*Research Article*
# LUIFT: LUminance Invariant Feature Transform

**Julia Diaz-Escobar** [iD],[1] **Vitaly Kober** [iD],[1,2] **and Jose A. Gonzalez-Fraga**[3]

[1]*Department of Computer Science, CICESE, Carretera Tijuana-Ensenada, Playitas 3918, Ensenada, B.C., Mexico*
[2]*Department of Mathematics, Chelyabinsk State University, Russia*
[3]*Universidad Autónoma de Baja California, Carretera Tijuana-Ensenada, Playitas 3917, Ensenada, B.C., Mexico*

Correspondence should be addressed to Julia Diaz-Escobar; jdiaz@cicese.edu.mx

Illumination-invariant method for computing local feature points and descriptors, referred to as LUminance Invariant Feature Transform (LUIFT), is proposed. The method helps us to extract the most significant local features in images degraded by nonuniform illumination, geometric distortions, and heavy scene noise. The proposed method utilizes image phase information rather than intensity variations, as most of the state-of-the-art descriptors. Thus, the proposed method is robust to nonuniform illuminations and noise degradations. In this work, we first use the monogenic scale-space framework to compute the local phase, orientation, energy, and phase congruency from the image at different scales. Then, a modified Harris corner detector is applied to compute the feature points of the image using the monogenic signal components. The final descriptor is created from the histograms of oriented gradients of phase congruency. Computer simulation results show that the proposed method yields a superior feature detection and matching performance under illumination change, noise degradation, and slight geometric distortions comparing with that of the state-of-the-art descriptors.

## 1. Introduction

Feature detection and description are low-level tasks used in many computer vision and pattern recognition applications such as image classification and retrieval [1, 2], optical flow estimation [3], tracking [4], biometric systems [5], image registration [6], and 3D reconstruction [7].

The local feature detection task consists of finding *"feature points"* (points, lines, blobs, etc.) in the image. The points should satisfy certain properties such as distinctiveness, quantity, locality, accuracy, and more important repeatability [8]. To represent each feature point in a distinctive way, a neighborhood around each feature is considered and encoded into a vector, known as *"feature descriptor."* The feature descriptors of different images are *"matched"* using either Euclidean or Mahalanobis distances.

It is desirable that the behavior of feature descriptors be invariant to viewpoint changes, blur effect, and affine transformations [9–13]; but also, it needs to be robust to noise and nonuniform illumination degradations. However, these last two conditions have not been completely solved, even when

they are common issues in real-world applications. Thus, the nonuniform illumination variations and noise degradations are still challenges that decrease the performance of the existing state-of-the-art methods.

Since Attneave research [14] about the importance of the image shape information, several techniques for feature detection have been developed [8, 15–17]. Many of the existing works are robust to affine transformations (scale and rotations), but they are not designed to work with complex illumination changes. Recently, to address the nonuniform illumination problem, different methods based on the order of the intensity values have been proposed [18–21]. However, these methods are only robust to monotonic intensity variations and are sensitive to heavy noise degradations.

On the other hand, the human visual system is able to recognize objects under different illumination conditions. The human eye perceives an amount of light energy passes through, reflected or emitted from an object surface, known as *luminance*. It converts the light energy into nerve impulses by the photoreceptor cells in the retina, where the information is encoded and sent to the primary visual cortex

(V1) [22]. Psychophysical evidence suggests that the human visual system decomposes the visual information in borders and lines components by using phase information. Besides, it is known that different groups of cells in V1 extract particular image features as frequency, orientation, and phase information [23].

In this work, to overcome the luminance variation problem inspired by the human visual system, a phase-based method for computing local feature points and descriptors, referred to as LUminance Invariant Feature Transform (LUIFT), is proposed. The LUIFT method helps us to extract the most significant local features in images degraded by nonuniform illumination, geometric distortions, and heavy scene noise. The proposed technique is suitable for recognition of rigid objects under real conditions. The LUIFT algorithm was extensively tested on common databases. The proposed method yields a competitive matching performance under slight scaling and in-plane rotation with that of the state-of-the-art algorithms. The LUIFT method shows improved performance regarding the feature points repeatability as well as the number of detected and matched feature descriptors under illumination changes and noise degradations.

The rest of this paper is organized as follows. In Section 2, the related works are recalled. In Section 3, the phase-based approach is described. In Section 4, the proposed LUIFT detector and descriptor are presented. In Section 5, computer simulation results are provided and discussed. Finally, Section 6 summarizes our conclusions.

## 2. Related Work

Early works on image feature points began with the research of Attneave [14], showing that the most important shape information of an image is concentrated at the contour points with high curvature values, such as corners and junctions. Since then, several techniques for features detection have been developed, such as contour curvature based methods [8, 24], blob-like detector techniques [16], differential approaches [8, 17], intensity variations based techniques [25, 26], and recently learning-based methods [27–29].

The Harris corner detector [30], which is an improvement of the Moravec approach [31], is one of the first and most used corner detectors, which describes the gradient distribution in a local neighborhood of a point based on the second-moment matrix. The feature points are obtained at the points where the local gradient varies significantly in two directions. Similarly to the Harris matrix, the Hessian matrix [32] is constructed by the second-order Taylor expansion of the intensity surface and encodes the shape information of the image. Recently, a Harris-based (HarrisZ) corner detector was proposed [33]. The HarrisZ corner detector considers a z-score to adapt the corner response function, searching the corners near to edges by a coarse gradient mask.

SUSAN (Smallest Univalue Segment Assimilating Nucleus) [25] and, more recently, FAST (Features from Accelerated Segment Test) [26] corner detectors are also intensity-based techniques. They obtain fast feature points associating to image points in a local area with similar brightness. The FAST detector is based on the SUSAN detector, but it uses more efficient decision trees to evaluate intensity pixel values.

