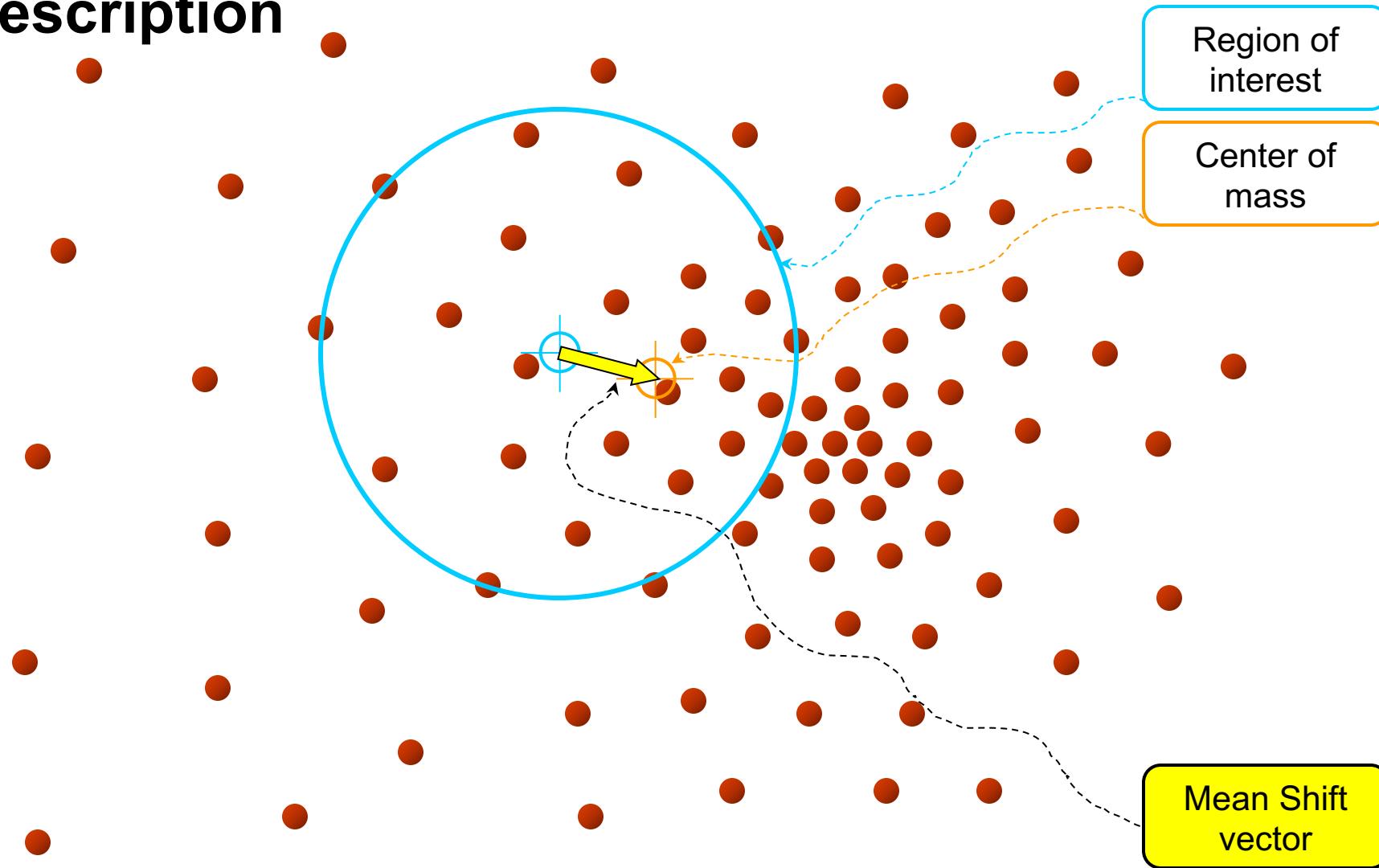
Intuitive Description



Objective : Find the densest region
Distribution of identical billiard balls



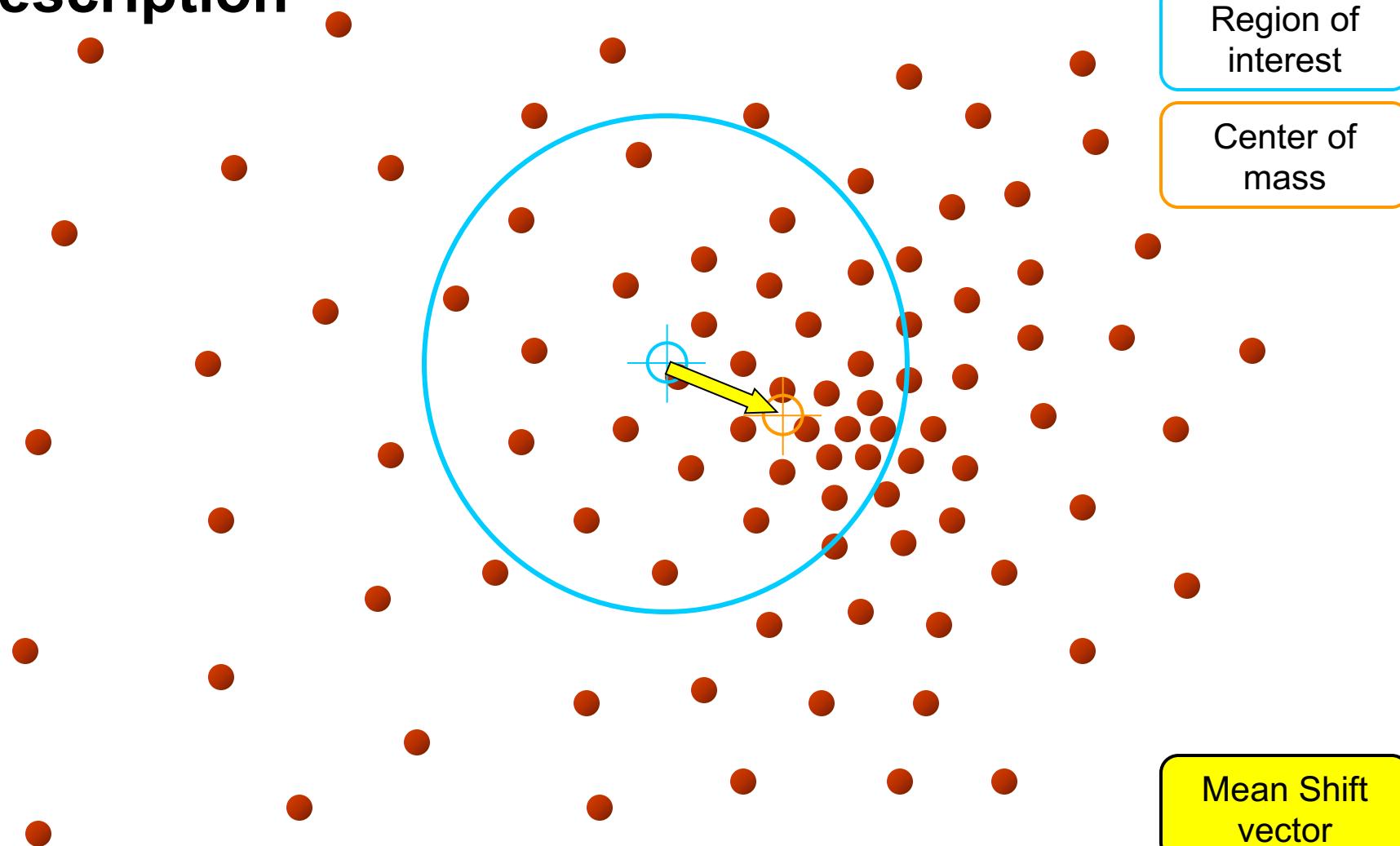
Intuitive Description



Objective : Find the densest region
Distribution of identical billiard balls



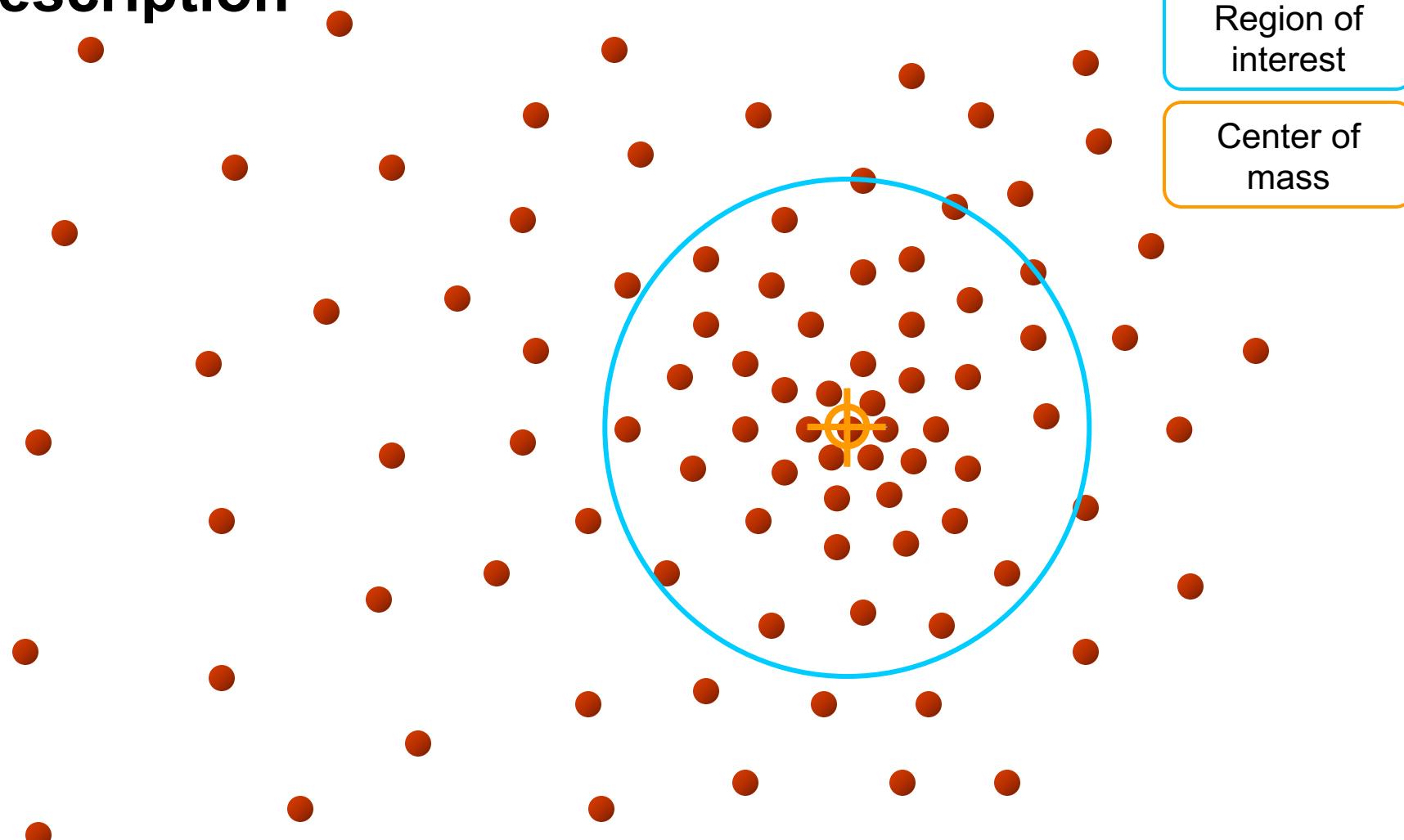
Intuitive Description



Objective : Find the densest region
Distribution of identical billiard balls



Intuitive Description



Objective : Find the densest region
Distribution of identical billiard balls



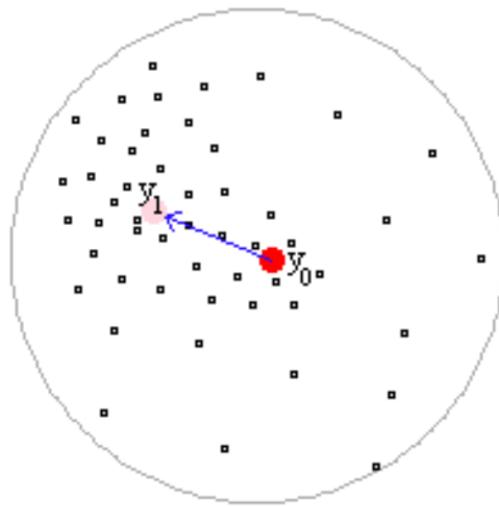
Mean Shift Vector

- Given:
 - Data points and approximate location of the mean of this data:
- Task:
 - Estimate the exact location of the mean of the data by determining the shift vector from the initial mean.



Mean Shift Vector

- An example



$$M_h(\mathbf{y}) = \left[\frac{1}{n_x} \sum_{i=1}^{n_x} \mathbf{x}_i \right] - \mathbf{y}_0$$

Mean shift vector always points towards the direction of the maximum increase in the density.



Mean Shift (Weighted)

$$M_h(\mathbf{y}_0) = \left[\frac{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0) \mathbf{x}_i}{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0)} \right] - \mathbf{y}_0$$

n_x : number of points in the kernel
 \mathbf{y}_0 : initial mean location
 \mathbf{x}_i : data points
 h : kernel radius

Weights are determined using kernels (masks):
Uniform, Gaussian or Epanechnikov



Properties of Mean Shift

- Mean shift vector has the direction of the gradient of the density estimate.
- It is computed iteratively for obtaining the maximum density in the local neighborhood.



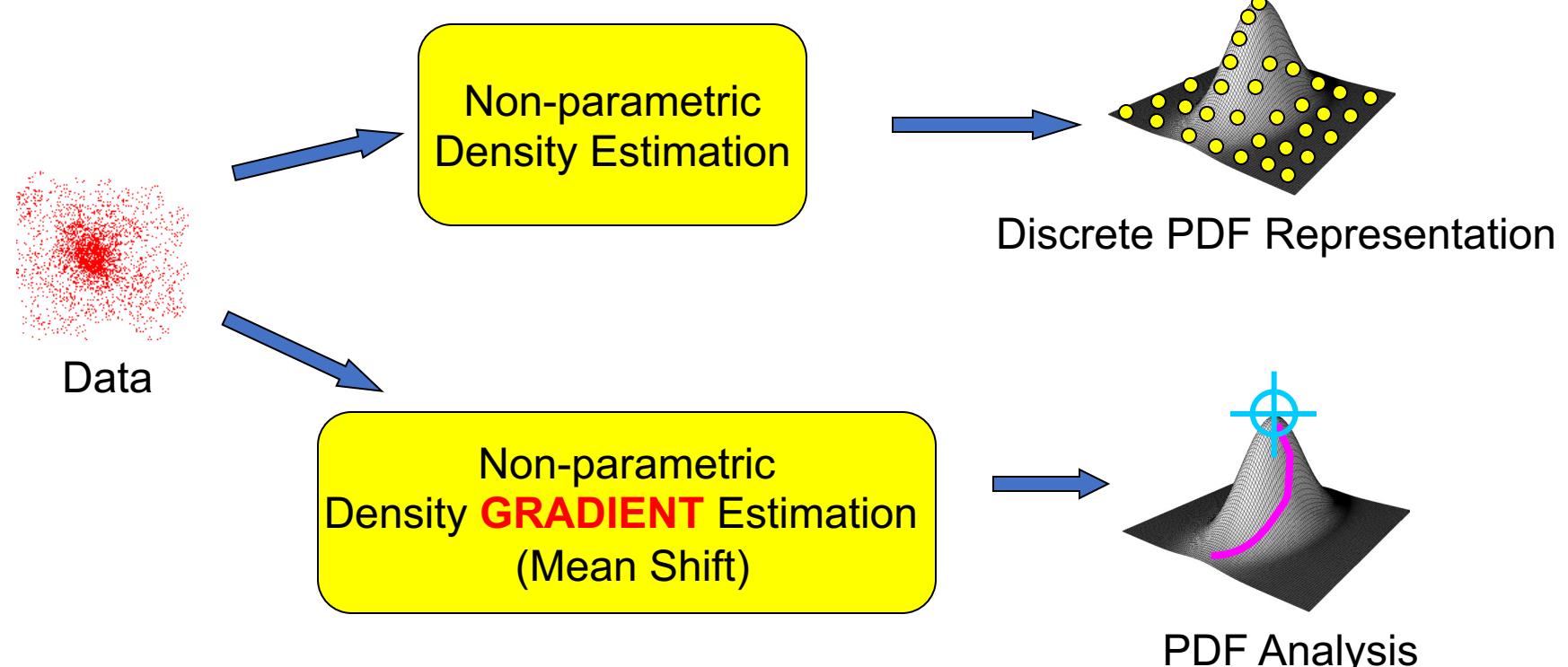
What is Mean Shift ?

