

# TSBB15 Computer Vision

#### Lecture 8 Local features

#### Per-Erik Forssén



# Today's lecture

- What are local features used for?
- The local (invariant) features paradigm
- Invariances: Geometric, Photometric
- Examples: SIFT, MSER/MSCR...
- Feature matching



#### What are local features used for?

KLT tracking and block matching are useful when matching between consecutive frames in a video sequence.



## What are local features used for?

- KLT tracking and block matching are useful when matching between consecutive frames in a video sequence.
- Images are from the same camera
- small changes in scale, rotation and illumination



# What are local features used for?

- KLT tracking and block matching are useful when matching between consecutive frames in a video sequence.
- Images are from the same camera
- small changes in scale, rotation and illumination

# Local invariant features work when these conditions are violated.



### Wide-baseline stereo

Problem 1: wide-baseline stereo

 Matching images of the same scene, captured at different positions.





## Wide-baseline stereo

Problem 1: wide-baseline stereo

 Matching images of the same scene, captured at different positions.





# Object instance recognition and pose estimation

- Problem 2: bin picking
  - identity and pose estimation under partial occlusion
  - training set
  - test set
  - 6dof pose





February 9, 2022





# Object recognition

#### • Example: Eddie the embodied



#### See webpage for details https://www.cvl.isy.liu.se/research/objrec/EVOR/



• In lecture 2 we discussed how to match across scale and translation. How?



• In lecture 2 we discussed how to match across scale and translation. How?





Example: face detection



- In lecture 2 we discussed how to match across scale and translation. How?
- Another option is to use interest points e.g. Harris points [Z. Zhang et al. 95]:
  - 1. Detect interest points
  - 2. Cut out image patches around each point
  - 3. Matches can now be found by comparing patches+epipolar geometry constraints.



 Correspondences from block matching at Harris points (assignment problem:LE7).





• After applying the Epipolar constraint (You will test this in CE3).





• The epipolar constraint:  $\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$ 







- The epipolar constraint:  $\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$
- x<sub>1</sub> and x<sub>2</sub> are projections of the same 3D point in two views.
- Scene is static, i.e. no motion has taken place (except the change of camera position).
- F can be estimated from 7 or more correspondences. E.g. 8-pt algorithm.



- The epipolar constraint:  $\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$
- See the compendium, Introduction to Representations and Estimation in Geometry (IREG), Klas Nordberg



- Zhang's interest point method. (repeat)
  - 1. Detect interest points
  - 2. Cut out image patches around each point
  - 3. Find matches, by comparing patch descriptors and epipolar geometry constraints.



 Zhang's method is invariant to translation (and partially to scale).

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = s \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \mathbf{t}$$

 2 degrees-of-freedom (DOF) of invariance (transl. only) (3 if scale is also counted)



• Zhang's method is invariant to translation (and partially to scale).

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = s \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \mathbf{t}$$

- 2 degrees-of-freedom (DOF) of invariance (transl. only) (3 if scale is also counted)
- We will now add invariance to image rotations and view changes.



- In general, the *local invariant feature approach* can be described as three steps:
  - Detection: Use a *detector* to find a local, canonical frame (a coordinate system)



- In general, the *local invariant feature approach* can be described as three steps:
  - Detection: Use a *detector* to find a local, canonical frame (a coordinate system)
  - Description: Compute a *descriptor*, by sampling the image in the canonical frame



- In general, the *local invariant feature approach* can be described as three steps:
  - Detection: Use a *detector* to find a local, canonical frame (a coordinate system)
  - Description: Compute a *descriptor*, by sampling the image in the canonical frame
  - Matching: Find correspondences, by comparing descriptors from two images



# Canonical frame example

Resampling to canonical frame





#### Geometric invariances

Robustness to view changes





Photometric invariances

Robustness to illumination changes





- Geometric invariances can be obtained by choosing a frame that is equivariant to rotations, scalings, and image skews
- Photometric invariances can be obtained by computing the descriptor in a more advanced way than direct sampling.



- The geometric invariances used in local features make a locally planar assumption.
- They can thus be described using homographies (See IREG, TSBB06).



