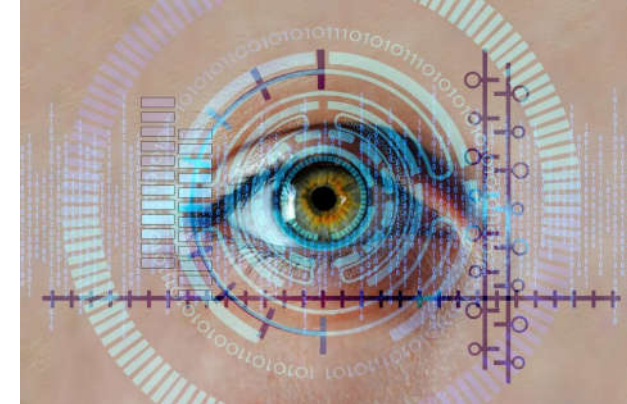




**北京理工大学**  
**BEIJING INSTITUTE OF TECHNOLOGY**

**Computer Vision**



# Lecture 16 Tracking

---

**School of Computer Science and Technology**

**Ying Fu**



[fuying@bit.edu.cn](mailto:fuying@bit.edu.cn)





# Outline

---

- Feature tracking
- Simple KLT tracker
- 2D transformations (recap)
- Iterative KLT tracker

# Outline

---



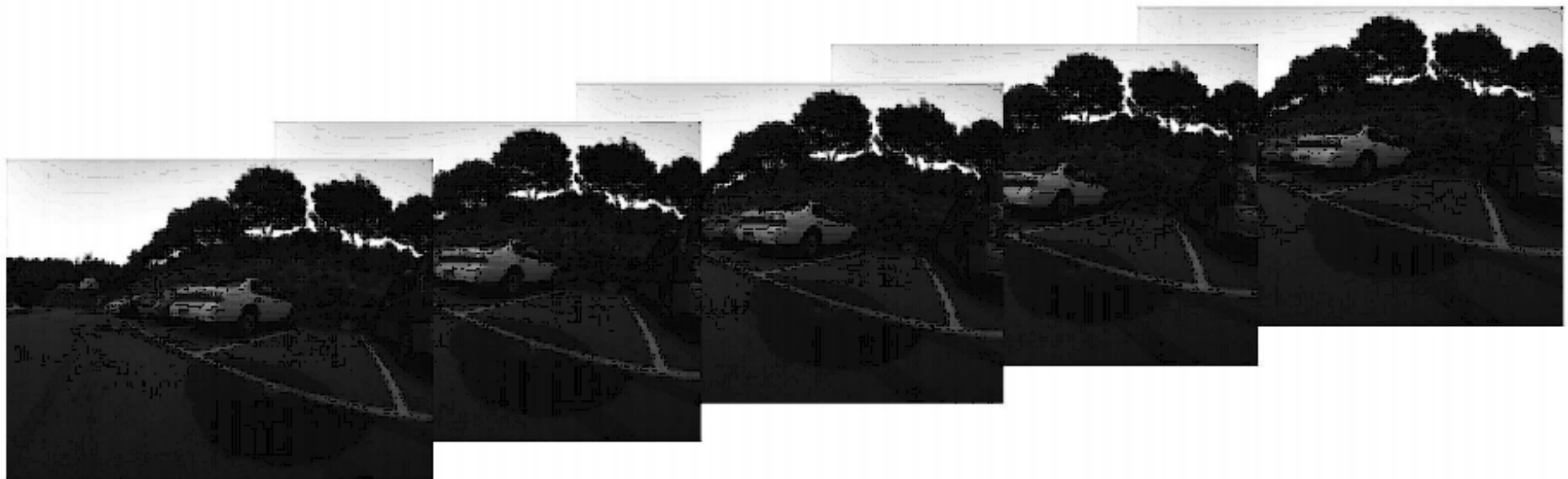
- Feature tracking
  - Problem statement
  - Overview