A tool for:

Finding modes in a set of data samples, manifesting an underlying probability density function (PDF) in \mathbb{R}^N

PDF in feature space

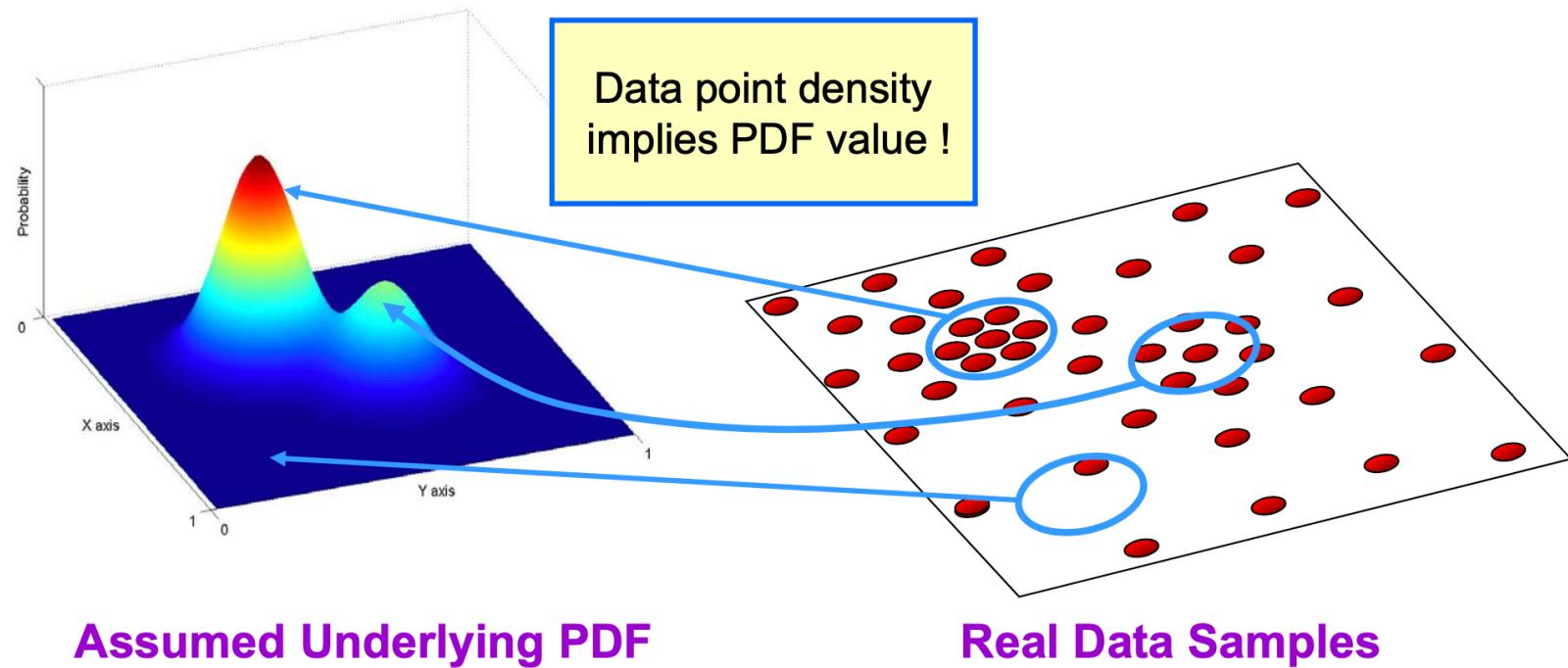
- Color space
- Scale space
- Actually, any feature space you can conceive





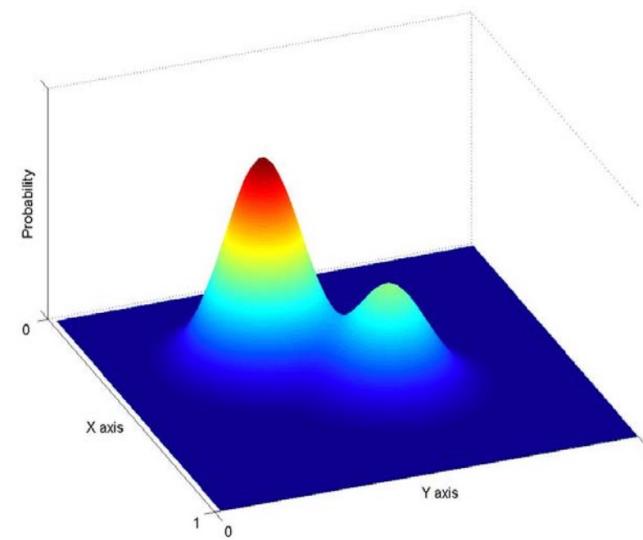
Non-Parametric Density Estimation

- Assumption : The data points are sampled from an underlying PDF

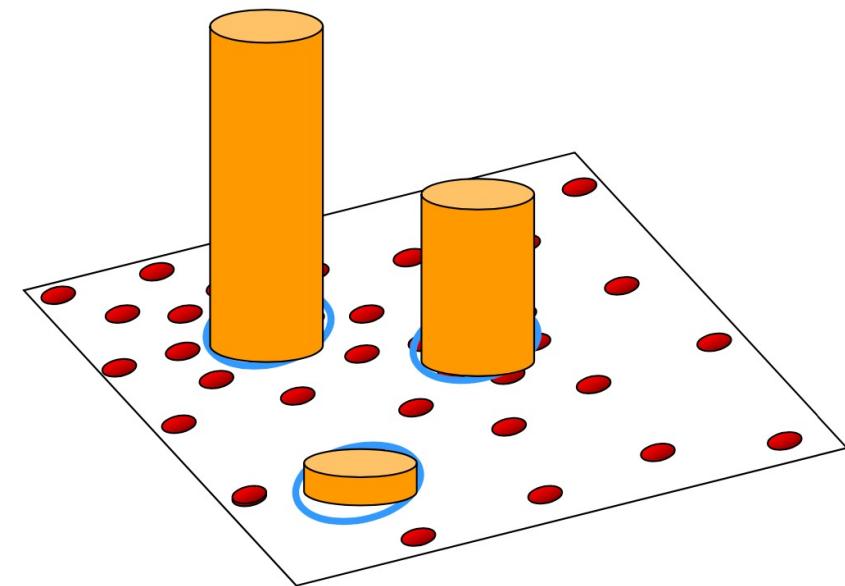




Non-Parametric Density Estimation



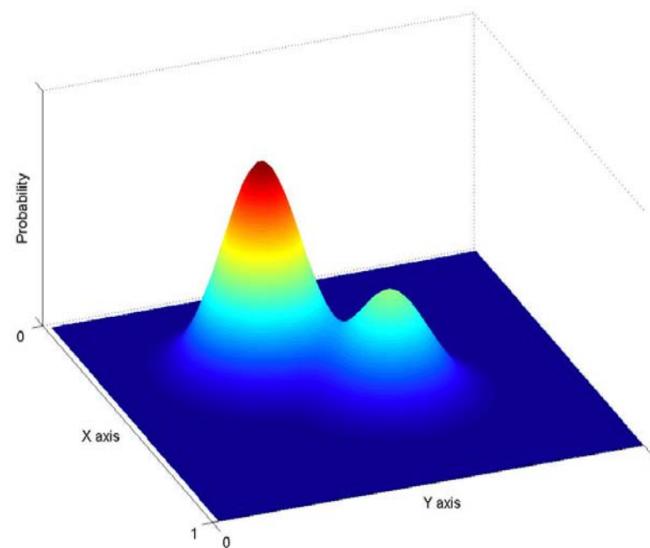
Assumed Underlying PDF



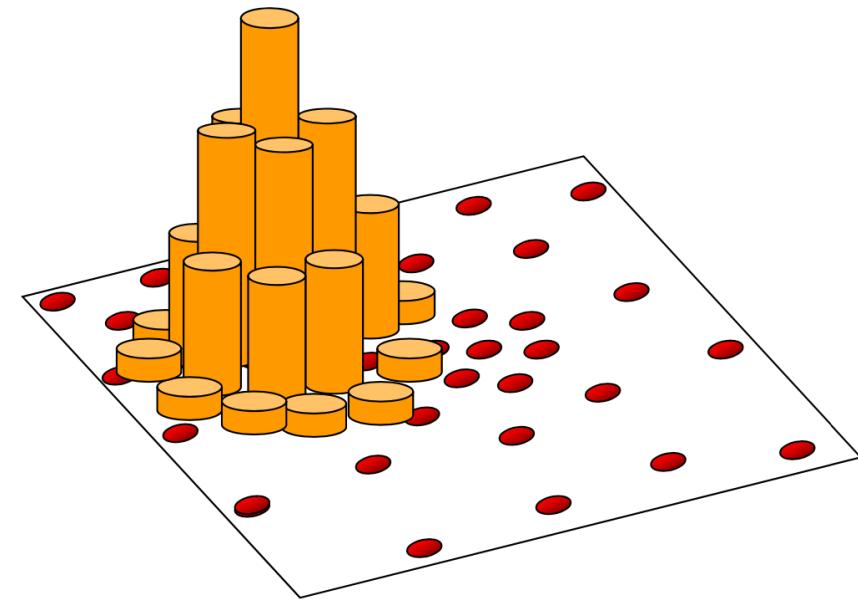
Real Data Samples



Non-Parametric Density Estimation



Assumed Underlying PDF

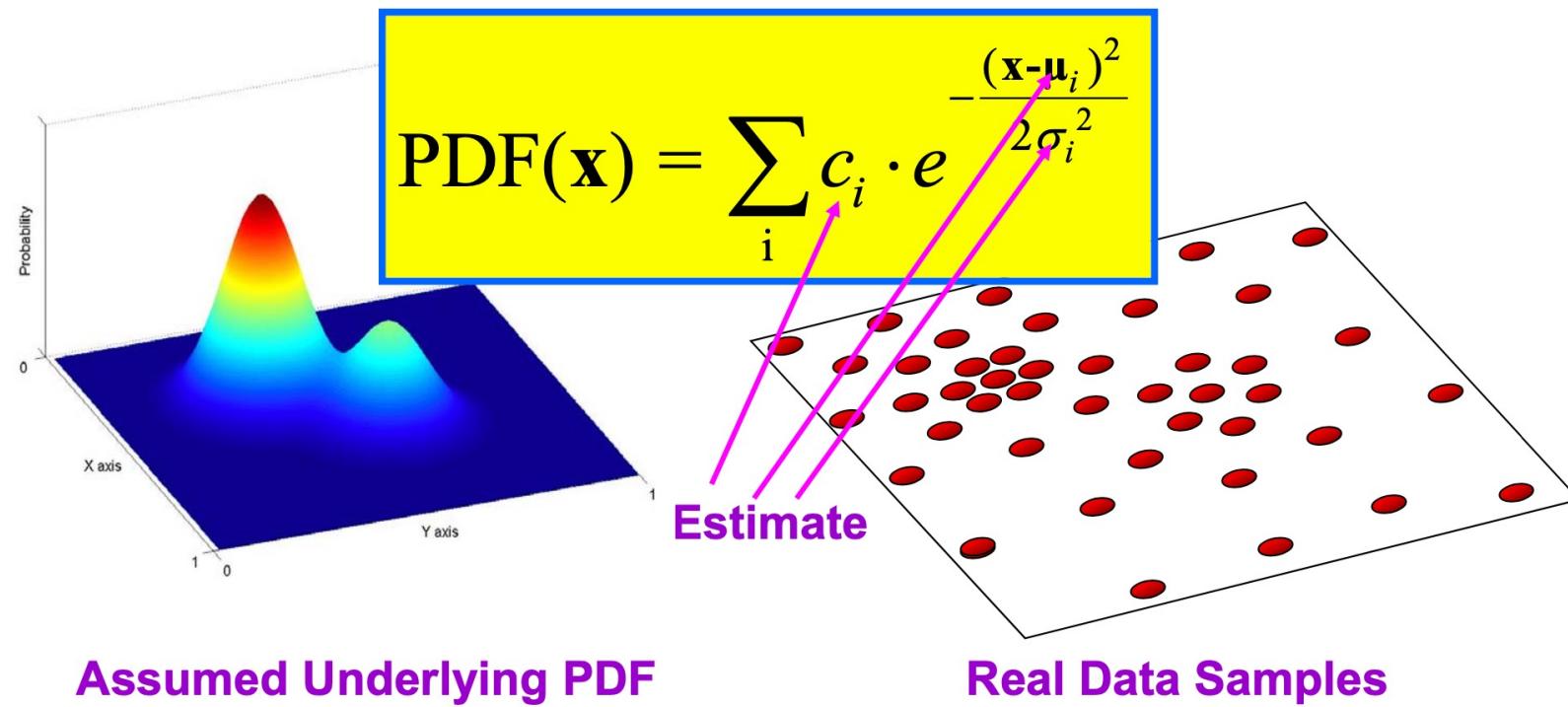


Real Data Samples



Parametric Density Estimation

- Assumption: The data points are samples from an underlying PDF



Assumed Underlying PDF

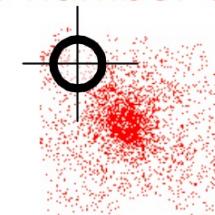
Real Data Samples



Kernel Density Estimation - Function Forms

$$P(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K(\mathbf{x} - \mathbf{x}_i)$$

A function of some finite number of data points
 $\mathbf{x}_1, \dots, \mathbf{x}_n$



Data

In practice one uses the forms:

$$K(\mathbf{x}) = c \prod_{i=1}^d k(x_i)$$

or

$$K(\mathbf{x}) = ck(\|\mathbf{x}\|)$$

Same function on each dimension

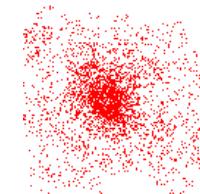
Function of vector length only



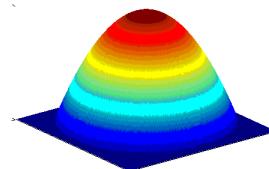
Kernel Density Estimation - Various Kernels

$$P(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K(\mathbf{x} - \mathbf{x}_i)$$

A function of some finite number of data points
 $x_1 \dots x_n$



Data



Examples:

- Epanechnikov Kernel

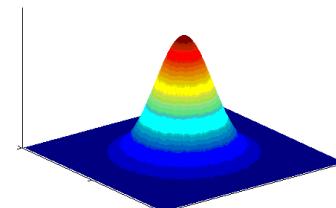
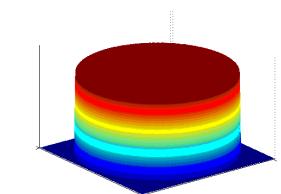
$$K_E(\mathbf{x}) = \begin{cases} c(1 - \|\mathbf{x}\|^2) & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- Uniform Kernel

$$K_U(\mathbf{x}) = \begin{cases} c & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- Normal Kernel

$$K_N(\mathbf{x}) = c \cdot \exp\left(-\frac{1}{2} \|\mathbf{x}\|^2\right)$$

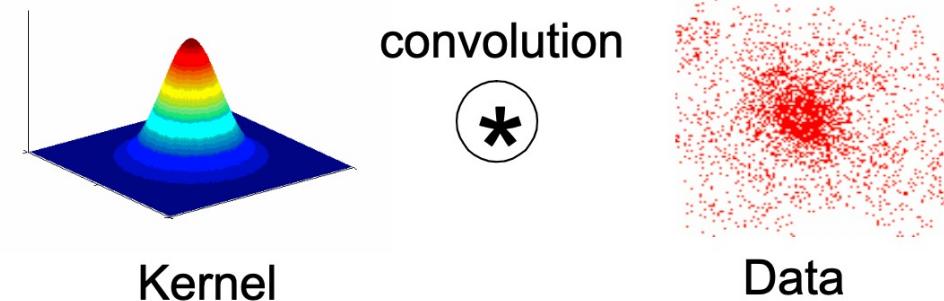




Kernel Density Estimation - Key Idea

$$P(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K(\mathbf{x} - \mathbf{x}_i)$$

Superposition of kernels, centered at each data point is equivalent to convolving the data points with the kernel.





