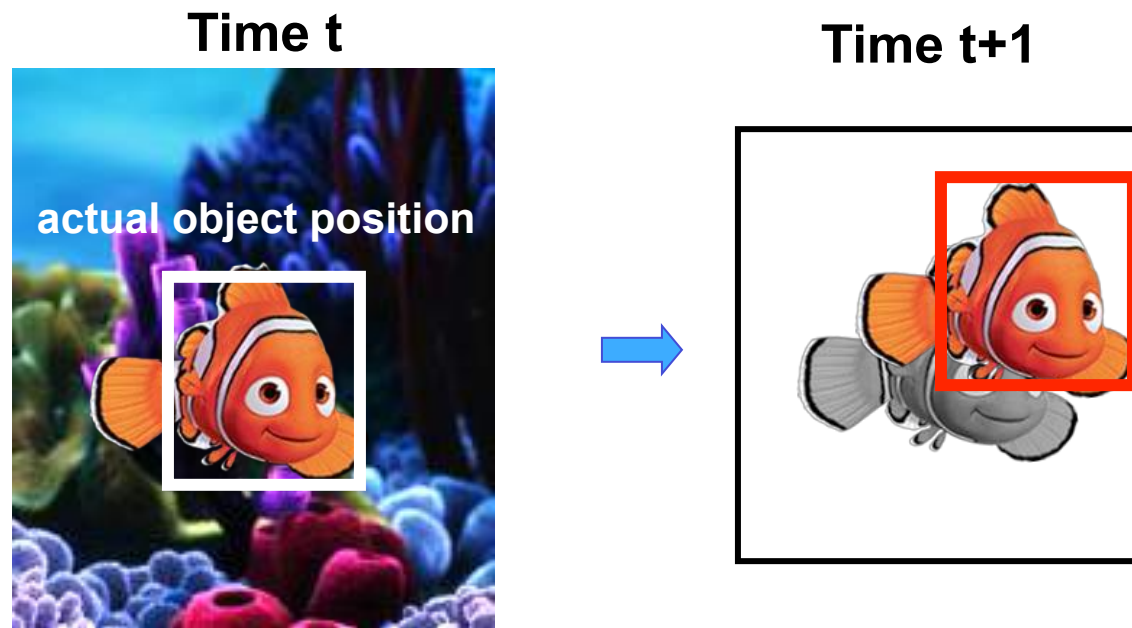# Tracking



Dictionary:

- [noun] "The pursuit (of a person or animal) by following tracks or marks they left behind"
- [verb] "Observe or plot the moving path of something (e.g., to track a missile)"

What does it mean in Computer Vision?

*Many thanks to:* H. Grabner, L. van Gool, and V. Ferrari for some of the slides & videos.
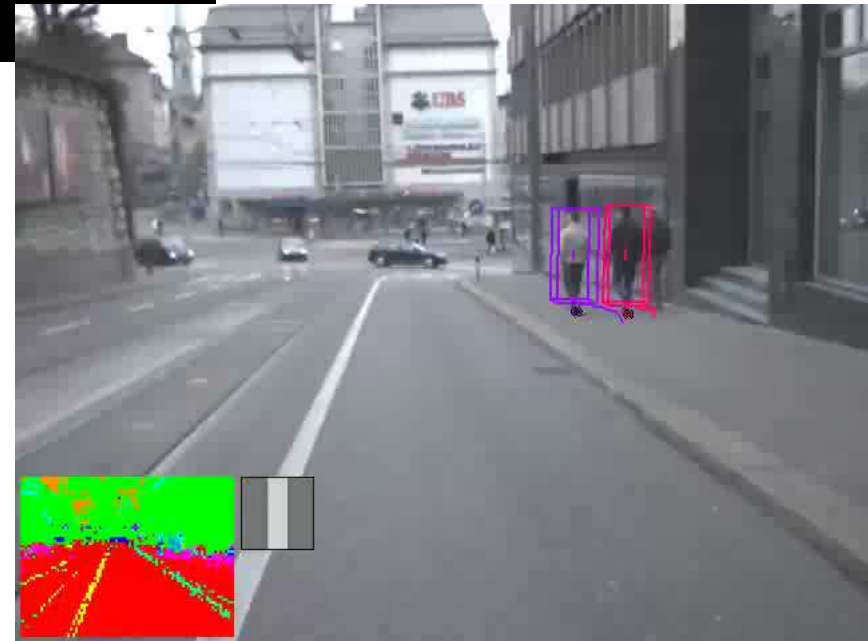
# What is Tracking

# Why do we need it

What is tracking for you? Why do you think it is relevant and may be important?

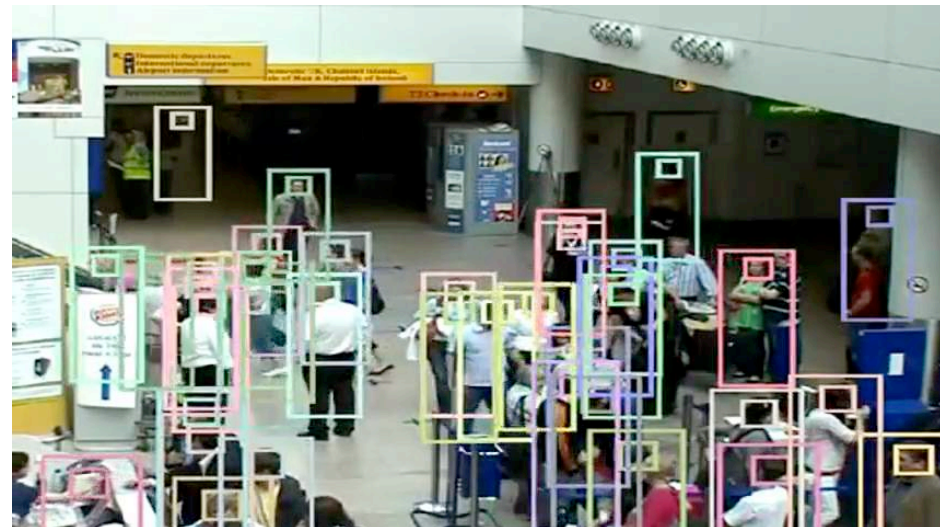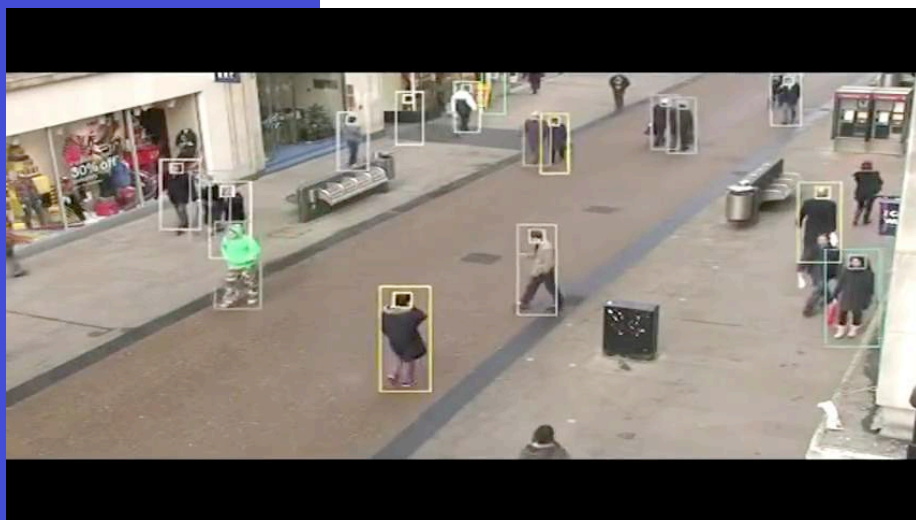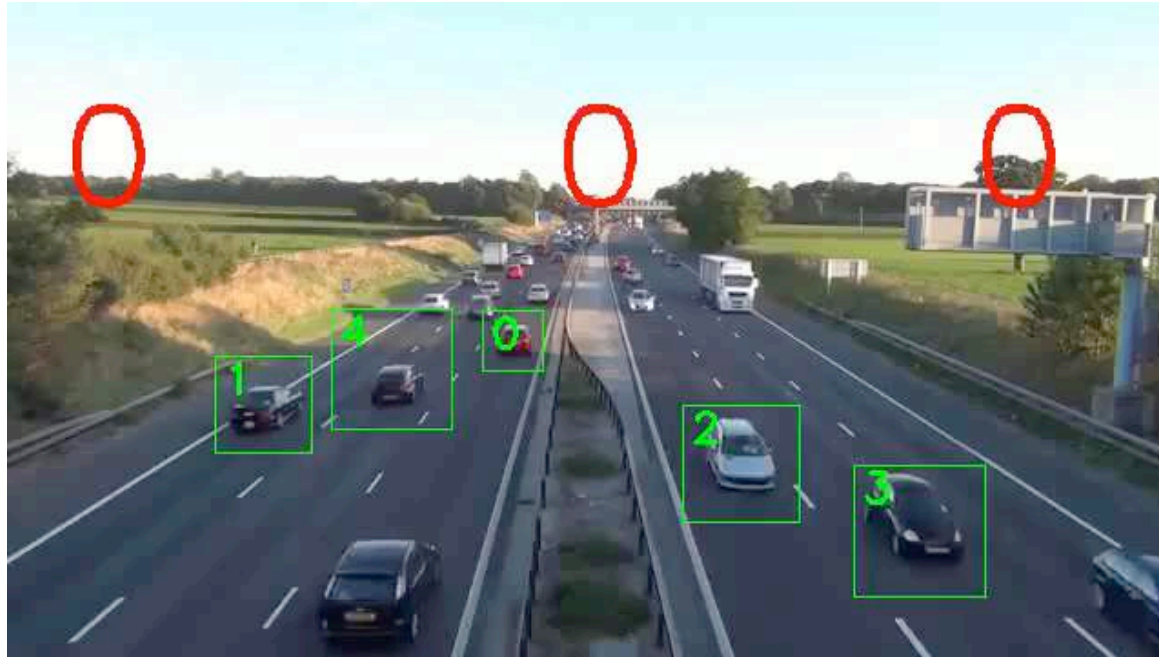Where could it be useful, in real-life application/engineering scenarios?

