**WIKIPEDIA**
The Free Encyclopedia

# Scale-invariant feature transform

The **scale-invariant feature transform** (**SIFT**) is a computer vision algorithm to detect, describe, and match local *features* in images, invented by David Lowe in 1999.[1] Applications include object recognition, robotic mapping and navigation, image stitching, 3D modeling, gesture recognition, video tracking, individual identification of wildlife and match moving.

SIFT keypoints of objects are first extracted from a set of reference images[1] and stored in a database. An object is recognized in a new image by individually comparing each feature from the new image to this database and finding candidate matching features based on Euclidean distance of their feature vectors. From the full set of matches, subsets of keypoints that agree on the object and its location, scale, and orientation in the new image are identified to filter out good matches. The determination of consistent clusters is performed [...] implementation of the generalised Hough transform. Each cluster of 3 or m[...] pose is then subject to further detailed model verification and subsequently [...] that a particular set of features indicates the presence of an object is compu[...] probable false matches. Object matches that pass all these tests can be identi[...]

## Overview

For any object in an image, interesting points on the object can be extract[...] object. This description, extracted from a training image, can then be use[...] locate the object in a test image containing many other objects. To perform[...] features extracted from the training image be detectable even under changes[...] points usually lie on high-contrast regions of the image, such as object edges.

Another important characteristic of these features is that the relative po[...] shouldn't change from one image to another. For example, if only the four c[...] would work regardless of the door's position; but if points in the frame wer[...] door is opened or closed. Similarly, features located in articulated or flexi[...] change in their internal geometry happens between two images in the set [...]

**Gesture recognition** is a topic in computer science and language technology with the goal of interpreting human gestures via mathematical algorithms. It is a subdiscipline of computer vision. Gestures can originate from any bodily motion or

detects and uses a much larger number of features from the images, which reduces the contribution of the errors caused by these local variations in the average error of all feature matching errors.

SIFT[3] can robustly identify objects even among clutter and under partial occlusion, because the SIFT feature descriptor is invariant to uniform scaling, orientation, illumination changes, and partially invariant to affine distortion.[1] This section summarizes the original SIFT algorithm and mentions a few competing techniques available for object recognition under clutter and partial occlusion.

The SIFT descriptor is based on image measurements in terms of *receptive fields*[4][5][6][7] over which *local scale invariant reference frames*[8][9] are established by *local scale selection*.[10][11][9] A general theoretical explanation about this is given in the Scholarpedia article on SIFT.[12]

| Problem | Technique | Advantage |
|---|---|---|
| key localization / scale / rotation | Difference of Gaussians / scale-space pyramid / orientation assignment | accuracy, stability, scale & rotational invariance |
| geometric distortion | blurring / resampling of local image orientation planes | affine invariance |
| indexing and matching | nearest neighbor / Best Bin First search | Efficiency / speed |
| Cluster identification | Hough Transform voting | reliable pose models |
| Model verification / outlier detection | Linear least squares | better error tolerance with fewer matches |
| Hypothesis acceptance | Bayesian Probability analysis | reliability |

## Types of features

The detection and description of local image features can help in object recognition. The SIFT features are local and based on the appearance of the object at particular interest points, and are invariant to image scale and rotation. They are also robust to changes in illumination, noise, and minor changes in viewpoint. In addition to these properties, they are highly distinctive, relatively easy to extract and allow for correct object identification with low probability of mismatch. They are relatively easy to match against a (large) database of local features but, however, the high dimensionality can be an issue, and generally probabilistic algorithms such as k-d trees with best bin first search are used. Object description by set of SIFT features is also robust to partial occlusion; as few as 3 SIFT features from an object are enough to compute its location and pose. Recognition can be performed in close-to-real time, at least for small databases and on modern computer hardware.

# Main Stages

## Scale-invariant feature detection

Lowe's method for image feature generation transforms an image into a large collection of feature vectors, each of which is invariant to image translation, scaling, and rotation, partially invariant to illumination changes, and robust to local geometric distortion. These features share similar properties with neurons in the primary visual cortex that encode basic forms, color, and movement for object detection in primate vision.[13] Key locations are defined as maxima and minima of the result of difference of Gaussians function applied in scale space to a series of smoothed and resampled images. Low-contrast candidate points and edge response points along an edge are discarded. Dominant orientations are assigned to localized key points. These steps ensure that the key points are more stable for matching and recognition. SIFT descriptors robust to local affine distortion are then obtained by considering pixels around a radius of the key location, blurring, and resampling local image orientation planes.

## Feature matching and indexing

Indexing consists of storing SIFT keys and identifying matching keys from the new image. Lowe used a modification of the k-d tree algorithm called the **best-bin-first search** method[14] that can identify the nearest neighbors with high probability using only a limited amount of computation. The BBF algorithm uses a modified search ordering for the k-d tree algorithm so that bins in feature space are searched in the order of their closest distance from the query location. This search order requires the use of a heap-based priority queue for efficient determination of the search order. We obtain a candidate for each keypoint by identifying its nearest neighbor in the database of keypoints from training images. The nearest neighbors are defined as the keypoints with minimum Euclidean distance from the given descriptor vector. The way Lowe[2] determined whether a given candidate should be kept or 'thrown out' is by checking the ratio between the distance from this given candidate and the distance from the closest keypoint which is not of the same object class as the candidate at hand (candidate feature vector / closest different class feature vector), the idea is that we can only be sure of candidates in which features/keypoints from distinct object classes don't "clutter" it (not geometrically clutter in the feature space necessarily but more so clutter along the right half (>0) of the real line), this is an obvious consequence of using Euclidean distance as our nearest neighbor measure. The ratio threshold for rejection is whenever it is above 0.8. This method eliminated 90% of false matches while discarding less than 5% of correct matches. To further improve the efficiency of the best-bin-first algorithm search was cut off after checking the first 200 nearest neighbor candidates. For a database of 100,000 keypoints, this provides a speedup over exact nearest neighbor search by about 2 orders of magnitude, yet results in less than a 5% loss in the number of correct matches.

## Cluster identification by Hough transform voting

Hough transform is used to cluster reliable model hypotheses to search for keys that agree upon a particular model pose. Hough transform identifies clusters of features with a consistent interpretation by using each feature to vote for all object poses that are consistent with the feature. When clusters of features are found to vote for the same pose of an object, the probability of the interpretation being correct is much higher than for any single feature. An entry in a hash table is created predicting the model location, orientation, and scale from the match hypothesis. The hash table is searched to identify all clusters of at least 3 entries in a bin, and the bins are sorted into decreasing order of size.

Each of the SIFT keypoints specifies 2D location, scale, and orientation, and each matched keypoint in the database has a record of its parameters relative to the training image in which it was found. The similarity transform implied by these 4 parameters is only an approximation to the full 6 degree-of-freedom pose space for a 3D object and also does not account for any non-rigid deformations. Therefore, Lowe[2] used broad bin sizes of 30 degrees for orientation, a factor of 2 for

scale, and 0.25 times the maximum projected training image dimension (using the predicted scale) for location. The SIFT key samples generated at the larger scale are given twice the weight of those at the smaller scale. This means that the larger scale is in effect able to filter the most likely neighbors for checking at the smaller scale. This also improves recognition performance by giving more weight to the least-noisy scale. To avoid the problem of boundary effects in bin assignment, each keypoint match votes for the 2 closest bins in each dimension, giving a total of 16 entries for each hypothesis and further broadening the pose range.