Profile and Kernel

Radially symmetric Kernel

$$K(x) = ck(\|x\|^2)$$

Profile

$$P(x) = \frac{1}{n} \sum_{i=1}^n K(x - x_i) = \frac{1}{n} c \sum_{i=1}^n k(\|x - x_i\|^2)$$

Kernel Density *Gradient* Estimation

$$P(\mathbf{x}) = \frac{1}{n} c \sum_{i=1}^n k(\|\mathbf{x} - \mathbf{x}_i\|^2)$$

$$\nabla P(\mathbf{x}) = \frac{1}{n} c \sum_{i=1}^n \nabla k(\|\mathbf{x} - \mathbf{x}_i\|^2)$$

$$\nabla P(\mathbf{x}) = \frac{1}{n} 2c \sum_{i=1}^n (\mathbf{x} - \mathbf{x}_i) k'(\|\mathbf{x} - \mathbf{x}_i\|^2)$$



Kernel Density *Gradient* Estimation

$$\nabla P(\mathbf{x}) = \frac{1}{n} 2c \sum_{i=1}^n (\mathbf{x} - \mathbf{x}_i) k'(\|\mathbf{x} - \mathbf{x}_i\|^2)$$

$$\nabla P(\mathbf{x}) = \frac{1}{n} 2c \sum_{i=1}^n (\mathbf{x}_i - \mathbf{x}) g(\|\mathbf{x} - \mathbf{x}_i\|^2)$$

$$\nabla P(\mathbf{x}) = \frac{1}{n} 2c \sum_{i=1}^n \mathbf{x}_i g(\|\mathbf{x} - \mathbf{x}_i\|^2) - \frac{1}{n} 2c \sum_{i=1}^n \mathbf{x} g(\|\mathbf{x} - \mathbf{x}_i\|^2)$$

$$\nabla P(\mathbf{x}) = \frac{1}{n} 2c \sum_{i=1}^n g(\|\mathbf{x} - \mathbf{x}_i\|^2) \left[\frac{\sum_{i=1}^n \mathbf{x}_i g(\|\mathbf{x} - \mathbf{x}_i\|^2)}{\sum_{i=1}^n g(\|\mathbf{x} - \mathbf{x}_i\|^2)} - \mathbf{x} \right]$$
$$g(\mathbf{x}) = k'(\mathbf{x})$$



Kernel Density *Gradient* Estimation

$$\nabla P(\mathbf{x}) = -\frac{1}{n} 2c \sum_{i=1}^n g(\|\mathbf{x} - \mathbf{x}_i\|^2) \left[\frac{\sum_{i=1}^n \mathbf{x}_i g(\|\mathbf{x} - \mathbf{x}_i\|^2)}{\sum_{i=1}^n g(\|\mathbf{x} - \mathbf{x}_i\|^2)} - \mathbf{x} \right]$$

$$\nabla P(\mathbf{x}) = -\frac{1}{n} 2c \sum_{i=1}^n g_i \left[\frac{\sum_{i=1}^n \mathbf{x}_i g_i}{\sum_{i=1}^n g_i} - \mathbf{x} \right]$$



Kernel Density *Gradient* Estimation

$$\nabla P(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \nabla K(\mathbf{x} - \mathbf{x}_i)$$

Give up estimating the PDF !
Estimate ONLY the gradient

Using the
Kernel form:

$$K(\mathbf{x} - \mathbf{x}_i) = ck \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h} \right)$$

We get :

Size of window

$$\nabla P(\mathbf{x}) = \frac{c}{n} \sum_{i=1}^n \nabla k_i = \frac{c}{n} \left[\sum_{i=1}^n g_i \right] \begin{bmatrix} \frac{\sum_{i=1}^n \mathbf{x}_i g_i}{n} - \mathbf{x} \\ \sum_{i=1}^n g_i \end{bmatrix}$$



Computing The Mean Shift

$$\nabla P(\mathbf{x}) = \frac{c}{n} \sum_{i=1}^n \nabla k_i = \frac{c}{n} \left[\sum_{i=1}^n g_i \right]$$

$\frac{\sum_{i=1}^n \mathbf{x}_i g_i}{\sum_{i=1}^n g_i} - \mathbf{x}$

Yet another Kernel density estimation !

Simple Mean Shift procedure:

- Compute mean shift vector

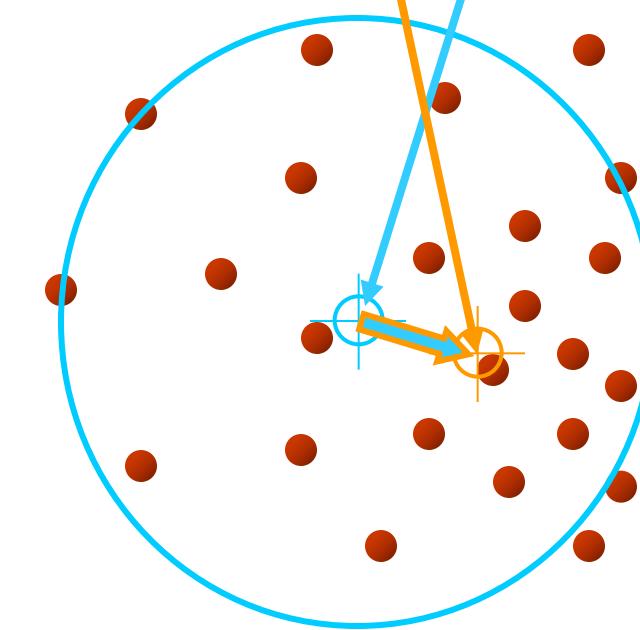
$$\mathbf{m}(\mathbf{x}) = \left[\frac{\sum_{i=1}^n \mathbf{x}_i g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h}\right)}{\sum_{i=1}^n g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h}\right)} - \mathbf{x} \right]$$

- Translate the Kernel window by $\mathbf{m}(\mathbf{x})$

$$M_h(\mathbf{y}_0) = \left[\frac{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0) \mathbf{x}_i}{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0)} \right] - \mathbf{y}_0$$

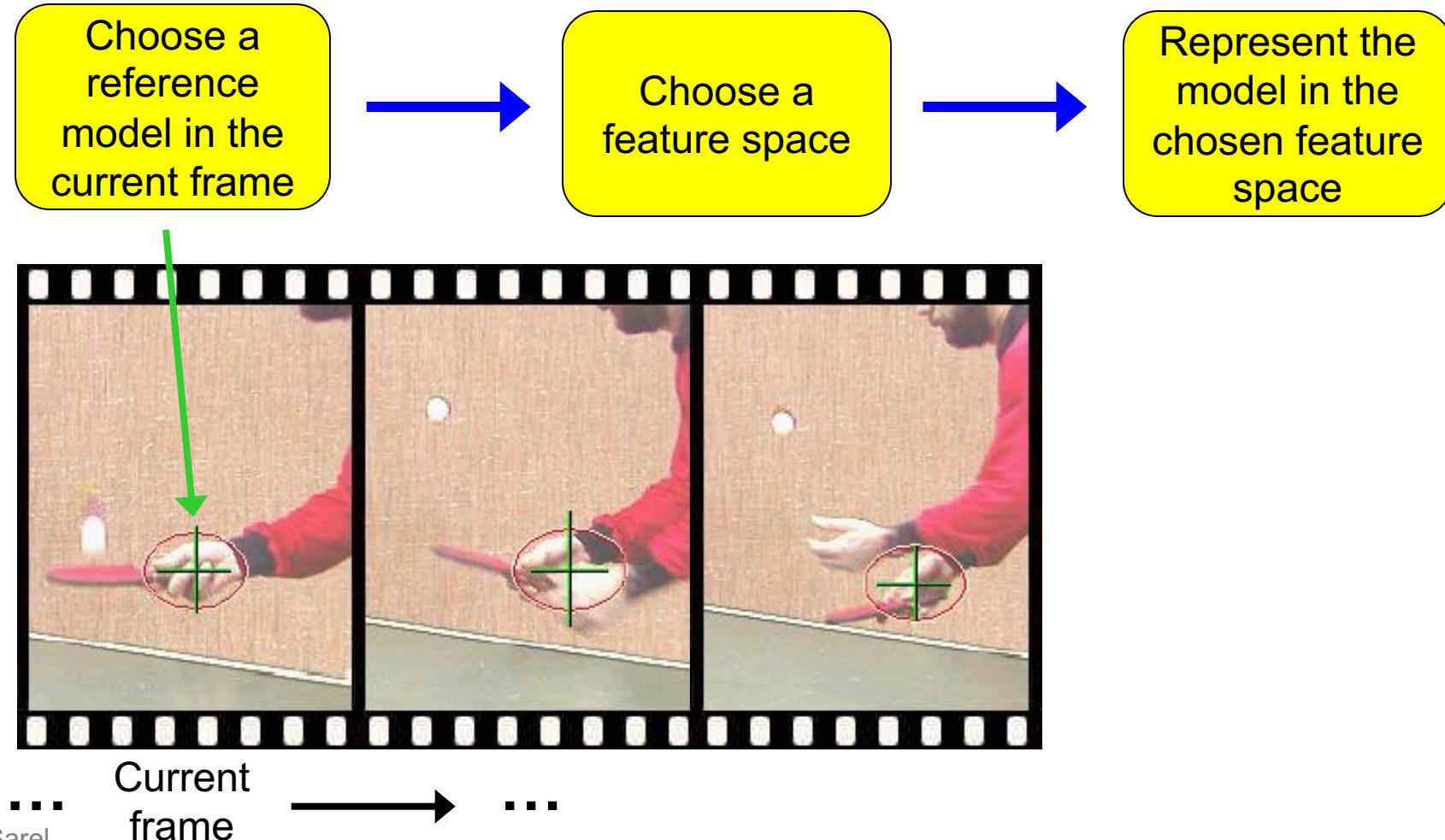
n_x : number of points in the kernel
 \mathbf{y}_0 : initial mean location
 \mathbf{x}_i : data points
 h : kernel radius

