

## SIFT (Scale Keypoint descriptor)

~~Autonomous Intelligent~~

- ① Take a  $16 \times 16$  window around the detected point.
- ② Compute the edge orientation for each pixel.
- ③ Divide into a  $4 \times 4$  grid of cells.
- ④ Compute the histogram in each cell.

### Orientation assignment (for rotation invariance)

- ① Create histogram of gradient directions, within a region around a keypoint, at selected scale.

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y)$$

magnitude  $\rightarrow m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2}$

angle / phase  $\rightarrow \theta(x, y) = \tan^{-1} \frac{(L(x, y+1) - L(x, y-1))}{(L(x+1, y) - L(x-1, y))}$

The reason for orientation assignment is to compute the main direction among the pixels. For this, we will build a histogram (from  $0-2\pi$ )

→ Histogram entries are weighted (i) gradient magnitude, (ii) a Gaussian function with  $\sigma$  equal to 1.5 times the scale of the keypoint (max at centre) + more weight.

→ We can calculate the keypoint by visualising the histogram bins. We set a threshold of  $\approx 80\%$  and we can have a main edge orientation of the point.

If there are multiple orientations in the point of interest, we will generate duplicate keypoints. limit to 2 peak.

→ To create a local histogram (Keypoint descriptors)

- ① Each histogram entry is weighted by
  - (i) a gradient magnitude
  - (ii) a Gaussian function with  $\sigma$  equal to 0.5 times the width of the descriptor window.

→ matching SIFT features :-

$$\text{SSD}(f_1, f_2) / \text{SSD}(f_1, f'_2) < t \quad \text{if } t = 0.8$$

### SIFT main steps

- ① Scale-space extrema detection → Extract scale and rotation invariant interest points (i.e. Keypoints)
- ② Keypoint localization → Determine location and scale for each interest point  
→ Eliminate weak keypoints.
- ③ Orientation assignment → Assign one or more orientation to each Keypoint
- ④ Keypoint descriptor → Use local image gradients at the selected scale.

## Feature matching, Tracking and optical flow.

→ **Feature matching** → If we have 2 images and we have to match the features we follow a simple approach.

- ✓ Step 1 → Detect corners in the image with for eg. Harris detection.
- ✓ Step 2 → Compute similarity / correlation scores. Take a point in the original image and find a point at the same location or in a proximity with the help of disparity window. With the help of all the points, draw a correlation curve on the intensity neighborhoods.
- ✓ Step 3 → Select the point with the highest score.

→ usually during template matching, many wrong matches (10-50%) are also done. But even with this error rate, we can start computation "motion". All thanks to outliers rejection techniques such as RANSAC.

→ With the help of invariant descriptors, we can improve the "similarity measure" of features.

invariant (to light, scale, viewpoint)  
features ( $\text{transform}(\text{image})$ ) = feature ( $\text{image}$ )

→ **Scale Invariant Feature Transform (SIFT)** → Advantages

- ① Locality → features are local, robust to occlusion and clutter.
- ② Distinctiveness → individual features can be matched to a large database of objects.
- ③ Quantity → many features can be generated for even small objects.
- ④ Efficiency → close to real-time performance.

Algorithm → ① **Scale-space peak selection** → Potential location for finding features.

The scale-space of an image is produced from the convolution of Gaussian kernel (blurring) at different scales with the original image. Scale space is separated into octaves and the no. of octaves and scale depends on the size of the original image.

Once the image is blurred, the pixel is compared with 8 neighbours. Also compared with the nine pixels in previous and next scales. (Total 26)

If it is local extrema, it is a potential keypoint.

② **Keypoint localisation** → If the intensity at this extrema is less than the threshold value, they are rejected.

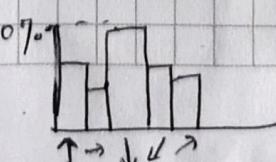
→ Edges are removed using eigen values and their ratios. We only need features.

→ Removal of low contrast keypoints and edge keypoints. Left with only strong interest points.