## Model verification by linear least squares

Each identified cluster is then subject to a verification procedure in which a linear least squares solution is performed for the parameters of the affine transformation relating the model to the image. The affine transformation of a model point [x y]$^T$ to an image point [u v]$^T$ can be written as below

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix}$$

where the model translation is [t$_x$ t$_y$]$^T$ and the affine rotation, scale, and stretch are represented by the parameters $m_1$, $m_2$, $m_3$ and $m_4$. To solve for the transformation parameters the equation above can be rewritten to gather the unknowns into a column vector.

$$\begin{bmatrix} x & y & 0 & 0 & 1 & 0 \\ 0 & 0 & x & y & 0 & 1 \\ \dots & & & & & \\ \dots & & & & & \end{bmatrix} \begin{bmatrix} m1 \\ m2 \\ m3 \\ m4 \\ tx \\ ty \end{bmatrix} = \begin{bmatrix} u \\ v \\ . \\ . \end{bmatrix}$$

This equation shows a single match, but any number of further matches can be added, with each match contributing two more rows to the first and last matrix. At least 3 matches are needed to provide a solution. We can write this linear system as

$$A\hat{\mathbf{x}} \approx \mathbf{b},$$

where $A$ is a known $m$-by-$n$ matrix (usually with $m > n$), $\mathbf{x}$ is an unknown $n$-dimensional parameter vector, and $\mathbf{b}$ is a known $m$-dimensional measurement vector.

Therefore, the minimizing vector $\hat{\mathbf{x}}$ is a solution of the **normal equation**

$$A^T A\hat{\mathbf{x}} = A^T \mathbf{b}.$$

The solution of the system of linear equations is given in terms of the matrix $(A^T A)^{-1} A^T$, called the pseudoinverse of $A$, by

$$\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}.$$

which minimizes the sum of the squares of the distances from the projected model locations to the corresponding image locations.

## Outlier detection

Outliers can now be removed by checking for agreement between each image feature and the model, given the parameter solution. Given the linear least squares solution, each match is required to agree within half the error range that was used for the parameters in the Hough transform bins. As outliers are discarded, the linear least squares solution is re-solved with the remaining points, and the process iterated. If fewer than 3 points remain after discarding outliers, then the match is rejected. In addition, a top-down matching phase is used to add any further matches that agree with the projected model position, which may have been missed from the Hough transform bin due to the similarity transform approximation or other errors.

The final decision to accept or reject a model hypothesis is based on a detailed probabilistic model.[15] This method first computes the expected number of false matches to the model pose, given the projected size of the model, the number of features within the region, and the accuracy of the fit. A Bayesian probability analysis then gives the probability that the object is present based on the actual number of matching features found. A model is accepted if the final probability for a correct interpretation is greater than 0.98. Lowe's SIFT based object recognition gives excellent results except under wide illumination variations and under non-rigid transformations.

# Algorithm

## Scale-space extrema detection

We begin by detecting points of interest, which are termed *keypoints* in the SIFT framework. The image is convolved with Gaussian filters at different scales, and then the difference of successive Gaussian-blurred images are taken. Keypoints are then taken as maxima/minima of the Difference of Gaussians (DoG) that occur at multiple scales. Specifically, a DoG image $D(x, y, \sigma)$ is given by

$$D(x, y, \sigma) = L(x, y, k_i\sigma) - L(x, y, k_j\sigma),$$

where $L(x, y, k\sigma)$ is the convolution of the original image $I(x, y)$ with the Gaussian blur $G(x, y, k\sigma)$ at scale $k\sigma$, i.e.,

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

Hence a DoG image between scales $k_i\sigma$ and $k_j\sigma$ is just the difference of the Gaussian-blurred images at scales $k_i\sigma$ and $k_j\sigma$. For scale space extrema detection in the SIFT algorithm, the image is first convolved with Gaussian-blurs at different scales. The convolved images are grouped by octave (an octave corresponds to doubling the value of $\sigma$), and the value of $k_i$ is selected so that we obtain a fixed number of convolved images per octave. Then the Difference-of-Gaussian images are taken from adjacent Gaussian-blurred images per octave.

Once DoG images have been obtained, keypoints are identified as local minima/maxima of the DoG images across scales. This is done by comparing each pixel in the DoG images to its eight neighbors at the same scale and nine corresponding neighboring pixels in each of the neighboring scales. If the pixel value is the maximum or minimum among all compared pixels, it is selected as a candidate keypoint.

This keypoint detection step is a variation of one of the blob detection methods developed by Lindeberg by detecting scale-space extrema of the scale normalized Laplacian;[10][11] that is, detecting points that are local extrema with respect to both space and scale, in the discrete case by comparisons with the nearest 26 neighbors in a discretized scale-space volume. The difference of Gaussians operator can be seen as an approximation to the Laplacian, with the implicit normalization in the pyramid also constituting a discrete approximation of the scale-normalized Laplacian.[12] Another real-time implementation of scale-space extrema of the Laplacian operator has been presented by Lindeberg and Bretzner based on a hybrid pyramid representation,[16] which was used for human-computer interaction by real-time gesture recognition in Bretzner et al. (2002).[17]

## Keypoint localization

Scale-space extrema detection produces too many keypoint candidates, some of which are unstable. The next step in the algorithm is to perform a detailed fit to the nearby data for accurate location, scale, and ratio of principal curvatures. This information allows the rejection of points which are low contrast (and are therefore sensitive to noise) or poorly localized along an edge.

### Interpolation of nearby data for accurate position

First, for each candidate keypoint, interpolation of nearby data is used to accurately determine its position. The initial approach was to just locate each keypoint at the location and scale of the candidate keypoint.[1] The new approach calculates the interpolated location of the extremum, which substantially improves matching and stability.[2] The interpolation is done using the quadratic Taylor expansion of the Difference-of-Gaussian scale-space function, $D(x, y, \sigma)$ with the candidate keypoint as the origin. This Taylor expansion is given by:

$$D(\mathbf{x}) = D + \frac{\partial D}{\partial \mathbf{x}}^T \mathbf{x} + \frac{1}{2}\mathbf{x}^T \frac{\partial^2 D}{\partial \mathbf{x}^2}\mathbf{x}$$

where D and its derivatives are evaluated at the candidate keypoint and $\mathbf{x} = (x, y, \sigma)^T$ is the offset from this point. The location of the extremum, $\hat{\mathbf{x}}$, is determined by taking the derivative of this function with respect to $\mathbf{x}$ and setting it to zero. If the offset $\hat{\mathbf{x}}$ is larger than 0.5 in any dimension, then that's an indication that the extremum lies closer to another candidate keypoint. In this case, the candidate keypoint is changed and the interpolation performed instead about that point. Otherwise the offset is added to its candidate keypoint to get the interpolated estimate for the location of the extremum. A similar subpixel determination of the locations of scale-space extrema is performed in the real-time implementation based on hybrid pyramids developed by Lindeberg and his co-workers.[16]

### Discarding low-contrast keypoints

To discard the keypoints with low contrast, the value of the second-order Taylor expansion $D(\mathbf{x})$ is computed at the offset $\hat{\mathbf{x}}$. If this value is less than 0.03, the candidate keypoint is discarded. Otherwise it is kept, with final scale-space location $\mathbf{y} + \hat{\mathbf{x}}$, where $\mathbf{y}$ is the original location of the keypoint.