Weights are determined using kernels (masks):
 Uniform, Gaussian or Epanechnikov



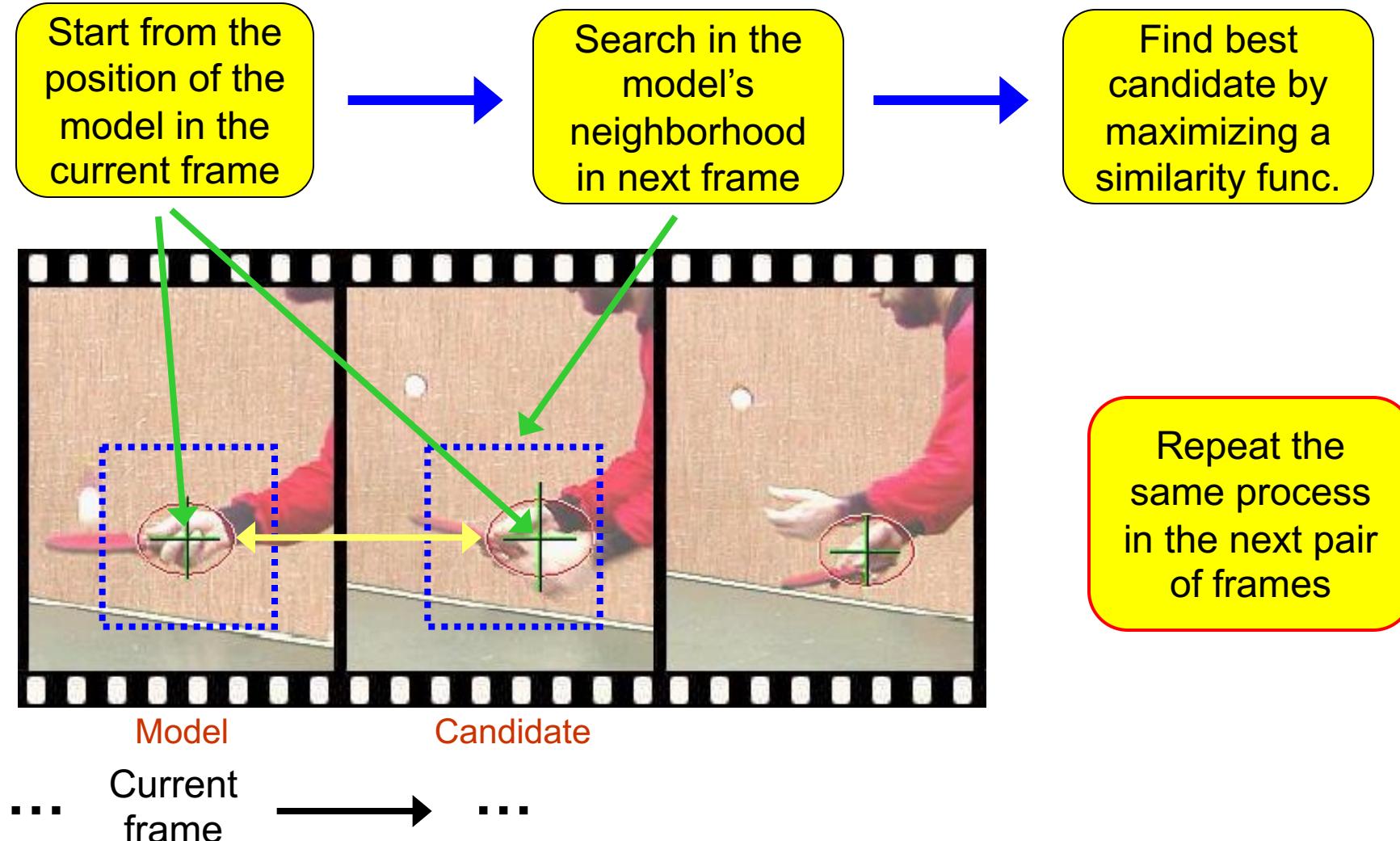


Mean-Shift Object Tracking - General Framework: Target Representation



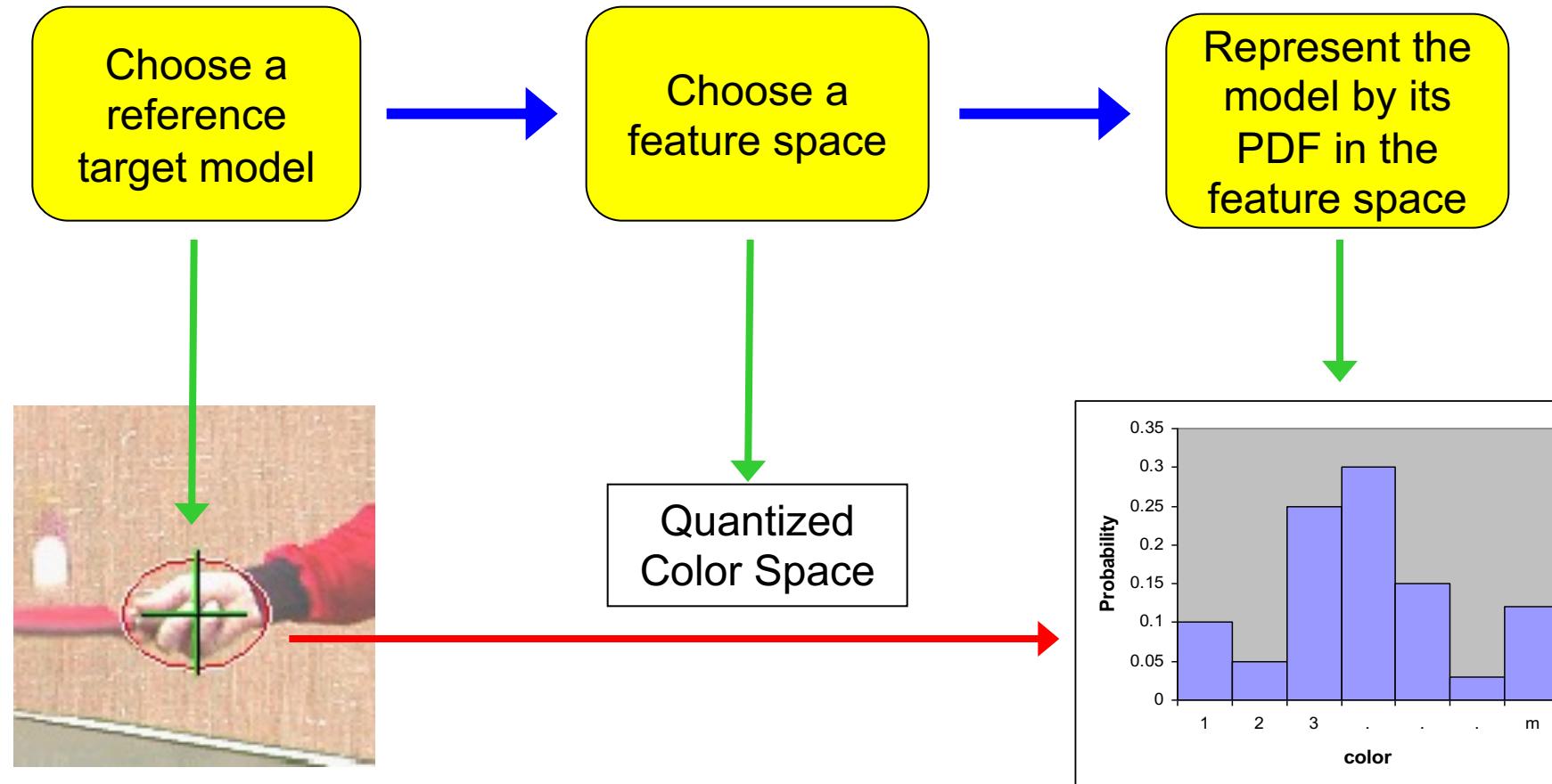


Mean-Shift Object Tracking - General Framework: Target Localization





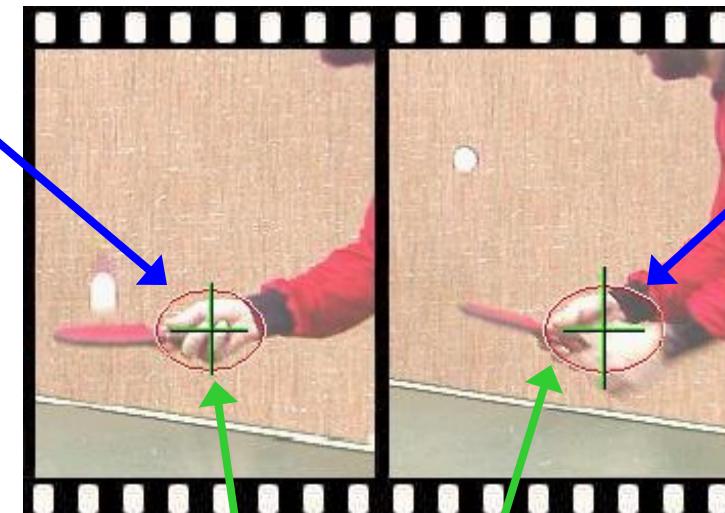
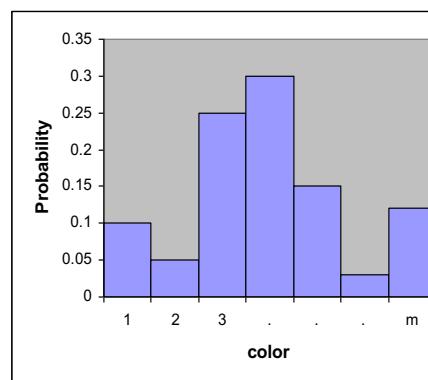
Mean-Shift Object Tracking - Target Representation



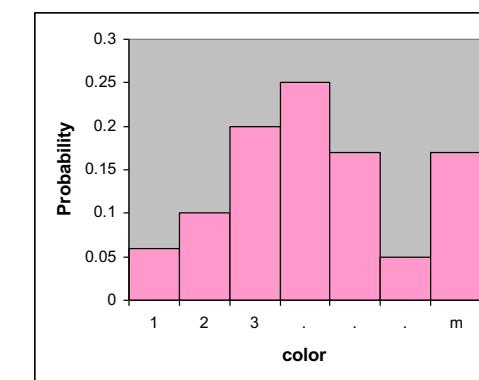


Mean-Shift Object Tracking - PDF Representation

Target Model
(centered at 0)



Target Candidate
(centered at y)



$$\vec{q} = \{q_u\}_{u=1..m} \quad \sum_{u=1}^m q_u = 1$$

$$\vec{p}(y) = \{p_u(y)\}_{u=1..m} \quad \sum_{u=1}^m p_u = 1$$

Similarity
Function:

$$f(y) = f[\vec{q}, \vec{p}(y)]$$



Mean-Shift Object Tracking - Similarity Function

Target model:

$$\vec{q} = (q_1, \dots, q_m)$$

Target candidate:

$$\vec{p}(y) = (p_1(y), \dots, p_m(y))$$

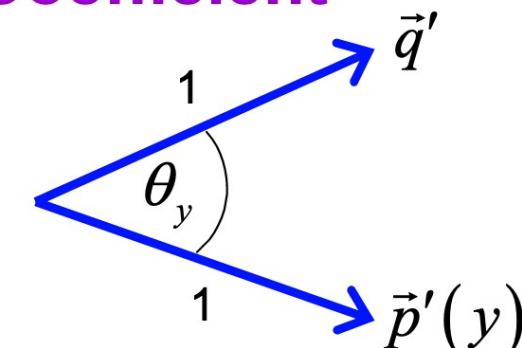
Similarity function:

$$f(y) = f[\vec{p}(y), \vec{q}] = ?$$

The Bhattacharyya Coefficient

$$\vec{q}' = (\sqrt{q_1}, \dots, \sqrt{q_m})$$

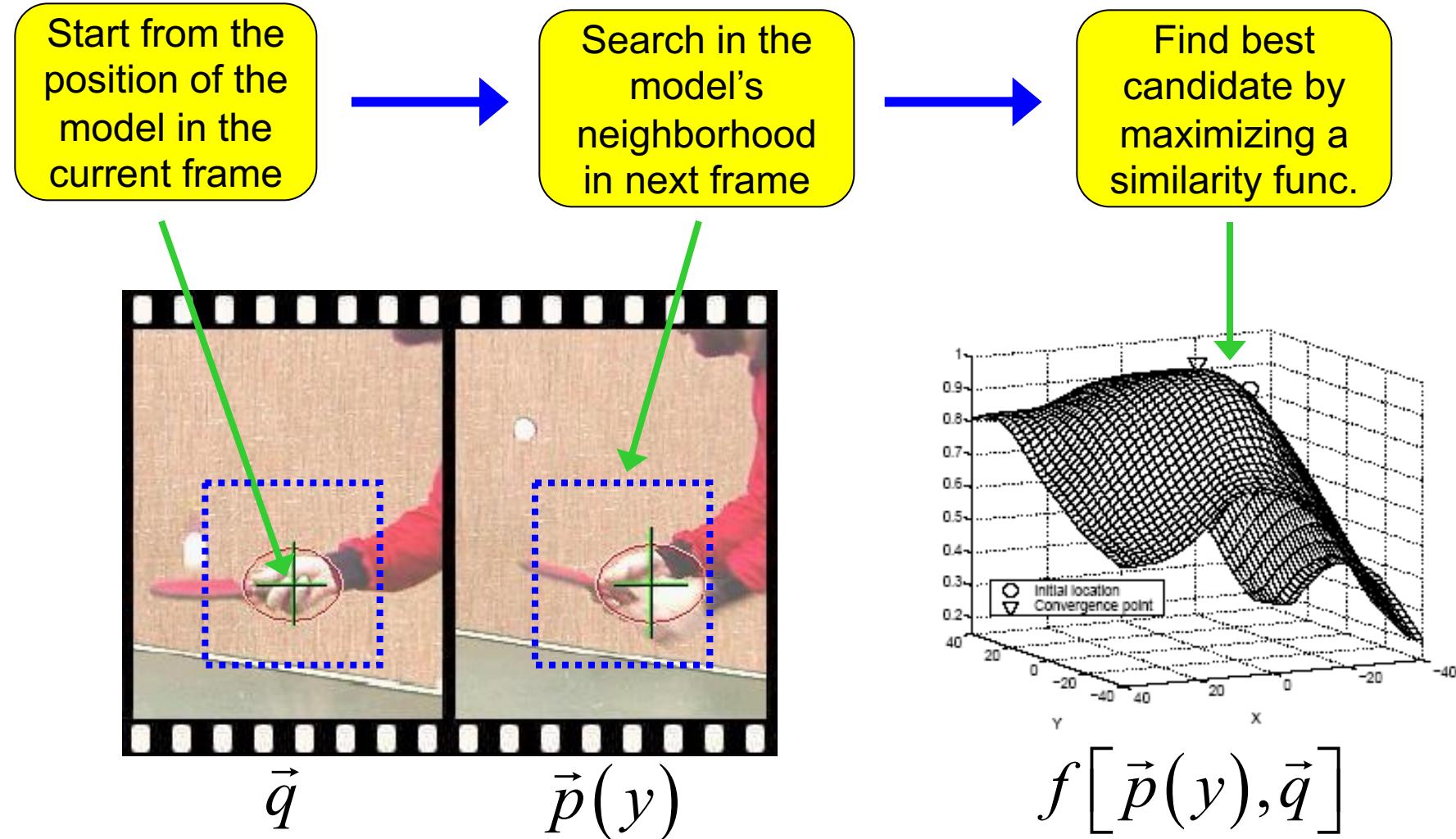
$$\vec{p}'(y) = (\sqrt{p_1(y)}, \dots, \sqrt{p_m(y)})$$



$$f(y) = \cos \theta_y = \frac{\vec{p}'(y)^T \vec{q}'}{\|\vec{p}'(y)\| \cdot \|\vec{q}'\|} = \sum_{u=1}^m \sqrt{p_u(y) q_u}$$



Mean-Shift Object Tracking - Target Localization Algorithm





Mean-Shift Object Tracking - Approximating the Similarity Func.

$$f(y) = \sum_{u=1}^m \sqrt{p_u(y)q_u}$$

Model location: y_0

Candidate location: y

Linear
approx.
(around y_0)

$$f(y) \approx \underbrace{\frac{1}{2} \sum_{u=1}^m \sqrt{p_u(y_0)q_u}}_{\text{Independent of } y} + \underbrace{\frac{1}{2} \sum_{u=1}^m p_u(y)}_{\text{Density estimate!}} \sqrt{\frac{q_u}{p_u(y_0)}}$$

Independent
of y

$$p_u(y) = C_h \sum_{b(x_i)=u} k\left(\frac{\|y-x_i\|^2}{h}\right)$$

$$\frac{C_h}{2} \sum_{i=1}^n w_i k\left(\frac{\|y-x_i\|^2}{h}\right)$$

Density
estimate!
(as a function
of y)

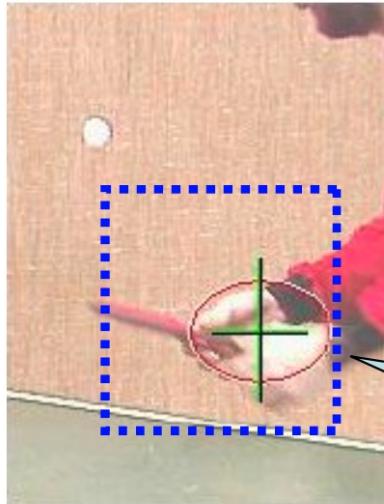


Mean-Shift Object Tracking - Maximizing the Similarity Function

The mode of

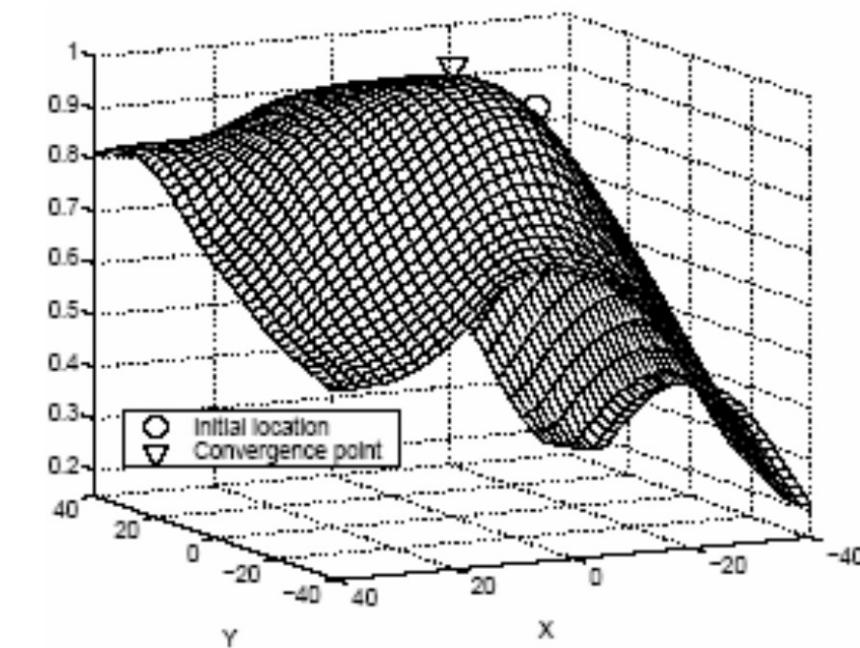
$$\frac{C_h}{2} \sum_{i=1}^n w_i k\left(\frac{\|y-x_i\|^2}{h}\right) = \text{sought maximum}$$

Important Assumption:



The target representation provides sufficient discrimination

One mode in the searched neighborhood





Mean-Shift Object Tracking - Applying Mean-Shift

The mode of

$$\frac{C_h}{2} \sum_{i=1}^n w_i k\left(\left\|\frac{y-x_i}{h}\right\|^2\right) = \text{sought maximum}$$

Original
Mean-Shift:

Find mode of $c \sum_{i=1}^n k\left(\left\|\frac{y-x_i}{h}\right\|^2\right)$ using

$$y_1 = \frac{\sum_{i=1}^n x_i g\left(\left\|\frac{y_0-x_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{y_0-x_i}{h}\right\|^2\right)}$$

Extended
Mean-Shift:

Find mode of $c \sum_{i=1}^n w_i k\left(\left\|\frac{y-x_i}{h}\right\|^2\right)$ using

$$y_1 = \frac{\sum_{i=1}^n x_i w_i g\left(\left\|\frac{y_0-x_i}{h}\right\|^2\right)}{\sum_{i=1}^n w_i g\left(\left\|\frac{y_0-x_i}{h}\right\|^2\right)}$$

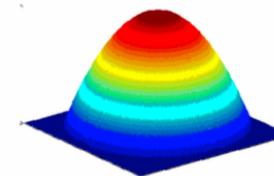


Mean-Shift Object Tracking - Choosing the Kernels

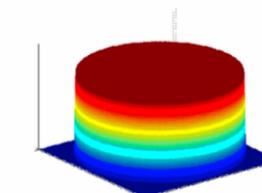
A special class of radially symmetric kernels:

$$K(x) = ck(\|x\|^2)$$

Epanechnikov kernel:



Uniform kernel:



$$k(x) = \begin{cases} 1-x & \text{if } \|x\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$g(x) = -k'(x) = \begin{cases} 1 & \text{if } \|x\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$y_1 = \frac{\sum_{i=1}^n x_i w_i g\left(\left\|\frac{y_0 - x_i}{h}\right\|^2\right)}{\sum_{i=1}^n w_i g\left(\left\|\frac{y_0 - x_i}{h}\right\|^2\right)}$$
$$y_1 = \frac{\sum_{i=1}^n x_i w_i}{\sum_{i=1}^n w_i}$$



Likelihood Maximization Using Mean Shift Vector

Maximization of the likelihood of target and candidate depends on the weights:

$$w_i(\mathbf{y}_o) = \sum_{u=1}^m \delta[S(\mathbf{x}_i) - u] \sqrt{\frac{q_u}{\hat{p}_u(\mathbf{y}_o)}}$$