③ **Orientation Assignment** → Orientation is assigned to each keypoint.

→ orientation histogram is created. Highest peak in the histogram and any peak above 80% is considered.

→ creates keypoints with same location and scale but different directions.



- ④ Keypoint descriptor → We start with a  $16 \times 16$  window around detected interest point. We compute edge orientation for each pixel. Sub-group  $16 \times 16$  window to  $4 \times 4$  sub-blocks. Compute histograms. So,  $4 \times 4$  descriptors over  $16 \times 16$  sample array.  $4 \times 4 \times 8 \rightarrow 128$  values.
- Histogram entries are weighted by: (i) gradient magnitude (ii) A gaussian function with  $\sigma$  equal to 1.5 times the scale of the keypoint.
- Partial Voting → distribute histogram entries into adjacent bins. Each entry is added to all bins, multiplied by a weight of  $1-d$ , ( $d$  is the distance from the bin it belongs to).
- Matching SIFT features → accept the match if  $SSD(F_1, F_2) / SSD(F_1, F'_2) < t$ .  $t = 0.8$  has given good results in object recognition.
- Advantages/Properties of SIFT
- ① can handle changes in viewpoint → upto  $60^\circ$  out of plane rotation.
  - ② can " " " " in illumination → Sift even day-night
  - ③ fast and efficient → can run in real-time.
- ### Optical flow
- A Vector field function of the spatio-temporal image brightness variations  
It basically shows how does every pixel move from one image to another. (over time and space).
- ✓ Motion field → It is the 2d representation/reprojection of 3d-motion of the object in the scene into the image. True motion of the object.
- ✓ Optical flow → apparent motion of brightness patterns in the image.
- optical flow ≠ motion field: (sphere-sum) example.
- ↳ Feature Tracking → Extract visual features and track them over multiple frames (e.g. corners, textured areas)
- ↳ Optical flow → Recover image motion at each pixel from spatio-temporal image brightness variations.
- Solving these 2 problems with one method
- Basic Assumptions: Brightness consistency, Spatial coherence, Temporal persistence.
- Brightness consistency → Image measurement in a small region remain the same although their location may change.
- Spatial coherence → Neighbouring points in the scene typically belong to the same surface and hence typically have similar motions. (consider a small patch)
- Since they also project to the nearby pixels in the image we expect spatial coherence in the image flow.
- Temporal consistency → The image motion of a surface patch changes gradually over time.

→ The brightness constancy constraint :-

$$\begin{array}{|c|} \hline (x, y) \\ \text{displac.} \\ = (u, v) \\ I(x, y, t) \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline (x+u, y+v) \\ I(x, y, t+1) \\ \hline \end{array}$$

→ The brightness constancy equation:-

$$I(x, y, t) = I(x+u, y+v, t+1)$$

→ Taking the ~~higher~~<sup>Taylor</sup> expansion of  $I(x+u, y+v, t+1)$  at  $(x, y, t)$  to minimize the right side.

$$I(x+u, y+v, t+1) \approx I(x, y, t) + I_x \cdot u + I_y \cdot v + I_t \rightarrow \text{Difference over frames.}$$

Image derivative along  $x$

$$I(x+u, y+v, t+1) - I(x, y, t) = I_x \cdot u + I_y \cdot v + I_t \approx 0$$

Since it is a scalar equation with 2 unknowns, it is not solvable

Solvable?  
NO

Brightness constancy equation

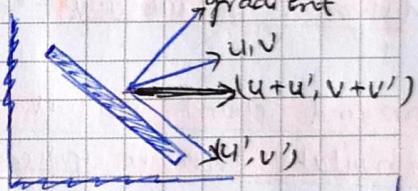
$$\approx \nabla I \cdot [u \ v]^T + I_t \approx 0$$

Aperture problem. → This problem states that if we are looking only at a single pixel, we cannot differentiate b/w the motion i.e. Different motions classified as similar or similar motions classified as different.