**Task: "List <u>applications</u> you can think of on a piece of paper"**

Discuss in groups of 3-4

# Autonomous Driving

NVIDIA GTC Europe

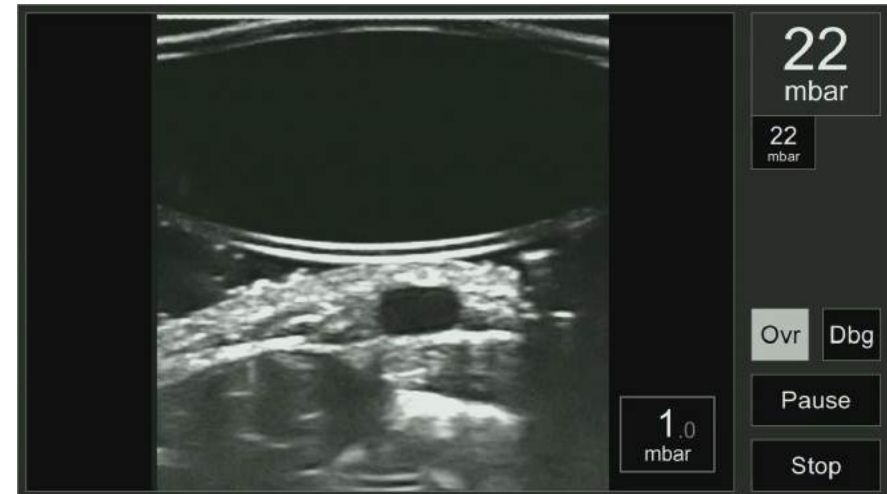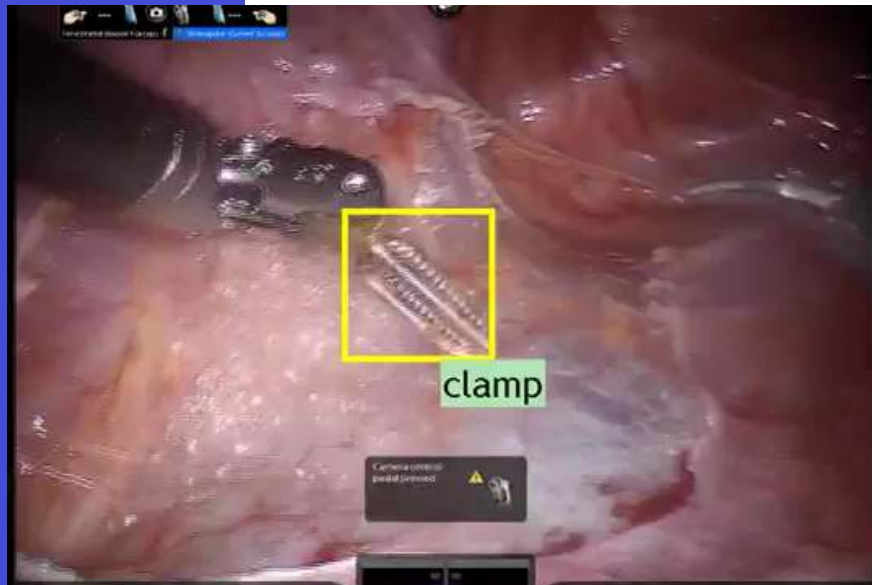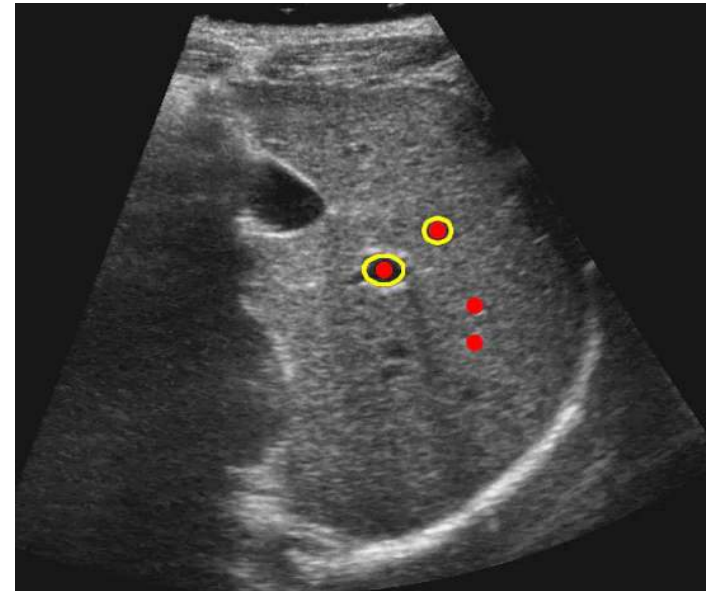Orcun Goksel,  ETH Zurich

# Surveillance, Safety, Security
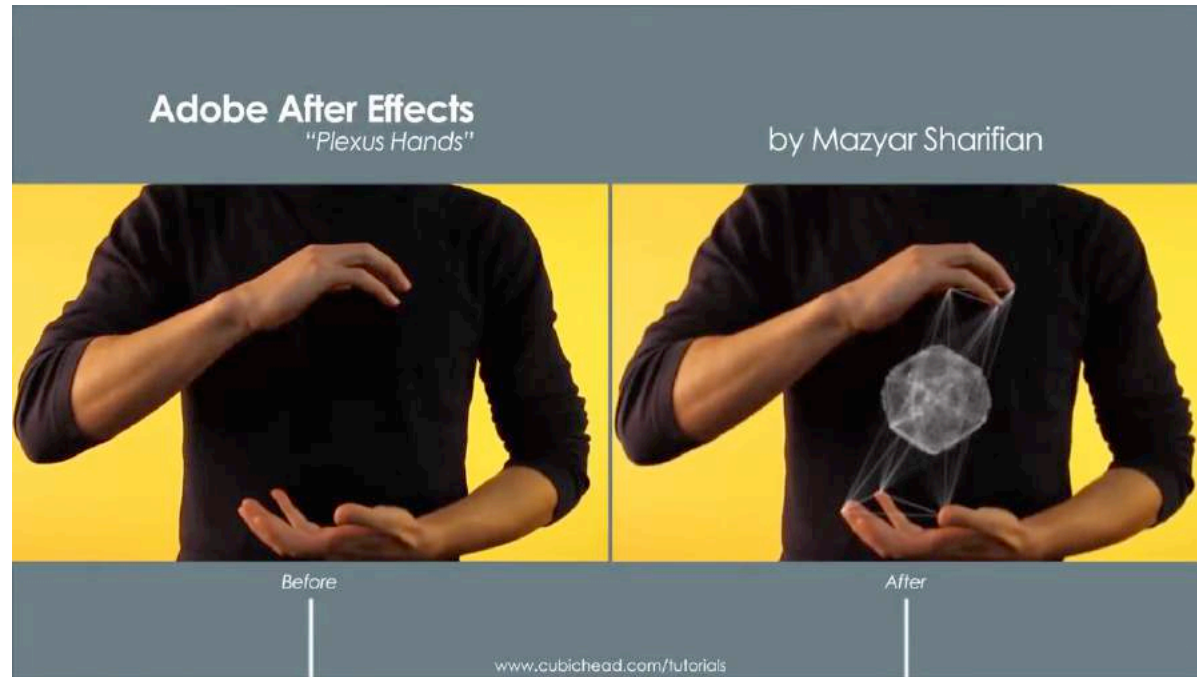
# Applications: VR/AR glasses



Microsoft HoloLens

# Medical Guidance

Orcun Goksel,  ETH Zurich

# Sports

# Video Editing

# SfM: Structure from Motion

- Tracked Points gives correspondences

# Defense

"Top Gun"



Orcun Goksel, ETH Zurich

# Of course, "very importantly" The Cow Tracker

Orcun Goksel,  ETH Zurich

# Applications

- Structure-from-Motion

- Autonomous Driving

- Gesture/Action Recognition

- Augmented Reality

- Navigation

- Safety and Security

- Medical Targeting / Guidance
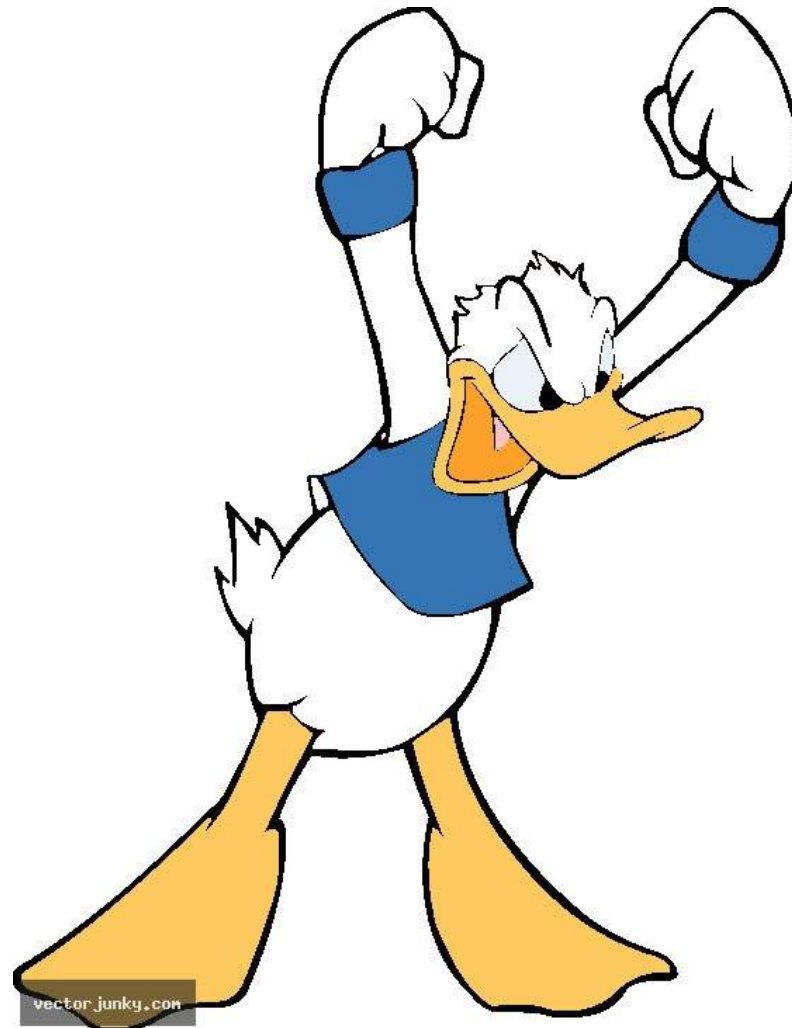
- Motion Compensation

- …

# You will be able to:

1. Determine applications of tracking and identify problems solvable by tracking

2. Analyze what methods could work in a practical scenario / situation

3. Assess potential limitations / pitfalls of particular approaches and scenarios

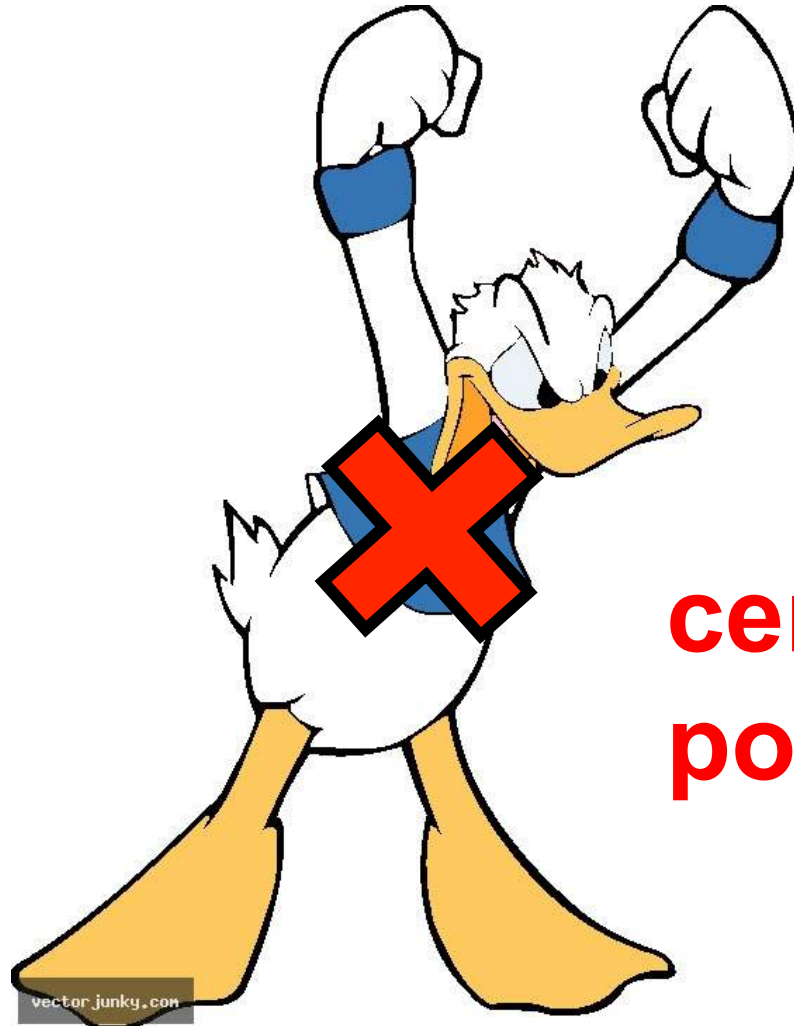4. **Propose an optimal tracking solution**

## How will we get there:

- (some) common tracking methods

- Few particular keywords & implementation

- What not: details of all individual implementations; cf. "how to google"
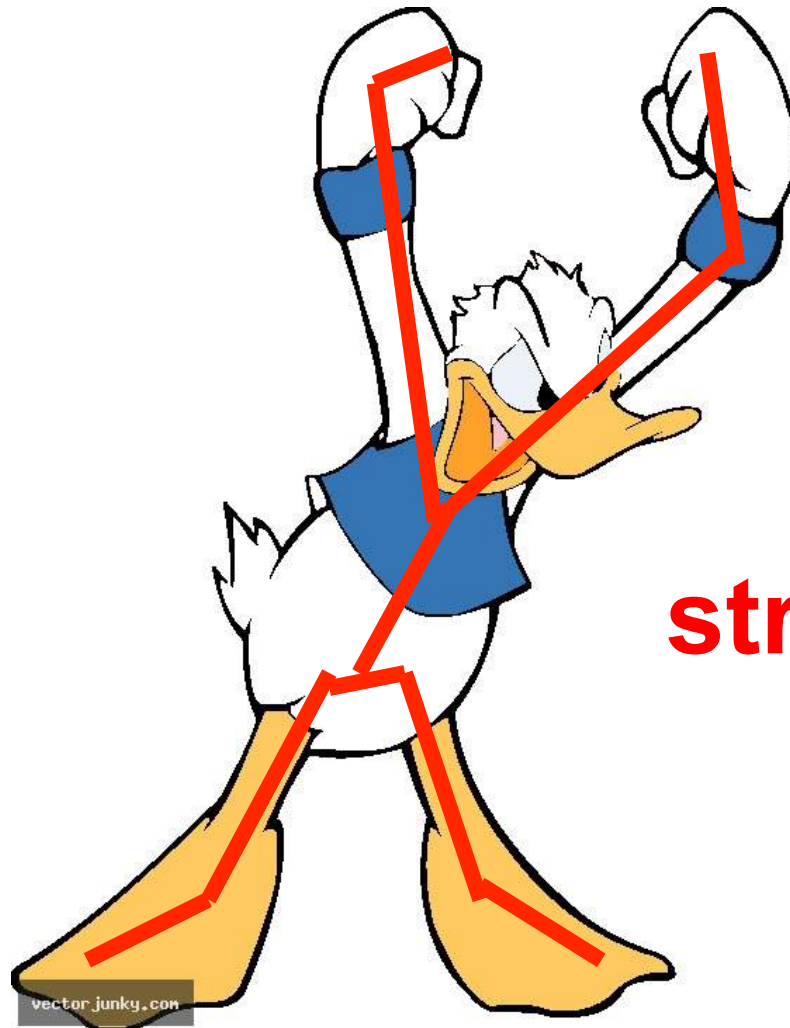
# What to track?