# Problem statement

---



- Image sequence



# Problem statement

---



- Feature point detection

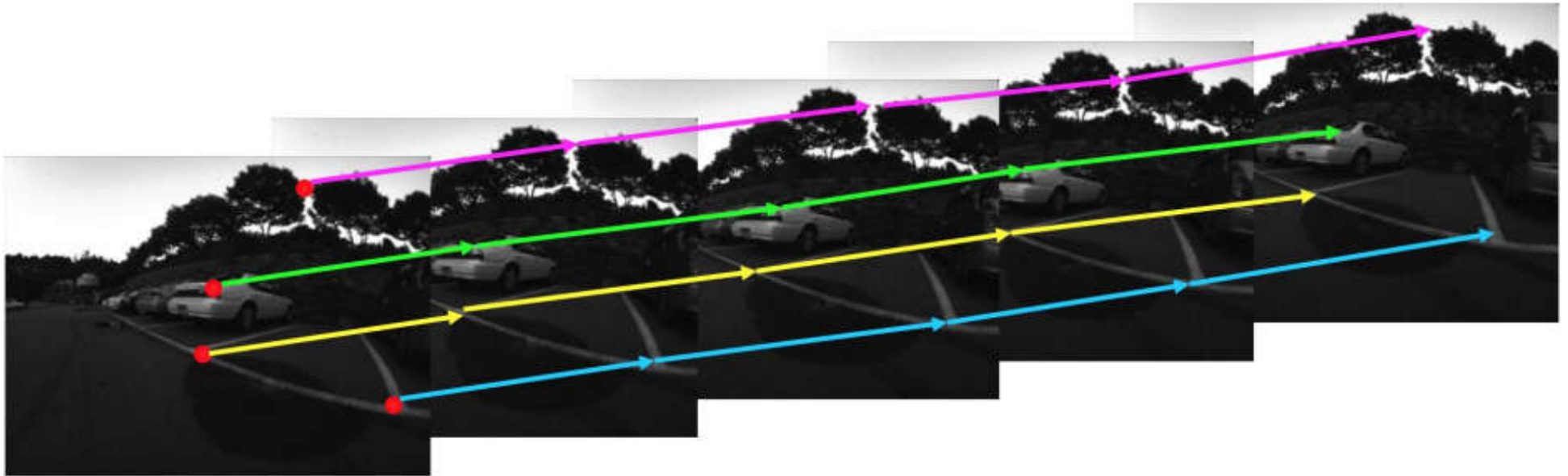




# Problem statement

---

- Feature point tracking



# Single object tracking

---



# Multiple object tracking

---





# Tracking with a fixed camera

---

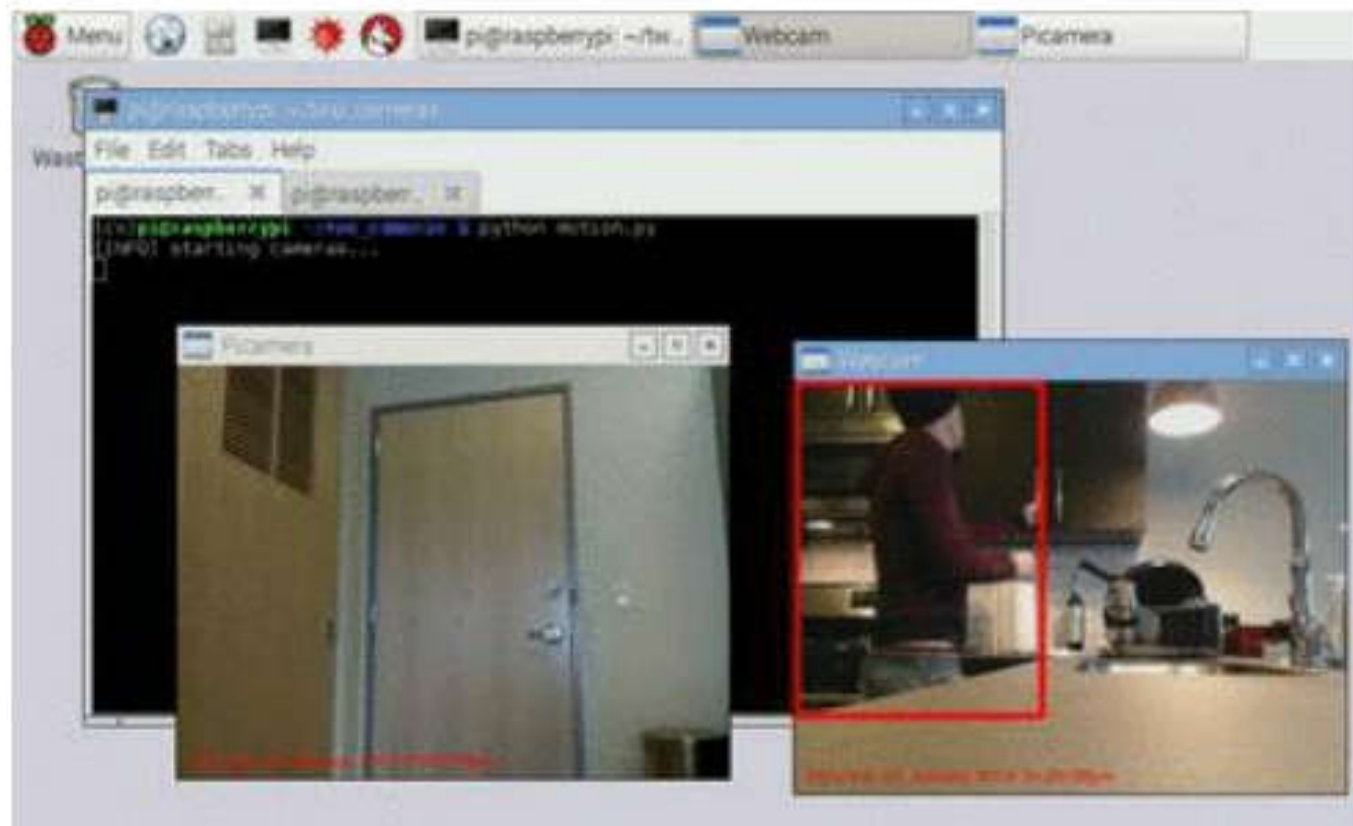


# Tracking with a moving camera

---



# Tracking with multiple cameras





# Challenges in Feature tracking

---

- Figure out which features can be tracked
  - Efficiently track across frames
- Some points may change appearance over time
  - e.g., due to rotation, moving into shadows, etc.
- Drift: small errors can accumulate as appearance model is updated
- Points may appear or disappear.
  - need to be able to add/delete tracked points.

# What are good features to track?



- 
- Intuitively, we want to avoid smooth regions and edges. But is there a more principled way to define good features?
  - **What kinds of image regions can we detect easily and consistently?** Think about what you learnt earlier in the class.

# What are good features to track?



- 
- Can measure “quality” of features from just a single image.
  - Hence: tracking Harris corners (or equivalent) guarantees small error sensitivity!



# Motion estimation techniques

---

- Optical flow
  - Recover image motion at each pixel from spatio-temporal image brightness variations (optical flow)
- Feature-tracking
  - Extract visual features (corners, textured areas) and “track” them over multiple frames