The SIFT (Scale Invariant Feature Transform) descriptor [9, 34] utilizes an approximation of the LoG (Laplacian of Gaussian) and HOG (Histograms of Oriented Gradient) [35] for scale and rotation invariance, respectively. Until now, the SIFT descriptor is the most popular state-of-the-art descriptor due to its effectiveness in the feature detection and matching under scale and rotation image changes. That is why different variations of the SIFT descriptor have been proposed. The SURF [11, 36] (Speed Up Robust Features) and the KAZE [12] descriptors are a couple of examples. Unlike the SIFT method, the SURF descriptor uses Haar-like filters and integral images to improve the processing time at the expense of the method performance; meanwhile, the KAZE descriptor is based on nonlinear scale space improving the locally adaptive blurring on the nonlinear scale-space construction. The CenSurE [37] (Center Surround Extremas) feature detector is based on the estimation of the LoG (Laplacian of Gaussian) using simple center-surround filters and integral images for real-time tasks. The Daisy descriptor [10] is inspired by the SIFT and GLOH [17] descriptors but computed more efficiently replacing weighted sums by sums of convolutions.

Binary descriptors have also been suggested. FREAK (Fast Retina Keypoints) [38], BRIEF (Binary Robust Independent Elementary Features) [39], and BRISK (Binary Robust Invariant Scalable Keypoints) [40] are some of them. Basically, they carry out pairwise intensity comparisons within an image patch and use the Hamming distance for fast feature matching.

Although all mentioned methods provide satisfactory results for affine image transformations (rotation and scale), they are usually constructed on the base of differences between the pixel intensities of the image, which makes them sensitive to nonuniform illumination variation and noise degradation. To obtain robust descriptors to intensity variations, new methods have been proposed. The DaLI [27] (Deformation and Light Invariant descriptor) descriptor was developed for nonrigid transformations and illumination changes. The 2D image patches are considered as 3D surfaces and described in terms of a heat kernel signature. Then, for descriptor dimensional reduction a Principal Component Analysis (PCA) is applied. However, DaLI descriptor is not invariant to scale and rotation distortions and has a high complexity due to the computation of eigenvalues for the heat diffusion equation. The TILDE [13] (Temporally Invariant Learned DEtector) and the LIFT [28] (Learned Invariant Feature Transform) methods consider a learned method for feature detection and description. Basically, the detector uses training to obtain those features that remain stable under different conditions. However, a prestage of training and a collection of image patches are needed. The LIOP [21] (Local Intensity Order Pattern) descriptor is based on the intensity values order, assuming the principle that the relative order of pixel intensities remains unchanged with monotonic intensity changes. However, nonuniform illumination variations are not considered.

In this work, we propose a phase-based feature detector and descriptor. Unlike the mentioned above methods, the proposed technique utilizes the image local phase information instead of relying on the image pixels intensities changes. So, there are two main contributions of the proposed work: first, since the local phase contains the most important image information and it is invariant to image pixel intensities [41], the proposed method is robust to nonuniform illumination variations; second, since the proposed method utilizes the local phase congruency approach rather than only image gradients, it is robust to heavy noise degradations.

## 3. Phase-Based Signal Model

Ever since the Hubel and Wiesel work [42], it has been known that different groups of neurons in the biological visual cortex, called simple cells, respond selectively to bars and edges at particular orientation and location. Furthermore, psychophysical evidence suggests the existence of the frequency-selective V1 neurons operating as bandpass filters and the computation of complex cells energies as a sum of squared responses of simple cells (see [23]).

Morrone and Owens proposed a model of feature perception such as edges, lines, and shadows called the *local energy model* [43–45]. According to this model, the human visual system is capable to determinate a square waveform and a trapezoid by using phase information, and it can be proved that the maximum of the energy function occurs at the points of the maximum *phase congruency* [46]. Continuing with this approach, Kovesi [47–49] proposed a dimensionless measure of phase congruency at each point of an image, where the phase congruency value indicates the significance of the current feature; that is, unity means the most significant feature, and zero indicates the lowest significance.

Felsberg et al. [50] provided a framework to obtain features based on the phase of an image. Unlike other works, they did not use steerable filters, such as Gabor filters, to get the image features. Instead, they proposed a new concept of a two-dimensional analytic signal, referred to as *the monogenic signal* [51].

*3.1. Local Energy Model and Phase Congruency Approach.* The local energy model [44, 45] establishes that the visual system could locate features by searching for maxima of local energy and identifies the feature type by evaluating the argument at that point.

Formally, let the pair of filters $H_e \in L^2$ and $H_o \in L^2$ be the basic operators of the model with equal magnitude spectra but with orthogonal phases (here $H_o$ denotes the Hilbert Transform of $H_e$). The local energy function is defined as

$$E(x) = \sqrt{\left(H_e(x) * f(x)\right)^2 + \left(H_o(x) * f(x)\right)^2}, \quad (1)$$

where $f(x) \in L^2$ is a periodic signal, and $(*)$ is the convolution operator.

The local energy function locates the position of image features but it has no information about the feature type. To determine the feature type, it is necessary to consider the argument defined as follows:

$$\phi(x) = \tan^{-1}\left(H_e(x) * f(x), H_o(x) * f(x)\right). \quad (2)$$

On the other hand, a periodic function, $f(x) \in L^2$, can be expanded in its Fourier components as follows:

$$f(x) = \sum_n A_n \cos\left(\varphi_n(x)\right), \quad (3)$$

where $A_n$ and $\varphi_n(x) = n\omega x + \phi(x)$ represent the magnitude and the local phase of the $n$th Fourier component, respectively. *The phase congruency function* is defined as follows [44]:

$$PC(x) = \max_{\overline{\varphi}(x) \in [0, 2\pi]} \frac{\sum_n A_n \cos\left(\varphi_n(x) - \overline{\varphi}(x)\right)}{\sum_n A_n}, \quad (4)$$

where $\overline{\varphi}(x)$ is the weighted mean local phase angle of all Fourier components at the point $x$ and $0 \leq PC(x) \leq 1$. The congruency of phase at any angle produces a local feature. A phase congruency value of one means that most of the Fourier components phases are similar ($\varphi_n(x) - \overline{\varphi}(x) \cong 0$) and; therefore, there exists a local feature (edge or line), while a phase congruency of value zero indicates the lack of structure. Besides, the value of $\overline{\varphi}(x)$ determines the nature of the feature: values near to zero and $\pi$ correspond to a line feature, and values near to $\pi/2$ and $3\pi/2$ correspond to an edge feature.