# What to track?

**center point**
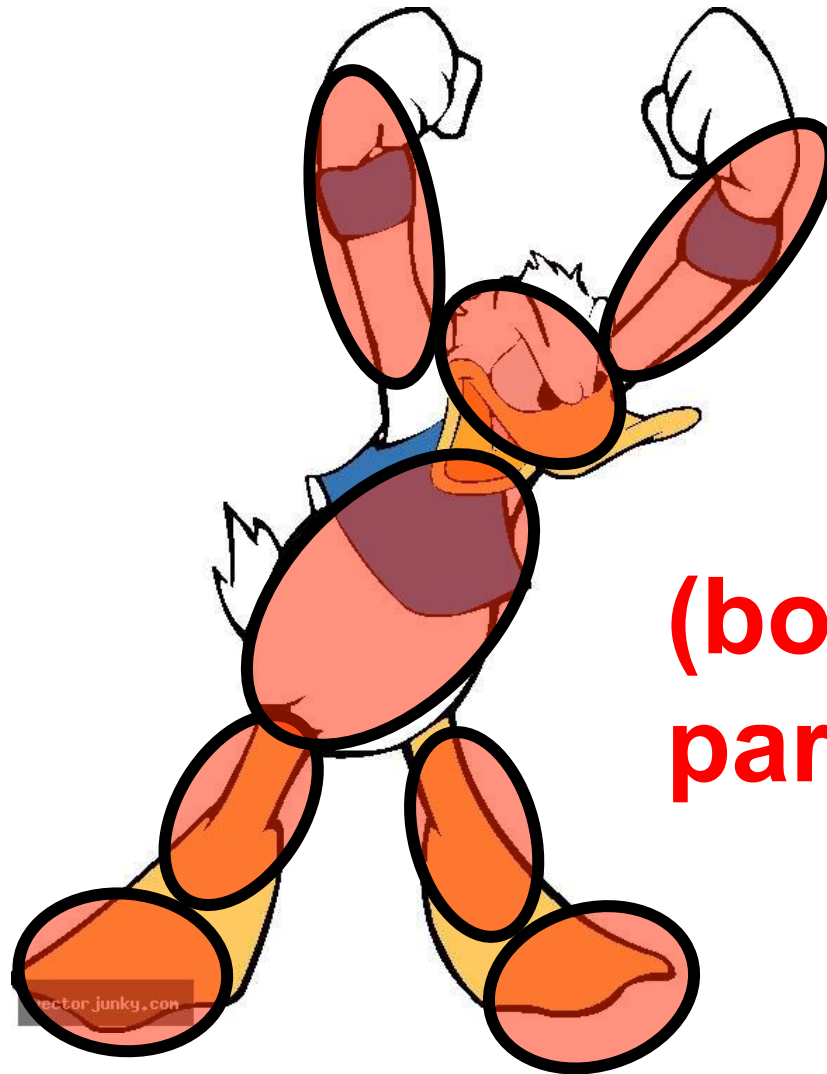
# What to track?

**multiple points**

# What to track?

**structure**

# What to track?



**(body) parts**

# What to track?

**region**

# What to track?



**outline**

# Approaches

**(i) Feature tracking**
**generic**

corners, blob/contours, regions, …

**(ii) Model-based tracking**
**application-specific**

face, human body, …

# Tracking Requirements

- Strongly depends on the **application**!

<span style="color:red">Robust, Accurate, Fast,…</span>

- Constrain the tracking task!

<span style="color:blue">Information about the object, dynamics, …</span>

Orcun Goksel,  ETH Zurich

# Tracking Cues

**Object**

**Saliency**

**Scene**

**Model/ Tracking History**

# Motion as a Cue

# Motion as a Cue



Orcun Goksel,  ETH Zurich

# Motion as a Cue



- Eye perceptive to temporal changes (gradients)
- "Event based camera"

Orcun Goksel,  ETH Zurich

# General Tracking Loop

**predict to t+1**

**measure at t+1**

**time t**



**update model**

**update location**

# Which strategy to use?
## Depends, No single solution
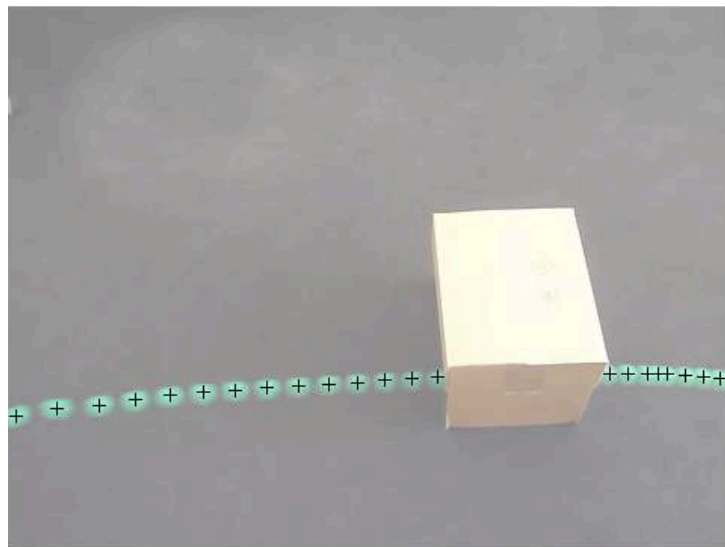
Some rule-of-thumb suggestions:

- If you can alter the "object" to be tracked,

➔ <u>modify/add tracking info</u>

e.g. optical IR markers, mark with patterns, etc

- If object is fixed/known, but modification not possible/ desired    ➔    <u>Utilize known info</u>

  e.g. use a template image and/or known object features

- If object unknown/variable object, but
  resides in a known (static) environment ➔ utilize this!

- If none above, simply follow from initial image/location

Tracking v.s. segmentation/localization:
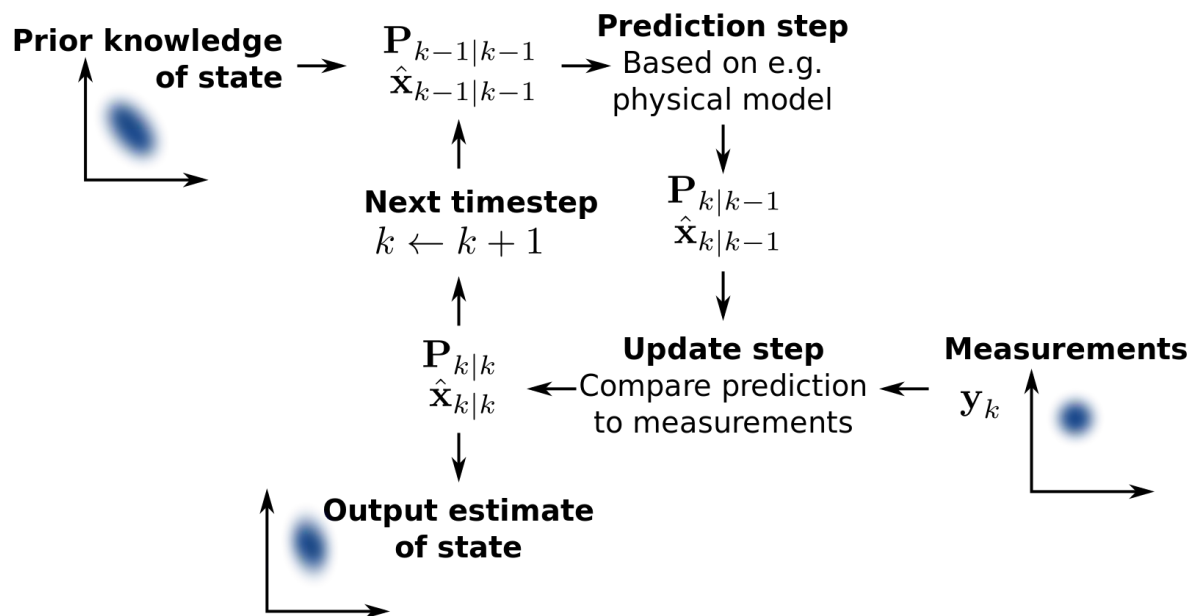Key difference is TEMPORAL consistency

Orcun Goksel,  ETH Zurich

# **Trajectory**

# **(Temporal Filtering)**

# Temporal Filtering/Predictions



- To predict location
- To reduce noise
- To disambiguate multiple objects

## Kalman Filtering

**Prior knowledge of state** → $\mathbf{P}_{k-1|k-1}$ $\hat{\mathbf{x}}_{k-1|k-1}$ → **Prediction step** Based on e.g. physical model

**Next timestep** $k \leftarrow k+1$

$\mathbf{P}_{k|k-1}$ $\hat{\mathbf{x}}_{k|k-1}$

$\mathbf{P}_{k|k}$ $\hat{\mathbf{x}}_{k|k}$ ← **Update step** Compare prediction to measurements ← **Measurements** $\mathbf{y}_k$

**Output estimate of state**

# Steps of Tracking

predict     correct

- Recap: Particle filtering
  - Tracking can be seen as the process of propagating the posterior distribution of state given measurements across time.

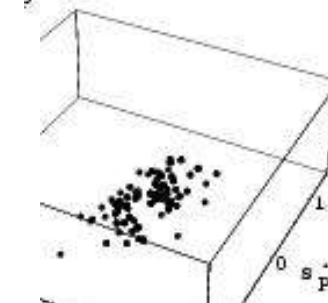**Particle Filter**

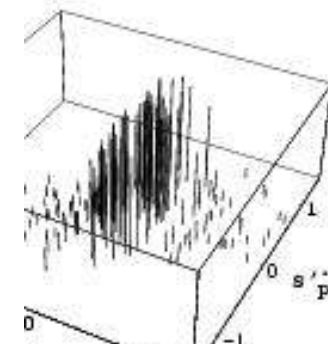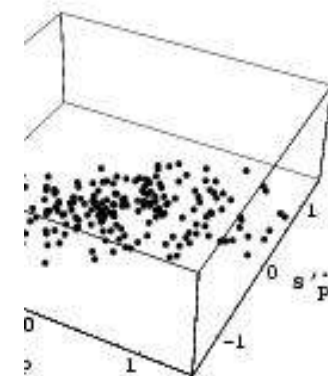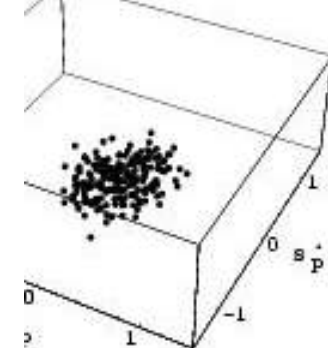$$p(p_{t-1}, \dot{p}_{t-1} \mid z_{t-1})$$
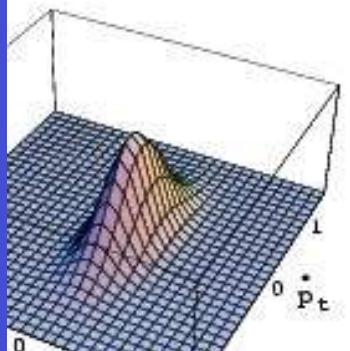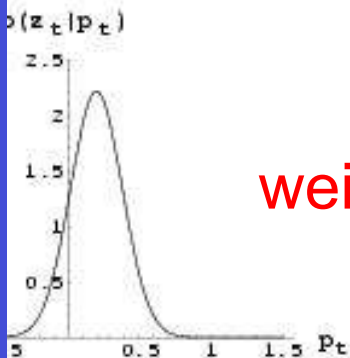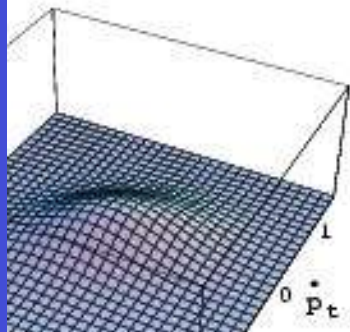
↓ prediction

$$p(p_t, \dot{p}_t \mid z_{t-1})$$

weighing with $\boxed{p(z_t \mid p_t)}$

↓ update

$$p(p_t, \dot{p}_t \mid z_t)$$

C O N D E N S A T I O N

# Traditional/Simple Tracking



**t=1**

**initialization**

**t=2**

**position in prev. frame**

**candidate new positions (e.g., dynamics)**

**best new position (e.g., max color similarity)**

# Tracking-by-Detection



...