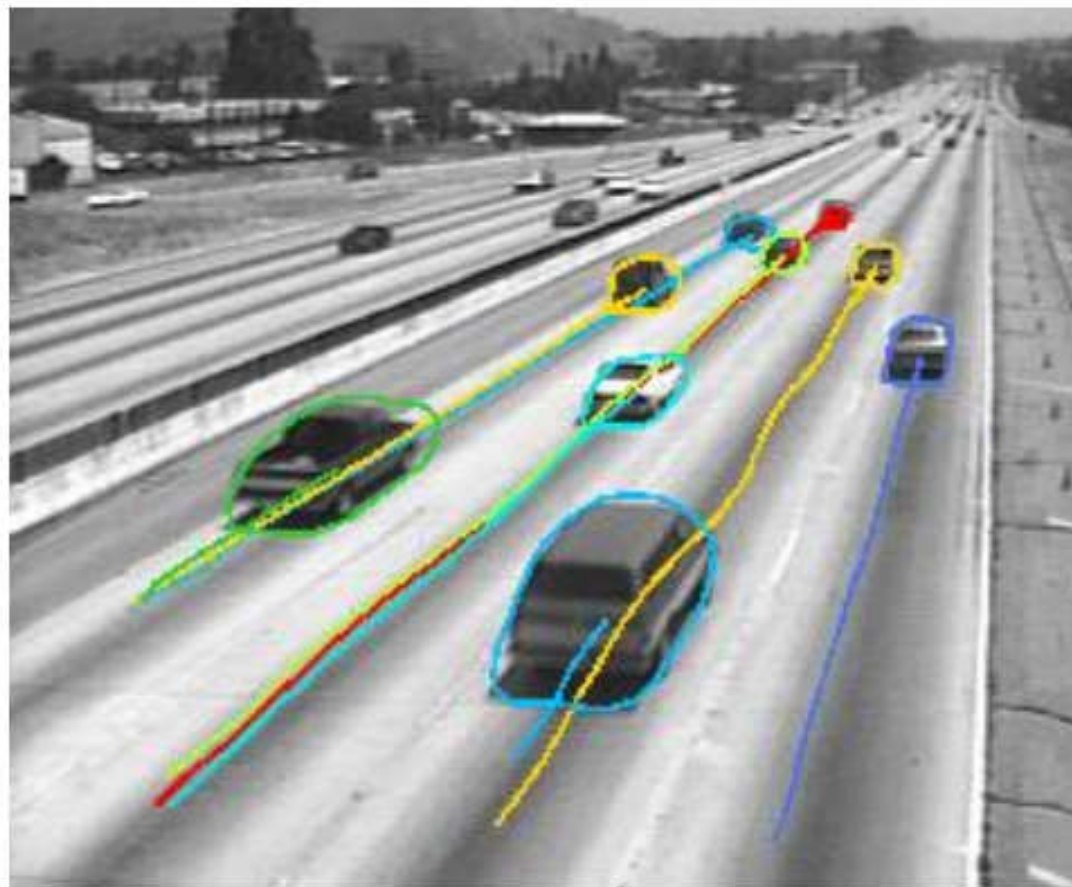




# Optical flow can help track features



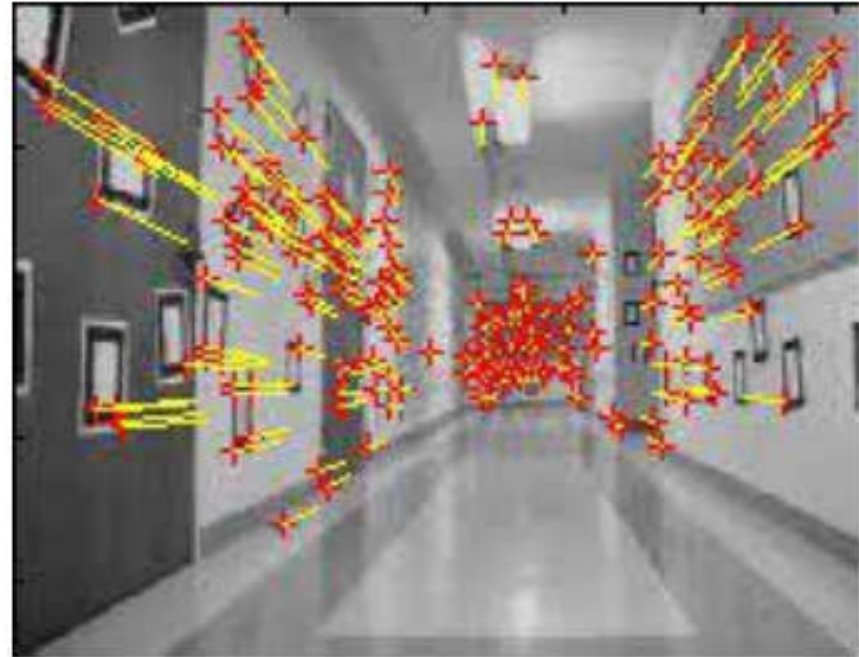
- Once we have the features we want to track, Lucas-Kanade or other optical flow algorithms can help track those features





# Feature-tracking

---





# Outline

---

- Feature tracking
- Simple KLT tracker  
(Kanade–Lucas–Tomasi feature tracker)
  - Pipeline
  - Results

# Simple KLT tracker

---



- 1. Find a good point to track (Harris corner)
- 2. For each Harris corner compute motion(translation or affine) between consecutive frames.
- 3. Link motion vectors in successive frames to get a track for each Harris point
- 4. Introduce new Harris points by applying Harris detector at every  $m$  (10 or 15) frames.
- 5. Track new and old Harris points using steps 1-3.