- Recap: A Homography is a transformation between points **x** on one plane, and points **y** on another.  $(y_1)$   $(x_1)$   $(h_{11}$   $h_{12}$   $h_{13})$   $(x_1)$ 

$$\lambda \begin{pmatrix} y_1 \\ y_2 \\ 1 \end{pmatrix} = \mathbf{H} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} = \begin{pmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$



- Recap: A **Homography** is a transformation between points **x** on one plane, and points **y** on another.  $(y_1)$   $(x_1)$   $(h_{11}$   $h_{12}$   $h_{13})$   $(x_1)$ 

$$\lambda \begin{pmatrix} y_1 \\ y_2 \\ 1 \end{pmatrix} = \mathbf{H} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} = \begin{pmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

- Degrees of freedom: minimal number of parameters needed in H.
  - at most 8dof (for plane projective case), as **H** and  $k\mathbf{H}$ ,  $k \in \mathbb{R} \setminus 0$  give the same output



- A hierarchy of transformations:
   scale+translation (3dof)
  - similarity (4dof) (scale+translation+rotation)
  - affine (6dof)(similarity+skew)
  - plane projective (8dof) (affine+foreshortening)





- We can find the canonical frame by using more than one point [Brown&Lowe 02] aka. *interest-point* groups
- We will now give some examples...



 Scale+translation: Useful if we know that there is no rotation. E.g. for a camera mounted in a car, looking at upright pedestrians.





 Similarity: Full invariance in image plane, none outside image plane.
 Useful e.g. for pose estimation.





#### Affine: Deals with most common projective distortions. Good if patch size is small relative to distance to patch.





 Plane projective: Full modelling of a plane in 3D. Requires more image measurements, but is better for extreme view angles.




#### **Geometric Invariance**

- Resampling to canonical frame
   results in
  - geometric invariance:





# **Geometric Invariance**

Problems with interest-point groups:
 – Sensitive to missing points:

If e=P(point-detected|present) then

#### P(frame-is-detected|present)=e<sup>N</sup>

where N is number of points in frame.



# **Geometric Invariance**

- Problems with interest-point groups:
  - Sensitive to missing points:
     If e=P(point-detected|present) then
     P(frame-is-detected|present)=e<sup>N</sup>
     where N is number of points in frame.
  - Combinatorics: if K points in image, we have  $\binom{N}{K}$  possible canonical frames.
- We will introduce other ways to find the frame soon.



- Image intensity is linear in radiance
- the sensor activation,  $a(\mathbf{x})$ :

$$a(\mathbf{x}) = \int s(\lambda) e(\lambda) d\lambda$$

- $s(\lambda)$  sensor absorption spectrum
- $e(\lambda)$  spectrum of incoming light (attenuated by the atmosphere)



- Image intensity is linear in radiance
- the sensor activation,  $a(\mathbf{x})$ :

$$a(\mathbf{x}) = \int s(\lambda) e(\lambda) d\lambda$$

- $s(\lambda)$  sensor absorption spectrum
- $e(\lambda)$  spectrum of incoming light (attenuated by the atmosphere)
- Adding a second, identical light source will double  $e(\lambda)$ , and thus also  $a(\mathbf{x})$ .
- After gamma correction this is perfect linearity is broken.
   But we still have approximate linearity.



• If illumination changes, image matching fails:

$$\begin{aligned}
I(\mathbf{x}) &= I_0(\mathbf{x})k_1 \\
J(\mathbf{x}) &= I_0(\mathbf{x})k_2 \end{aligned} \Rightarrow \sum_{x \in \Omega} (I(\mathbf{x}) - J(\mathbf{x}))^2 = \text{non-zero}
\end{aligned}$$

- We want a function that is invariant to scalings:  $\sum_{x \in \Omega} (f(I(\mathbf{x})) - f(J(\mathbf{x})))^2 = \text{small number}$
- How should we choose the invariant f()?



• For cameras with gamma correction, or if two different cameras are used we may use the affine model:

$$I(\mathbf{x}) = I_0(\mathbf{x})k_1 + k_2$$

• How should we choose *f* ()? we want:

$$\sum_{x \in \Omega} (f(I(\mathbf{x})) - f(J(\mathbf{x})))^2 = \text{small number}$$



- Invariance to intensity offsets:
   Mean subtraction, and any DC free linear filters, e.g. derivatives.
- Scaling invariance:
  - Normalising a patch by an  $L_p$ -norm, e.g. the  $L_2$ -norm or the standard deviation
- Affine invariance by combining both:

$$\hat{I}(\mathbf{x}) = (I(\mathbf{x}) - \mu_I) / \sigma_I$$







### Local Invariant Features

- There are many examples of features that fit the descriptor+detector paradigm.
- The two most widely used are:
  - SIFT Scale Invariant Feature Transform (Lowe 99)
  - MSER Maximally Stable Extremal Regions (Matas et al. 02)
- We will look at these two in more detail.