Unfortunately, the $PC(x)$ function is highly sensitive to noise and frequency spread. To overcome this problem, the following definition of the phase congruency function was proposed [47]:

$$PC(x) = \frac{W(x) \lfloor E(x) - T \rfloor}{\sum_n A_n(x) + \varepsilon}, \quad (5)$$

where $W(x)$ is a weight for the frequency spread, $E(x)$ represents the signal energy, $T$ is a noise threshold parameter, and $\varepsilon$ is a small constant to avoid division by zero. We refer to the following papers [47–49] for more details.

In practice, local frequency information is obtained via banks of oriented 2D Gabor filters, but this procedure is computationally expensive. Recently, Felsberg and Sommer [52] proposed the monogenic signal, which is a generalization of the 1D analytic signal. It gives us a theoretical framework to obtain local frequency information.

*3.2. The Monogenic Signal.* The monogenic signal [52] is defined as a combination of the 2D signal and its first-order Riesz transform, defined as follows.

Let $R_x$, $R_y$ be the transfer functions of the first-order Riesz transform in the frequency domain:

$$R_x(u, v) = i\frac{u}{\sqrt{u^2 + v^2}} = \mathscr{F}\left\{\frac{x}{2\pi\left(x^2 + y^2\right)^{3/2}}\right\}, \quad (6)$$

$$R_y(u, v) = i\frac{v}{\sqrt{u^2 + v^2}} = \mathscr{F}\left\{\frac{y}{2\pi\left(x^2 + y^2\right)^{3/2}}\right\}. \quad (7)$$

The monogenic signal $F_M(u, v)$ in the frequency domain is defined as follows:

$$F_M(u, v) = F(u, v) + i\mathbf{R} \cdot F(u, v), \qquad (8)$$

where $F(u, v) = \mathscr{F}\{f(x, y)\}$ is the Fourier transform of $f(x, y)$ and $\mathbf{R} = (\mathrm{R_x}, \mathrm{R_y})$.

In order to perform scale decomposition of a signal into a set of partial signals, it is necessary to calculate the monogenic signal for narrow bandwidths. A good approximation of the scale decomposition can be done by using appropriate bandpass filters to obtain localization in the spatial and frequency domains.

### 3.3. Scale-Space Monogenic Signal.

Felsberg and Sommer [53] defined the linear Poisson scale-space representation as an alternative to the well-known Gaussian scale-space, because it is related to the monogenic signal. The Poisson scale-space $f_p(x, y, s)$ is defined as the convolution of the image $f(x, y)$ with the Poisson kernel, as follows:

$$f_p(x_1, x_2, h) = \frac{s}{2\pi(x^2 + y^2 + s^2)} * f(x, y)$$
$$= \mathscr{F}^{-1}\left\{e^{-2\pi s\sqrt{u^2+v^2}} \cdot F(u, v)\right\}, \qquad (9)$$

where $s$ is the scale parameter that controls the degree of image resolution. The combination of two lowpass filters with a fixed ratio of scale parameters gives us a family of bandpass filters with a constant relative bandwidth, defined as

$$B_{s_0, \lambda, k}(u, v) = \left(e^{-2\pi s_o \lambda^k \sqrt{u^2+v^2}} - e^{-2\pi s_o \lambda^{k-1} \sqrt{u^2+v^2}}\right), \qquad (10)$$

where $\lambda \in (0, 1)$ indicates the relative bandwidth, $s_0$ is the coarsest scale, and $k \in \mathbb{N}$ denotes the bandpass number [54]. The Poisson scale-space representation in the frequency domain of the image $F(u, v)$ filtered by the bandpass filter $B_{s_0, \lambda, k}(u, v)$ is given by

$$F_{bp}(u, v) = \left(e^{-2\pi s_o \lambda^k \sqrt{u^2+v^2}} - e^{-2\pi s_o \lambda^{k-1} \sqrt{u^2+v^2}}\right) \\ \cdot F(u, v). \qquad (11)$$

Then, the Poisson scale-space monogenic signal representation is formed by

$$F_{Mbp}(u, v) = F_{bp}(u, v) + i\mathbf{R} \cdot F_{bp}(u, v), \qquad (12)$$

where

$$f_p(x, y) = \mathscr{F}^{-1}\left\{F_{bp}(u, v)\right\}, \qquad (13)$$

$$f_x(x, y) = \mathscr{F}^{-1}\left\{\mathrm{R}_x(u, v) \cdot F_{bp}(u, v)\right\}, \qquad (14)$$

and

$$f_y(x, y) = \mathscr{F}^{-1}\left\{\mathrm{R}_y(u, v) \cdot F_{bp}(u, v)\right\}, \qquad (15)$$

in the spatial domain.

Therefore, the local energy $E(x, y)$, local orientation $\theta_{or}(x, y)$, local direction $\theta_{dir}(x, y)$, and local phase $\varphi(x, y)$ (Note that the function $\mathrm{atan2}(|y|/x) = \mathrm{sign}(y) \cdot \tan^{-1}(|y|/x)$, where the factor $\mathrm{sign}(y)$ indicates the direction of rotation) can be computed as follows:

$$\begin{aligned} &E(x, y) \\ &= \sqrt{\left(f_p(x, y)\right)^2 + (f_x(x, y))^2 + \left(f_y(x, y)\right)^2}, \end{aligned} \qquad (16)$$

$$\theta_{or}(x, y) = \tan^{-1}\left(\frac{f_y(x, y)}{f_x(x, y)}\right), \qquad (17)$$

$$\theta_{dir}(x, y) = \mathrm{atan2}\left(\frac{f_y(x, y)}{f_x(x, y)}\right), \qquad (18)$$

$$\varphi(x, y) = \tan^{-1}\left(\frac{\sqrt{(f_x(x, y))^2 + \left(f_y(x, y)\right)^2}}{f_p(x, y)}\right). \qquad (19)$$

Figure 1 shows a block diagram for computing the monogenic scale-space signal.

## 4. Proposed Feature Detector and Descriptor

In this section, the proposed LUIFT feature detector and descriptor are described. The feature detector is constructed using a modified Harris corner detector and the phase congruency approach, while the feature descriptor is constructed using a modified HOG-based method.