⇒ To solve this problem :- We have 2 methods :-

(A) Local method → Lukas and Kanade

(B) Global method → Horn and Schunck



Local method → we try to get more equations for 1 pixel. How? By spatial coherence constraint.  
→ Assuming the pixel's neighbours to have the same displacement  $(u, v)$   
→ for e.g. if we use a  $5 \times 5$  window, that gives us 25 equations/pixel.

$$\Rightarrow 0 = I_t(p_i) + \nabla I(p_i) \cdot [u \ v] \Rightarrow \begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ I_x(p_3) & I_y(p_3) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

(A) 25x2      (d) 2x1

$$(b) \quad 25 \times 1$$

Now the Problem: we have more equations than unknowns.

Solution: minimize  $\|Ad - b\|^2$

Solution → Solve least square problem  
 $(x = (A^T A)^{-1} A^T b)$   
in MATLAB  $x = A \backslash b$

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_y I_x & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix} \quad \xrightarrow{\text{ATA}} \begin{bmatrix} d \\ A^T b \end{bmatrix}$$

→ conditions for solvability:

- ① ATA should be invertible
- ② ATA should not be too small due to noise  
i.e. the eigen values of  $A^T A$  should not be too small.
- ③  $A^T A$  should be well-conditioned i.e.  $\lambda_1 / \lambda_2$  should not be too large.

**Some as Horn's corner detector**

- Observation:**
- ① Can measure sensitivity by just looking at one of the images
  - ② This tells us which pixels are easy to track and which are hard.

**Errors in Lucas-Kanade:** Potential causes of errors :-

- ① Suppose  $A^T A$  is easily invertible
- ② Suppose there is not much noise in the image.

When are our assumptions violated?

- ✗ Brightness constancy is not satisfied.
- ✗ The motion is not small.
- ✗ A point does not move like its neighbours.
  - window size is too large
  - what is the ideal window size?

→ Implementation issues

- ① Window Size
  - ✓ Small window size more sensitive to noise and may miss larger motions.
  - ✓ Large window more likely to cross an occlusion boundary (and is also slower)
  - ✓ Typical size  $15 \times 15 - 31 \times 31$

② Weighting the window → common to apply weights so that the center matters more

→ Horn and Schunk global method (Dense optical flow)

③ Formulate error in optical flow constraint:

$E_C = \iint_{\text{image}} (I_x u + I_y v + I_t)^2 dx dy$

should be as small as possible

④ Adding the smoothness constraint:

$E_S = \iint_{\text{image}} (u_x^2 + u_y^2) + (v_x^2 + v_y^2) dx dy$

→ usually motion field varies smoothly in the image. Penalize departure from smoothness.

⑤ find  $(u, v)$  at each point that minimizes:

$E = E_C + \lambda E_S$

global error      Regulation term.

What happens if the motion is not small? (i.e. much larger than 1 pixel)

→ and order terms dominate. Aliasing problem occurs.

Aliasing → Temporal aliasing causes ambiguities in optical flow because images can have many pixels with the same intensity.

Solution → coarse to fine estimation i.e. Reducing the resolution

→ Dealing with larger movements: Iterative refinement

① Initialize  $(x', y') = (x, y)$

$$I_t = I(x', y', t+1) - I(x, y, t)$$

② Compute  $(u, v)$  by:  $\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_y I_x & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$   
↳ displacement

③ Shift the window by  $(u, v)$ :  $x' = x + u$   $y' = y + v$

④ Recalculate  $I_t$

⑤ Repeat the steps 2-4 until small change

Iterative Refinement → Iterative KLT algorithm.

- ① Estimate displacement at each pixel by solving L-K equations
- ② Wrap  $I_t$  towards  $I_{t+1}$  using estimated flow field (Interpolation)
- ③ Repeat until convergence.

X

X

→ Tracking over many frames:-

① Select features in the first frame

② For each frame:

① update the position of tracked features:

By Discrete search or Lukas-Kanade method

② Terminate inconsistent tracks:

compute similarity with the corresponding feature in the previous frame or in the first frame where it's visible.

③ Find more features to track