# SIFT

- Scale Invariant Feature Transform [Lowe'99]. In brief:
  - The SIFT detector finds points using Difference-of-Gaussians in a pyramid Gives: position x,y and scale s
  - Rotation is found from a gradient histogram
  - This gives a frame for the SIFT descriptor, which is computed from gradient orientation histograms.



- Scale space (recap. from LE2)
  - The image is extended with an extra dimension for scale/blur:

$$f(x, y, s) = (f_0 * g(s))(x, y)$$

– The blurring kernel g(s) is typically a Gaussian:

$$g(\mathbf{x},s) = \frac{1}{2\pi s} e^{-\mathbf{x}^T \mathbf{x}/2s^2}$$



- Scale selection [Lindeberg'93]
  - Find a characteristic point (e.g. local max) on a function of position and scale:

$$(\hat{\mathbf{x}}, \hat{s}) = rg \max h(f(\mathbf{x}, s))$$

– Example: Maximum of normalised Laplacian:

$$h(f(\mathbf{x},s)) = s^2(f * \nabla^2 g(s))(\mathbf{x})$$







 In SIFT, scale selection is done using difference-of-Gaussians:

 $h_{\text{SIFT}}(f(\mathbf{x},\sigma)) = (f * (g(\sigma) - g(k\sigma)))(\mathbf{x})$ 

- Efficient implementation using pyramids [Lowe'99]
- Sampling in scale space with  $\Delta \sigma = 1/\sqrt{2}$









Non-max suppression in (x,y,s)



 Finally we find one or more reference directions using a gradient orientation histogram *h* at the found location in scale

space.





- The SIFT detector gives us a similarity frame. What is this?
  - We now want to convert the image patch at the frame to a 128-byte *descriptor vector*.
  - The purpose of this is to add photometric invariance, and some extra translation and scale robustness.



Compute x- and y-gradients through convolution:

$$\nabla \mathbf{f}(\mathbf{x}) = \begin{bmatrix} (d_x * f)(\mathbf{x}) \\ (d_y * f)(\mathbf{x}) \end{bmatrix}$$

- Rotate gradient map to direction from orient-hist:  $\nabla \hat{\mathbf{f}}(\mathbf{x}) = \mathbf{R} \nabla \mathbf{f}(\mathbf{R^T} \mathbf{x})$
- Compute gradient orientation histograms in 4x4 spatial regions:

$$h_{kl} = \sum_{\mathbf{x} \in \text{patch}_l} |\nabla \hat{\mathbf{f}}(\mathbf{x})| w(\mathbf{x} + \mathbf{d}_l) B_k(\tan^{-1} \nabla \hat{\mathbf{f}}(\mathbf{x}))$$



Compute gradient orientation histograms in 4x4 spatial regions :

$$h_{kl} = \sum_{\mathbf{x} \in \text{patch}_l} |\nabla \hat{\mathbf{f}}(\mathbf{x})| w(\mathbf{x} + \mathbf{d}_l) B_k(\tan^{-1} \nabla \hat{\mathbf{f}}(\mathbf{x}))$$

- B<sub>k</sub>(x) linear interpolation kernel
   Quadratic is better (Jonsson&Felsberg)
- Subwindows  $l \in [1 \dots 16]$  directions  $k \in [1 \dots 8]$
- Spatial weight  $w(\mathbf{x} + \mathbf{d}_l)$  (Gaussian decay)



 Implementation with source code in both VLfeat and OpenCV.



Note that 4x4 regions are actually used, with 8 orientations -> 128 elements



- Affine illumination invariance by using gradients and normalising descriptor  $\hat{\mathbf{h}} = \mathbf{h}/||\mathbf{h}||$
- Some robustness by truncating and normalising again  $\hat{\hat{\mathbf{h}}} = \min(\mathbf{t}, \hat{\mathbf{h}}) / \|\min(\mathbf{t}, \hat{\mathbf{h}})\|$
- The spatial histogramming gives robustness to scale/rotation/translation errors.