### Eliminating edge responses

The DoG function will have strong responses along edges, even if the candidate keypoint is not robust to small amounts of noise. Therefore, in order to increase stability, we need to eliminate the keypoints that have poorly determined locations but have high edge responses.

For poorly defined peaks in the DoG function, the principal curvature across the edge would be much larger than the principal curvature along it. Finding these principal curvatures amounts to solving for the eigenvalues of the second-order Hessian matrix, $\mathbf{H}$:
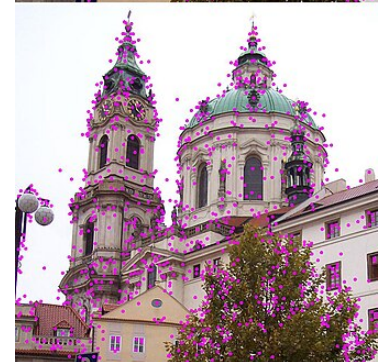
$$\mathbf{H} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix}$$

The eigenvalues of $\mathbf{H}$ are proportional to the principal curvatures of D. It turns out that the ratio of the two eigenvalues, say $\alpha$ is the larger one, and $\beta$ the smaller one, with ratio $r = \alpha/\beta$, is sufficient for SIFT's purposes. The trace of $\mathbf{H}$, i.e., $D_{xx} + D_{yy}$, gives us the sum of the two eigenvalues, while its determinant, i.e., $D_{xx}D_{yy} - D_{xy}^2$, yields the product. The ratio $\mathrm{R} = \mathrm{Tr}(\mathbf{H})^2/\mathrm{Det}(\mathbf{H})$ can be shown to be equal to $(r+1)^2/r$, which depends only on the ratio of the eigenvalues rather than their individual values. R is minimum when the eigenvalues are equal to each other.



After scale space extrema are detected (their location being shown in the uppermost image) the SIFT algorithm discards low-contrast keypoints (remaining points are shown in the middle image) and then filters out those located on edges. Resulting set of keypoints is shown on last image.

Therefore, the higher the absolute difference between the two eigenvalues, which is equivalent to a higher absolute difference between the two principal curvatures of D, the higher the value of R. It follows that, for some threshold eigenvalue ratio $r_{\mathrm{th}}$, if R for a candidate keypoint is larger than $(r_{\mathrm{th}} + 1)^2/r_{\mathrm{th}}$, that keypoint is poorly localized and hence rejected. The new approach uses $r_{\mathrm{th}} = 10$.[2]

This processing step for suppressing responses at edges is a transfer of a corresponding approach in the Harris operator for corner detection. The difference is that the measure for thresholding is computed from the Hessian matrix instead of a second-moment matrix.

## Orientation assignment

In this step, each keypoint is assigned one or more orientations based on local image gradient directions. This is the key step in achieving invariance to rotation as the keypoint descriptor can be represented relative to this orientation and

therefore achieve invariance to image rotation.

First, the Gaussian-smoothed image $L(x, y, \sigma)$ at the keypoint's scale $\sigma$ is taken so that all computations are performed in a scale-invariant manner. For an image sample $L(x, y)$ at scale $\sigma$, the gradient magnitude, $m(x, y)$, and orientation, $\theta(x, y)$, are precomputed using pixel differences:

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

$$\theta(x, y) = \text{atan2}(L(x, y + 1) - L(x, y - 1), L(x + 1, y) - L(x - 1, y))$$

The magnitude and direction calculations for the gradient are done for every pixel in a neighboring region around the keypoint in the Gaussian-blurred image L. An orientation histogram with 36 bins is formed, with each bin covering 10 degrees. Each sample in the neighboring window added to a histogram bin is weighted by its gradient magnitude and by a Gaussian-weighted circular window with a $\sigma$ that is 1.5 times that of the scale of the keypoint. The peaks in this histogram correspond to dominant orientations. Once the histogram is filled, the orientations corresponding to the highest peak and local peaks that are within 80% of the highest peaks are assigned to the keypoint. In the case of multiple orientations being assigned, an additional keypoint is created having the same location and scale as the original keypoint for each additional orientation.

## Keypoint descriptor

Previous steps found keypoint locations at particular scales and assigned orientations to them. This ensured invariance to image location, scale and rotation. Now we want to compute a descriptor vector for each keypoint such that the descriptor is highly distinctive and partially invariant to the remaining variations such as illumination, 3D viewpoint, etc. This step is performed on the image closest in scale to the keypoint's scale.

First a set of orientation histograms is created on 4×4 pixel neighborhoods with 8 bins each. These histograms are computed from magnitude and orientation values of samples in a 16×16 region around the keypoint such that each histogram contains samples from a 4×4 subregion of the original neighborhood region. The image gradient magnitudes and orientations are sampled around the keypoint location, using the scale of the keypoint to select the level of Gaussian blur for the image. In order to achieve orientation invariance, the coordinates of the descriptor and the gradient orientations are rotated relative to the keypoint orientation. The magnitudes are further weighted by a Gaussian function with $\sigma$ equal to one half the width of the descriptor window. The descriptor then becomes a vector of all the values of these histograms. Since there are 4 × 4 = 16 histograms each with 8 bins the vector has 128 elements. This vector is then normalized to unit length in order to enhance invariance to affine changes in illumination. To reduce the effects of non-linear illumination a threshold of 0.2 is applied and the vector is again normalized. The thresholding process, also referred to as clamping, can improve matching results even when non-linear illumination effects are not present.[18] The threshold of 0.2 was empirically chosen, and by replacing the fixed threshold with one systematically calculated, matching results can be improved.[18]

Although the dimension of the descriptor, i.e. 128, seems high, descriptors with lower dimension than this don't perform as well across the range of matching tasks[2] and the computational cost remains low due to the approximate BBF (see below) method used for finding the nearest neighbor. Longer descriptors continue to do better but not by much and there is an additional danger of increased sensitivity to distortion and occlusion. It is also shown that feature matching accuracy is above 50% for viewpoint changes of up to 50 degrees. Therefore, SIFT descriptors are invariant to minor affine changes. To test the distinctiveness of the SIFT descriptors, matching accuracy is also measured against varying number of keypoints in the testing database, and it is shown that matching accuracy decreases only very slightly for very large database sizes, thus indicating that SIFT features are highly distinctive.

# Comparison of SIFT features with other local features

There has been an extensive study done on the performance evaluation of different local descriptors, including SIFT, using a range of detectors.[19] The main results are summarized below:

- SIFT and SIFT-like GLOH features exhibit the highest matching accuracies (recall rates) for an affine transformation of 50 degrees. After this transformation limit, results start to become unreliable.
- Distinctiveness of descriptors is measured by summing the eigenvalues of the descriptors, obtained by the Principal components analysis of the descriptors normalized by their variance. This corresponds to the amount of variance

captured by different descriptors, therefore, to their distinctiveness. PCA-SIFT (Principal Components Analysis applied to SIFT descriptors), GLOH and SIFT features give the highest values.