**detect object(s) independently in
each frame**

**associate detections over time into
*tracks***

# Outline

Computer Vision

Feature

- Region Tracking
- Point Tracking
- Template Tracking
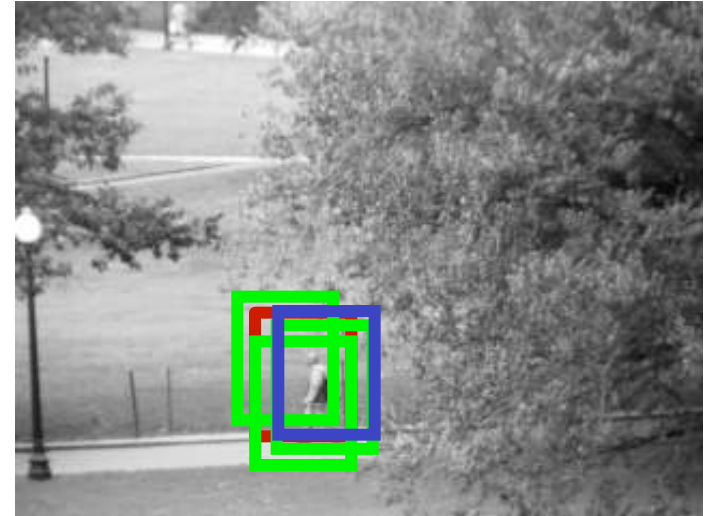
---

Model

- Tracking-by-Detection
  - of a specific target
  - of the object class
- Model-based Body Articulation
- On-line Learning

---

- Misc (preventing drift, context, issues)

# Region Tracking
## (and Mean Shift Algorithm)

# Background Modeling

**Input**

**Background Model**

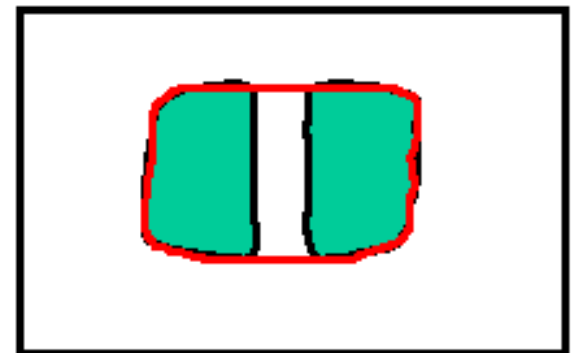**Moving Foreground Blobs (Objects)**

# Deformable models

- One option: Fit deformable curves



$(x_0, y_0)$

$(x_3, y_3)$

$(x_{17}, y_{17})$

# Mean Shift Tracking

- The mean shift tracker tracks a region, with a prescribed (color) distribution

- The similarity between the tracked region and the target region is maximized, through evolution towards higher density in a parameter space

- Typically this search only takes a few iterations

[Comaniciu and Meer, ICCV'99]

Orcun Goksel,  ETH Zurich

Computer Vision

# Meanshift Tracking

**Region of interest (Kernel)**

**Center of mass**

**Mean Shift vector**

**Measurements**

# Intuitive Description

# Intuitive Description

# Intuitive Description

# Intuitive Desciption

# Example: Safety Monitoring

# Outline

Computer Vision

**Feature**

- Region Tracking
- Point Tracking ←
- Template Tracking

**Model**

- Tracking-by-Detection
  - of a specific target
  - of the object class
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Point Tracking
# (and Aperture Problem)

# Estimate Optimal Transformation

# When can we estimate motion?

Q1. Which direction is the pattern behind the circular hole moving?

a) ↙   b) ←   c) ↓   d) **?**

gradients

Q2. 1D motion: We cannot determine the direction of motion from red to green line on the right.  Why not?

?

Q3. Any similarity/connection between Q1 & Q2?

# Sum of Squared Differences

$I_0(x)$

$I_1(x) = I_0(x+h)$

**h**

$$E(h) = [I_0(x+h) - I_1(x)]^2$$

# Displacement

$$E(h) = [\ I_0(x+h) - I_1(x)\ ]^2$$

$$E(h) \approx [\ I_0(x) + hI_0{'}(x) - I_1(x)\ ]^2$$

$$\frac{\partial E}{\partial h} \approx 2\ I_0{'}(x)\ [\ I_0(x) + hI_0{'}(x) - I_1(x)\ ] = 0$$

$$h \approx \frac{I_1(x) - I_0(x)}{I_0{'}(x)}$$

Computer Vision

$I_1(x)-I_0(x)$

$I_0(x)$

$I_1(x)$

$h$

$I_0'(x)$

$$h \approx \frac{I_1(x) - I_0(x)}{I_0'(x)}$$

# Problem 1: Zero Gradient

$$h \approx \frac{I_1(x) - I_0(x)}{\boxed{I_0{'}(x)}}$$

# Problem 1: "Aperture problem"

- Tracking needs gradients in all possible directions to be well defined

- If no gradient along one direction, we cannot determine relative motion in that axis

# Problem 2: Local Minima

**(a)**

**(b)**

- Motion to closest minimum has to be assumed

- Indirect result: Frame-rate should be faster than motion of half-wavelength (Nyquist rate)

- Nonconvex regions may indicate multiple sols

# Problem 2: Local Minima

# Recall: Optical Flow



$I_0$ (time t)  motion v  $I_1$ (time t+1)

- OF recovers (smooth) motion everywhere
- Least-squares regularization: Horn-Schunk makes smooth spatial change assumption
- In contrast, tracking seeks a single motion!

# Recall: Optical Flow

$$I_x u + I_y v + I_t = 0$$

$$I_x = \frac{\partial I}{\partial x}, \qquad I_y = \frac{\partial I}{\partial y}, \qquad I_t = \frac{\partial I}{\partial t}$$

$$u = \frac{dx}{dt}, \qquad v = \frac{dy}{dt}$$

**1 equation in 2 unknowns**

# Treating Aperture Problem in Tracking

- Get additional info to constrain motion:
  - OF: Smoothly regularize in space
  - Tracking: Assume single motion for a region

- Spatial coherence constraint: "A pixel's neighbours all move together"

$I_0(x)$

$I_1(x)$

# Least Squares Problem:
## Single motion with multiple equations

$$\begin{bmatrix} I_x(\mathbf{p_1}) & I_y(\mathbf{p_1}) \\ I_x(\mathbf{p_2}) & I_y(\mathbf{p_2}) \\ \vdots & \vdots \\ I_x(\mathbf{p_{25}}) & I_y(\mathbf{p_{25}}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p_1}) \\ I_t(\mathbf{p_2}) \\ \vdots \\ I_t(\mathbf{p_{25}}) \end{bmatrix}$$

**Over determined System of Equations**

$$\begin{array}{ccc} A & d & = b \\ \text{25x2} & \text{2x1} & \text{25x1} \end{array}$$

**Pseudo Inverse**

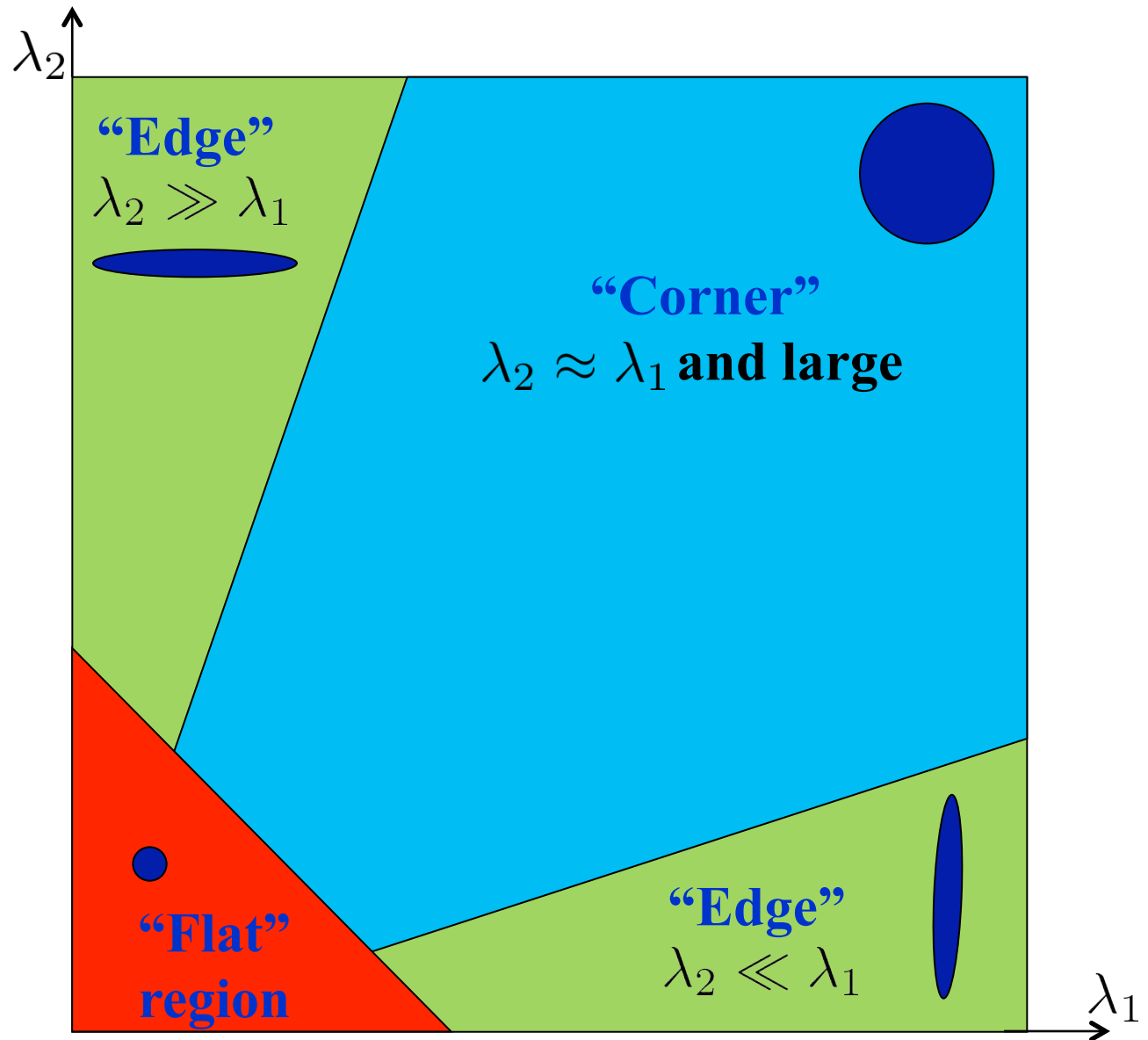$$\begin{array}{ccc} (A^T A) & d & = A^T b \\ \text{2x2} & \text{2x1} & \text{2x1} \end{array}$$