### 4.1. Feature Detector.

First, using the monogenic scale-space framework (see Figure 1) with a bandpass filter set $\{B_{s_0=3}, \lambda = 0.5, k = 3\}$, the scale-space monogenic signal $f_m = (f_p, f_x, f_y)$ and the sum of amplitudes $\sum_n A_n(x, y)$ are computed. Note that, by increasing the bandpass number $k$, more fine scale features are revealed. The phase congruency function in (5) can be calculated for each point of the image as follows:

$$PC(x, y) = \frac{W(x, y)\lfloor E(x, y) - T\rfloor}{\sum_n A_n(x, y) + \varepsilon}, \qquad (20)$$

where the energy $E(x, y) = \sqrt{f_p^2 + f_x^2 + f_y^2}$ and the sum of the amplitudes $\sum_n A_n(x, y)$ are obtained from the scale-space monogenic signal. The frequency spread weight $W(x, y)$ and the noise threshold $T$ are calculated as in [47].

Next, in order to obtain the feature point candidates, a modified Harris corner detector is utilized.

Let $H$ be the Harris matrix defined by

$$H = \sum_m \sum_n w(m, n)\begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}, \qquad (21)$$

where $I_x$ and $I_y$ are the partial derivatives of the image $I$. Considering the scale-space monogenic signal, the derivatives of the Harris matrix ($H$) are replaced by the monogenic signal components ($T_M$) as follows:

$$T_M = \sum_m \sum_n w(m, n)\begin{bmatrix} \left(f_x'\right)^2 & f_x' f_y' \\ f_x' f_y' & \left(f_y'\right)^2 \end{bmatrix}, \qquad (22)$$
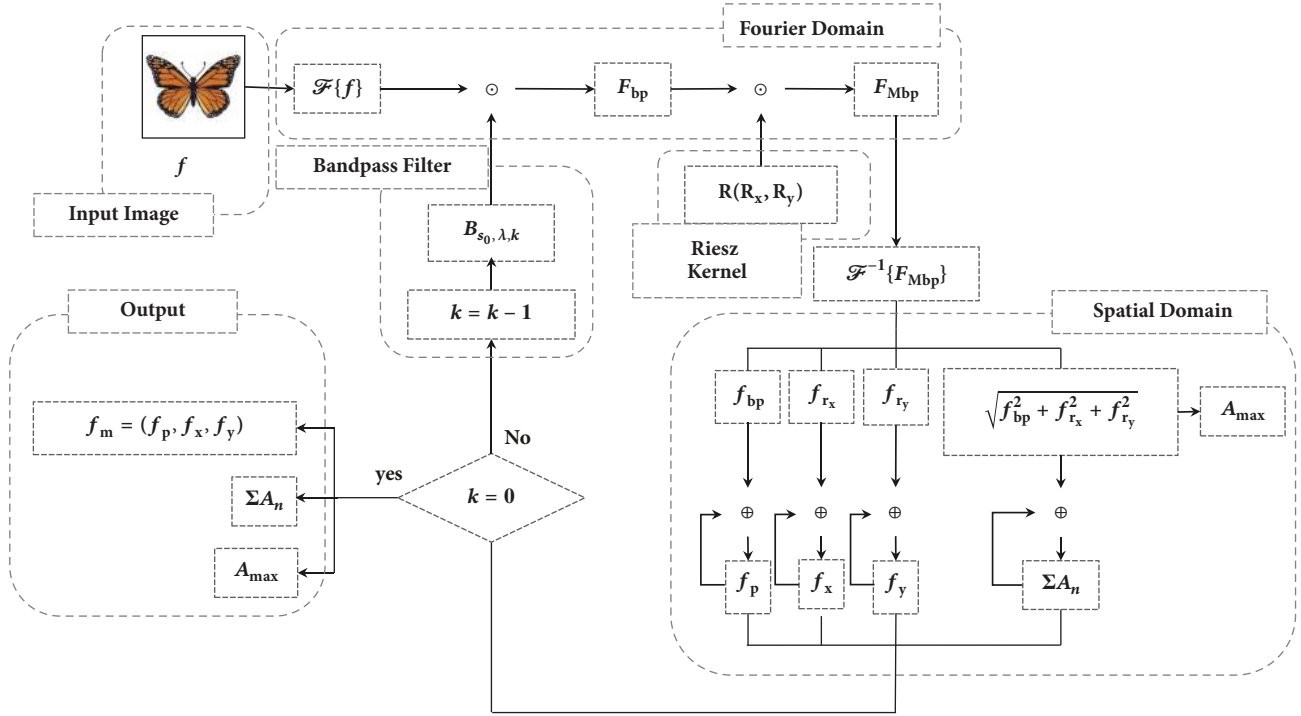
FIGURE 1: Block-diagram for forming the Scale-Space Monogenic Signal.

where

$$f_x'(x, y) = \frac{f_x(x, y)}{\sum_n A_n(x, y)} \qquad (23)$$

and

$$f_y'(x, y) = \frac{f_y(x, y)}{\sum_n A_n(x, y)} \qquad (24)$$

are normalized. Then, the corner detector function defined in [30] is utilized to obtain corner feature candidates,

$$M_c(x, y) = \det(T_M) - \beta \cdot \text{trace}^2(T_M), \qquad (25)$$

where $\beta$ is a sensitivity parameter, commonly used $\beta = 0.04$.

The obtained candidate features $M_c$ are weighted by its corresponding $PC(x, y)$ value, in order to extract feature points with high phase congruency; that is,

$$M_c'(x, y) = M_c(x, y) \cdot PC(x, y). \qquad (26)$$

Then, a thresholding followed by $3 \times 3$ nonmax suppression algorithm is applied to obtain the final feature points. Since the PC value indicates the significance of the detected features (see Section 3.1), the threshold value controls the number of features to be preserved or eliminated. A threshold close to one keeps only those features that belong to sharp lines or borders in the image. By changing the threshold value, important features belonging to borders, and lines with low contrast, high brightness or blur degradations could be preserved. For our experiments, a threshold of 0.3 was experimentally defined. Figure 2 illustrates the performance of the proposed feature detector.