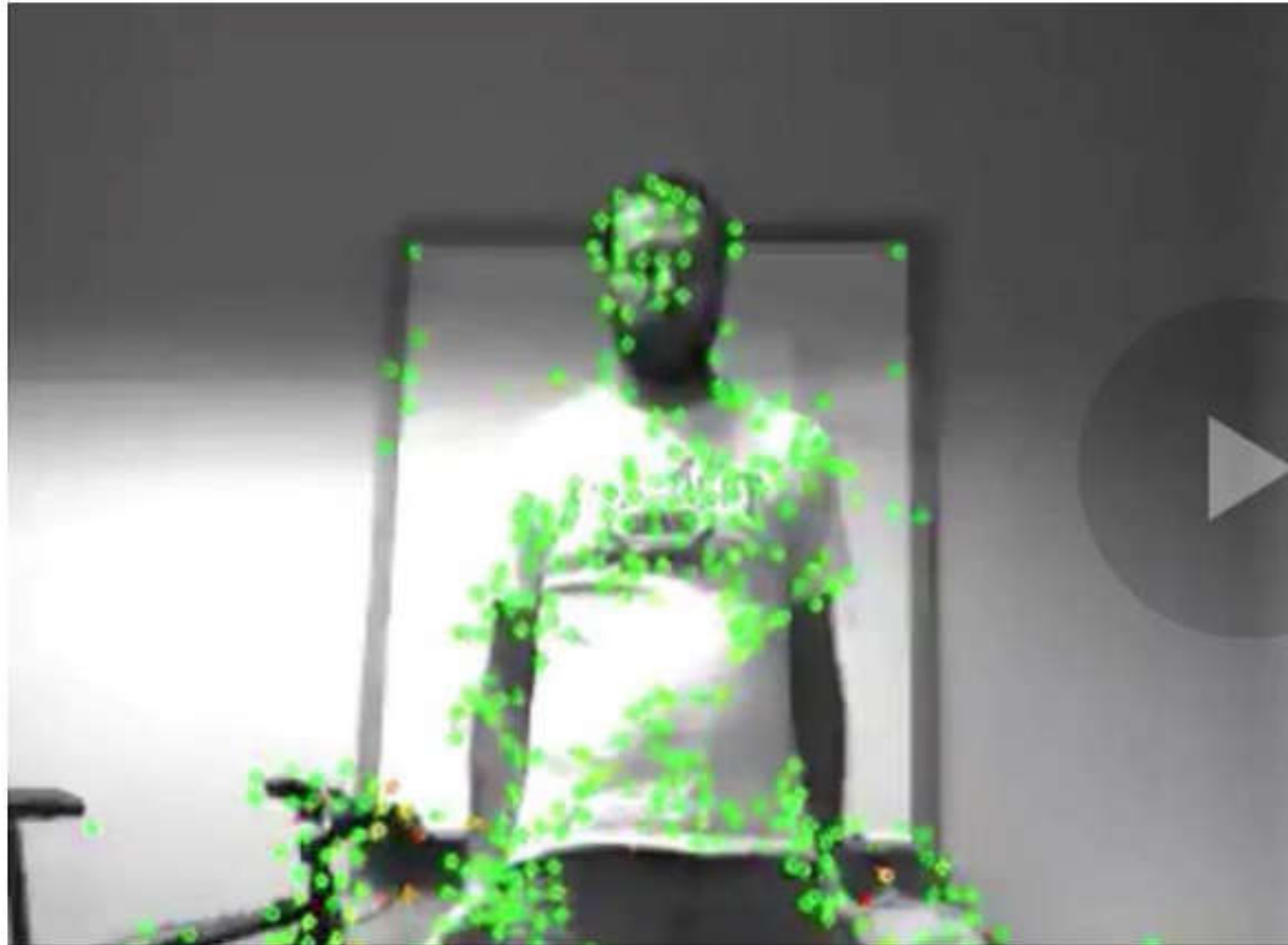
# Tracking cars

---



# Tracking movement

---



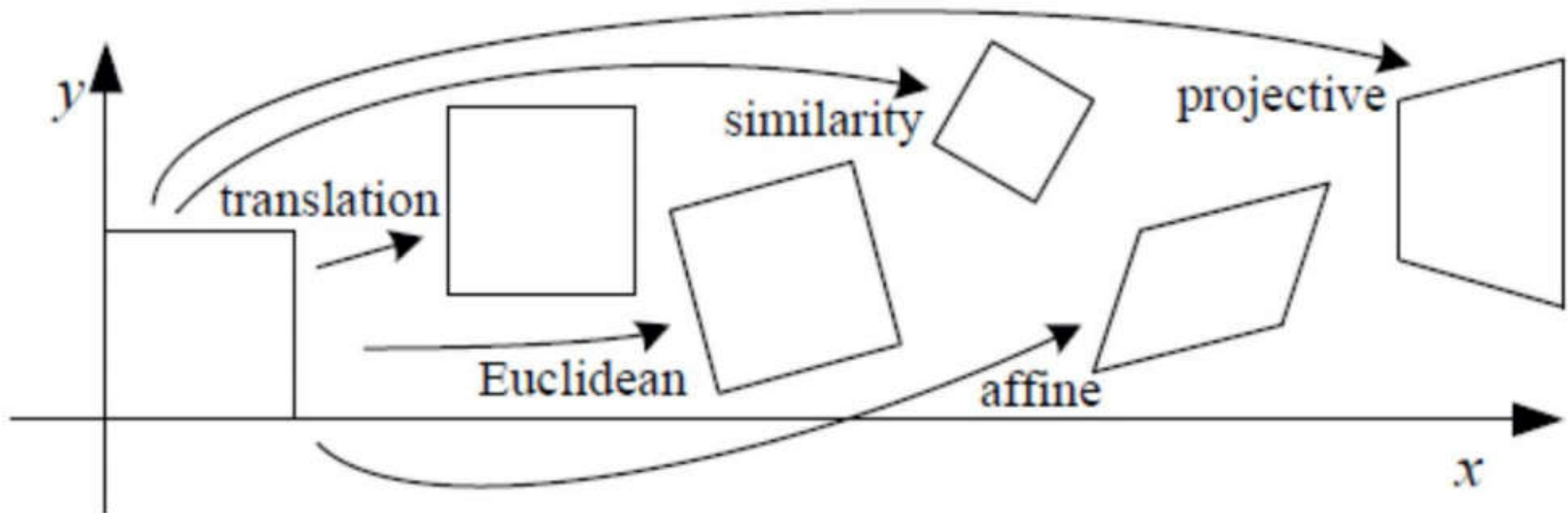
# Outline

---



- Feature tracking
- Simple KLT tracker
- 2D transformations (recap)

# Types of 2D transformations





# Translation

- Let the initial feature be located by  $(x, y)$ .
- In the next frame, it has translated to  $(x', y')$ .
- We can write the transformation as:

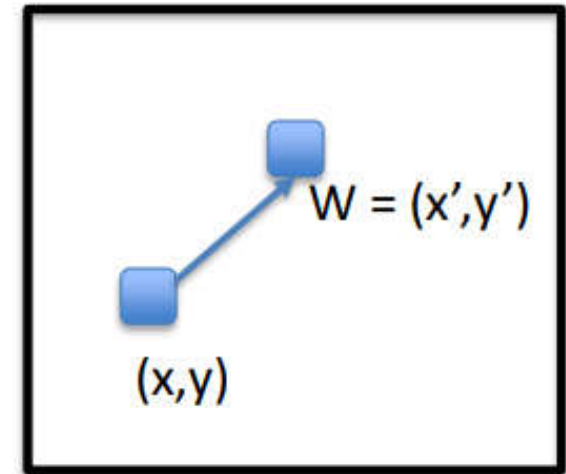
$$x' = x + b_1$$

$$y' = y + b_2$$

- We can write this as a matrix transformation using homogeneous coordinates:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} 1 & 0 & b_1 \\ 0 & 1 