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$
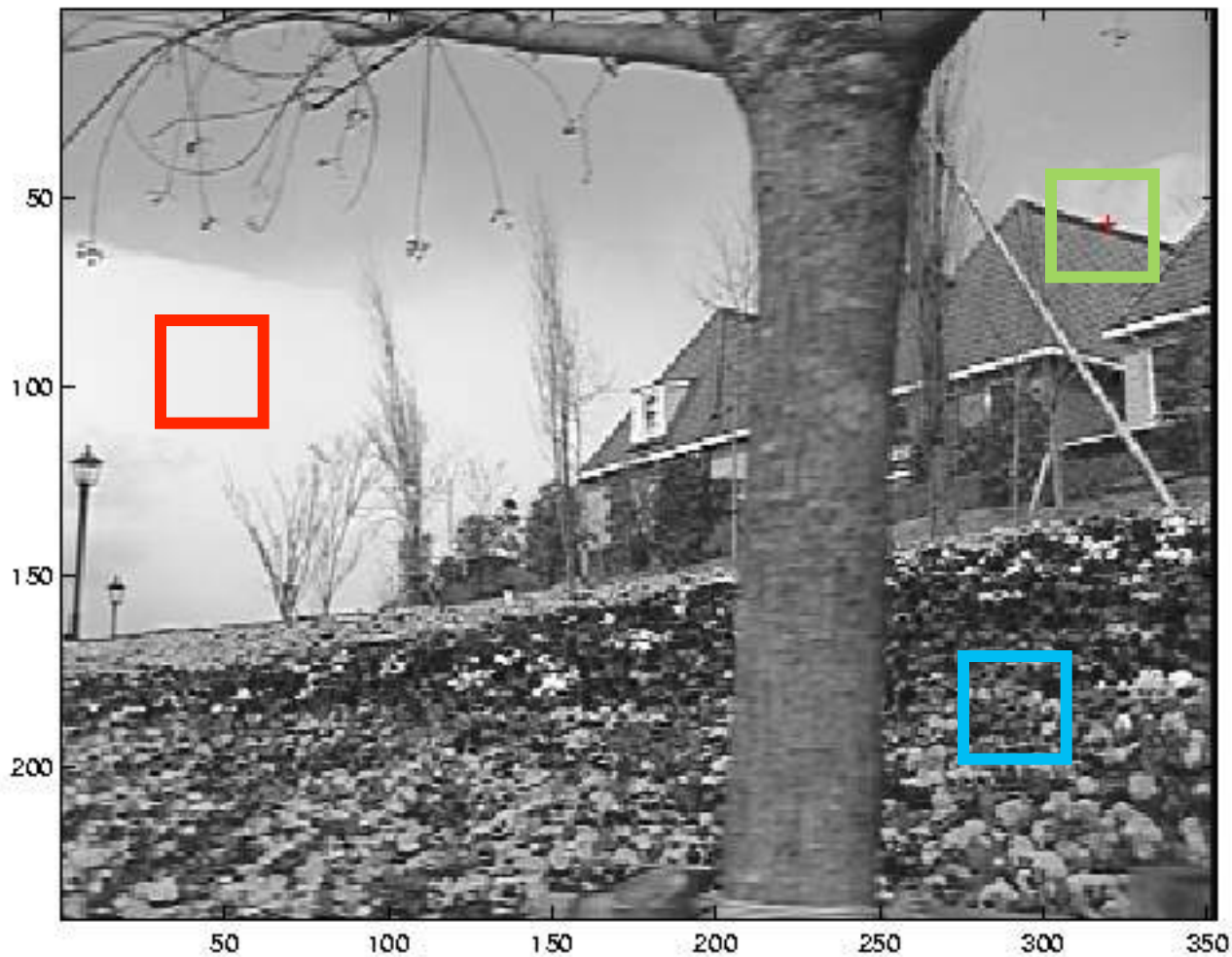
# Eigenvectors of A$^\mathrm{T}$A

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] = \sum \nabla I (\nabla I)^T$$

- Haven't we seen an equation like this before?

- Recall Harris corner detector!

- Thus, "good image features are also good for tracking"

# Interpreting the Eigenvalues



$\lambda_2$

**"Edge"**
$\lambda_2 \gg \lambda_1$

**"Corner"**
$\lambda_2 \approx \lambda_1$ **and large**

**"Flat" region**

**"Edge"**
$\lambda_2 \ll \lambda_1$

$\lambda_1$

Computer Vision

# Outline

Computer Vision

**Feature**

- Region Tracking

- Point Tracking

- Template Tracking ⬅

**Model**

- Tracking-by-Detection
  - of a specific target
  - of the object class
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Template Tracking

# Lucas-Kanade Template Tracker

- Lucas-Kanade is typically for small patches, e.g. 5x5
- Why not run it for the entire object (for a larger window)



- Locally, translation is sufficient to explain motion; but…

# Lucas-Kanade Template Tracker

- Motion is more complex in a larger window



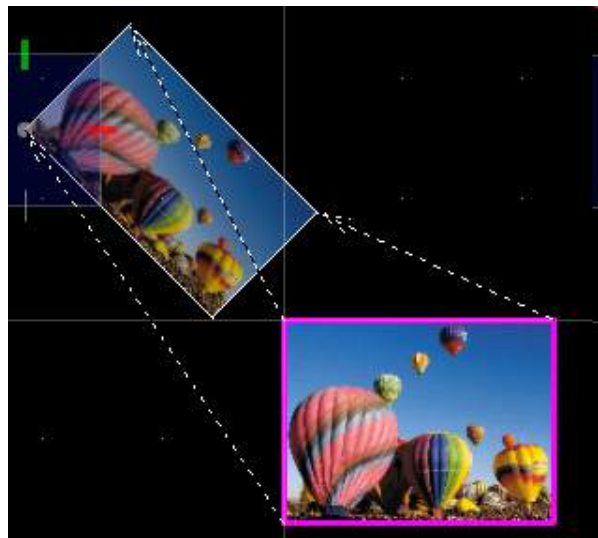- Nonetheless, we can easily generalize the motion model to other parametric models!

e.g., translation, affine, projective, "warp"

$$E(u,v) = \sum_{x,y} [I(x+u, y+v) - T(x,y)]^2$$

$$E(u,v) = \sum_{x,y} [I(W(x;p)) - T(x,y)]^2$$

# Lucas-Kanade Template Tracker
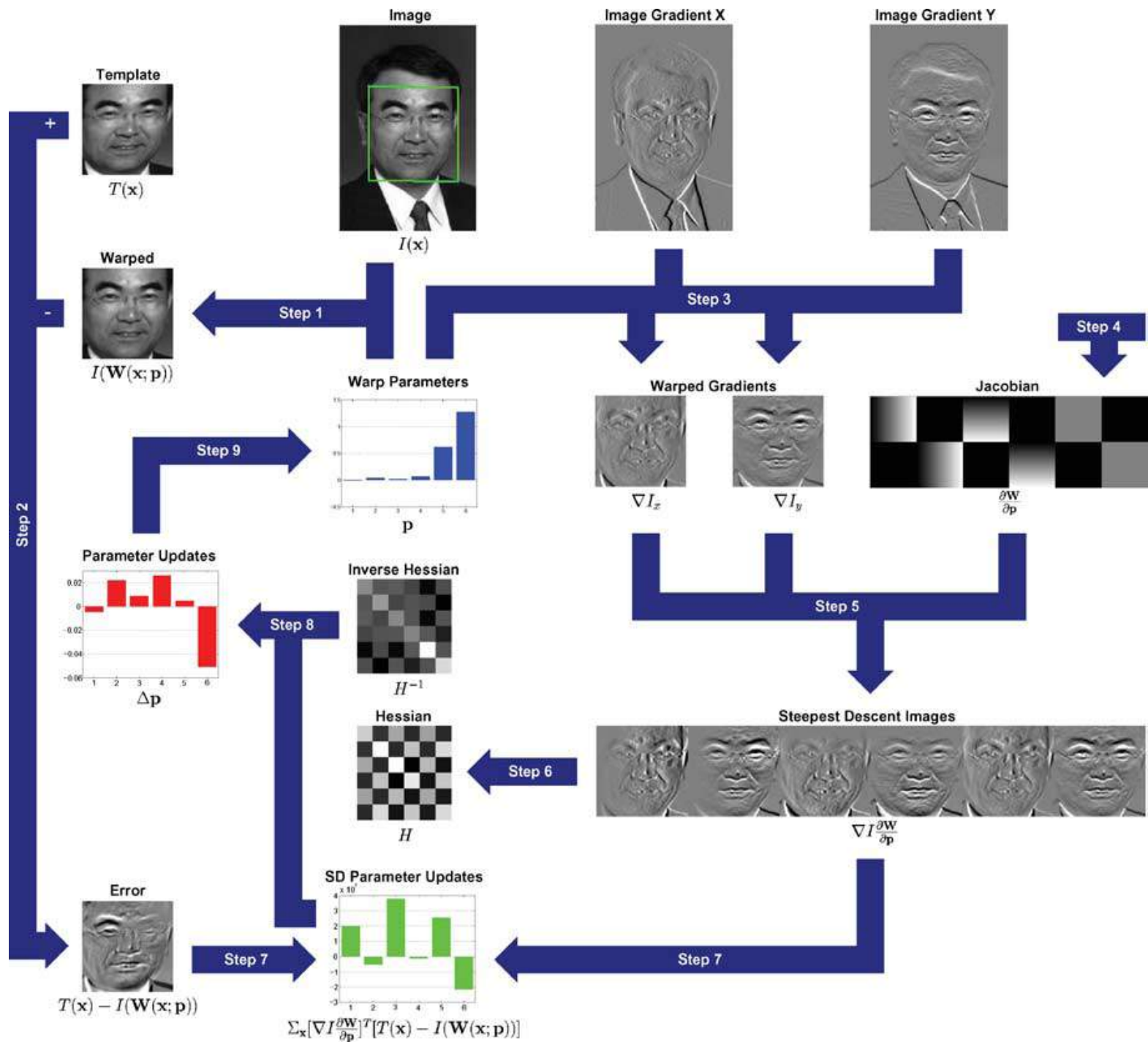
- From Points to templates

- Estimate „optimal" warp $W$

$$\sum_{\mathbf{x}} [I(\mathbf{W}(\mathbf{x}; \mathbf{p})) - T(\mathbf{x})]^2$$

$$\sum_{\mathbf{x}} [I(\mathbf{W}(\mathbf{x}; \mathbf{p} + \triangle\mathbf{p})) - T(\mathbf{x})]^2$$

[Baker & Matthews, IJCV'04, Lucas-Kanade 20 Years On: A Unifying Framework]

# Lucas-Kanade Template Tracker

Step 1. Warp $I$ to obtain $I(W([x\ y];P))$

Step 2. Compute the error image $T(x) - I(W([x\ y];\ P))$

Step 3. Warp the gradient $\nabla I$ with $W([x\ y];\ P)$

Step 4. Evaluate $\dfrac{\partial W}{\partial P}$ at $([x\ y];\ P)$ (Jacobian)

Step 5. Compute steepest descent images $\nabla I \dfrac{\partial W}{\partial P}$

Step 6. Compute Hessian matrix $\sum (\nabla I \dfrac{\partial W}{\partial P})^T (\nabla I \dfrac{\partial W}{\partial P})$