- SIFT-based descriptors outperform other contemporary local descriptors on both textured and structured scenes, with the difference in performance larger on the textured scene.
- For scale changes in the range 2–2.5 and image rotations in the range 30 to 45 degrees, SIFT and SIFT-based descriptors again outperform other contemporary local descriptors with both textured and structured scene content.
- Introduction of blur affects all local descriptors, especially those based on edges, like shape context, because edges disappear in the case of a strong blur. But GLOH, PCA-SIFT and SIFT still performed better than the others. This is also true for evaluation in the case of illumination changes.

The evaluations carried out suggests strongly that SIFT-based descriptors, which are region-based, are the most robust and distinctive, and are therefore best suited for feature matching. However, most recent feature descriptors such as SURF have not been evaluated in this study.

SURF has later been shown to have similar performance to SIFT, while at the same time being much faster.[20] Other studies conclude that when speed is not critical, SIFT outperforms SURF.[21][22] Specifically, disregarding discretization effects the pure image descriptor in SIFT is significantly better than the pure image descriptor in SURF, whereas the scale-space extrema of the determinant of the Hessian underlying the pure interest point detector in SURF constitute significantly better interest points compared to the scale-space extrema of the Laplacian to which the interest point detector in SIFT constitutes a numerical approximation.[21]

The performance of image matching by SIFT descriptors can be improved in the sense of achieving higher efficiency scores and lower 1-precision scores by replacing the scale-space extrema of the difference-of-Gaussians operator in original SIFT by scale-space extrema of the determinant of the Hessian, or more generally considering a more general family of generalized scale-space interest points.[21]

Recently, a slight variation of the descriptor employing an irregular histogram grid has been proposed that significantly improves its performance.[23] Instead of using a 4×4 grid of histogram bins, all bins extend to the center of the feature. This improves the descriptor's robustness to scale changes.

The SIFT-Rank[24] descriptor was shown to improve the performance of the standard SIFT descriptor for affine feature matching. A SIFT-Rank descriptor is generated from a standard SIFT descriptor, by setting each histogram bin to its rank in a sorted array of bins. The Euclidean distance between SIFT-Rank descriptors is invariant to arbitrary monotonic changes in histogram bin values, and is related to Spearman's rank correlation coefficient.

# Applications

## Object recognition using SIFT features

Given SIFT's ability to find distinctive keypoints that are invariant to location, scale and rotation, and robust to affine transformations (changes in scale, rotation, shear, and position) and changes in illumination, they are usable for object recognition. The steps are given below.

- First, SIFT features are obtained from the input image using the algorithm described above.
- These features are matched to the SIFT feature database obtained from the training images. This feature matching is done through a Euclidean-distance based nearest neighbor approach. To increase robustness, matches are rejected for those keypoints for which the ratio of the nearest neighbor distance to the second-nearest neighbor distance is greater than 0.8. This discards many of the false matches arising from background clutter. Finally, to avoid the expensive search required for finding the Euclidean-distance-based nearest neighbor, an approximate algorithm called the best-bin-first algorithm is used.[14] This is a fast method for returning the nearest neighbor with high probability, and can give speedup by factor of 1000 while finding nearest neighbor (of interest) 95% of the time.
- Although the distance ratio test described above discards many of the false matches arising from background clutter, we still have matches that belong to different objects. Therefore, to increase robustness to object identification, we want to cluster those features that belong to the same object and reject the matches that are left out in the clustering process. This is done using the Hough transform. This will identify clusters of features that vote for the same object pose. When clusters of features are found to vote for the same pose of an object, the probability of the interpretation being correct is much higher than for any single feature. Each keypoint votes for the set of object poses that are consistent with the keypoint's location, scale, and orientation. *Bins* that accumulate at least 3 votes are identified as candidate object/pose matches.
- For each candidate cluster, a least-squares solution for the best estimated affine projection parameters relating the

training image to the input image is obtained. If the projection of a keypoint through these parameters lies within half the error range that was used for the parameters in the Hough transform bins, the keypoint match is kept. If fewer than 3 points remain after discarding outliers for a bin, then the object match is rejected. The least-squares fitting is repeated until no more rejections take place. This works better for planar surface recognition than 3D object recognition since the affine model is no longer accurate for 3D objects.

- In this journal,[25] authors proposed a new approach to use SIFT descriptors for multiple object detection purposes. The proposed multiple object detection approach is tested on aerial and satellite images.

SIFT features can essentially be applied to any task that requires identification of matching locations between images. Work has been done on applications such as recognition of particular object categories in 2D images, 3D reconstruction, motion tracking and segmentation, robot localization, image panorama stitching and epipolar calibration. Some of these are discussed in more detail below.

## Robot localization and mapping

In this application,[26] a trinocular stereo system is used to determine 3D estimates for keypoint locations. Keypoints are used only when they appear in all 3 images with consistent disparities, resulting in very few outliers. As the robot moves, it localizes itself using feature matches to the existing 3D map, and then incrementally adds features to the map while updating their 3D positions using a Kalman filter. This provides a robust and accurate solution to the problem of robot localization in unknown environments. Recent 3D solvers leverage the use of keypoint directions to solve trinocular geometry from three keypoints[27] and absolute pose from only two keypoints,[28] an often disregarded but useful measurement available in SIFT. These orientation measurements reduce the number of required correspondences, further increasing robustness exponentially.

## Panorama stitching

SIFT feature matching can be used in image stitching for fully automated panorama reconstruction from non-panoramic images. The SIFT features extracted from the input images are matched against each other to find $k$ nearest-neighbors for each feature. These correspondences are then used to find $m$ candidate matching images for each image. Homographies between pairs of images are then computed using RANSAC and a probabilistic model is used for verification. Because there is no restriction on the input images, graph search is applied to find connected components of image matches such that each connected component will correspond to a panorama. Finally for each connected component bundle adjustment is performed to solve for joint camera parameters, and the panorama is rendered using multi-band blending. Because of the SIFT-inspired object recognition approach to panorama stitching, the resulting system is insensitive to the ordering, orientation, scale and illumination of the images. The input images can contain multiple panoramas and noise images (some of which may not even be part of the composite image), and panoramic sequences are recognized and rendered as output.[29]

## 3D scene modeling, recognition and tracking

This application uses SIFT features for 3D object recognition and 3D modeling in context of augmented reality, in which synthetic objects with accurate pose are superimposed on real images. SIFT matching is done for a number of 2D images of a scene or object taken from different angles. This is used with bundle adjustment initialized from an essential matrix or trifocal tensor to build a sparse 3D model of the viewed scene and to simultaneously recover camera poses and calibration parameters. Then the position, orientation and size of the virtual object are defined relative to the coordinate frame of the recovered model. For online match moving, SIFT features again are extracted from the current video frame and matched to the features already computed for the world model, resulting in a set of 2D-to-3D correspondences. These correspondences are then used to compute the current camera pose for the virtual projection and final rendering. A regularization technique is used to reduce the jitter in the virtual projection.[30] The use of SIFT directions have also been used to increase robustness of this process.[27][28] 3D extensions of SIFT have also been evaluated for true 3D object recognition and retrieval.[31][32]

## 3D SIFT-like descriptors for human action recognition

Extensions of the SIFT descriptor to 2+1-dimensional spatio-temporal data in context of human action recognition in video sequences have been studied.[31][33][34][35] The computation of local position-dependent histograms in the 2D SIFT algorithm are extended from two to three dimensions to describe SIFT features in a spatio-temporal domain. For

application to human action recognition in a video sequence, sampling of the training videos is carried out either at spatio-temporal interest points or at randomly determined locations, times and scales. The spatio-temporal regions around these interest points are then described using the 3D SIFT descriptor. These descriptors are then clustered to form a spatio-temporal Bag of words model. 3D SIFT descriptors extracted from the test videos are then matched against these *words* for human action classification.