*4.2. Feature Descriptor.* Because the histograms of oriented gradients [35] show robustness to small deformation such scale and rotations, a modified HOG-based descriptor is constructed. For each detected feature point, a $16 \times 16$ spatial neighborhood around each feature is constructed and weighted by a Gaussian kernel ($\sigma = 1.5$). Next, the neighborhood is split onto $4 \times 4$ subneighborhoods. For each sub-neighborhood, *nbins* Histogram of Oriented Phase Congruency (HOPC) is computed using the local direction $\theta_{\text{dir}}(x, y)$ (see (18)) between 0 and 360 degrees in such a manner that the amount added to each bin depends on the $PC(x, y)$ value of each point, as follows:

$$HOPC(bin_\theta) = HOPC(bin_\theta) + PC(x, y), \qquad (27)$$

where

$$bin_\theta = \left\lfloor \left(\frac{nbins}{360}\right) \cdot \theta_{\text{dir}}(x, y) \right\rfloor. \qquad (28)$$

Figure 3 illustrates the formation of the proposed feature descriptor.

Now, let $r_\theta$ be the remainder of the modulus (mod),

$$r_\theta = \theta_{\text{dir}}(x, y) \bmod \left(\frac{360}{nbins}\right). \qquad (29)$$

If either $r_\theta$ or $360/nbins - r_\theta$ are near to zero, it means that $\theta_{\text{dir}}(x, y)$ is near to the border between two adjacent bins. Therefore, $\theta_{\text{dir}}(x, y)$ could be assigned to one of the bins or divided between the bins. So, we assign the half of the $PC(x, y)$ value to each of the adjacent bins.
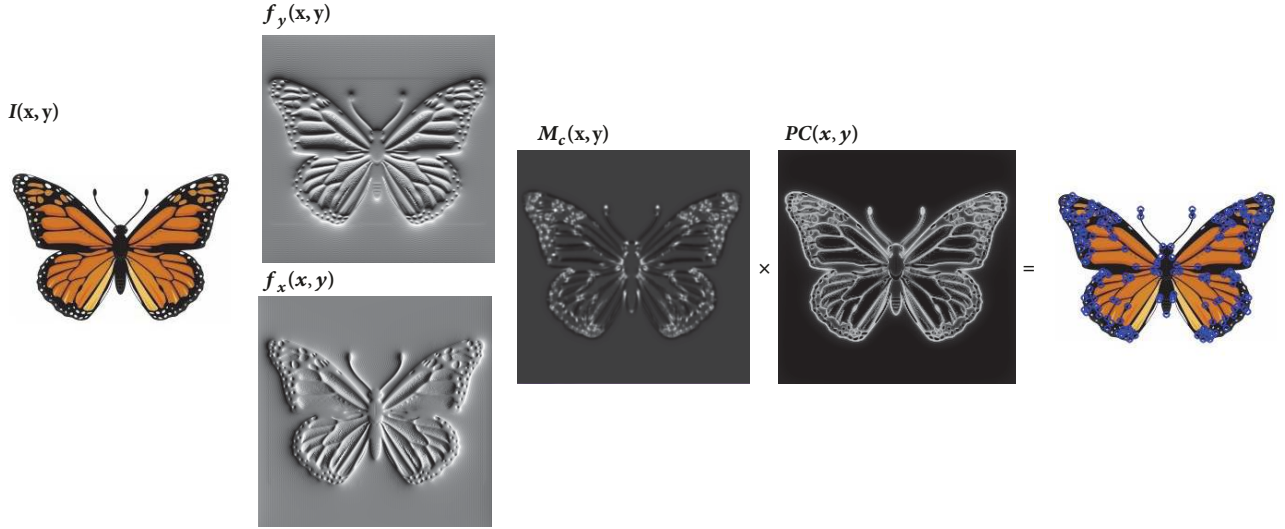
$f_y(x, y)$

$I(x, y)$

$M_c(x, y)$

$PC(x, y)$

$f_x(x, y)$

×

=

Figure 2: Proposed feature detector.

$PC(x, y) \in [0,1]$

$I(x, y) \in [0,255]$

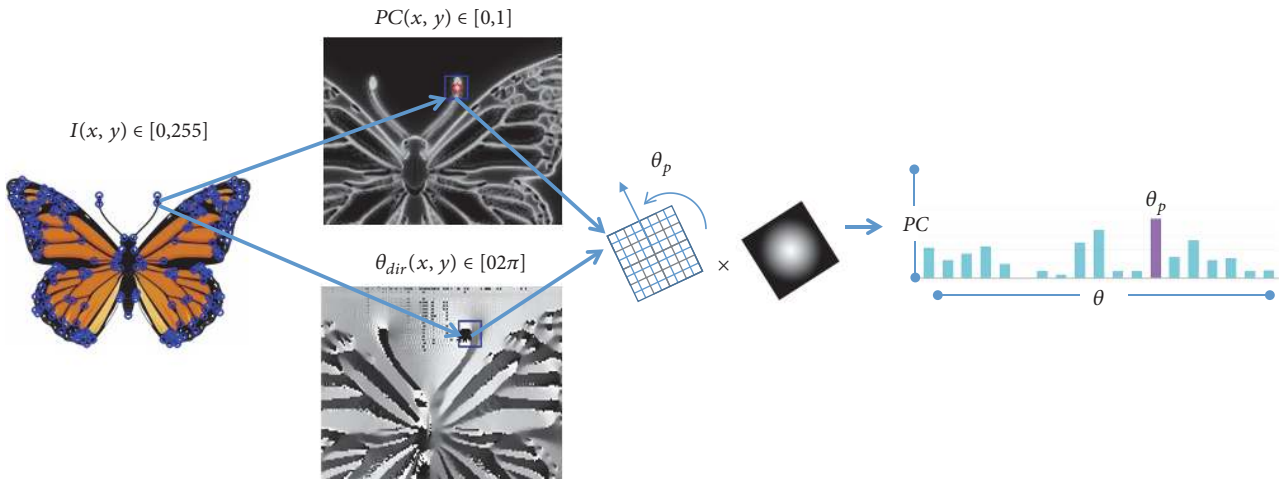$\theta_{dir}(x, y) \in [0 2\pi]$

$\theta_P$

×

$PC$

$\theta_P$

$\theta$

Figure 3: Proposed feature descriptor.

Besides, to provide invariance to rotation, each histogram is normalized using the prominent orientation ($\theta_p$) obtained as in [34], but taking into account the local direction $\theta_{dir}(x, y)$. Then, sixteen histograms are concatenated and normalized (using the $L^2$ norm) in order to form the final descriptor.

