# SCALE INVARIANT FEATURE TRANSFORM (SIFT)

#### OUTLINE

- SIFT Background
- SIFT Extraction
- Application in Content Based Image Search
- Conclusion

#### SIFT BACKGROUND

#### Scale-invariant feature transform

- **SIFT:** to **detect** and **describe** local features in an images.
- Proposed by *David Lowe* in ICCV1999.
- Refined in IJCV 2004.
- Cited more than 65953 times till now(March 2022).
- Wildly used in image search, object recognition, video tracking, gesture recognition, *etc*.

#### WHY SIFT IS SO POPULAR?

#### An instance of object matching



#### WHY SIFT IS SO POPULAR?

• Desired property of SIFT

- Invariant to scale change
- Invariant to rotation change
- Invariant to illumination change
- Robust to addition of noise
- Robust to substantial range of affine transformation
- Robust to 3D view point
- Highly distinctive for discrimination

# CLAIMED ADVANTAGES OF SIFT

- **Locality:** features are local, so robust to occlusion and clutter (no prior segmentation)
- **Distinctiveness:** individual features can be matched to a large database of objects
- **Quantity:** many features can be generated for even small objects
- Efficiency: close to real-time performance

#### MAJOR STAGES OF SIFT COMPUTATION

- Scale-space extrema detection
- Keypoint localization
- Orientation assignment
- Keypoint descriptor

#### HOW TO EXTRACT SIFT



Test image



**Detector**: where are the local features?



**Descriptor**: how to describe them?

## SIFT DETECTOR

• Desired properties for detector



Image 2

region size

- **Position**: Repeatable across different changes
- Scale: automatic scale estimation
- Intuition: Find scale that gives local maxima of some function *f* in both position and scale.



# WHAT CAN BE THE "SIGNATURE" FUNCTION F?

Laplacian-of-Gaussian = "**blob**" detector



#### AT A GIVEN POINT IN THE IMAGE:

• We define the *characteristic scale* as the scale that produces peak of Laplacian response























![](_page_14_Figure_1.jpeg)

![](_page_14_Figure_2.jpeg)

![](_page_15_Picture_0.jpeg)

![](_page_15_Figure_1.jpeg)

![](_page_15_Figure_2.jpeg)

![](_page_16_Figure_0.jpeg)

![](_page_16_Picture_1.jpeg)

![](_page_16_Picture_2.jpeg)

#### LOG VS. DOG

$$\nabla^2 G_{\sigma}(x, y) = \left(\frac{x^2 + y^2}{\sigma^4} - \frac{2}{\sigma^2}\right) G_{\sigma}(x, y)$$

$$G(x, y, k\sigma) - G(x, y, \sigma) \approx (k-1)\sigma \frac{\partial G}{\partial \sigma} = (k-1)\sigma^2 \nabla^2 G$$
  
DoG

#### TECHNICAL DETAIL

• We can approximate the Laplacian with a difference of Gaussians; more efficient to implement.

![](_page_19_Figure_0.jpeg)

For each octave of scale space, the initial image is repeatedly convolved with Gaussians to produce the set of scale images shown on the left. Adjacent Gaussian images are subtracted to produce the difference-of-Gaussian images on the right. After each octave, the Gaussian image is down-sampled by a factor 2, and the process repeated.

## Lowe's Pyramid Scheme

- Scale space is separated into octaves:
  - Octave 1 uses scale  $\sigma$
  - Octave 2 uses scale  $2\sigma$
  - etc.
- In each octave, the initial image is repeatedly convolved with Gaussians to produce a set of scale space images.
- Adjacent Gaussians are subtracted to produce the DOG
- After each octave, the Gaussian image is down-sampled by a factor of 2 to produce an image ¼ the size to start the next level.

#### DOG IMAGE PYRAMID

![](_page_21_Figure_1.jpeg)

Octave 1

 $\sigma_0 \rightarrow 2\sigma_0$ 

DoG Octave 1

#### LOCAL EXTREMA DETECTION

Maxima and minima
Compare x with its 26 neighbors at 3 scales

![](_page_22_Figure_2.jpeg)

#### FREQUENCY OF SAMPLING IN DOMAIN

- Trade-off between sampling frequency and rate of detection
- Sigma=1.6

![](_page_23_Figure_3.jpeg)

#### ELIMINATING EDGE RESPONSES

• Motivation

- DoG aims to detect "blob".
- DoG function have a strong response along edges.
- Remove such key points by Hessian Matrix analysis

#### • Hessian matrix

• Formulation

$$\mathbf{H} = \left[ \begin{array}{cc} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{array} \right]$$

# ELIMINATING EDGE RESPONSES $\mathbf{H} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix}$ $\mathrm{Tr}(\mathbf{H}) = D_{xx} + D_{yy} = \alpha + \beta,$

$$Det(\mathbf{H}) = D_{xx}D_{yy} - (D_{xy})^2 = \alpha\beta.$$

 $\alpha$ : larger eigenvalue  $\beta$ : smaller eigenvalue

$$\begin{aligned} \frac{\mathrm{Tr}(\mathbf{H})^2}{\mathrm{Det}(\mathbf{H})} &= \frac{(\alpha + \beta)^2}{\alpha\beta} = \frac{(r\beta + \beta)^2}{r\beta^2} = \frac{(r+1)^2}{r},\\ \frac{\mathrm{Tr}(\mathbf{H})^2}{\mathrm{Det}(\mathbf{H})} &\leq \frac{(r+1)^2}{r}. \end{aligned}$$

 $\alpha = r\beta$ 

r = 10 for this experimants

## ORIENTATION

#### Gradient and angle:

$$m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2}$$
  
$$\theta(x, y) = a \tan 2((L(x, y+1) - L(x, y-1))/(L(x+1, y) - L(x-1, y)))$$

#### Orientation selection

![](_page_26_Figure_4.jpeg)

# SIFT DESCRIPTOR

![](_page_27_Picture_1.jpeg)

#### • Making descriptor rotation invariant

![](_page_27_Picture_3.jpeg)

- Rotate patch according to its dominant gradient orientation
- This puts the patches into a canonical orientation.

## SIFT DESCRIPTOR

![](_page_28_Picture_1.jpeg)

• Use histograms to bin pixels within sub-patches according to their orientation.

![](_page_28_Figure_3.jpeg)

#### SUMMARY OF SIFT FEATURE

- Descriptor: 128-D
  - 4 by 4 patches, each with 8-D gradient angle histogram:
    - $4 \times 4 \times 8 = 128$
  - Normalized to reduce the effects of illumination change.
- Position: (x, y)
  - Where the feature is located at.
- Scale
  - Control the region size for descriptor extraction.
- Orientation
  - To achieve rotation-invariant descriptor.

#### LOCAL FEATURES

- Detection: Identify the interest points.
- Description: Extract vector feature descriptor around each interest point.
- Matching: Determine correspondence between descriptors in two views.

![](_page_30_Figure_4.jpeg)

#### MATCHING LOCAL FEATURES

- To generate candidate matches, find patches that have the most similar appearance (e.g., lowest SSD)
- Simplest approach: compare them all, take the closest (or closest k, or within a thresholded distance)

![](_page_31_Picture_3.jpeg)

#### AMBIGUOUS MATCHES

- At what SSD value do we have a good match?
- To add robustness, consider ratio of distance to best match to distance to second best match
  - If low, first match looks good.
  - If high, could be ambiguous match.

![](_page_32_Picture_5.jpeg)

#### Closest / NextClosest< threshold(0.6 for example)