The authors report much better results with their 3D SIFT descriptor approach than with other approaches like simple 2D SIFT descriptors and Gradient Magnitude.[36]

### Analyzing the Human Brain in 3D Magnetic Resonance Images

The Feature-based Morphometry (FBM) technique[37] uses extrema in a difference of Gaussian scale-space to analyze and classify 3D magnetic resonance images (MRIs) of the human brain. FBM models the image probabilistically as a collage of independent features, conditional on image geometry and group labels, e.g. healthy subjects and subjects with Alzheimer's disease (AD). Features are first extracted in individual images from a 4D difference of Gaussian scale-space, then modeled in terms of their appearance, geometry and group co-occurrence statistics across a set of images. FBM was validated in the analysis of AD using a set of ~200 volumetric MRIs of the human brain, automatically identifying established indicators of AD in the brain and classifying mild AD in new images with a rate of 80%.[37]

# Competing methods

Competing methods for scale invariant object recognition under clutter / partial occlusion include the following.

RIFT[38] is a rotation-invariant generalization of SIFT. The RIFT descriptor is constructed using circular normalized patches divided into concentric rings of equal width and within each ring a gradient orientation histogram is computed. To maintain rotation invariance, the orientation is measured at each point relative to the direction pointing outward from the center.

RootSIFT[39] is a variant of SIFT that modifies descriptor normalization. Since SIFT descriptors are histograms (and as such probability distributions), employing Euclidean distance to determine their similarity is not a natural choice. Comparing such descriptors using similarity measures tailored to probability distributions such as Bhattacharyya coefficient (also known as Hellinger kernel) turns out to be more beneficial. For this purpose, the originally $\ell^2$ normalized descriptor is first $\ell^1$ normalized and the square root of each element is computed followed by $\ell^2$ renormalization. After these algebraic manipulations, RootSIFT descriptors can be normally compared using Euclidean distance which is equivalent to using the Hellinger kernel on the original SIFT descriptors. This normalization scheme termed "L1-sqrt" was previously introduced for the block normalization of HOG features whose rectangular block arrangement descriptor variant (R-HOG) is conceptually similar to the SIFT descriptor.

G-RIF:[40] Generalized Robust Invariant Feature is a general context descriptor which encodes edge orientation, edge density and hue information in a unified form combining perceptual information with spatial encoding. The object recognition scheme uses neighboring context based voting to estimate object models.

"SURF:[41] Speeded Up Robust Features" is a high-performance scale- and rotation-invariant interest point detector / descriptor claimed to approximate or even outperform previously proposed schemes with respect to repeatability, distinctiveness, and robustness. SURF relies on integral images for image convolutions to reduce computation time, builds on the strengths of the leading existing detectors and descriptors (using a fast Hessian matrix-based measure for the detector and a distribution-based descriptor). It describes a distribution of Haar wavelet responses within the interest point neighborhood. Integral images are used for speed and only 64 dimensions are used reducing the time for feature computation and matching. The indexing step is based on the sign of the Laplacian, which increases the matching speed and the robustness of the descriptor.

PCA-SIFT[42] and GLOH[19] are variants of SIFT. PCA-SIFT descriptor is a vector of image gradients in x and y direction computed within the support region. The gradient region is sampled at 39×39 locations, therefore the vector is of dimension 3042. The dimension is reduced to 36 with PCA. Gradient location-orientation histogram (GLOH) is an extension of the SIFT descriptor designed to increase its robustness and distinctiveness. The SIFT descriptor is computed for a log-polar location grid with three bins in radial direction (the radius set to 6, 11, and 15) and 8 in angular direction, which results in 17 location bins. The central bin is not divided in angular directions. The gradient orientations are quantized in 16 bins resulting in 272-bin histogram. The size of this descriptor is reduced with PCA. The covariance matrix for PCA is estimated on image patches collected from various images. The 128 largest eigenvectors are used for

description.

Gauss-SIFT[21] is a pure image descriptor defined by performing all image measurements underlying the pure image descriptor in SIFT by Gaussian derivative responses as opposed to derivative approximations in an image pyramid as done in regular SIFT. In this way, discretization effects over space and scale can be reduced to a minimum allowing for potentially more accurate image descriptors. In Lindeberg (2015)[21] such pure Gauss-SIFT image descriptors were combined with a set of generalized scale-space interest points comprising the Laplacian of the Gaussian, the determinant of the Hessian, four new unsigned or signed Hessian feature strength measures as well as Harris-Laplace and Shi-and-Tomasi interests points. In an extensive experimental evaluation on a poster dataset comprising multiple views of 12 posters over scaling transformations up to a factor of 6 and viewing direction variations up to a slant angle of 45 degrees, it was shown that substantial increase in performance of image matching (higher efficiency scores and lower 1-precision scores) could be obtained by replacing Laplacian of Gaussian interest points by determinant of the Hessian interest points. Since difference-of-Gaussians interest points constitute a numerical approximation of Laplacian of the Gaussian interest points, this shows that a substantial increase in matching performance is possible by replacing the difference-of-Gaussians interest points in SIFT by determinant of the Hessian interest points. Additional increase in performance can furthermore be obtained by considering the unsigned Hessian feature strength measure $D_1 L = \det HL - k \operatorname{trace}^2 HL$ if $\det HL - k \operatorname{trace}^2 HL > 0$ or 0 otherwise. A quantitative comparison between the Gauss-SIFT descriptor and a corresponding Gauss-SURF descriptor did also show that Gauss-SIFT does generally perform significantly better than Gauss-SURF for a large number of different scale-space interest point detectors. This study therefore shows that disregarding discretization effects the pure image descriptor in SIFT is significantly better than the pure image descriptor in SURF, whereas the underlying interest point detector in SURF, which can be seen as numerical approximation to scale-space extrema of the determinant of the Hessian, is significantly better than the underlying interest point detector in SIFT.

Wagner et al. developed two object recognition algorithms especially designed with the limitations of current mobile phones in mind.[43] In contrast to the classic SIFT approach, Wagner et al. use the FAST corner detector for feature detection. The algorithm also distinguishes between the off-line preparation phase where features are created at different scale levels and the on-line phase where features are only created at the current fixed scale level of the phone's camera image. In addition, features are created from a fixed patch size of 15×15 pixels and form a SIFT descriptor with only 36 dimensions. The approach has been further extended by integrating a Scalable Vocabulary Tree in the recognition pipeline.[44] This allows the efficient recognition of a larger number of objects on mobile phones. The approach is mainly restricted by the amount of available RAM.