Again, the mean shift vector M is:

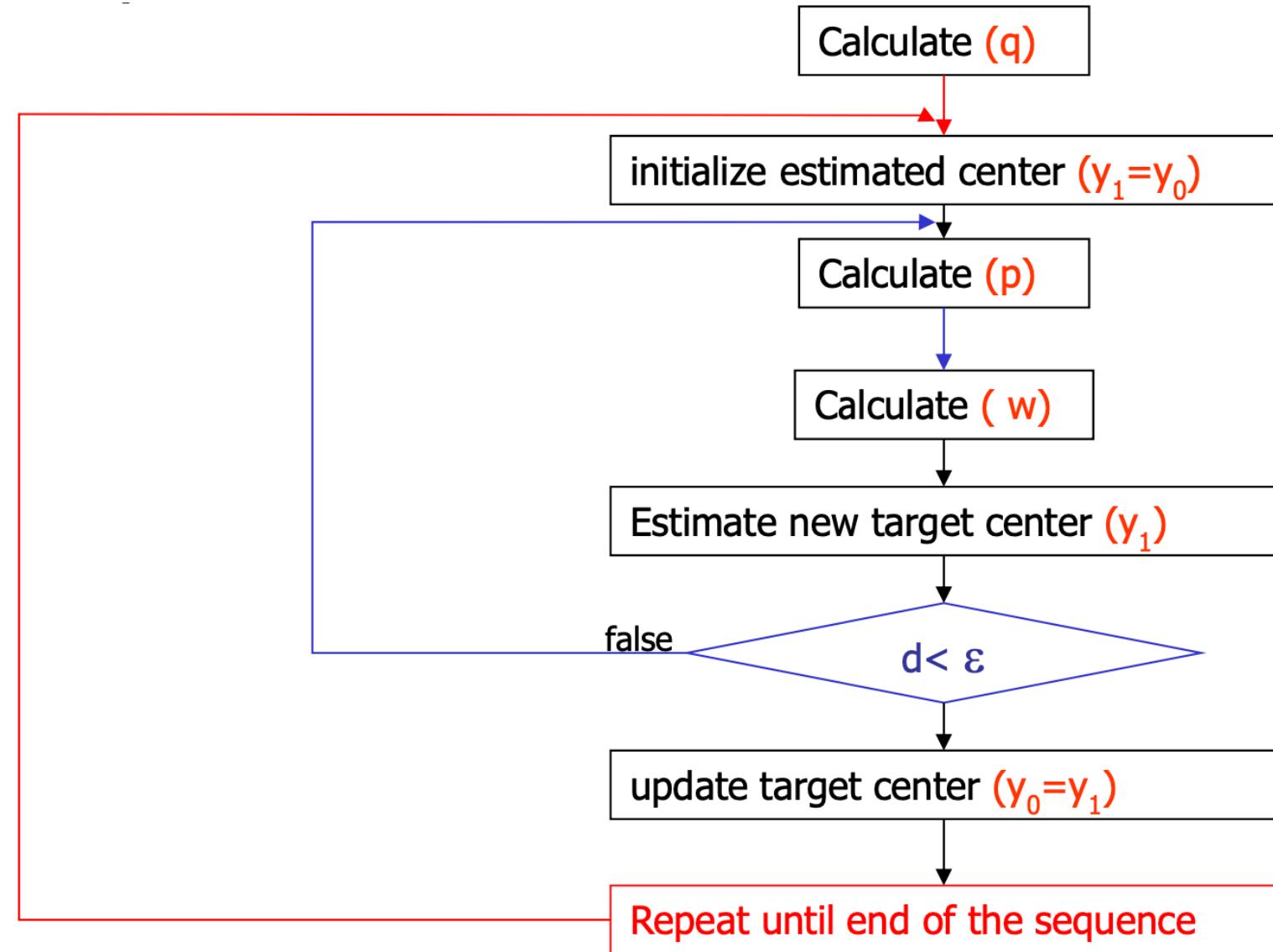
$$M_h(\mathbf{y}_0) = \frac{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0) \mathbf{x}_i}{\sum_{i=1}^{n_x} w_i(\mathbf{y}_0)} - \mathbf{y}_0$$

Thus, new target center is

$$\hat{\mathbf{y}} = \mathbf{y}_0 + M_h(\mathbf{y}_0)$$



Algorithm





Mean-Shift Pros/Cons

- Low computational cost (easily real-time)
- Surprisingly robust
 - Invariant to pose and viewpoint
 - Often no need to update reference color model
- Invariance comes at a price
 - Position estimate prone to fluctuation
 - Scale and orientation not well captured
 - Sensitive to color clutter (e.g., teammates in team sports)
- Local search by gradient descent
- Problems:
 - abrupt moves
 - occlusions



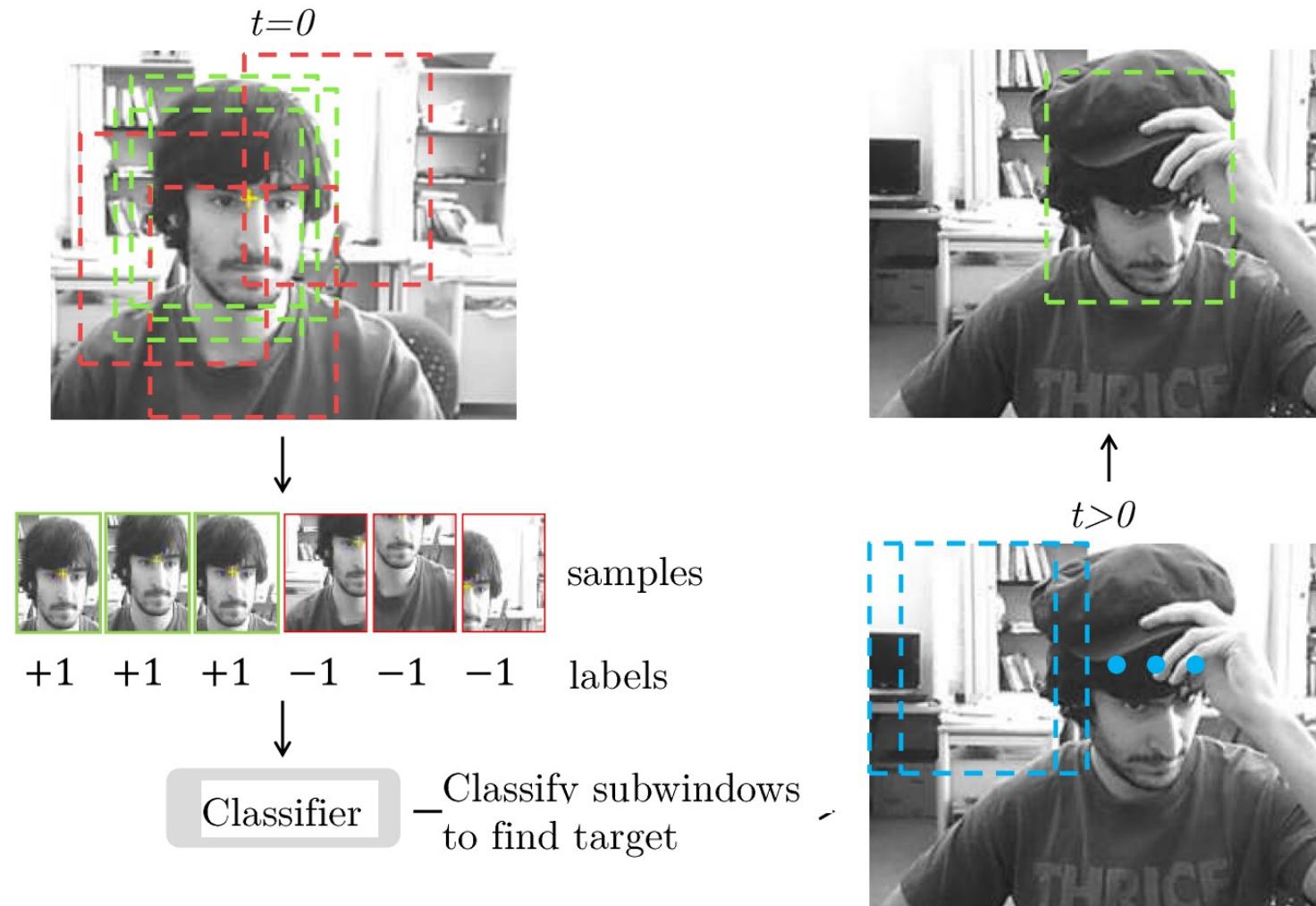
Mean-Shift Example



https://youtu.be/RG5uV_h50b0



Discriminative Tracking: Tracking by Detection





Connection to Correlation

- Let's have a linear classifier with weights \mathbf{w}

$$y = \mathbf{w}^T \mathbf{x}$$

$i = 3$
 $i = 2$
 $i = 1$



- During tracking we want to evaluate the classifier at subwindows \mathbf{x}_i :

$$y_i = \mathbf{w}^T \mathbf{x}_i$$

- Then we can concatenate y_i into a vector \mathbf{y} (i.e. response map)

- This is equivalent to **cross-correlation** formulation which can be computed **efficiently** in Fourier domain

$$\mathbf{y} = \mathbf{x} \odot \mathbf{w}$$

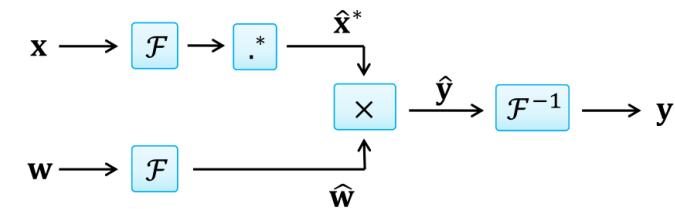
- Note: Convolution is related; it is the same as cross-correlation, but with the flipped image of \mathbf{w} ($\mathbf{P} \rightarrow \mathbf{d}$).

The Convolution Theorem

“Cross-correlation is **equivalent** to an **element-wise product** in Fourier domain”

$$\mathbf{y} = \mathbf{x} \odot \mathbf{w} \quad \Leftrightarrow \quad \hat{\mathbf{y}} = \hat{\mathbf{x}}^* \times \hat{\mathbf{w}}$$

- In practice:





How to Decide Classifier Weights, w ?

Training image: $\mathbf{x} =$

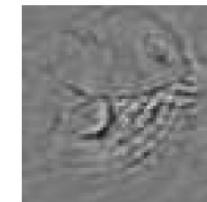


Classifier
(\mathbf{w})

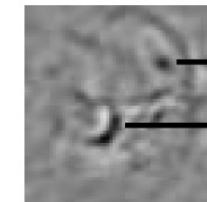
Naïve filter
($\mathbf{w} = \mathbf{x}$)



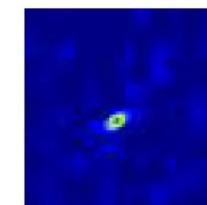
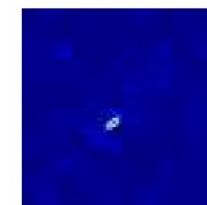
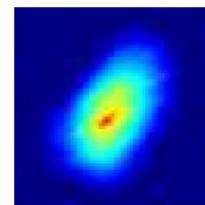
Sharp peak
(UMACE)



Gaussian peak
(GMACE)



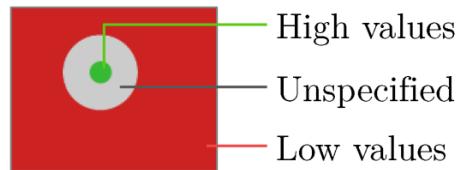
Output
($\mathbf{w} * \mathbf{x}$)



Objective



$$\odot \mathbf{w} =$$



- A good compromise.
- Tiny details are ignored.
- focuses on **larger, more robust structures**.



Minimum Output Sum of Squared Error (MOSSE) based Correlation Filter



A Nice Summary for Tracking API in OpenCV

- <https://www.learnopencv.com/object-tracking-using-opencv-cpp-python/>



Deep Learning based Tracking



Supervised Tracking

- Paper : Detect to Track and Track to Detect
- A ConvNet architecture for simultaneous detection and tracking, using a multi-task objective for frame-based object detection and across-frame track regression
- Introduces correlation features that represent object co-occurrences across time to aid the ConvNet during tracking.
- Links the frame level detections based on our across-frame tracklets to produce high accuracy detections at the video level



Backbone network : R-FCN

- R-FCNs are convolutional networks used for object detection.
- It is faster than traditional object detection network such as R-CNN as its region-based feature maps are independent of region of interests.
- The network is lighter because it is fully convolutional and has no FC layers.
- The output of the network is a set of bounding coordinates and class labels for detected objects.

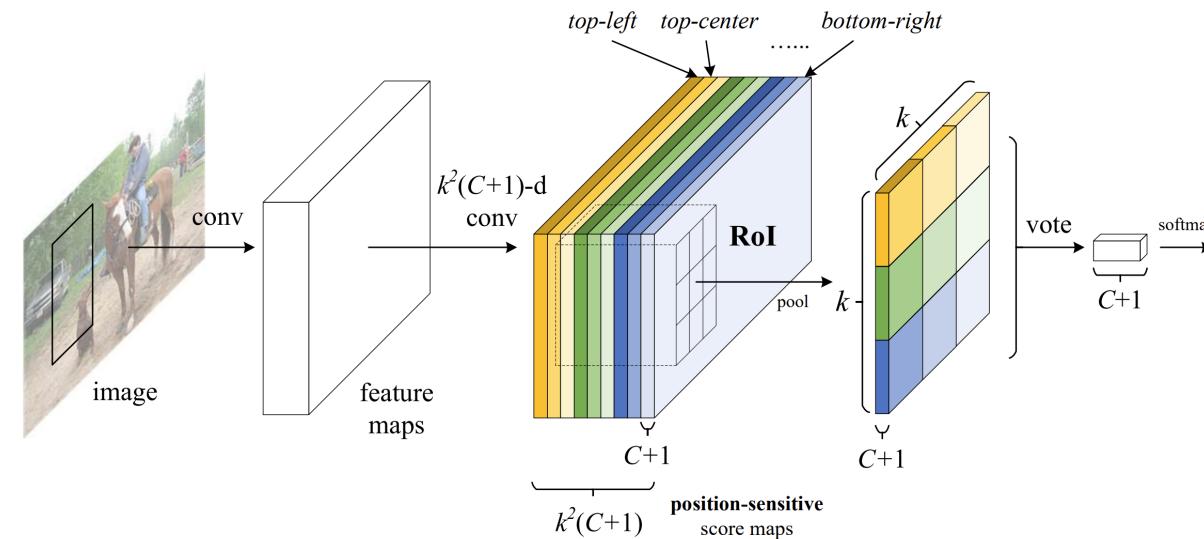
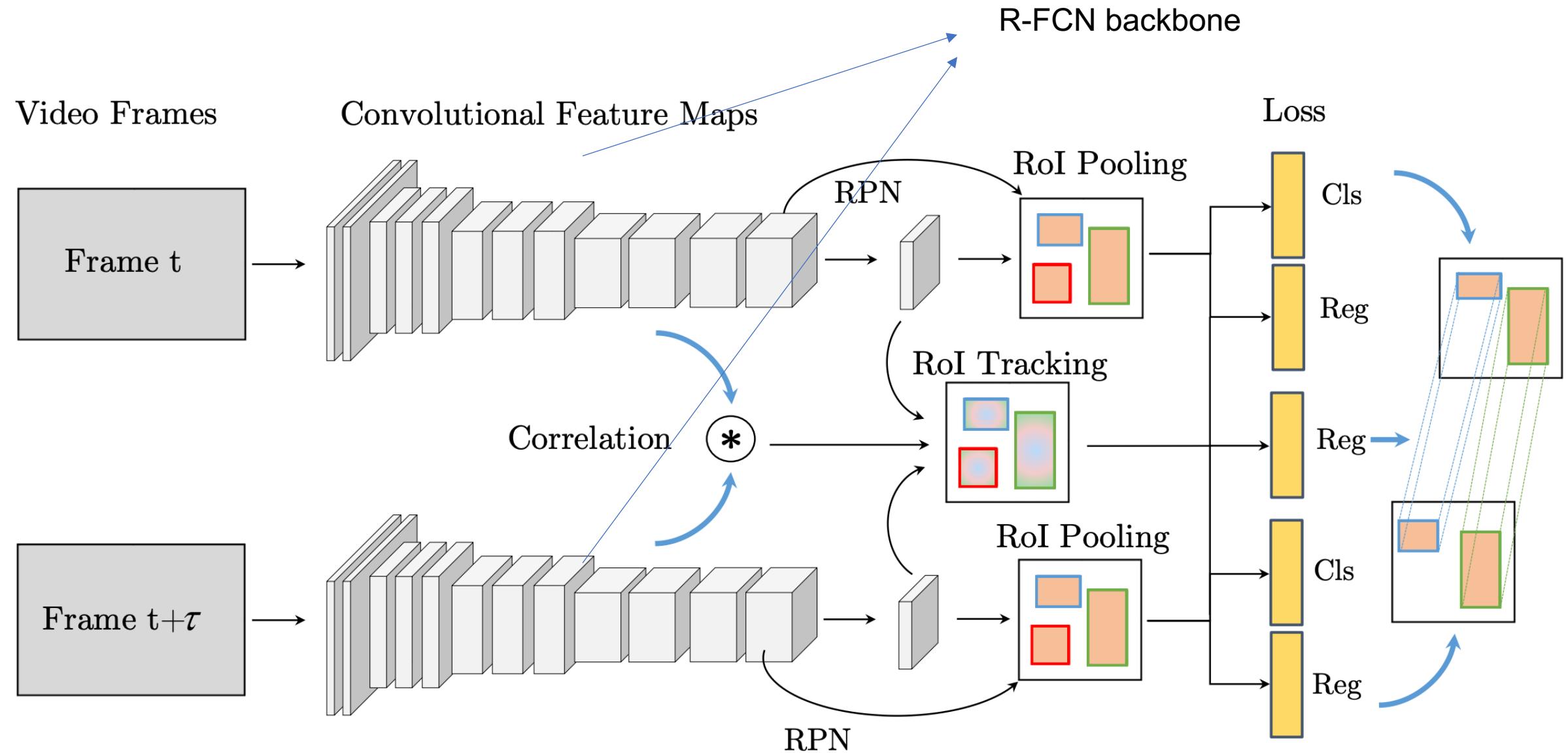