Step 7. Compute $\sum (\nabla I \dfrac{\partial W}{\partial P})^T (T(x,y) - I(W([x,y];P)))$
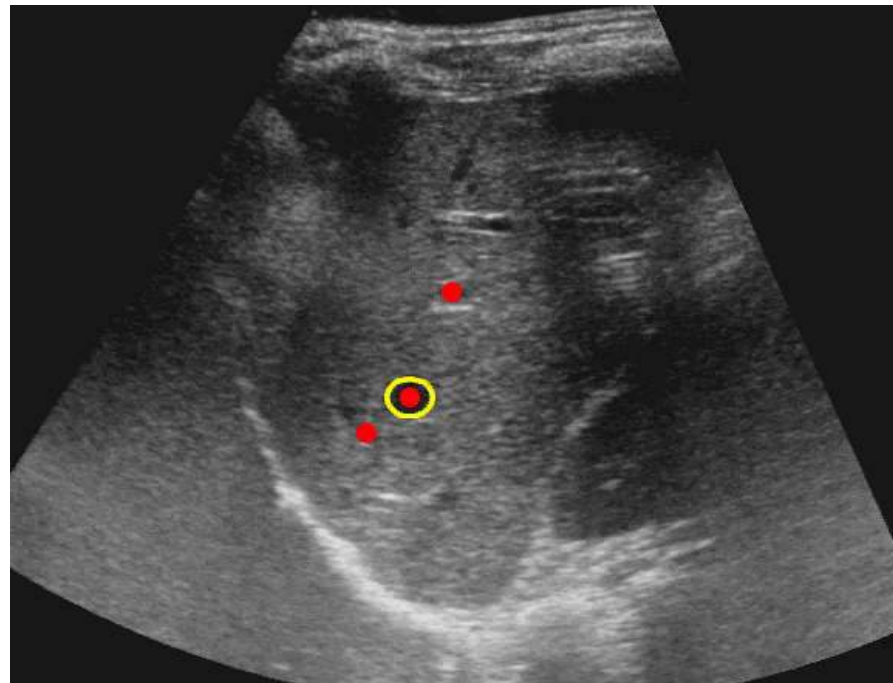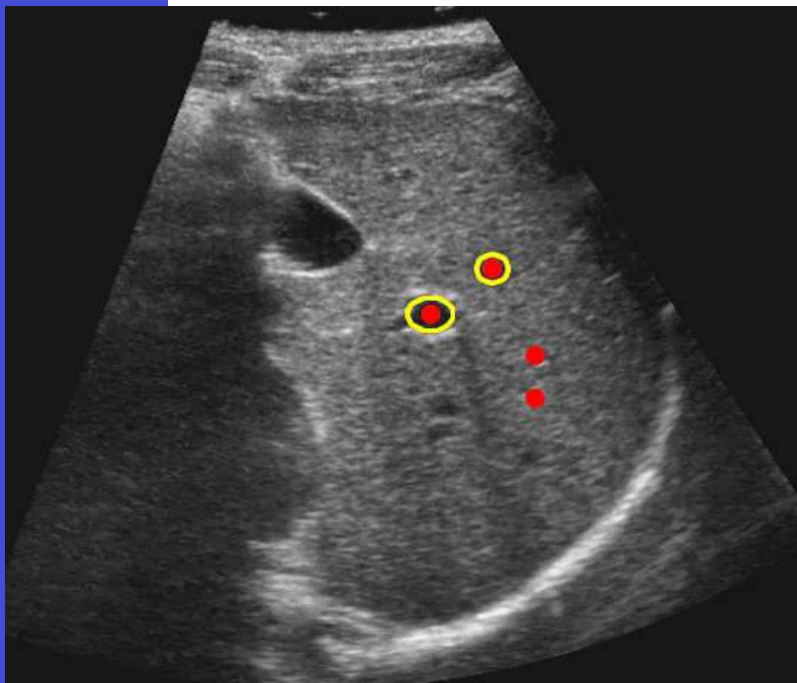
Step 8. Compute $\Delta P$

Step 9. Update $P \longleftarrow P + \Delta P$

● Our tracking
+ Manual annotation

Orcun Goksel, ETH Zurich
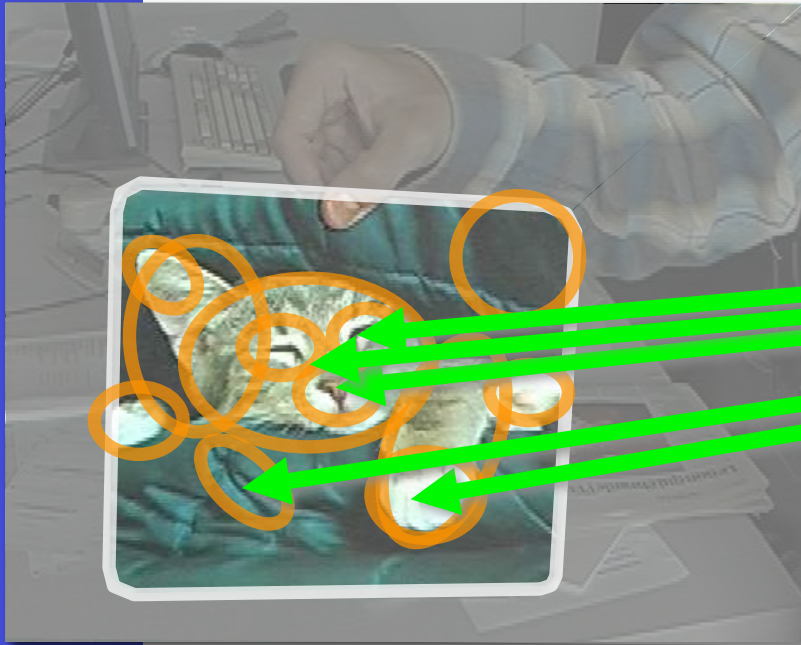
# Outline

**Feature**

- Region Tracking
- Point Tracking
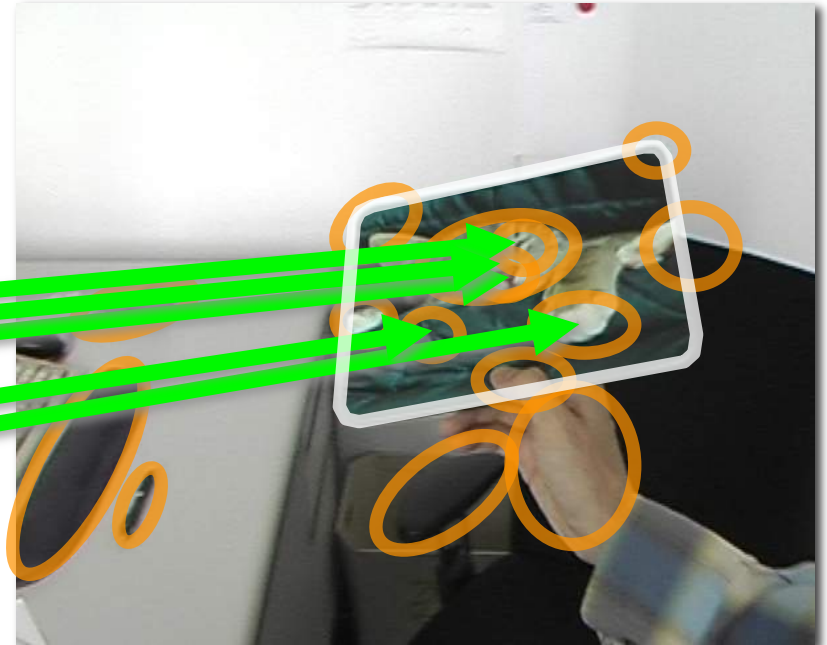- Template Tracking

**Model**

- Tracking-by-Detection
  – of a specific target ⬅
  – of the object class
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Tracking by Detection

## (of a specific target)

# 3D Object Detection



**Reference image(s) of
the object to detect**

**Test image**

# 3D Object Detection

MathWorks



**Reference image(s) of
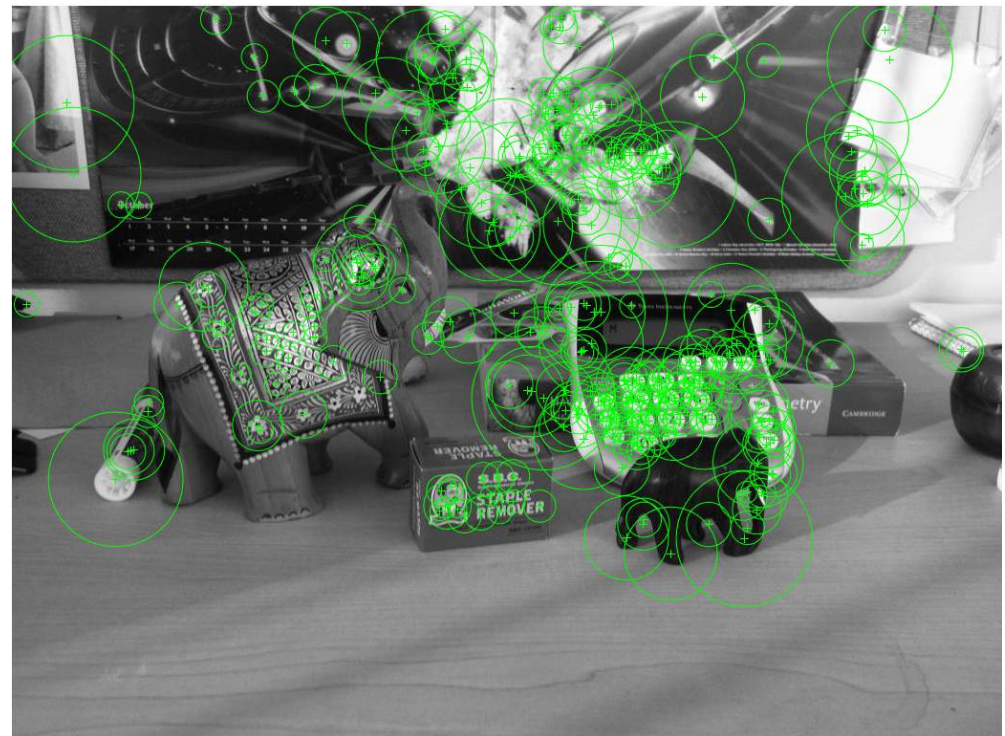the object to detect**

**Test image**

# 1. Detect Keypoints
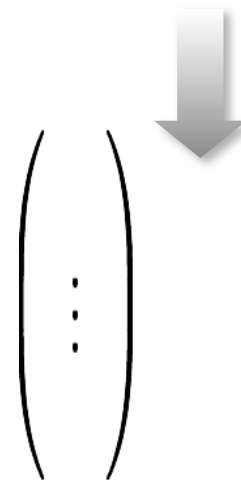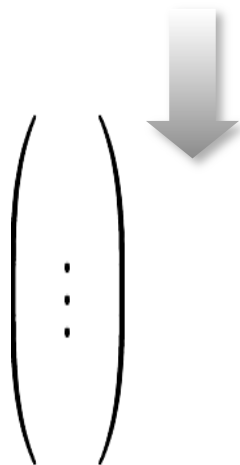
– invariant to scale, rotation, or perspective