KAZE and A-KAZE *(KAZE Features and Accelerated-Kaze Features)* is a new 2D feature detection and description method that perform better compared to SIFT and SURF. It gains a lot of popularity due to its open source code. KAZE was originally made by Pablo F. Alcantarilla, Adrien Bartoli and Andrew J. Davison.[45]

## See also

- Convolutional neural network
- Image stitching
- Scale space
- Scale space implementation
- Simultaneous localization and mapping
- Structure from motion

## References

1. Lowe, David G. (1999). "Object recognition from local scale-invariant features" (http://www.cs.ubc.ca/~lowe/papers/iccv99.pdf) (PDF). *Proceedings of the International Conference on Computer Vision*. Vol. 2. pp. 1150–1157. doi:10.1109/ICCV.1999.790410 (https://doi.org/10.1109%2FICCV.1999.790410).

2. Lowe, David G. (2004). "Distinctive Image Features from Scale-Invariant Keypoints" (http://citeseer.ist.psu.edu/lowe04distinctive.html). *International Journal of Computer Vision*. **60** (2): 91–110. CiteSeerX 10.1.1.73.2924 (https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.73.2924). doi:10.1023/B:VISI.0000029664.99615.94 (https://doi.org/10.1023%2FB%3AVISI.0000029664.99615.94). S2CID 221242327 (https://api.semanticscholar.org/CorpusID:221242327).

3. U.S. Patent 6,711,293 (https://patents.google.com/patent/US6711293), "Method and apparatus for identifying scale invariant features in an image and use of same for locating an object in an image", David Lowe's patent for the SIFT algorithm, March 23, 2004

4. Koenderink, Jan and van Doorn, Ans: "Representation of local geometry in the visual system (http://www-prima.inrial pes.fr/perso/Tran/Documents/Articles/Divers/koenderink87.pdf)", Biological Cybernetics, vol 3, pp 383-396, 1987

5. Koenderink, Jan and van Doorn, Ans: "Generic neighbourhood operators", IEEE Transactions on Pattern Analysis and Machine Intelligence, vol 14, pp 597-605, 1992

6. Lindeberg, Tony (December 1, 2013). "A computational theory of visual receptive fields" (https://doi.org/10.1007/s004 22-013-0569-z). *Biological Cybernetics*. **107** (6): 589–635. doi:10.1007/s00422-013-0569-z (https://doi.org/10.1007% 2Fs00422-013-0569-z). PMC 3840297 (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3840297). PMID 24197240 (h ttps://pubmed.ncbi.nlm.nih.gov/24197240) – via Springer Link.

7. Lindeberg, Tony (2013). "T. Generalized axiomatic scale-space theory" (https://www.sciencedirect.com/science/articl e/pii/B9780124077010000017). In Hawkes, Peter W. (ed.). *Advances in Imaging and Electron Physics* (http://urn.kb.s e/resolve?urn=urn:nbn:se:kth:diva-118695). Vol. 178. Elsevier. pp. 1–96. doi:10.1016/b978-0-12-407701-0.00001-7 (https://doi.org/10.1016%2Fb978-0-12-407701-0.00001-7). ISBN 9780124077010 – via ScienceDirect.

8. Lindeberg, Tony (July 19, 2013). "Invariance of visual operations at the level of receptive fields" (https://www.ncbi.nl m.nih.gov/pmc/articles/PMC3716821). *PLOS ONE*. **8** (7): e66990. arXiv:1210.0754 (https://arxiv.org/abs/1210.0754). Bibcode:2013PLoSO...866990L (https://ui.adsabs.harvard.edu/abs/2013PLoSO...866990L). doi:10.1371/journal.pone.0066990 (https://doi.org/10.1371%2Fjournal.pone.0066990). PMC 3716821 (https://www.n cbi.nlm.nih.gov/pmc/articles/PMC3716821). PMID 23894283 (https://pubmed.ncbi.nlm.nih.gov/23894283).

9. T. Lindeberg (2014) "Scale selection", Computer Vision: A Reference Guide, (K. Ikeuchi, Editor), Springer, pages 701-713. (http://kth.diva-portal.org/smash/record.jsf?pid=diva2%3A1392833&dswid=-9445)

10. Lindeberg, T., Scale-Space Theory in Computer Vision, Kluwer Academic Publishers, 1994 (http://www.csc.kth.se/~to ny/book.html),ISBN 0-7923-9418-6

11. Lindeberg, Tony (1998). "Feature detection with automatic scale selection" (http://kth.diva-portal.org/smash/record.js f?pid=diva2%3A453064&dswid=5931). *International Journal of Computer Vision*. **30** (2): 79–116. doi:10.1023/A:1008045108935 (https://doi.org/10.1023%2FA%3A1008045108935). S2CID 723210 (https://api.sema nticscholar.org/CorpusID:723210).

12. Lindeberg, Tony (2012). "Scale invariant feature transform" (https://doi.org/10.4249%2Fscholarpedia.10491). *Scholarpedia*. **7** (5): 10491. Bibcode:2012SchpJ...710491L (https://ui.adsabs.harvard.edu/abs/2012SchpJ...710491L) . doi:10.4249/scholarpedia.10491 (https://doi.org/10.4249%2Fscholarpedia.10491).

13. Serre, T., Kouh, M., Cadieu, C., Knoblich, U., Kreiman, G., Poggio, T., "A Theory of Object Recognition: Computations and Circuits in the Feedforward Path of the Ventral Stream in Primate Visual Cortex (http://cbcl.mit.ed u/projects/cbcl/publications/ai-publications/2005/AIM-2005-036.pdf)", Computer Science and Artificial Intelligence Laboratory Technical Report, December 19, 2005 MIT-CSAIL-TR-2005-082.

14. Beis, J.; Lowe, David G. (1997). "Shape indexing using approximate nearest-neighbour search in high-dimensional spaces" (http://www.cs.ubc.ca/~lowe/papers/cvpr97.pdf) (PDF). *Conference on Computer Vision and Pattern Recognition, Puerto Rico: sn*. pp. 1000–1006. doi:10.1109/CVPR.1997.609451 (https://doi.org/10.1109%2FCVPR.19 97.609451).

15. Lowe, D.G., Local feature view clustering for 3D object recognition (http://www.cis.rit.edu/~cnspci/references/dip/feat ure_extraction/lowe2001.pdf). IEEE Conference on Computer Vision and Pattern Recognition, Kauai, Hawaii, 2001, pp. 682-688.

16. Lindeberg, Tony & Bretzner, Lars (2003). *Real-time scale selection in hybrid multi-scale representations* (http://kth.div a-portal.org/smash/record.jsf?pid=diva2%3A440700&dswid=7232). *Proc. Scale-Space'03, Springer Lecture Notes in Computer Science*. Vol. 2695. pp. 148–163. doi:10.1007/3-540-44935-3_11 (https://doi.org/10.1007%2F3-540-44935 -3_11). ISBN 978-3-540-40368-5.

17. Lars Bretzner, Ivan Laptev, Tony Lindeberg "Hand gesture recognition using multi-scale colour features, hierarchical models and particle filtering" (http://kth.diva-portal.org/smash/record.jsf?pid=diva2%3A462620&dswid=608), Proceedings of the Fifth IEEE International Conference on Automatic Face and Gesture Recognition, Washington, DC, USA, 21–21 May 2002, pages 423-428. ISBN 0-7695-1602-5, doi:10.1109/AFGR.2002.1004190 (https://doi.org/ 10.1109%2FAFGR.2002.1004190)

18. Kirchner, Matthew R. "Automatic thresholding of SIFT descriptors (https://arxiv.org/abs/1811.03173)." In *Image Processing (ICIP), 2016 IEEE International Conference on*, pp. 291-295. IEEE, 2016.