## 5. Experimental Results

In this section, the performance of the proposed LUIFT algorithm is experimentally presented and analyzed. Three versions of the LUIFT descriptor are evaluated, that is, LUIFT_8, LUIFT_36, and LUIFT_64 which utilize 8, 36, and 64 bins, respectively. The performance of the proposed LUIFT method is compared with FAST [26], STAR[37], SIFT [9], SURF [11], KAZE [12], HARRISZ[33], DAISY [10], and

LIOP [21] detectors and descriptors. All simulations were performed using C$^{++}$ and openCV (http://opencv.org/) library, with the exception of the LIOP descriptor, which was performed in Matlab using the VLFeat (http://www.vlfeat.org/) library.

*5.1. Evaluation Setup.* To evaluate the performance of the tested methods, the *repeatability score*, *matching score* and the *overlap error* are considered.

Let be $P_i = \{fp_i \mid i = 1, 2, \cdots, N. \ N \in \mathbb{N}\}$ a set of feature points $fp_i = I(x, y)$ detected in the original image $I$, $\mathbf{T}$ be a transformation matrix, and $P_j = \{fp_j \mid j = 1, 2, \cdots, M. \ M \in \mathbb{N}\}$ be the set of feature points $fp_j = I'(x, y)$ detected in the test image $I'$. A *correspondence* is considered if $\|\mathbf{T} \cdot fp_i - fpj\| \leq \epsilon$, where $\|\cdot\|$ denotes the Euclidean distance,
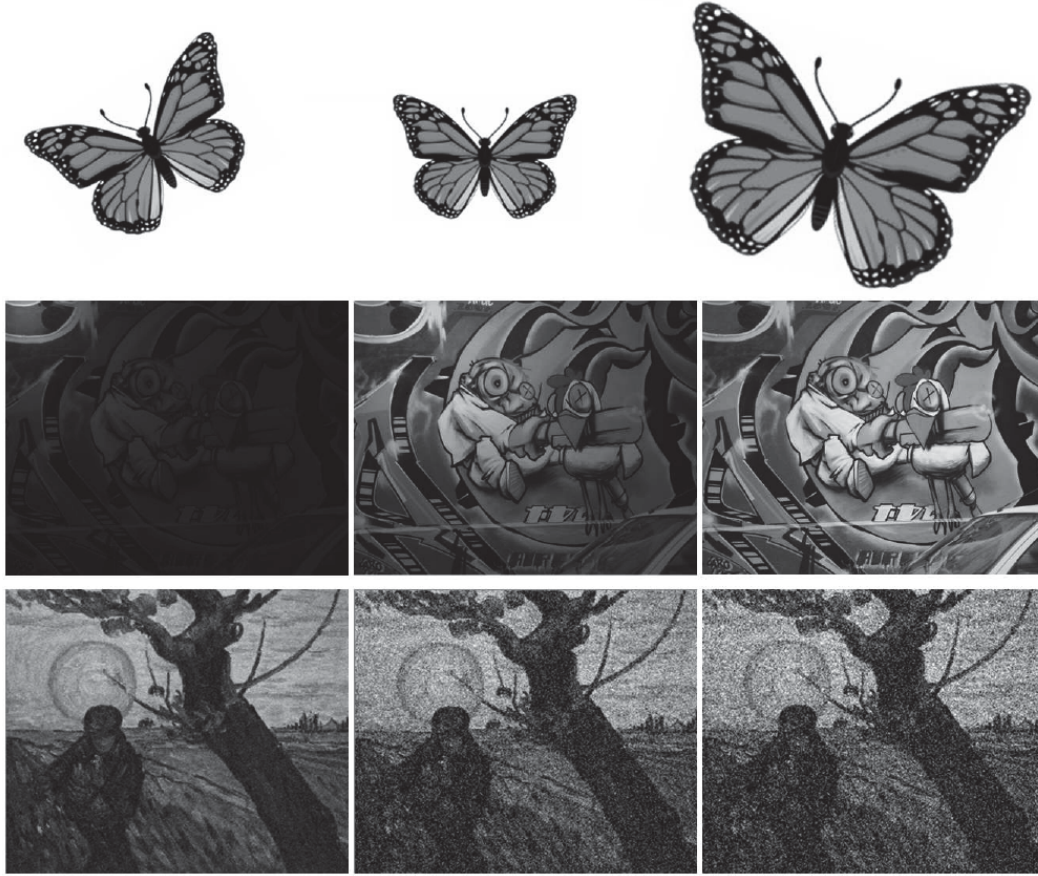
FIGURE 4: Example of synthetic dataset images. From top to bottom: butterfly scene under rotation and scale distortions; graffiti scene under nonuniform illumination variations; and gogh scene under additive noise degradations.

and $\epsilon = \sqrt{2}$ pixels [55]. The feature detector performance is evaluated using the *repeatability score* [15] defined as the ratio between the number of point-to-point correspondences and the minimum number of points detected in both images.

For the descriptor matching performance, two descriptors are matched if the distance between the descriptors is below a threshold $t$. According to [34], if the ratio is less or equal to 0.9, then a *correspondence* is considered. To find the nearest neighbors, the Fast Approximate Nearest Neighbor Search algorithm (FLANN) [56] is exploited.

The results are presented by the *recall-vs-1-precision* curve. *Recall* and *1-precision* are defined as follows [17]:

$$recall = \frac{correct\ matches}{correspondences}, \tag{30}$$

$$1 - precision = \frac{false\ matches}{correct\ matches + false\ matches}. \tag{31}$$

The *correct matches* are determined with the overlap error ($\epsilon < 0.5$) [15]. Basically, the overlap error measure (also called surface error) indicates how well two detected feature regions intersect. The overlap error is defined as the ratio of the intersection of the regions ($\mu_A \cap (H_L^T \mu_B H_L)$) and their union ($\mu_A \cup (H_L^T \mu_B H_L)$) as follows:

$$\epsilon = 1 - \frac{\left(\mu_A \cap \left(H_L^T \mu_B H_L\right)\right)}{\left(\mu_A \cup \left(H_L^T \mu_B H_L\right)\right)}, \tag{32}$$

where $\mu_A$ and $\mu_B$ are the elliptic regions defined by the second moment matrix that satisfy $x^T \mu x = 1$ and $H_L$ is the locally linearized homography $\mathbf{H}$ in the point $x_B$.