100 strongest feature points
in the reference image

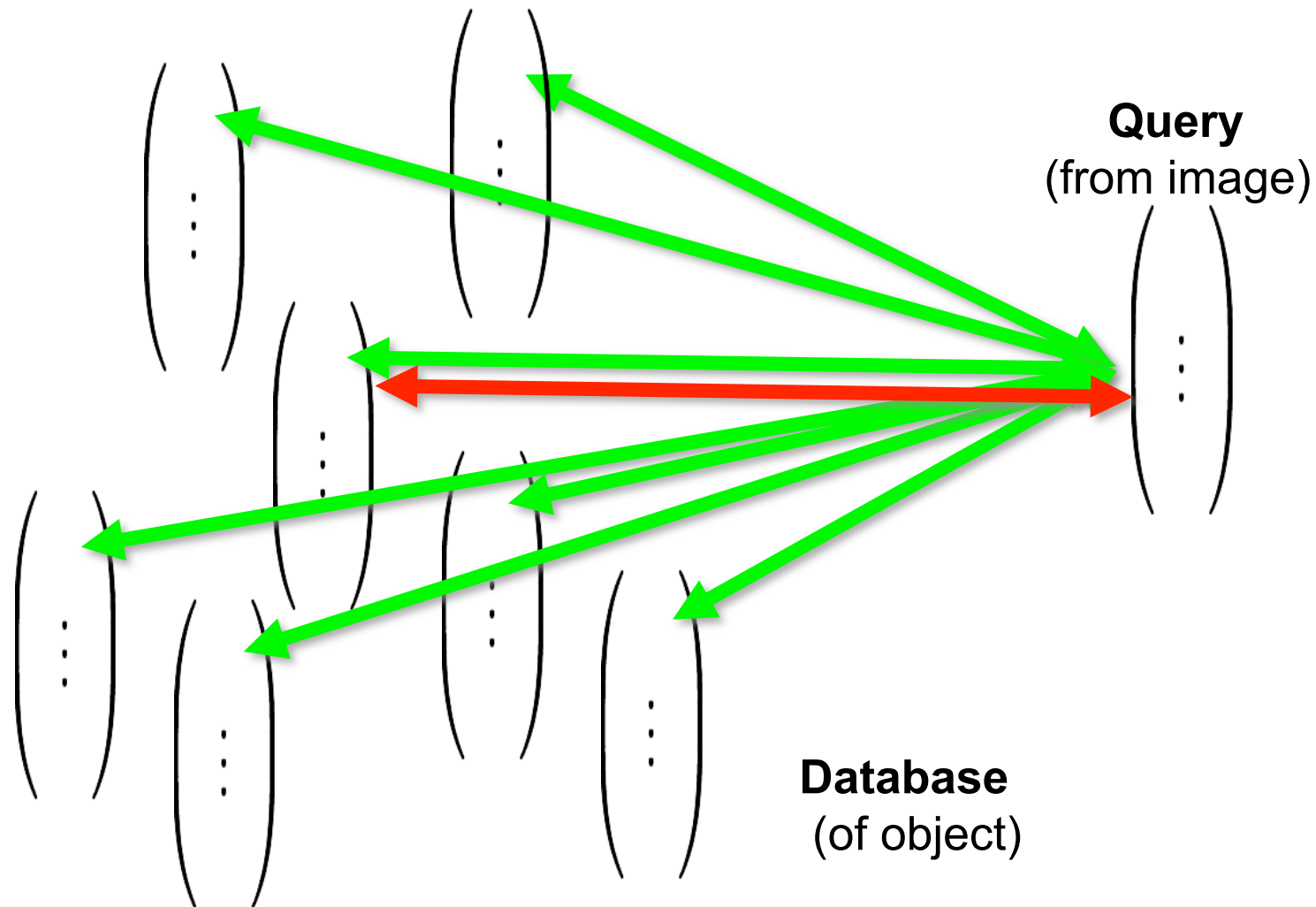300 strongest feature points
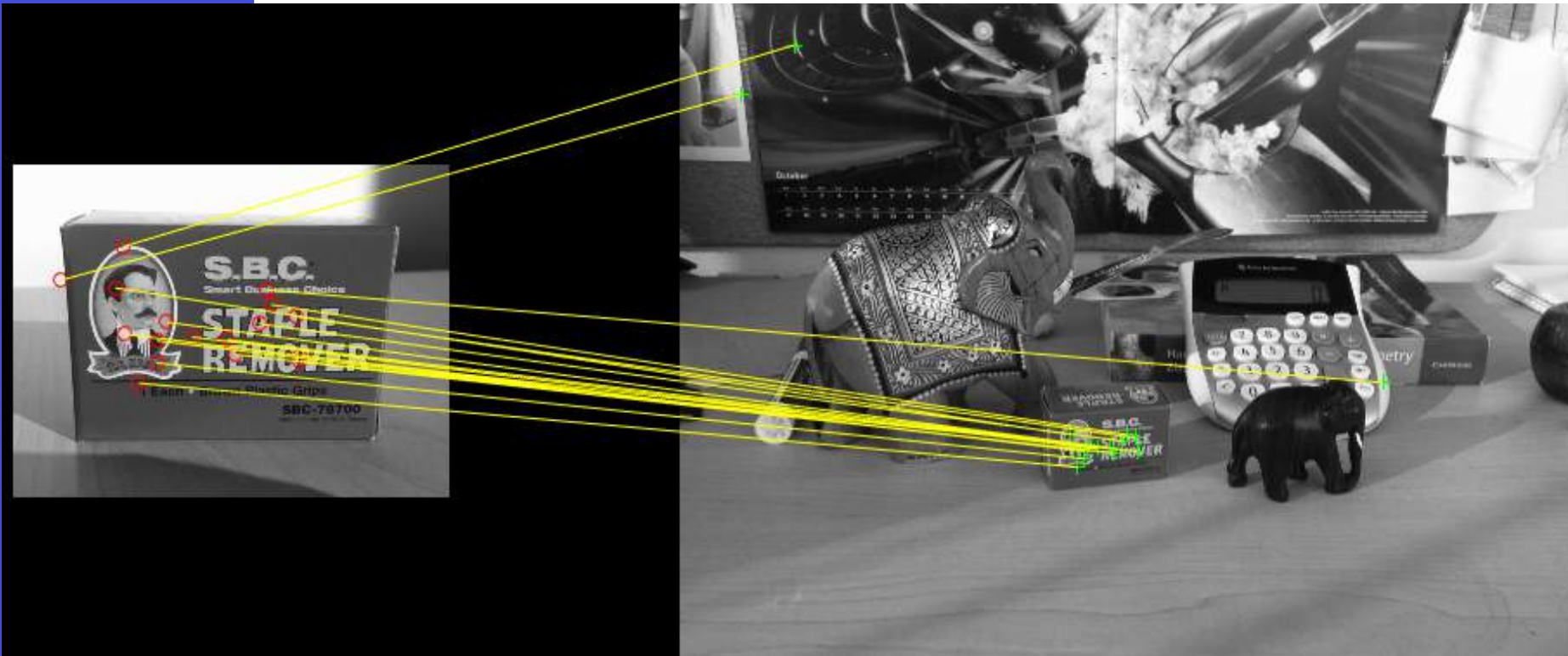in the test image

Computer
Vision

# 2. Build Feature Descriptors

# 3. Match Keypoint Descriptors

- Search in the Database



**Query**
(from image)

**Database**
(of object)

# 3. Search in the Database

# 4. Outlier Elimination

# Summary

**Keypoint Detection**

**Keypoint Recognition**

Search in the
Database

$$\left( \vdots \right)$$

Robust 3D Pose
Calculation

(RANSAC)

**Geometric
verification**

# Computer Vision

Overall: 3.42 ms
Find Pts: 1.25 ms
Track Pts: 0.32 ms
Features: 1.16 ms
Outliers: 0.50 ms
Pose: 0.19 ms

Corners: 166
Matched Features: 29
 Wrong Rotation: 0
 Bad Linetest: 0
 Bad Homographytest: 0
Correct: 29
 From Cache: 0
 From ActiveSearch: 20
 Levels: 0 0 0 0 0 0 0 0 0
Rotation: 6
Avg. Reproj. Err: 1.31

www.topmodelgossip.com

Orcun Goksel, ETH Zurich

Orcun Goksel,  ETH Zurich

# Outline

**Feature**

- Region Tracking

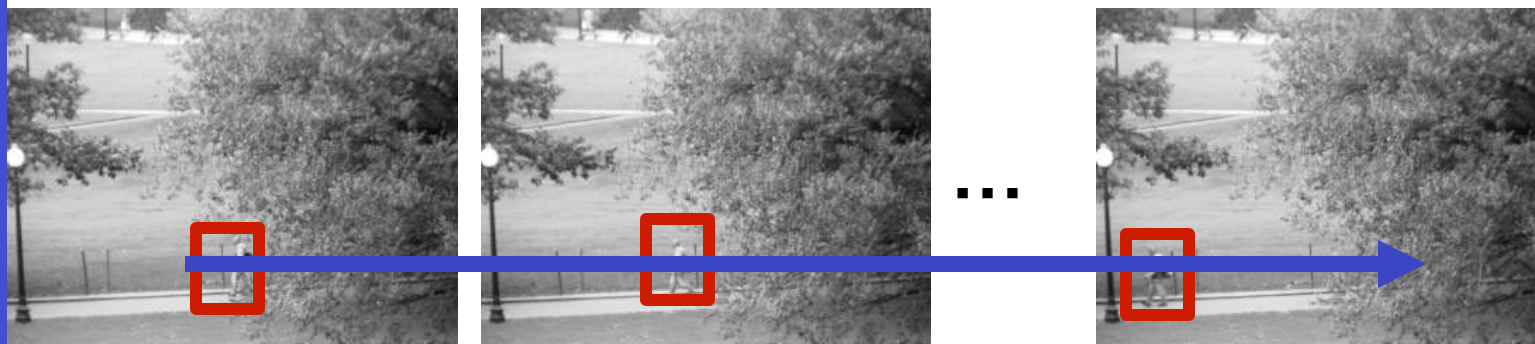- Point Tracking

- Template Tracking

**Model**

- Tracking-by-Detection
  - of a specific target
  - of the object class ⬅

- Model-based Body Articulation

- On-line Learning

- Misc (preventing drift, context, issues)

# Tracking by Detection
## (of the object class)

## "Multiple Object Tracking"

# Tracking-by-Detection

**detect object(s) independently in each frame**
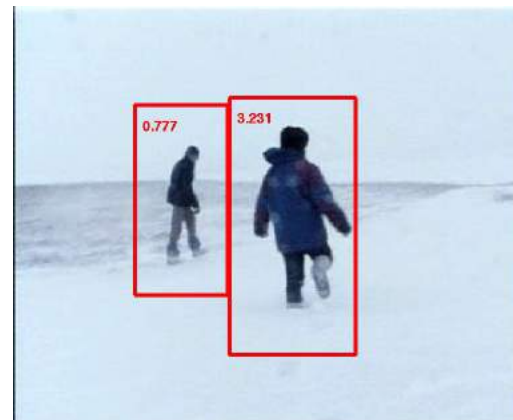
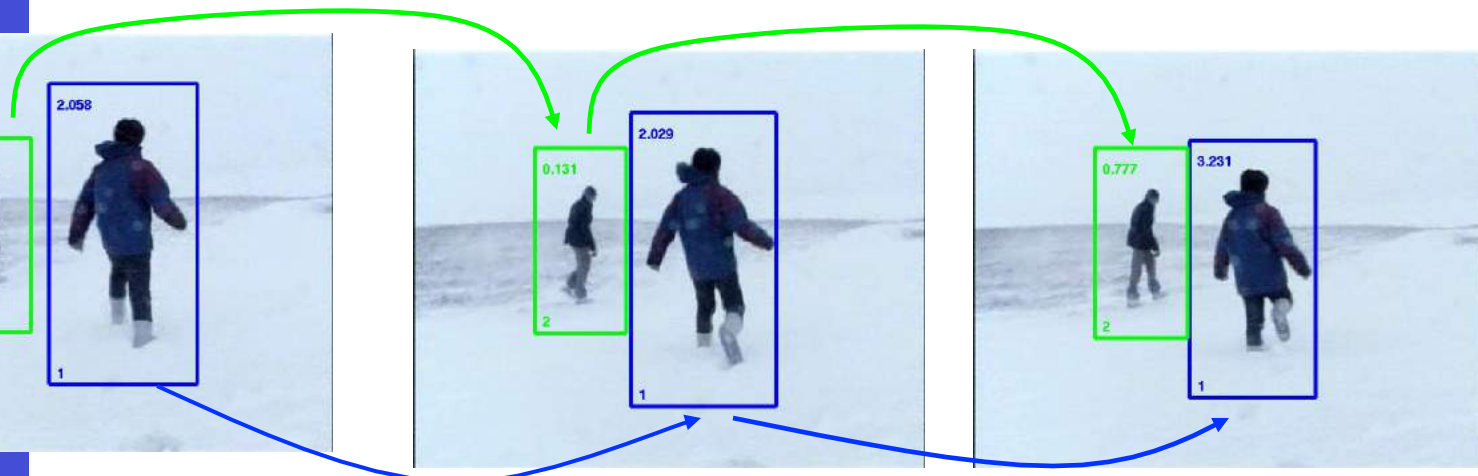**associate detections over time into *tracks***

# Multiple Objects

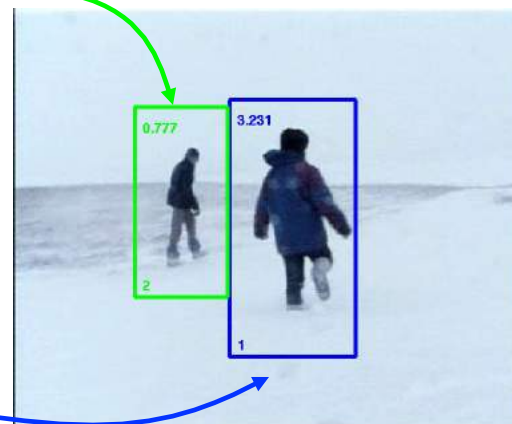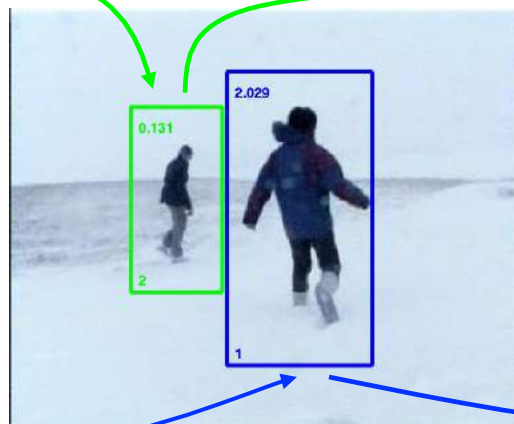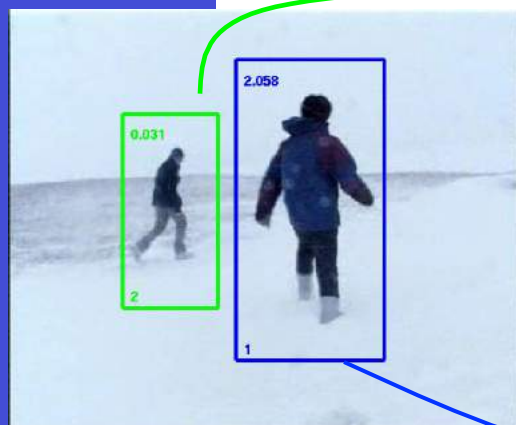Frame 1                    Frame 5                    Frame 9

# Examples:
# Multiple Object Tracking

# How to get the detections?



**Persons**

**Background**

**Supervised Learning**

(Support Vector Machines,
Random Forests,
Neural Networks, ...)