- Affine illumination invariance by using gradients and normalising descriptor  $\hat{\mathbf{h}} = \mathbf{h}/||\mathbf{h}||$
- Some robustness by truncating and normalising again  $\hat{\hat{\mathbf{h}}} = \min(\mathbf{t}, \mathbf{\hat{h}}) / \|\min(\mathbf{t}, \mathbf{\hat{h}})\|$
- The spatial histogramming gives robustness to scale/rotation/translation errors.
- SIFT is used commercially in many places. (The Sony AIBO anno 1999, was an early example.) Patent has now expired.





- Maximally Stable Extremal Regions
  [Matas et al.'02]
- Consider the set of all possible thresholdings of an image...

[Movie clip]







- Connected regions form segments.
  - Cf. Watershed algorithm (similar idea but different output)
  - Look at stability of a function of segment across image evolution. e.g.

area(component(t))

 MSERs are components that are maximally stable, i.e., have a local minimum of the rate of change: ∂area(component(t))



- compare: Maximal Stability, Scale Selection
- Stability measure: Range of stable thresholds t<sub>2</sub>-t<sub>1</sub> around min is called the margin of the region.



– Two possible thresholdings:  $I(\mathbf{x}) < t$  ,  $I(\mathbf{x}) > t$ 



Input image

64 MSER- (total 272) 64 MSER+ (total 294)

Very fast (using union/find+path compression).
MSER type (+/-) is useful for matching How?



- MSER is invariant to monotonic changes of intensity. i.e. I(x) and f(I(x)) have the same output if  $f(x+k) > f(x) \forall k > 0$
- Wide range of sizes obtained without a scale pyramid.
   Better still with a pyramid (Forssén&Lowe ICCV'07)
- Colour objects can be tracked by computing MSERs on the Mahalanobis distance to a colour distribution. (Donoser&Bischof CVPR'06)
- Colour regions by looking at gradients.
   Called MSCR (Forssén CVPR'07)





#### **MSCR**







#### MSCR





Reference directions from extremal points
 along ellipse-normalized contour.





• Approximating ellipse

– from moments of binary mask  $v:\Omega\mapsto\{0,1\}$ 

$$\mu_{k,l}(v) = \sum_{x} \sum_{y} x^{k} y^{l} v(x, y)$$
$$\mathbf{m} = \frac{1}{\mu_{0,0}} \begin{bmatrix} \mu_{1,0} \\ \mu_{0,1} \end{bmatrix} \quad \mathbf{C} = \frac{1}{\mu_{0,0}} \begin{bmatrix} \mu_{2,0} & \mu_{1,1} \\ \mu_{1,1} & \mu_{0,2} \end{bmatrix} - \mathbf{m} \mathbf{m}^{T}$$
$$\mathcal{R}(\mathbf{m}, \mathbf{C}) = \{ \mathbf{x} : (\mathbf{x} - \mathbf{m})^{T} \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}) \le 4 \}$$



Normalisation to a circle (axis aligned)
 Compute the eigenfactorisation:

 $\mathbf{C} = \mathbf{R} \mathbf{D} \mathbf{R}^T, \quad \det \mathbf{R} > 0$ 

The circle normalisation can now be performed as:

$$\mathbf{x} = \mathbf{A}\hat{\mathbf{x}} + \mathbf{m}$$
, for  $\mathbf{A} = 2\mathbf{R}\mathbf{D}^{1/2}$ 

- $\hat{\mathbf{x}}$  canonical coordinates
- x image coordinates



- Ellipse+extrema of distance to centre is just one frame construction option.
- Other (affine covariant) choices:
  - Points of maximum curvature.
  - Bi-tangens.
  - See Obdrzalek&Matas BMVC'02
- Implementation w. source: in both VLfeat and OpenCV



# MSER descriptor

• The MSER detector originally used normalized colour patches as descriptor vectors:

$$\hat{I}_r(\mathbf{x}) = (I_r(\mathbf{x}) - \mu_r) / \sigma_r$$
$$\hat{I}_g(\mathbf{x}) = (I_g(\mathbf{x}) - \mu_g) / \sigma_g$$
$$\hat{I}_b(\mathbf{x}) = (I_b(\mathbf{x}) - \mu_b) / \sigma_b$$

• Nowadays other descriptors, e.g. the SIFT descriptor are used.