Finally, the *matching score* is computed as

$$matching\ score = \frac{correct\ matches}{total\ matches}. \tag{33}$$

*5.2. Synthetic Dataset Evaluation.* In order to evaluate the performance of the proposed LUIFT detector and descriptor, a synthetic grayscale (range from 0 to 255) dataset was created. The synthetic dataset contains 7164 images, of which 2,106 ones correspond to three different scenes (butterfly, gogh, and graffiti) scaled (6 scales) and rotated (13 rotations) under nonuniform illumination (9 variations); 2,106 ones correspond to three different scenes scaled and rotated under additive Gaussian noise (9 variations); and 3042 images correspond to three different scenes scaled and rotated under brightness and contrast (13 variations) changes. Figure 4 shows examples of the synthetic dataset images.

TABLE 1: Parameters used to generate synthetic image dataset.

| Degradation | Step | Range |
|---|---|---|
| Illumination ($\rho$) | 10 | $[10, 50]$ |
| Additive noise ($\sigma$) | 5 | $[0, 40]$ |
| Brightness ($b$) | 30 | $[-90, 90]$ |
| Contrast ($c$) | 0.3 | $[0.5, 2]$ |
| Distortion | | |
| Rotation | 5 | $[-30, 30]$ |
| Scale | 0.1 | $[0.8, 1.3]$ |

The test images are corrupted by zero-mean additive white Gaussian noise, varying the standard deviation $\sigma$.

Nonuniform illumination is simulated using the Lambertian model [57] defined as

$$d(x, y) = \cos\left(\frac{\phi}{2}\right.$$
$$\left. - \tan^{-1}\left(\frac{\rho}{\cos(\phi)}\left[(s_x - x)^2 + (s_y - y)^2\right]^{-1/2}\right)\right), \quad (34)$$

where

$$s_x = \rho \cdot \tan(\phi)\cos(\psi) \quad (35)$$

and

$$s_y = \rho \cdot \tan(\phi)\sin(\psi). \quad (36)$$

The multiplicative function $d(x, y)$ depends on the parameter $\rho$; that is the distance between a point in the surface and the light source, and the parameters $\phi$ and $\psi$ are tilt and slang angles, respectively. In our experiments the following parameters were used: $\phi = 45$ and $\psi = 90$, varying the distance parameter $\rho$.

Brightness and contrast are simulated by

$$f'(x, y) = c \cdot f(x, y) + b. \quad (37)$$

where $b$ and $c$ represent the brightness and contrast parameters, respectively.

Table 1 summarizes the parameters used to generate the synthetic images.

*5.2.1. Simulation Results.* Using the synthetic dataset (Section 5.2), four experiments were conducted in order to evaluate the performance of the proposed LUIFT method under nonuniform illumination, noise, brightness and contrast variations. The performance of the proposed LUIFT method performance is compared with that of the common methods SIFT [34] and SURF [36], in terms of *repeatability* and the *matching score*.

Our first experiment for nonuniform illumination conditions is carried out by varying the distance parameter $\rho$ in test images (rotated and scaled scenes). Figure 5 shows the obtained simulation results for nonuniform illumination
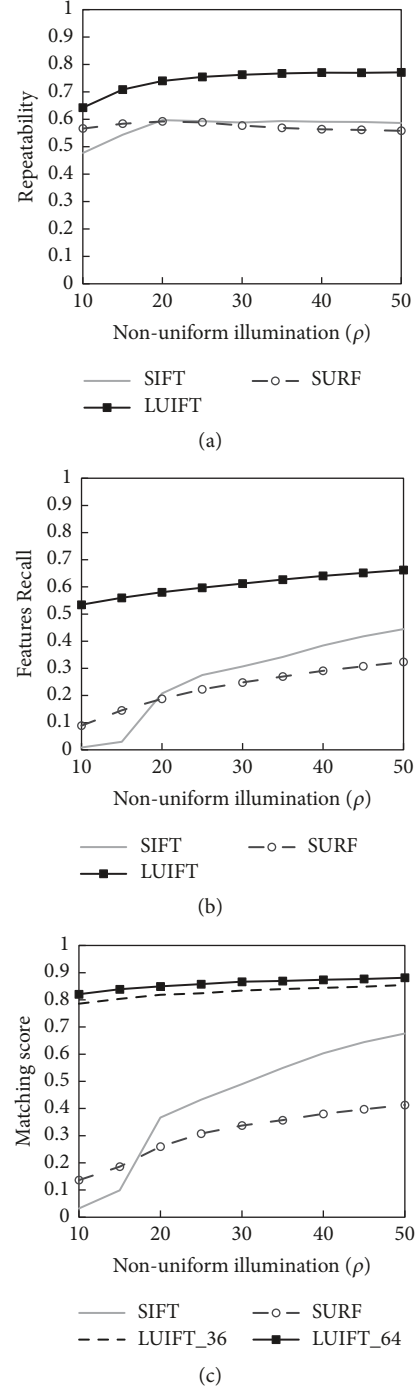


(a)



(b)



(c)

FIGURE 5: Performance of the tested methods on synthetic images rotated and scaled under the nonuniform illumination variations. (a) Feature percentage that remains stable under illumination variations; (b) percentage of correct detected features with respect the original image; (c) feature descriptor performance.

in terms of the *repeatability* and *matching score*. It can be observed that all the tested methods are capable to detect and match feature points of the synthetic test images. However, the feature detection performance, as well as the matching

performance of the SIFT and the SURF methods, decreases considerably when illumination becomes more nonuniform. Note that the proposed method significantly outperforms the tested methods on low-illuminated scenes, reaching up to 50% of improvement.

The next experiment consists in testing of the method performance under Gaussian noise degradations carried out by varying the standard deviation value $\sigma$ in test images (rotated and scaled scenes). Figure 6 shows the simulation results for noise degradation in terms of the *repeatability* and the *matching score*.

The performance of the SIFT method decreases as the noise variance increases, meanwhile the performance of the SURF detector remains stable. In terms of the *repeatability score*, the performance of the SIFT and SURF detectors is worse by almost 20% than that of the proposed LUIFT method, whereas the SURF method shows the worst performance with respect to the *matching score* among all tested descriptors.