# Using the classifier

# Space-Time Analysis

- Collect detections in space-time volume

**Detections**

**Space Time Volume**



Orcun Goksel, ETH Zurich

# Trajectory Estimation

- Trajectory growing and selection



**Space Time Volume**

Orcun Goksel,  ETH Zurich

# Trajectory Estimation

- Trajectory growing and selection

# Driving

**Input (Object Detections)**



**"Tracking" Result**

# Outline

**Feature**

- Region Tracking  (and Mean Shift Algorithm)
- Point Tracking  (and Aperture Problem)
- Template Tracking  (Lucas-Kanade)

**Model**

- Tracking-by-Detection
  - a specific target (e.g., keypoints + Ransac)
  - object class  (multiple object tracking)
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Model based Tracking

# Articulated Tracking: Part-Based Models





- Intuitive model of an object
- Model has two components
  1. parts (2D image fragments)
  2. structure (configuration of parts)
- Dates back to Fischler & Elschlager 1973

# Parts-based analysis

Objective: detect human and determine upper body pose (layout)



## Model as a graph labelling problem

- Vertices $\mathcal{V}$ are parts, $a_i, i = 1, \cdots, n$

- Edges $\mathcal{E}$ are pairwise linkages between parts

- For each part there are $h$ possible poses $\mathbf{p}_j = (x_j, y_j, \phi_j, s_j)$

- Label each part by its pose: $f : \mathcal{V} \longrightarrow \{1, \cdots, h\}$, i.e. part $a$ takes pose $\mathbf{p}_{f(a)}$.

# Parts-based analysis

## Pictorial structure model – CRF



- Each labelling has an energy (cost):

$$E(f) = \underbrace{\sum_{a \in \mathcal{V}} \theta_{a;f(a)}}_{\text{unary terms (appearance)}} + \underbrace{\sum_{(a,b) \in \mathcal{E}} \theta_{ab;f(a)f(b)}}_{\text{pairwise terms (configuration)}}$$

Features for unary:
- colour
- HOG

for limbs/torso

- Fit model (inference) as labelling with lowest energy

# Walking

- What temporal info can we use for tracking?

- What variation would we expect in population?

# Articulation Space

## Tracking Articulated Motion as High-Dimensional Inference

- Walking cycles have one main (periodic) DOF

- Regressors to learn this (latent) space, and its variation (Gaussian Process regression, **PCA**, etc)

- (Pose,Silhouette) training data can be obtained by 3D rendering

Multiple manual orientations (ω)

Motion Capture

Pose Data (p)          3D Render          Silhouettes (s)

# Articulation Space

## Tracking Articulated Motion as High-Dimensional Inference

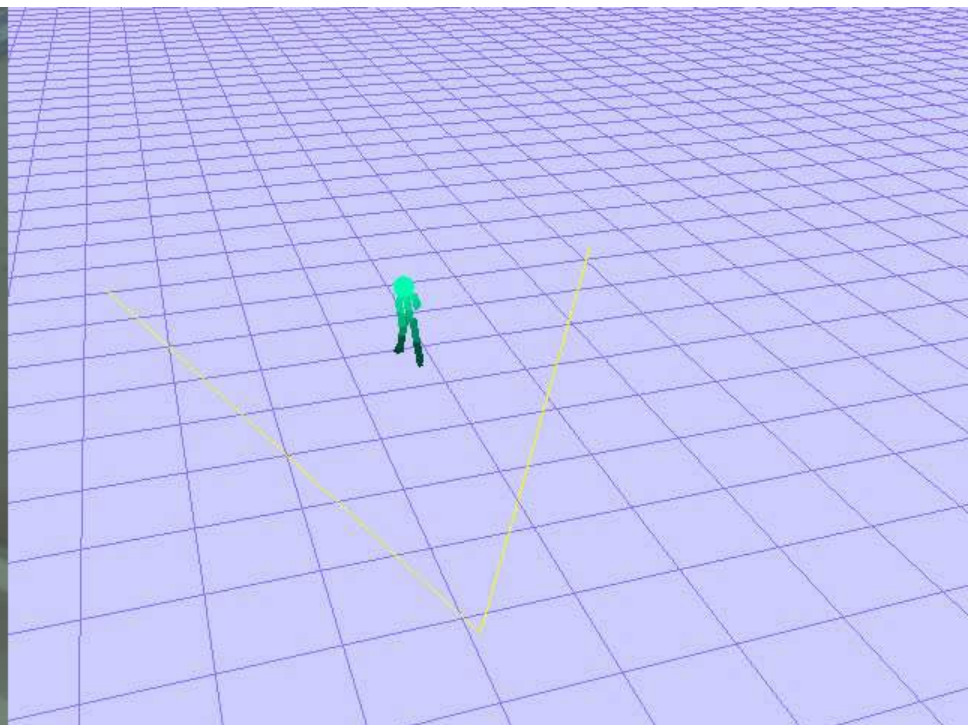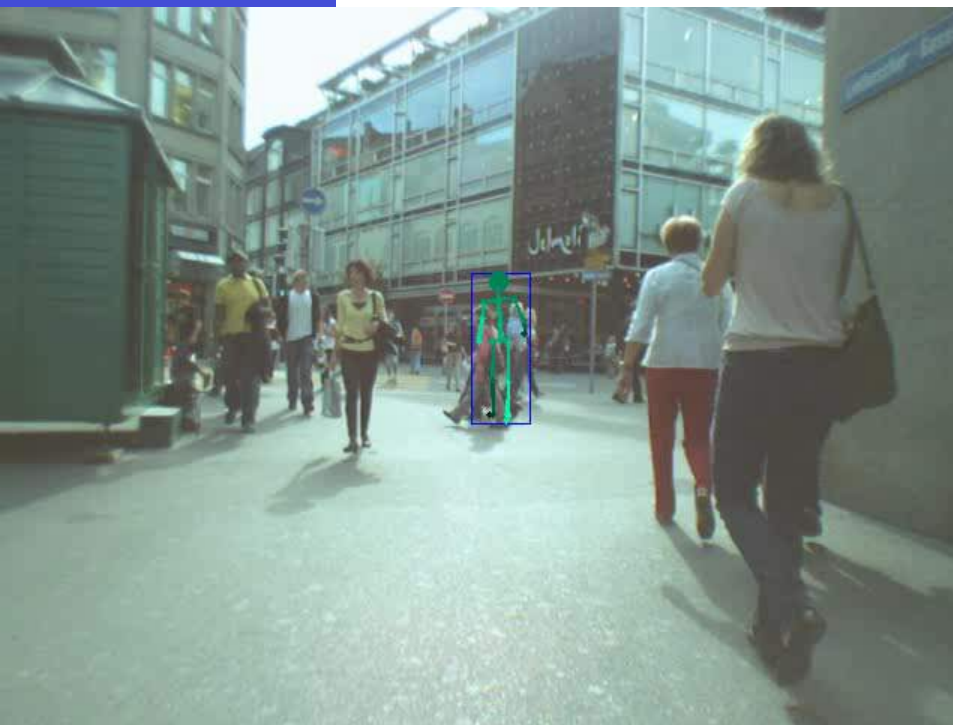- Walking cycles have one main (periodic) DOF

- Regressors to learn this (latent) space, and its variation (Gaussian Process regression, PCA, etc)

- (Pose,Silhouette) training data can be obtained by 3D rendering



**P(Silhouette|k)**
perform inference
on silhouettes

**P(Pose|k)**
recover pose
from latent space

# Articulation Space Tracking

# Outline

**Computer Vision**

**Feature**

- Region Tracking  (and Mean Shift Algorithm)
- Point Tracking  (and Aperture Problem)
- Template Tracking  (Lucas-Kanade)

**Model**

- Tracking-by-Detection
  - a specific target (e.g., keypoints + Ransac)
  - object class  (multiple object tracking)
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Tracking
# as On-line learning
## (updating tracking models)

# Tracking as Classification

- Learning current object appearance vs. local background.



**current
background**

**current
obj. appearance**

# Tracking as Classification



**object**

**vs.**

**background**

# Tracking as Classification



**object**

**vs.**

**background**

# Tracking Loop



**from time t to t+1**

**evaluate classifier on sub-patches**

**actual object posit**

**search Region**

**analyze map and set new object position**

**update classifier (tracker)**

**create confidence map**

# "Tracking the Invisible"

# When does it fail…

# When does it fail…

Computer Vision

actual object position

from time t to t+1

evaluate classifier on sub-patches

search Region

update classifier (tracker)

analyze map and set new object position

create confidence map

**Self-learning!**

# Drift

**Tracked Patches**

**Confidence**

# Drift

# Outline

**Feature**

- Region Tracking  (and Mean Shift Algorithm)
- Point Tracking  (and Aperture Problem)
- Template Tracking  (Lucas-Kanade)

**Model**

- Tracking-by-Detection
  - a specific target (e.g., keypoints + Ransac)
  - object class  (multiple object tracking)
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

Computer
Vision

# Combining Tracking and Detection

## (to avoid drift)

# Refining an object model

- Only thing we are sure about the object is its initial model (e.g. appearance in first frame)

- We can "anchor" / correct our model with that

- This can limit drift

**Current Model**

**Fix (initial) Model**

# Recover from Drift
## using a fixed/anchor model (e.g. first frame)

# Context
# in Tracking

# Humans use context to track

- … objects which change there appearance very quickly.

- … occluded objects or object outside the image.

- … small and/or low textured objects or even "virtual points".

# Computer Vision

Orcun Goksel, ETH Zurich

# Using Supporters

# Assumptions should hold

**With Supporters**

# Problems in Tracking

# Tracking Issues

- Initialization

**Time t = 0**



**object position**

# Tracking Issues

- Obtaining observation…
  - <u>Generative</u>: "render" the state on top of the image and compare
  - <u>Discriminative</u>: classifier or detector score

- …and dynamics model
  - specify using domain knowledge
  - learn (very difficult)

# Tracking Issues

- Model- vs. Model-free-Tracking

# Tracking Issues

- Nonlinear dynamics
  - Sometimes needed to keep multiple trackers in parallel
  - E.g., for abrupt direction changes („Persons")

Correct prediction

Wrong prediction

# Tracking Issues

- Prediction vs. Correction  (cf. Kalman Filtering)

    - If the <u>dynamics</u> model is <u>too strong</u>, tracking will end up <u>ignoring the data</u>.

    - If the <u>observation</u> model is <u>too strong</u>, tracking is <u>reduced to repeated detection</u>.



08.10.2009
<< Rudolf Kalman, ETH-Zurich emeritus professor of mathematics, is awarded the National Medal of Science by Barack Obama – one of the highest accolades for researchers in the USA.

In January 2008, Hungarian-born Kalman received the Charles Draper Prize, which is regarded as the "Nobel Prize" of the engineering world. >>

**http://www.ethlife.ethz.ch/archive_articles/091008_kalman_per/index**

# Tracking Issues

- Data Association –
  Multiple Object Tracking
  - What if we don't know which measurements
    to associate with which tracks?

# Tracking Issues

- Data Association –
  Occlusions / Self Occlusions

# Tracking Issues

- Data Association – Fast Motion

# Tracking Issues

- Data Association –
  Background / Appearance Change
  - Cluttered Background
  - Changes in shape, orientation, color,…

# Tracking Issues

- Drift

  – Errors caused by dynamical model, observation model, and data association tend to accumulate over time

Computer Vision

# Summary

Feature

- Region Tracking  (and Mean Shift Algorithm)
- Point Tracking  (and Aperture Problem)
- Template Tracking  (Lucas-Kanade)

Model

- Tracking-by-Detection
  - a specific target (e.g., keypoints + Ransac)
  - object class  (multiple object tracking)
- Model-based Body Articulation
- On-line Learning

- Misc (preventing drift, context, issues)

# Let's apply

Q. What tracking method would you use in each following application scenario?

What limitations you may expect?

**Task: "Discuss one (or more) in groups"**

App1. Safety: In a lumbar mill, you wish to use CV to stop the blade if a hand reaches nearby.

App2. Medical: You wish to track the motion of an ultrasound probe, to relate images in space.

App3. Autonomous driving: Tracking other nearby vehicles to adjust speed and course

AppX. Your favorite tracking app

Orcun Goksel, ETH Zurich