19. Mikolajczyk, K.; Schmid, C. (2005). "A performance evaluation of local descriptors" (http://research.microsoft.com/us ers/manik/projects/trade-off/papers/MikolajczykPAMI05.pdf) (PDF). *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **27** (10): 1615–1630. CiteSeerX 10.1.1.230.255 (https://citeseerx.ist.psu.edu/viewdoc/summar y?doi=10.1.1.230.255). doi:10.1109/TPAMI.2005.188 (https://doi.org/10.1109%2FTPAMI.2005.188). PMID 16237996 (https://pubmed.ncbi.nlm.nih.gov/16237996).

20. "TU-chemnitz.de" (http://www.tu-chemnitz.de/etit/proaut/rsrc/iav07-surf.pdf) (PDF).

21. Lindeberg, Tony (May 1, 2015). "Image Matching Using Generalized Scale-Space Interest Points" (https://doi.org/10.1007/s10851-014-0541-0). *Journal of Mathematical Imaging and Vision*. **52** (1): 3–36. doi:10.1007/s10851-014-0541-0 (https://doi.org/10.1007%2Fs10851-014-0541-0). S2CID 254657377 (https://api.semanticscholar.org/CorpusID:254657377) – via Springer Link.

22. Edouard Oyallon, Julien Rabin, "An Analysis and Implementation of the SURF Method, and its Comparison to SIFT (http://www.ipol.im/pub/pre/69/)", Image Processing On Line

23. Cui, Y.; Hasler, N.; Thormaehlen, T.; Seidel, H.-P. (July 2009). "Scale Invariant Feature Transform with Irregular Orientation Histogram Binning" (http://www.mpi-inf.mpg.de/~hasler/download/CuiHasThoSei09igSIFT.pdf) (PDF). *Proceedings of the International Conference on Image Analysis and Recognition (ICIAR 2009)*. Halifax, Canada: Springer.

24. Matthew Toews; William M. Wells III (2009). "SIFT-Rank: Ordinal Descriptors for Invariant Feature Correspondence" (http://www.matthewtoews.com/papers/cvpr09-matt.final.pdf) (PDF). *IEEE International Conference on Computer Vision and Pattern Recognition*. pp. 172–177. doi:10.1109/CVPR.2009.5206849 (https://doi.org/10.1109%2FCVPR.2009.5206849).

25. Beril Sirmacek & Cem Unsalan (2009). "Urban Area and Building Detection Using SIFT Keypoints and Graph Theory". *IEEE Transactions on Geoscience and Remote Sensing*. **47** (4): 1156–1167. Bibcode:2009ITGRS..47.1156S (https://ui.adsabs.harvard.edu/abs/2009ITGRS..47.1156S). doi:10.1109/TGRS.2008.2008440 (https://doi.org/10.1109%2FTGRS.2008.2008440). S2CID 6629776 (https://api.semanticscholar.org/CorpusID:6629776).

26. Se, S.; Lowe, David G.; Little, J. (2001). "Vision-based mobile robot localization and mapping using scale-invariant features" (http://citeseer.ist.psu.edu/425735.html). *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*. Vol. 2. p. 2051. doi:10.1109/ROBOT.2001.932909 (https://doi.org/10.1109%2FROBOT.2001.932909).

27. Fabbri, Ricardo; Duff, Timothy; Fan, Hongyi; Regan, Margaret; de Pinho, David; Tsigaridas, Elias; Wampler, Charles; Hauenstein, Jonathan; Kimia, Benjamin; Leykin, Anton; Pajdla, Tomas (23 Mar 2019). "Trifocal Relative Pose from Lines at Points and its Efficient Solution". arXiv:1903.09755 (https://arxiv.org/abs/1903.09755) [cs.CV (https://arxiv.org/archive/cs.CV)].

28. Fabbri, Ricardo; Giblin, Peter; Kimia, Benjamin (2012). "Camera Pose Estimation Using First-Order Curve Differential Geometry" (https://rfabbri.github.io/stuff/fabbri-giblin-kimia-eccv2012-final-ext.pdf) (PDF). *Lecture Notes in Computer Science (ECCV 2012)*. Lecture Notes in Computer Science. **7575**: 231–244. doi:10.1007/978-3-642-33765-9_17 (https://doi.org/10.1007%2F978-3-642-33765-9_17). ISBN 978-3-642-33764-2. S2CID 15402824 (https://api.semanticscholar.org/CorpusID:15402824).

29. Brown, M.; Lowe, David G. (2003). "Recognising Panoramas" (http://graphics.cs.cmu.edu/courses/15-463/2005_fall/www/Papers/BrownLowe.pdf) (PDF). *Proceedings of the ninth IEEE International Conference on Computer Vision*. Vol. 2. pp. 1218–1225. doi:10.1109/ICCV.2003.1238630 (https://doi.org/10.1109%2FICCV.2003.1238630).

30. Iryna Gordon and David G. Lowe, "What and where: 3D object recognition with accurate pose (http://www.cs.ubc.ca/labs/lci/papers/docs2006/lowe_gordon.pdf)," in Toward Category-Level Object Recognition, (Springer-Verlag, 2006), pp. 67-82

31. Flitton, G.; Breckon, T. (2010). "Object Recognition using 3D SIFT in Complex CT Volumes" (http://www.durham.ac.uk/toby.breckon/publications/papers/flitton10baggage.pdf) (PDF). *Proceedings of the British Machine Vision Conference*. pp. 11.1–12. doi:10.5244/C.24.11 (https://doi.org/10.5244%2FC.24.11).

32. Flitton, G.T., Breckon, T.P., Megherbi, N. (2013). "A Comparison of 3D Interest Point Descriptors with Application to Airport Baggage Object Detection in Complex CT Imagery". *Pattern Recognition*. **46** (9): 2420–2436. Bibcode:2013PatRe..46.2420F (https://ui.adsabs.harvard.edu/abs/2013PatRe..46.2420F). doi:10.1016/j.patcog.2013.02.008 (https://doi.org/10.1016%2Fj.patcog.2013.02.008). hdl:1826/15213 (https://hdl.handle.net/1826%2F15213).

33. Laptev, Ivan & Lindeberg, Tony (2004). "Local descriptors for spatio-temporal recognition" (http://kth.diva-portal.org/smash/record.jsf?pid=diva2%3A445261&dswid=2390). *ECCV'04 Workshop on Spatial Coherence for Visual Motion Analysis, Springer Lecture Notes in Computer Science, Volume 3667*. pp. 91–103. doi:10.1007/11676959_8 (https://doi.org/10.1007%2F11676959_8).

34. Ivan Laptev, Barbara Caputo, Christian Schuldt and Tony Lindeberg (2007). "Local velocity-adapted motion events for spatio-temporal recognition" (http://kth.diva-portal.org/smash/record.jsf?pid=diva2%3A335153&dswid=8365). *Computer Vision and Image Understanding*. **108** (3): 207–229. CiteSeerX 10.1.1.168.5780 (https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.168.5780). doi:10.1016/j.cviu.2006.11.023 (https://doi.org/10.1016%2Fj.cviu.2006.11.023).