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$W(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} 1 & 0 & b_1 \\ 0 & 1 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



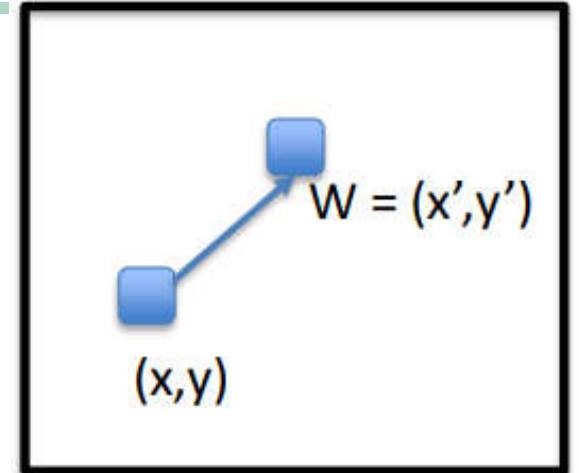


# Displacement Model for Translation



$$W(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} 1 & 0 & b_1 \\ 0 & 1 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

- There are only two parameters:  $\mathbf{p} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$



- The derivative of the transformation w.r.t.  $\mathbf{p}$  :

$$\frac{\partial W}{\partial \mathbf{p}}(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- This is called the Jacobian.