Figure 1: Key idea of **R-FCN** for object detection. In this illustration, there are $k \times k = 3 \times 3$ position-sensitive score maps generated by a fully convolutional network. For each of the $k \times k$ bins in an ROI, pooling is only performed on one of the k^2 maps (marked by different colors).

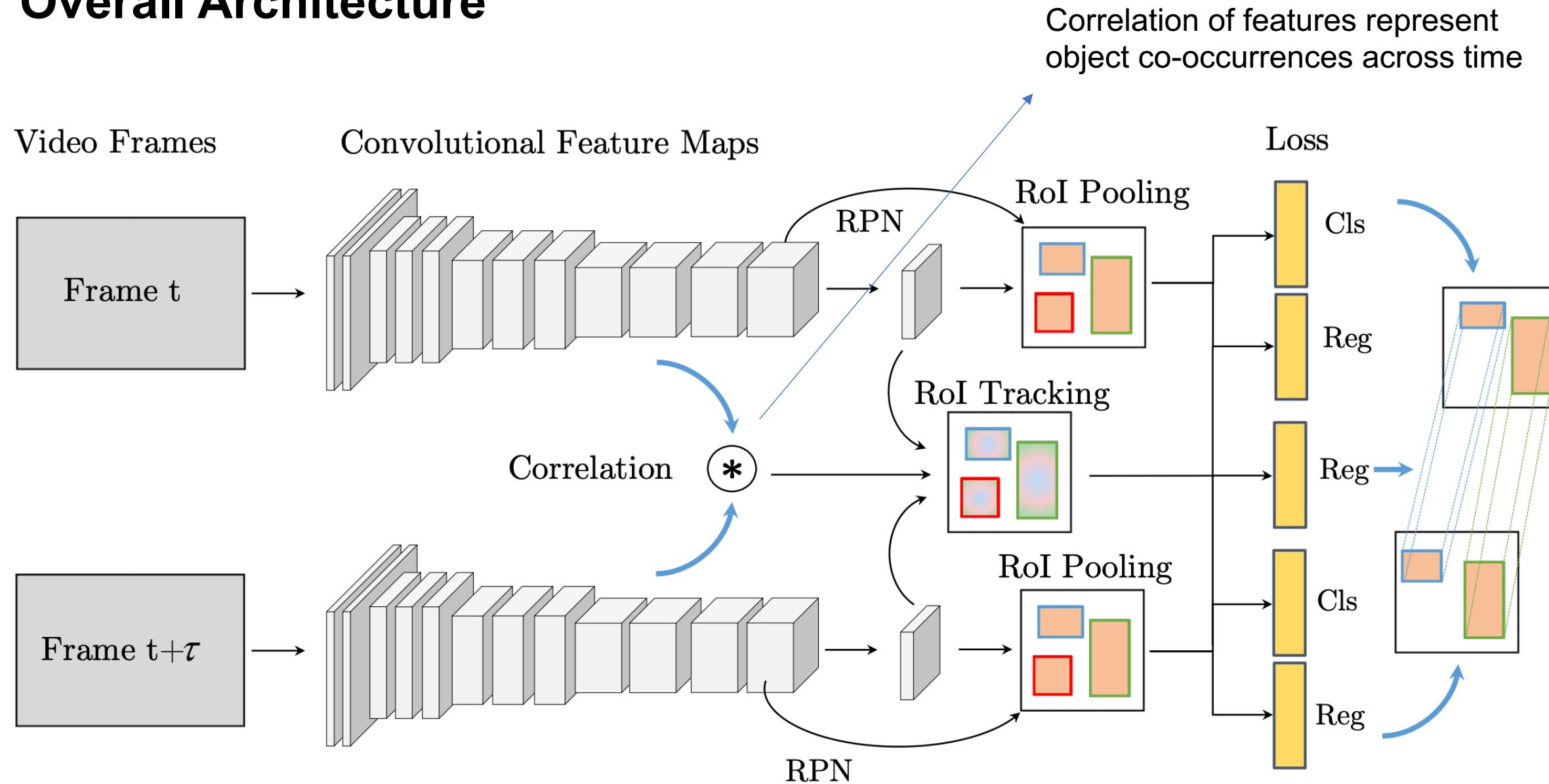


Overall Architecture



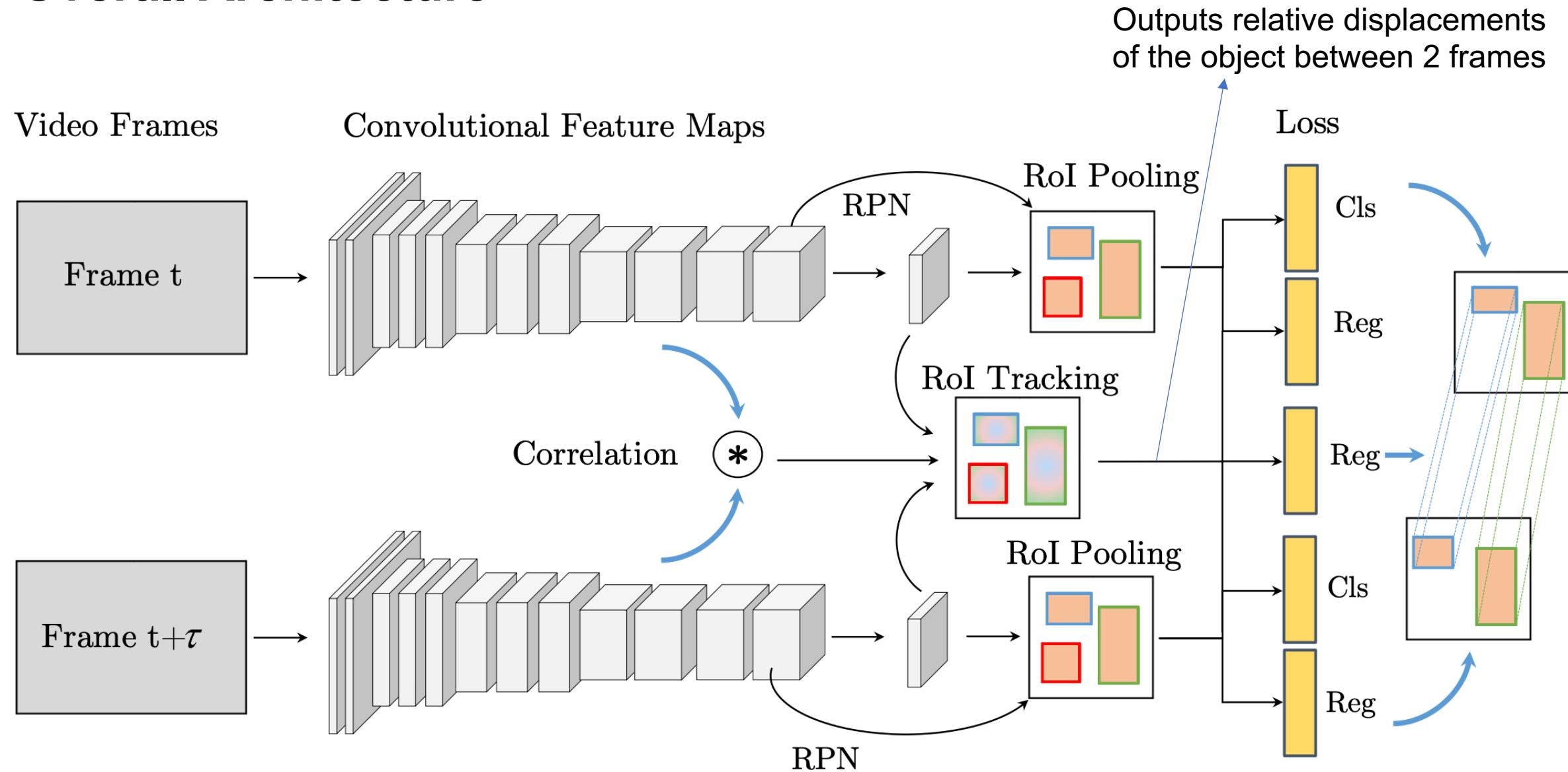


Overall Architecture





Overall Architecture





Objective function

$$L(\{p_i\}, \{b_i\}, \{\Delta_i\}) = \frac{1}{N} \sum_{i=1}^N L_{cls}(p_i, c^*) \longrightarrow$$

Classification loss
(foreground/background)

$$+ \lambda \frac{1}{N_{fg}} \sum_{i=1}^N [c_i^* > 0] L_{reg}(b_i, b_i^*) \longrightarrow$$

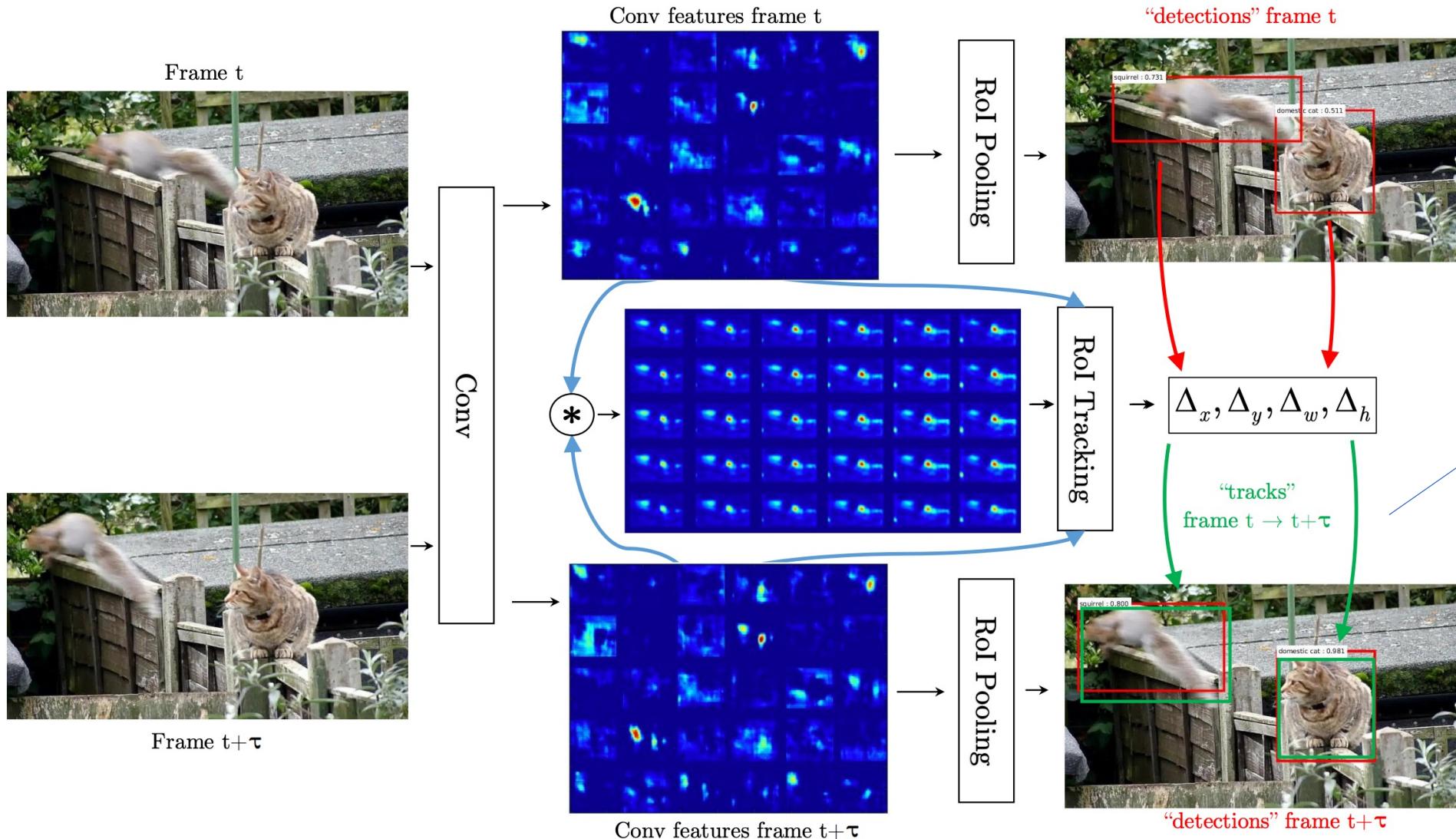
Bounding box regression
cost (for both frames)

$$+ \lambda \frac{1}{N_{tra}} \sum_{i=1}^{N_{tra}} L_{tra}(\Delta_i^{t+\tau}, \Delta_i^{*, t+\tau}). \longrightarrow$$

Regression offsets



Inference





Self-supervised Tracking

- To learn a feature extractor without any supervision
- Self-supervised learning via cycle consistency

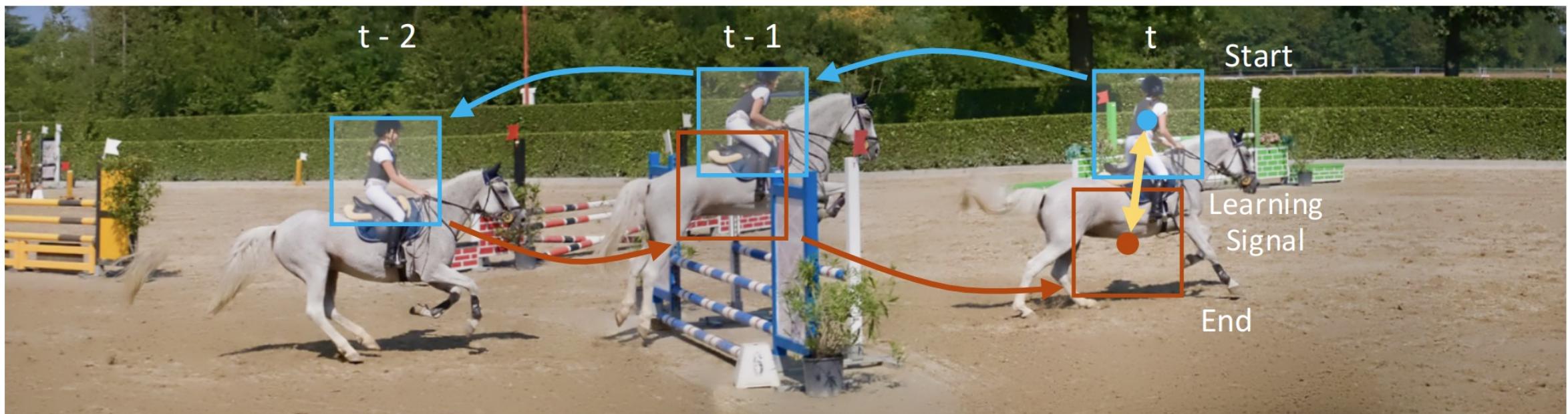


Figure 2: **A Cycle in Time.** Given a video, tracking along the sequence formed by a cycle in time can be self-supervised: the target is simply the beginning of the cycle. The yellow arrow between the start and end represents the differentiable learning signal.