# Other local invariant features

• SFOP

http://www.ipb.uni-bonn.de/sfop/

- BRISK
   Source Code+description
   <a href="http://www.asl.ethz.ch/people/lestefan/personal/BRISK">http://www.asl.ethz.ch/people/lestefan/personal/BRISK</a>
- FREAK, ORB
   In OpenCV features2d
- SURF and SIFT

in OpenCV xfeatures2d (in LiU installation)



## Binary descriptors

• To save memory and time, many descriptors use local binary patterns:



Image from Alexandre et al. CVPR 2012 10110
 sign of intensity difference has monotonic illumination invariance



## Binary descriptors

• To save memory and time, many descriptors use local binary patterns:



#### E.g. BRIEF (ECCV'10), BRISK (ICCV'11), ORB (ICCV'11), FREAK (CVPR'12)



# Deep learning descriptors

Examples:

- DeCAF (ArXiv'13) descriptors
- TILDE (CVPR'15) detector
- LIFT (ECCV'16) detector and descriptor
- SuperPoint (CVPRw'18) detector + descriptor
- LF-Net (NIPS'18) detector+descriptor

# Better matching performance at the price of more expensive computations.



#### A note on invariance

Always strive to limit amount of invariance

- For hand-coded features: use knowledge on imaging situation
  - e.g. a car mounted camera may not need rotation invariance for pedestrians.
  - e.g. in a video with smooth illumination changes, affine illumination invariance is not necessary
- Learned local features do this based on the training set
  - Knowing the training set is important!



- The *Local Invariant Feature* method:
- Detection
- Description
- Matching



For a descriptor *q* in a query image. Which prototype in memory (*p*<sub>1</sub>,*p*<sub>2</sub>,...,*p*<sub>N</sub>) is most likely to correspond to the same world object?



- For a descriptor *q* in a query image. Which prototype in memory (*p*<sub>1</sub>,*p*<sub>2</sub>,...,*p*<sub>N</sub>) is most likely to correspond to the same world object?
- Assuming additive i.i.d. Gaussian noise on all elements:  $D = 5(n_1 - n_2)^2/\sigma^2$

$$p(\mathbf{q}|\mathbf{p}_k) \propto \prod_{l=1}^{D} e^{-.5(p_{kl}-q_l)^2/\sigma}$$
$$\max(J) \Leftrightarrow \min(-\log(J))$$
$$-\log(p(\mathbf{q}|\mathbf{p}_k)) \propto \sum_{l=1}^{D} (p_{kl}-q_l)^2$$



- So, the match with smallest distance is most likely correct, assuming i.i.d. Gaussian noise.
- What about the scalar product for normalised vectors/NCC?



- So, the match with smallest distance is most likely correct, assuming i.i.d. Gaussian noise.
- What about the scalar product for normalised vectors/NCC?

$$||\mathbf{p} - \mathbf{q}||^2 = \mathbf{p}^T \mathbf{p} + \mathbf{q}^T \mathbf{q} - 2\mathbf{p}^T \mathbf{q} = 2(1 - \mathbf{p}^T \mathbf{q})$$

But are all values identically distributed?
...are they all independent?



 For binary descriptors (e.g. BRIEF) the Hamming distance is used:

s = bitcnt(XOR(P,Q))

- Also makes i.i.d. assumption.
- Binomial distributed s~Bin(n,p)



#### **Distance** ratio

#### **Risk of mismatch**

can also be taken into account by looking at the ratio of the best and second best match.



$$p(r|\text{correct}) \text{ and } p(r|\text{incorrect})$$

$$r = d_{\min}/d_{\text{second\_smallest}}$$



#### **Distance** ratio

#### If we have a set of matches for descriptors $q_1$ and $q_2$ in the image. Which one is better?





#### Dense invariant features

- (semi-)dense flow for wide baseline problems can be obtained by matching invariant features
- at every pixel and at several scales
- e.g. SIFTflow, DSIFT, PHOW, DAISY
- Expensive to compute, unless GPGPU is used.



## Summary

- Use local invariant features: when KLT fails
- But use no more invariance than needed
- Two types of invariance: Photometric
   and Geometric invariance
- Recognition in three steps: Detection, Description and Matching



#### Upcoming course events

- CE2: Tomorrow 13-17 in Asgård. Checkup for CE1 at 13.00.
- Next Lecture (16/2, 10-12) Biological vision. Voluntary. Based on PhD course on Biological Vision Systems.