# Similarity motion

---

- Rigid motion includes scaling + translation.
- We can write the transformations as:

$$x' = ax + b_1$$

$$y' = ay + b_2$$

$$W(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} a & 0 & b_1 \\ 0 & a & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad \mathbf{p} = [a \quad b_1 \quad b_2]^T$$

$$\frac{\partial W}{\partial \mathbf{p}}(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} x & 1 & 0 \\ y & 0 & 1 \end{bmatrix}$$

# Affine motion includes scaling + rotation + translation

---



- Affine motion includes scaling + rotation + translation.

$$x' = a_1x + a_2y + b_1$$

$$y' = a_3x + a_4y + b_2$$

$$W(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} a_1 & a_2 & b_1 \\ a_3 & a_4 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad \mathbf{p} = [a_1 \quad a_2 \quad b_1 \quad a_3 \quad a_4 \quad b_2]^T$$

$$\frac{\partial W}{\partial \mathbf{p}}(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} x & y & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & x & y & 1 \end{bmatrix}$$



# Outline

---

- Feature tracking
- Simple KLT tracker
- 2D transformations (recap)
- **Iterative KLT tracker**



# Problem setting

---

- Given a video sequence, find all the features and track them across the video.
- First, use Harris corner detection to find features and their location  $x$ .
- For each feature at location  $x=[x \ y]^T$ :
  - Create an initial template for that feature:  $T(x)$ .
  - $T(x)$  is usually an image patch around  $x$ .
- Goal: find new location of feature  $x$  at the next frame.
- We will assume  $x$  undergoes a transformation (translation, affine, ...) parametrized by  $p$  to reach its new location  $W(x;p)$ .



# KLT objective

---

- Our aim is to find the  $\mathbf{p}$  that minimizes the difference between the template  $T(\mathbf{x})$  and the image region around the new location of  $\mathbf{x}$  after undergoing the transformation.

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

- $W(\mathbf{x}; \mathbf{p})$  is the new location of feature  $\mathbf{x}$ .
- $I(W(\mathbf{x}; \mathbf{p}))$  is image intensity at the new location.
- Recall that  $\mathbf{p}$  is our vector of parameters that define the transformation that took  $\mathbf{x}$  to its new location  $W(\mathbf{x}; \mathbf{p})$ .
- Sum is over an image patch around  $\mathbf{x}$ .



# KLT objective

---

- Since  $\mathbf{p}$  may be large, minimizing this function may be difficult:

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

- We will instead break down  $\mathbf{p} = \mathbf{p}_0 + \Delta\mathbf{p}$ 
  - Large + small/residual motion
  - Where  $\mathbf{p}_0$  is going to be fixed and we will solve for  $\Delta\mathbf{p}$ , which is a small value.
  - We can initialize  $\mathbf{p}_0$  with our best guess of what the motion is, and then calculate  $\Delta\mathbf{p}$ .
- We can substitute  $\mathbf{p}$  to get:

$$\sum_x [I(W(\mathbf{x}; \mathbf{p}_0 + \Delta\mathbf{p})) - T(\mathbf{x})]^2$$

# Taylor series

---



- Taylor series is defined as

$$f(x + \Delta x) = f(x) + \Delta x \frac{\partial f}{\partial x} + \Delta x^2 \frac{\partial^2 f}{\partial x^2} + \dots$$

- Assuming that  $\Delta x$  is small.
- We can apply this expansion to the KLT tracker and only use the first two terms:





# Expanded KLT objective

---

$$\sum_x [I(W(\mathbf{x}; \mathbf{p}_0 + \Delta \mathbf{p})) - T(\mathbf{x})]^2$$
$$\approx \sum_x \left[ I(W(\mathbf{x}; \mathbf{p}_0)) + \nabla I \frac{\partial W}{\partial \mathbf{p}} \Delta \mathbf{p} - T(\mathbf{x}) \right]^2$$

- It's a good thing we have already calculated what  $\frac{\partial W}{\partial \mathbf{p}}$  would look like for affine, translations and other transformations!