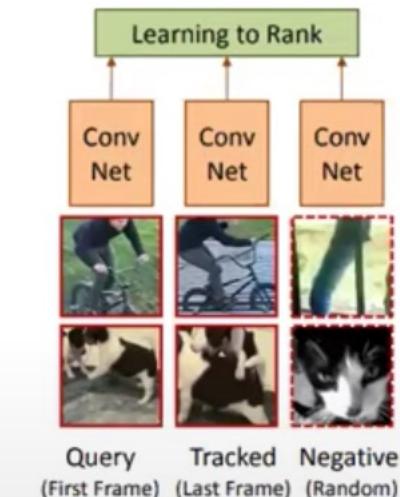
Previous Self-Supervised Correspondence Learning

- Tracking-based Similarity Learning

- (a) Prepare Query (Positive & Negative) ←
(b) Obtain Positive via Hand-craft tracker
(c) Use Siamese-Triplet Network



(a) Unsupervised Tracking in Videos



(b) Siamese-triplet Network



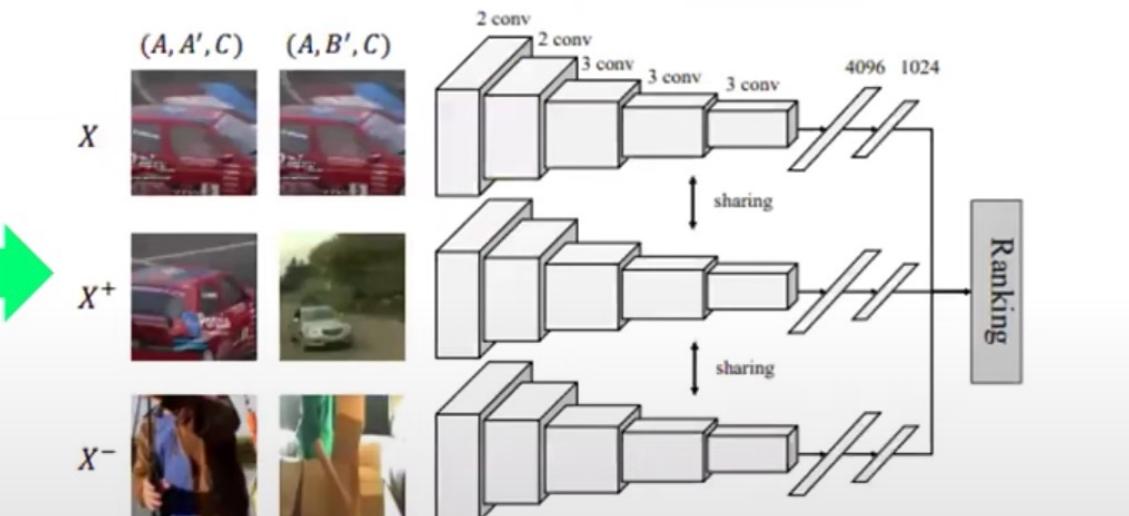
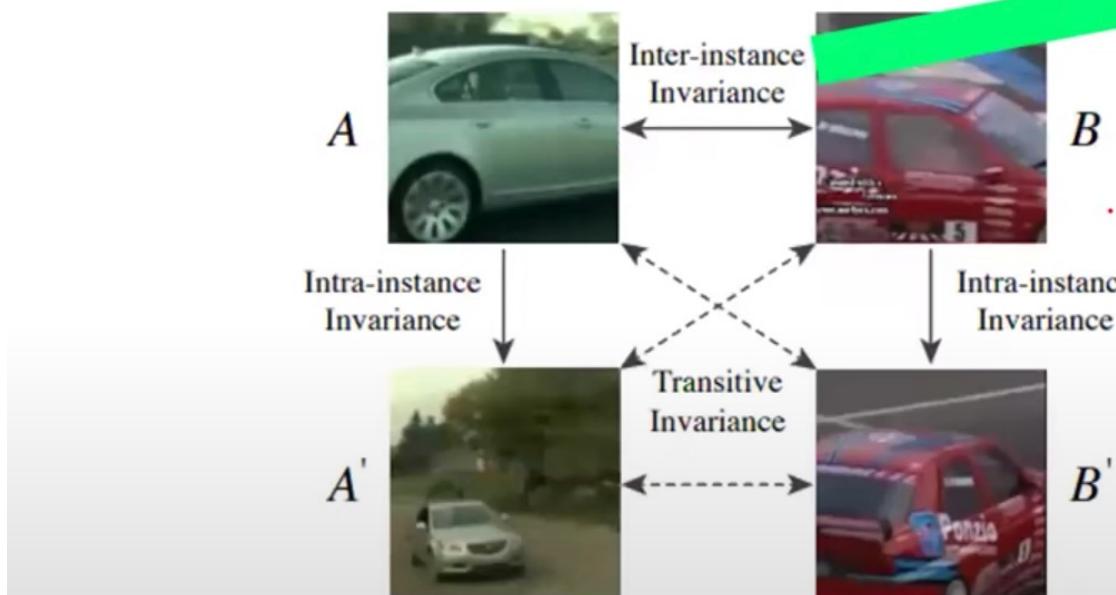
(c) Ranking Objective

Previous Self-Supervised Correspondence Learning

- Tracking-based Similarity Learning
- Transitive Invariance Learning



$X = (\text{cat}, \text{cat}), Y = 7$ $X = (\text{cat}, \text{dog}), Y = 7$ $X = (\text{dog}, \text{dog}), Y = 2$

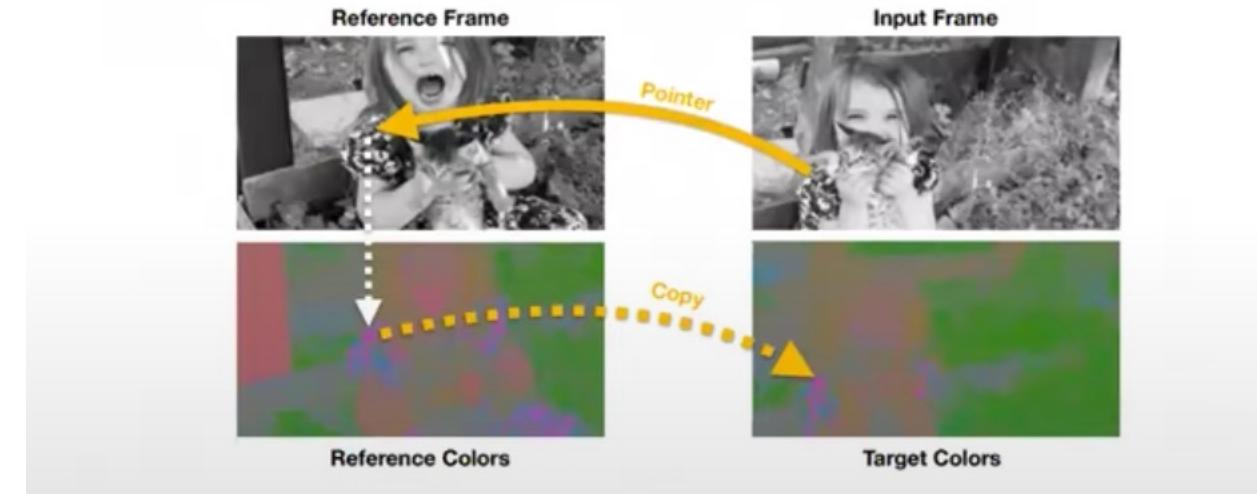
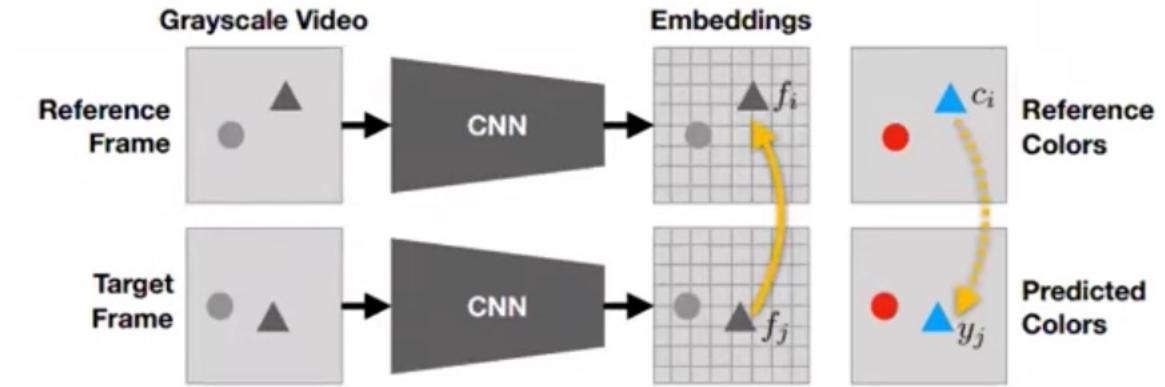


X. Wang, K. He, and A. Gupta. Transitive invariance for selfsupervised visual representation learning. In ICCV, 2017



Previous Self-Supervised Correspondence Learning

- Tracking-based Similarity Learning
- Transitive Invariance Learning
- Colorization
 - (a) Find similar features
 - (b) Copy colors from reference frame
 - (c) Learn CNN using colorization loss



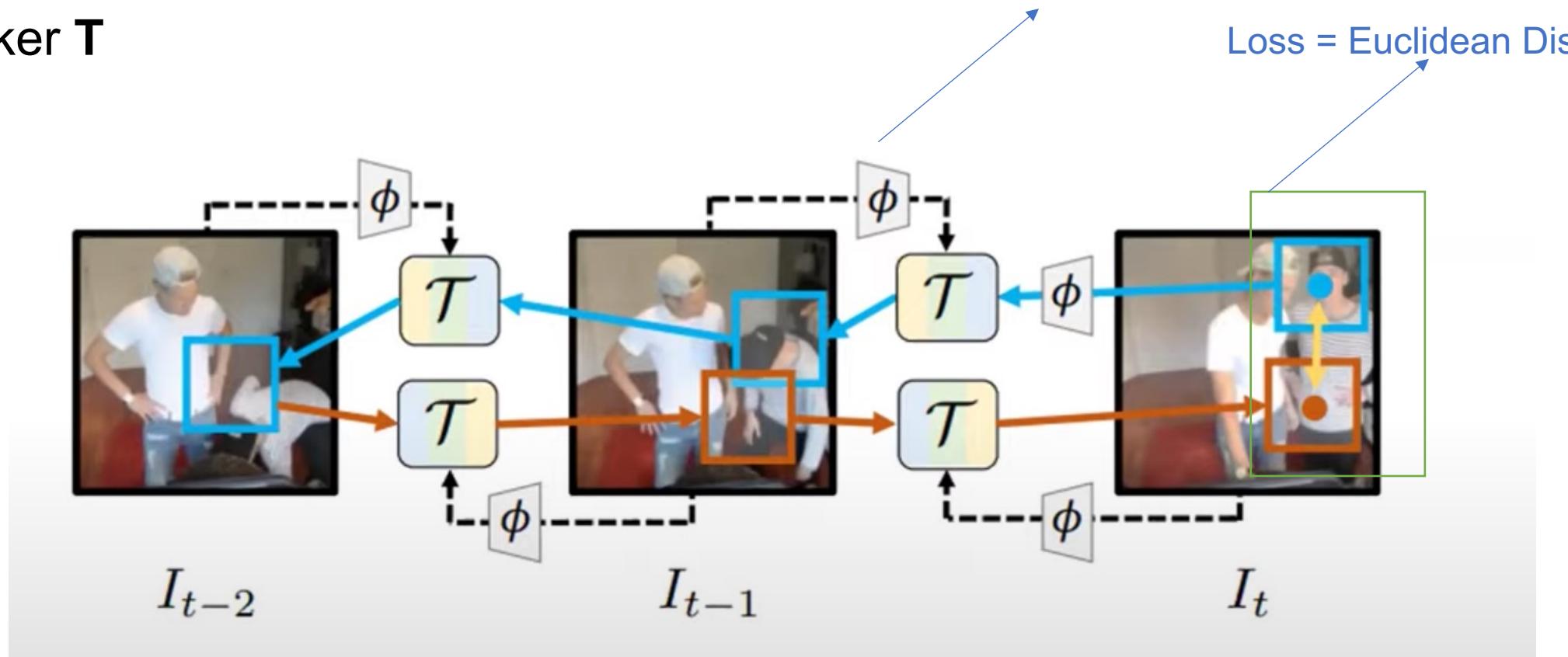


Cycle-in-Time

- Feature extractor Φ
- Tracker \mathcal{T}

Similar patches should be closer in embedding space

Loss = Euclidean Distance





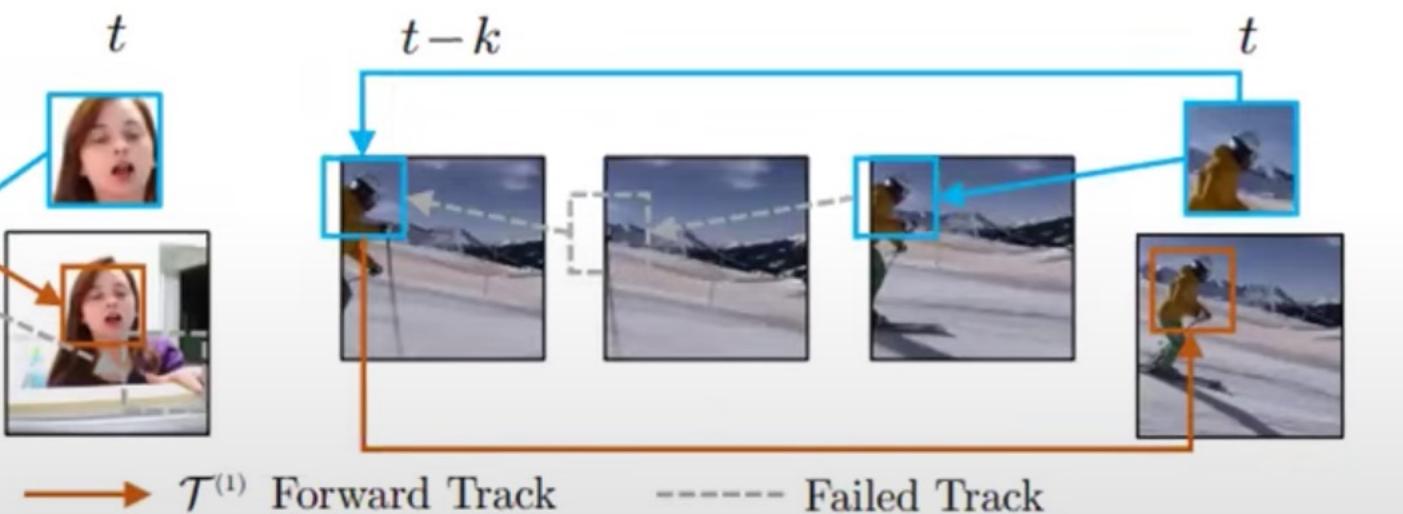
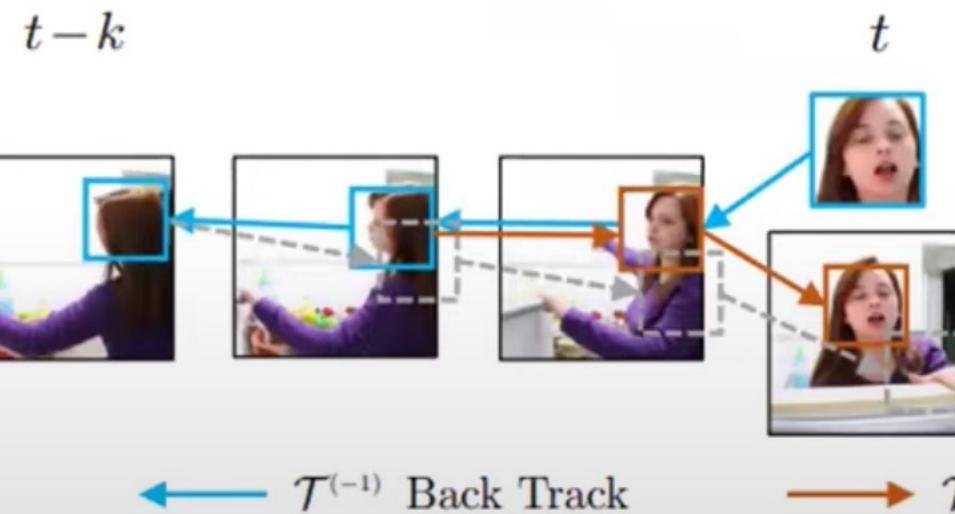
Multiple Cycles / Skip Cycle