35. Scovanner, Paul; Ali, S; Shah, M (2007). "A 3-dimensional sift descriptor and its application to action recognition". *Proceedings of the 15th International Conference on Multimedia*. pp. 357–360. doi:10.1145/1291233.1291311 (https://doi.org/10.1145%2F1291233.1291311).

36. Niebles, J. C. Wang, H. and Li, Fei-Fei (2006). "Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words" (http://vision.cs.princeton.edu/niebles/humanactions.htm). *Proceedings of the British Machine Vision Conference (BMVC)*. Edinburgh. Retrieved 2008-08-20.

37. Matthew Toews; William M. Wells III; D. Louis Collins; Tal Arbel (2010). "Feature-based Morphometry: Discovering Group-related Anatomical Patterns" (http://www.matthewtoews.com/papers/matt_neuroimage10.pdf) (PDF). *NeuroImage*. **49** (3): 2318–2327. doi:10.1016/j.neuroimage.2009.10.032 (https://doi.org/10.1016%2Fj.neuroimage.2009.10.032). PMC 4321966 (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4321966). PMID 19853047 (https://pubmed.ncbi.nlm.nih.gov/19853047).

38. Lazebnik, S., Schmid, C., and Ponce, J., "Semi-Local Affine Parts for Object Recognition (http://hal.archives-ouvertes.fr/docs/00/54/85/42/PDF/bmvc04.pdf)", Proceedings of the British Machine Vision Conference, 2004.

39. Arandjelović, Relja; Zisserman, Andrew (2012). "Three things everyone should know to improve object retrieval". *2012 IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2911–2918. doi:10.1109/CVPR.2012.6248018 (https://doi.org/10.1109%2FCVPR.2012.6248018).

40. Sungho Kim, Kuk-Jin Yoon, In So Kweon, "Object Recognition Using a Generalized Robust Invariant Feature and Gestalt's Law of Proximity and Similarity", Conference on Computer Vision and Pattern Recognition Workshop (CVPRW'06), 2006

41. Bay, H., Tuytelaars, T., Van Gool, L., "SURF: Speeded Up Robust Features (http://www.vision.ee.ethz.ch/~surf/eccv06.pdf)", Proceedings of the ninth European Conference on Computer Vision, May 2006.

42. Ke, Y., and Sukthankar, R., "PCA-SIFT: A More Distinctive Representation for Local Image Descriptors (https://www.cs.cmu.edu/~rahuls/pub/cvpr2004-keypoint-rahuls.pdf)", Computer Vision and Pattern Recognition, 2004.

43. D. Wagner, G. Reitmayr, A. Mulloni, T. Drummond, and D. Schmalstieg, "Pose tracking from natural features on mobile phones (http://mi.eng.cam.ac.uk/~gr281/docs/WagnerIsmar08NFT.pdf) Archived (https://web.archive.org/web/20090612124616/http://mi.eng.cam.ac.uk/~gr281/docs/WagnerIsmar08NFT.pdf) 2009-06-12 at the Wayback Machine" Proceedings of the International Symposium on Mixed and Augmented Reality, 2008.

44. N. Henze, T. Schinke, and S. Boll, "What is That? Object Recognition from Natural Features on a Mobile Phone (http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.149.9089&rep=rep1&type=pdf)" Proceedings of the Workshop on Mobile Interaction with the Real World, 2009.

45. "kaze" (http://www.robesafe.com/personal/pablo.alcantarilla/kaze.html). *www.robesafe.com*.

# External links

**Related studies:**

- The Invariant Relations of 3D to 2D Projection of Point Sets, Journal of Pattern Recognition Research (http://www.jprr.org/index.php/jprr/article/view/26)(JPRR) (http://www.jprr.org/), Vol. 3, No 1, 2008.

- Lowe, D. G., "Distinctive Image Features from Scale-Invariant Keypoints", International Journal of Computer Vision, 60, 2, pp. 91-110, 2004. (http://citeseer.ist.psu.edu/lowe04distinctive.html)

- Mikolajczyk, K., and Schmid, C., "A performance evaluation of local descriptors", IEEE Transactions on Pattern Analysis and Machine Intelligence, 10, 27, pp 1615--1630, 2005. (http://lear.inrialpes.fr/pubs/2005/MS05/)

- Andrea Maricela Plaza Cordero, Jorge Luis Zambrano-Martinez, " Estudio y Selección de las Técnicas SIFT, SURF y ASIFT de Reconocimiento de Imágenes para el Diseño de un Prototipo en Dispositivos Móviles (http://41jaiio.sadio.org.ar/sites/default/files/6_EST_2012.pdf)" , 15º Concurso de Trabajos Estudiantiles, EST 2012

- "PCA-SIFT: A More Distinctive Representation for Local Image Descriptors" (https://web.archive.org/web/20200126202031/https://www.cs.cmu.edu/~yke/pcasift/). Archived from the original (https://www.cs.cmu.edu/~yke/pcasift/) on 26 January 2020.

- Lazebnik, S., Schmid, C., and Ponce, J., Semi-Local Affine Parts for Object Recognition, BMVC, 2004. (http://www-cvr.ai.uiuc.edu/ponce_grp/publication/paper/bmvc04.pdf)

**Tutorials:**

- Scale-Invariant Feature Transform (SIFT) in Scholarpedia (http://www.scholarpedia.org/article/SIFT)

- A simple step by step guide to SIFT (http://www.aishack.in/tutorials/sift-scale-invariant-feature-transform-introduction/)

- "SIFT for multiple object detection" (https://web.archive.org/web/20150403120732/http://www.berilsirmacek.com/sift_multiple_object_detection.html). Archived from the original (http://www.berilsirmacek.com/sift_multiple_object_detection.html) on 3 April 2015.

- "The Anatomy of the SIFT Method (http://www.ipol.im/pub/pre/82/)" in Image Processing On Line, a detailed study of every step of the algorithm with an open source implementation and a web demo to try different parameters

**Implementations:**

- Rob Hess's implementation of SIFT (https://robwhess.github.com/opensift/) accessed 21 Nov 2012
- ASIFT (Affine SIFT) (http://www.ipol.im/pub/algo/my_affine_sift/): large viewpoint matching with SIFT, with source code and online demonstration
- VLFeat (http://www.vlfeat.org/api/sift.html), an open source computer vision library in C (with a MEX interface to MATLAB), including an implementation of SIFT
- LIP-VIREO (http://pami.xmu.edu.cn/~wlzhao/lip-vireo.htm), A toolkit for keypoint feature extraction (binaries for Windows, Linux and SunOS), including an implementation of SIFT
- (Parallel) SIFT in C# (https://sites.google.com/site/btabibian/projects/3d-reconstruction/code), SIFT algorithm in C# using Emgu CV and also a modified parallel version of the algorithm.
- DoH & LoG + affine (http://www.mathworks.com/matlabcentral/fileexchange/38782), Blob detector adapted from a SIFT toolbox
- ezSIFT: an easy-to-use standalone SIFT implementation in C/C++ (https://github.com/robertwgh/ezSIFT). A self-contained open-source SIFT implementation which does not require other libraries.
- A 3D SIFT implementation: detection and matching in volumetric images. (http://www.matthewtoews.com/fba/featExtract1.3.zip)

---

Retrieved from "https://en.wikipedia.org/w/index.php?title=Scale-invariant_feature_transform&oldid=1155989186"

-