# Expanded KLT objective

---

- So our aim is to find the  $\Delta \mathbf{p}$  that minimizes the following:

$$\operatorname{argmin}_{\Delta \mathbf{p}} \sum_{\mathbf{x}} \left[ I(W(\mathbf{x}; \mathbf{p}_0)) + \nabla I \frac{\partial W}{\partial \mathbf{p}} \Delta \mathbf{p} - T(\mathbf{x}) \right]^2$$

- Where  $\nabla I = [I_x \ I_y]$
- Differentiate wrt  $\Delta \mathbf{p}$  and setting it to zero

$$\sum_{\mathbf{x}} \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T \left[ I(W(\mathbf{x}; \mathbf{p}_0)) + \nabla I \frac{\partial W}{\partial \mathbf{p}} \Delta \mathbf{p} - T(\mathbf{x}) \right] = 0$$



# Solving for $\Delta \mathbf{p}$

---

- Solving for  $\Delta \mathbf{p}$  in

$$\sum_x \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T \left[ I(W(\mathbf{x}; \mathbf{p}_0)) + \nabla I \frac{\partial W}{\partial \mathbf{p}} \Delta \mathbf{p} - T(\mathbf{x}) \right] = 0$$

- We get

$$\Delta \mathbf{p} = H^{-1} \sum_x \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T [T(\mathbf{x}) - I(W(\mathbf{x}; \mathbf{p}_0))]$$

- where

$$H = \sum_x \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]$$

- **$H$**  must be invertible!

# Interpreting the H matrix for translation transformations



$$\mathbf{H} = \sum_x \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]$$

Recall that

- 1.  $\nabla I = [I_x \ I_y]$  and
- 2. for translation motion,  $\frac{\partial W}{\partial \mathbf{p}}(\mathbf{x}; \mathbf{p}) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Therefore,

$$\begin{aligned} \mathbf{H} &= \sum_x \left[ [I_x \ I_y] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right]^T \left[ [I_x \ I_y] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right] \\ &= \sum_x \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \end{aligned}$$

That's the matrix from the Harris corner detector we learnt in class! <sup>b6</sup>

# Interpreting the H matrix for affine transformations



$$H = \sum_{\mathbf{x}} \begin{bmatrix} I_x^2 & I_x I_y & x I_x^2 & y I_x I_y & x I_x I_y & y I_x I_y \\ I_x I_y & I_y^2 & x I_x I_y & y I_y^2 & x I_y^2 & y I_y^2 \\ x I_x^2 & y I_x I_y & x^2 I_x^2 & y^2 I_x I_y & xy I_x I_y & y^2 I_x I_y \\ y I_x I_y & y I_y^2 & xy I_x I_y & y^2 I_y^2 & xy I_y^2 & y^2 I_y^2 \\ x I_x I_y & x I_y^2 & x^2 I_x I_y & xy I_y^2 & x^2 I_y^2 & xy I_y^2 \\ y I_x I_y & y I_y^2 & xy I_x I_y & y^2 I_y^2 & xy I_y^2 & y^2 I_y^2 \end{bmatrix}$$



# Overall KLT $\Delta \mathbf{p} = H^{-1} \sum_x \left[ \nabla I \frac{\partial W}{\partial \mathbf{p}} \right]^T [T(\mathbf{x}) - I(W(\mathbf{x}; \mathbf{p}_0))]$

Given the features from Harris detector:

- 1. Initialize  $\mathbf{p}_0$ .
- 2. Compute the initial templates  $T(\mathbf{x})$  for each feature.
- 3. Transform the features in the image  $I$  with  $W(\mathbf{x}; \mathbf{p}_0)$ .
- 4. Measure the error:  $I(W(\mathbf{x}; \mathbf{p}_0)) - T(\mathbf{x})$ .
- 5. Compute the image gradients  $\nabla I = [I_x \ I_y]$ .
- 6. Evaluate the Jacobian  $\frac{\partial W}{\partial \mathbf{p}}$ .
- 7. Compute steepest descent  $\nabla I \frac{\partial W}{\partial \mathbf{p}}$ .
- 8. Compute Inverse Hessian  $H^{-1}$
- 9. Calculate the change in parameters  $\Delta \mathbf{p}$
- 10. Update parameters  $\mathbf{p}_0 = \mathbf{p}_0 + \Delta \mathbf{p}$
- 11. Repeat 3 to 10 until  $\Delta \mathbf{p}$  is small.



# KLT over multiple frames

---

- Once you find a transformation for two frames, you will repeat this process for every couple of frames.
- Run Harris detector every 15-20 frames to find new features.



# Challenges to consider

---

## Implementation issues

- Window size (size of neighborhood/template around  $x$ )
  - Small window more sensitive to noise and may miss larger motions (without pyramid)
  - Large window more likely to cross an occlusion boundary (and it's slower)
  - 15x15 to 31x31 seems typical
- Weighting the window
  - Common to apply weights so that center matters more (e.g., with Gaussian)



# References

---



- Basic reading:
  - Szeliski textbook, Section 9.3, 9.4