- Multiple cycles for variation
- Skip cycle for occlusion

Mean region
similarity

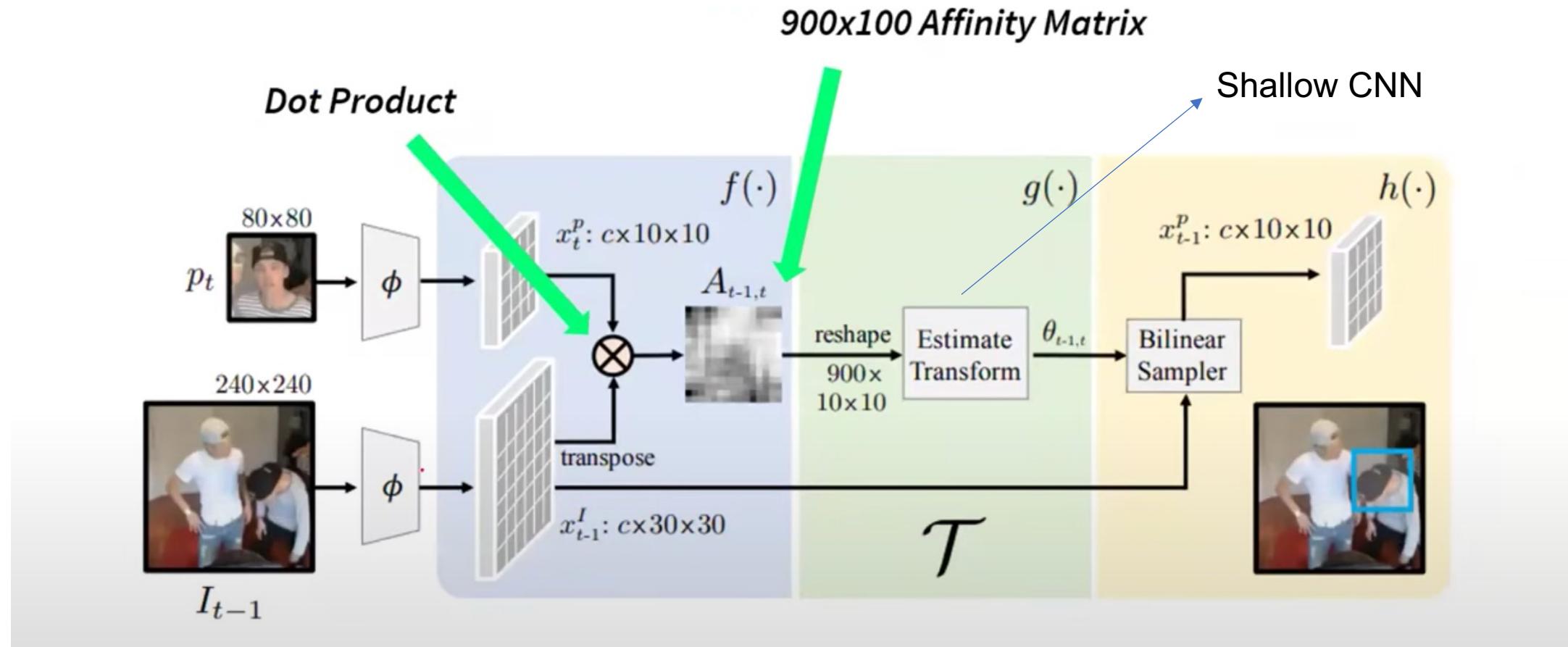
Contour based accuracy

Experiment	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
Ours	41.9	39.4
Ours without Skip Cycles	39.5	37.9



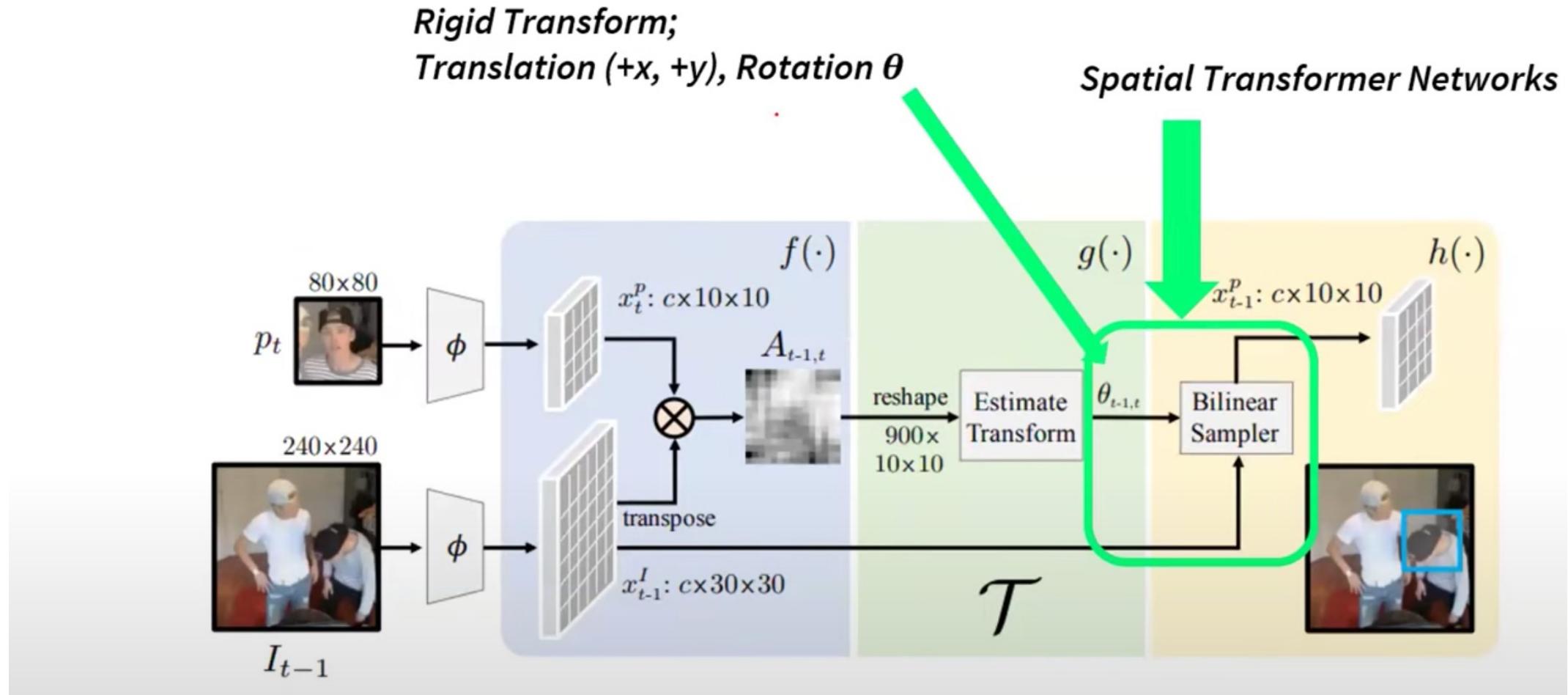


Differentiable Tracker



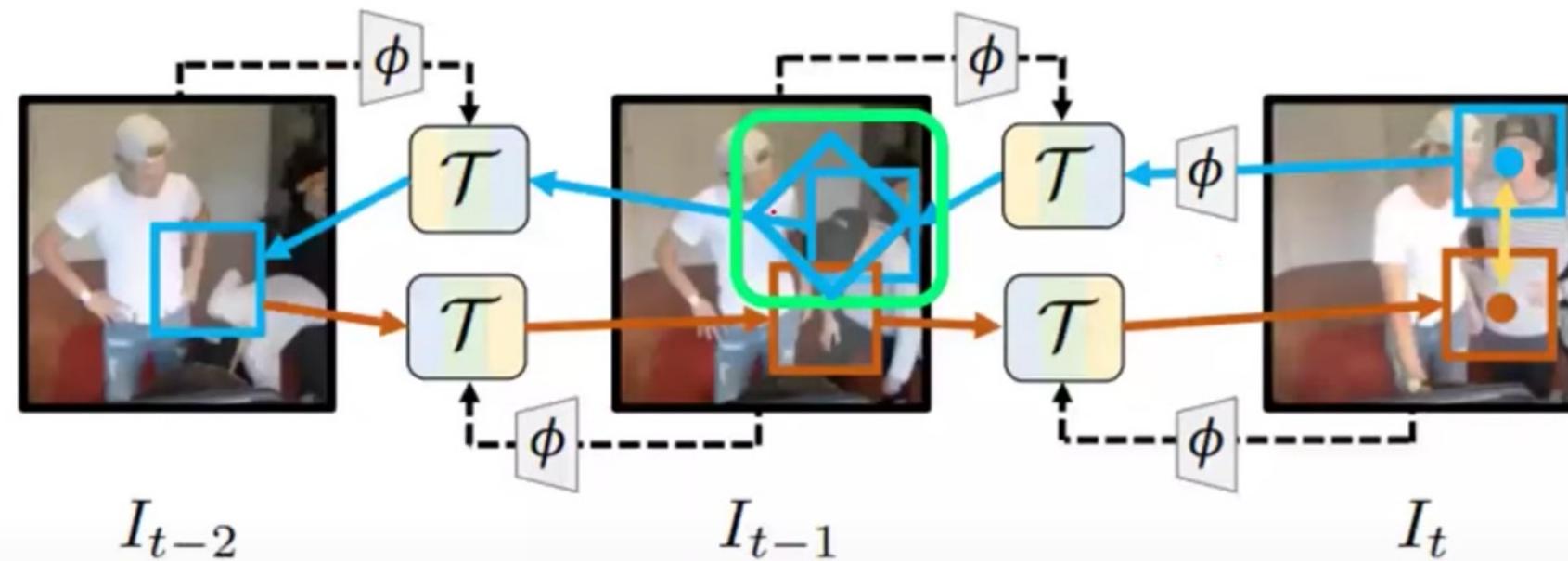


Differentiable Tracker





Training





Results : Video Object Segmentation



model	Supervised	\mathcal{J} (Mean)	\mathcal{F} (Mean)
Identity		22.1	23.6
Random Weights (ResNet-50)		12.4	12.5
Optical Flow (FlowNet2) [22]		26.7	25.2
SIFT Flow [39]		33.0	35.0
Transitive Inv. [74]		32.0	26.8
DeepCluster [8]		37.5	33.2
Video Colorization [69]		34.6	32.7
Ours (ResNet-18)		40.1	38.3
Ours (ResNet-50)		41.9	39.4
ImageNet (ResNet-50) [18]	✓	50.3	49.0
Fully Supervised [81, 7]	✓	55.1	62.1



Results: Pose Propagation

(b)
Pose

model	Supervised	PCK@.1	PCK@.2
Identity		43.1	64.5
Optical Flow (FlowNet2) [22]		45.2	62.9
SIFT Flow [39]		49.0	68.6
Transitive Inv. [74]		43.9	67.0
DeepCluster [8]		43.2	66.9
Video Colorization [69]		45.2	69.6
Ours (ResNet-18)		57.3	78.1
Ours (ResNet-50)		57.7	78.5
ImageNet (ResNet-50) [18]	✓	58.4	78.4
Fully Supervised [59]	✓	68.7	92.1

Probability of
Correct Keypoint

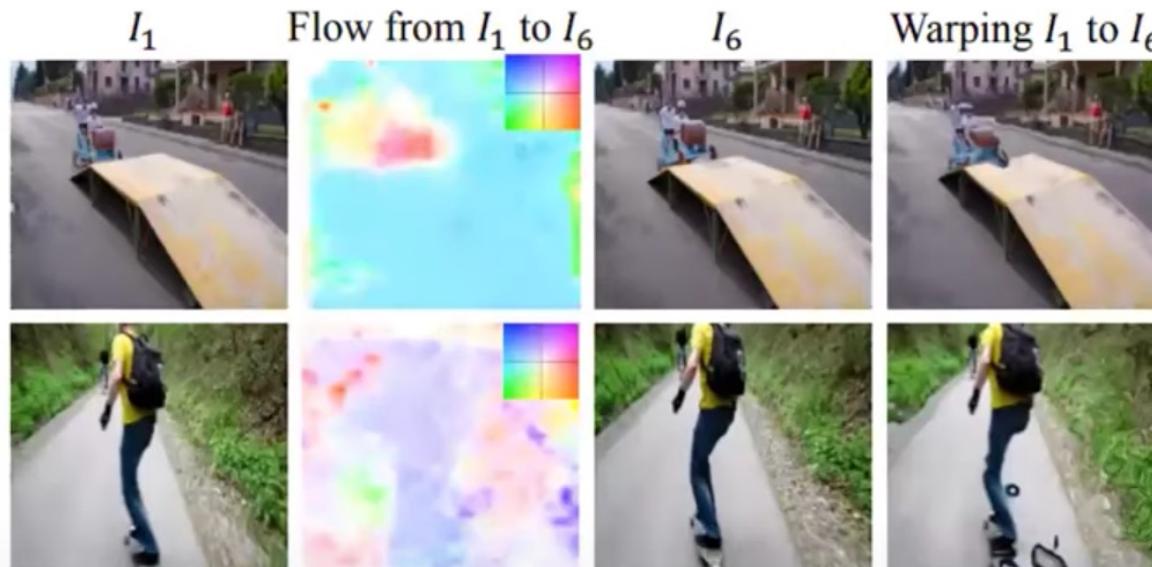
Results: Human Part Labels Propagation



model	Supervised	mIoU	AP_{vol}^r
Identity		13.6	4.0
Optical Flow (FlowNet2) [22]		16.1	8.3
SIFT Flow [39]		21.3	10.5
Transitive Inv. [74]		19.4	5.0
DeepCluster [8]		21.8	8.1
Ours (ResNet-50)		28.9	15.6
ImageNet (ResNet-50) [18]	✓	34.7	16.1
Fully Supervised [85]	✓	37.9	24.1



Results: Long-range Optical Flow



model	5-F	10-F
Identity	82.0	97.7
Optical Flow (FlowNet2) [22]	62.4	90.3
ImageNet (ResNet-50) [18]	64.0	79.2
Ours (ResNet-50)	60.4	76.4



Next week

- Segmentation
 - + Graph-based segmentation
 - + Semantic segmentation
 - + Instance segmentation
 - + Panoptic segmentation
- 3D Deep Learning (Recognition and Segmentation)
 - + VoxNet
 - + PointNet
 - + MVCNN
 - + FoldingNet
- + : know the concept



References for next week

- Sz: Ch 6.4
- Maturana, Daniel, and Sebastian Scherer. "Voxnet: A 3d convolutional neural network for real-time object recognition." *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2015.
- Su, Hang, et al. "Multi-view convolutional neural networks for 3d shape recognition." *Proceedings of the IEEE international conference on computer vision*. 2015.
- Qi, C. R., Su, H., Mo, K., & Guibas, L. J. (2017). Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 652-660).
- Yang, Y., Feng, C., Shen, Y., & Tian, D. (2018). Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 206-215).