The final experiments for brightness and contrast variations are carried out by varying the $c$ and $b$ parameters in test images (rotated and scaled scenes). Figures 7 and 8 show the simulation results for contrast and brightness variations in terms of the *repeatability* and *matching score*, respectively. The obtained results show that the SIFT method is less sensitive to monotonic illumination changes. However, the proposed method yields the best performance in terms of *repeatability* and *matching score*.

Next, in order to compare the performance of the proposed detector and descriptor to that of the state-of-the-art methods in real scene images, the OFFICE (http://www.zhwang.me/datasets.html) and the PHOS (http://www.computervisiononline.com/dataset/1105138614) datasets were utilized.

*5.3. Real Dataset Experiments.* The OFFICE dataset, proposed in [21], contains two different scenes called corridor and desktop. Each scene set contains 5 images with monotonic illumination variations (see Figure 9). For each image set, the performance of the proposed descriptor and the state-of-the-art methods are evaluated.

Figure 10 shows the performance of the tested methods in terms of *repeatability* for feature detector, and the *recall vs 1-precision* curve for the feature descriptor. It can be observed that the proposed descriptor obtain a superior performance compared with that of the state-of-the-art evaluated methods. Despite that the performance of the FAST feature detector looks to that of the proposed LUIFT detector for the corridor scene in terms of *repeatability* (Figure 10(a)), the number of correct feature points detected in all the images for the proposed detector is greater than for the FAST detector (Figure 10(b)). Furthermore, the number of features detected in the original image using the FAST detector decreases by more than 50% as the corridor scene is degraded (Figure 10(b)), and almost 75% for the desktop scene (Figure 10(e)). The main drawback of FAST detector is that the desired number of features detected by the method needs to be adjusted for each type of scene or task. Note that
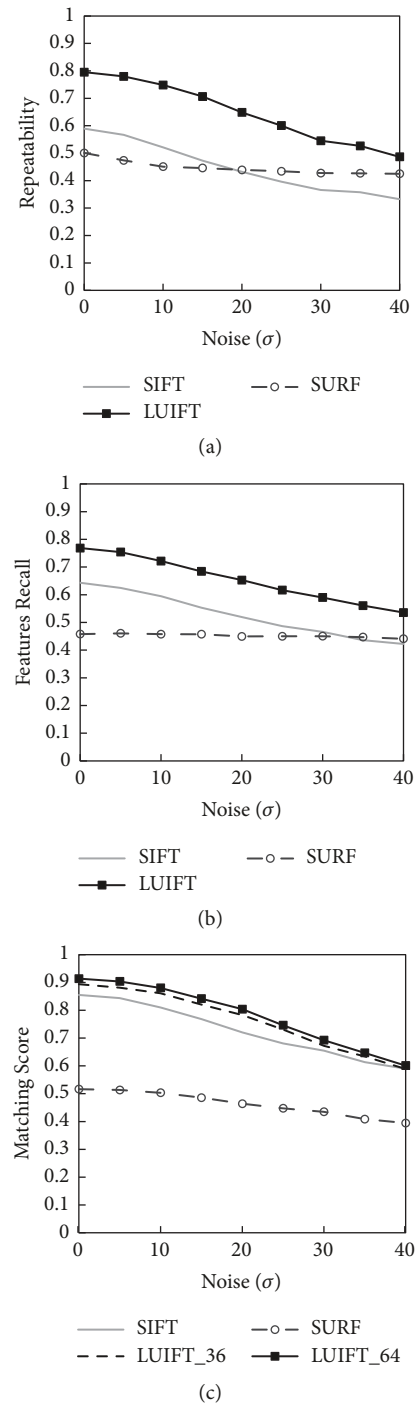


(a)



(b)



(c)

FIGURE 6: Performance of the tested methods on synthetic images rotated and scaled under the noise degradations. (a) Feature percentage that remains stable under noise degradations; (b) percentage of correct detected features with respect the original image; (c) feature descriptor performance.

it is important for the detector methods to have not only a high *repeatability* score, but also to obtain a high number of correct points.
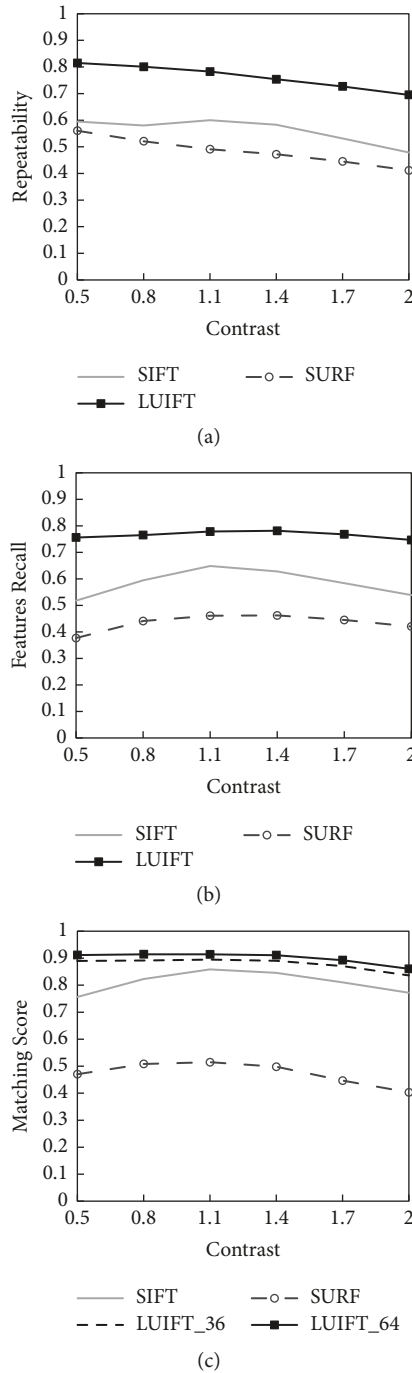
(a)

(b)

(c)

FIGURE 7: Performance of the tested methods on synthetic images rotated and scaled under the contrast variations. (a) Feature percentage that remains stable under contrast variations; (b) percentage of correct detected features with respect the original image; (c) feature descriptor performance.
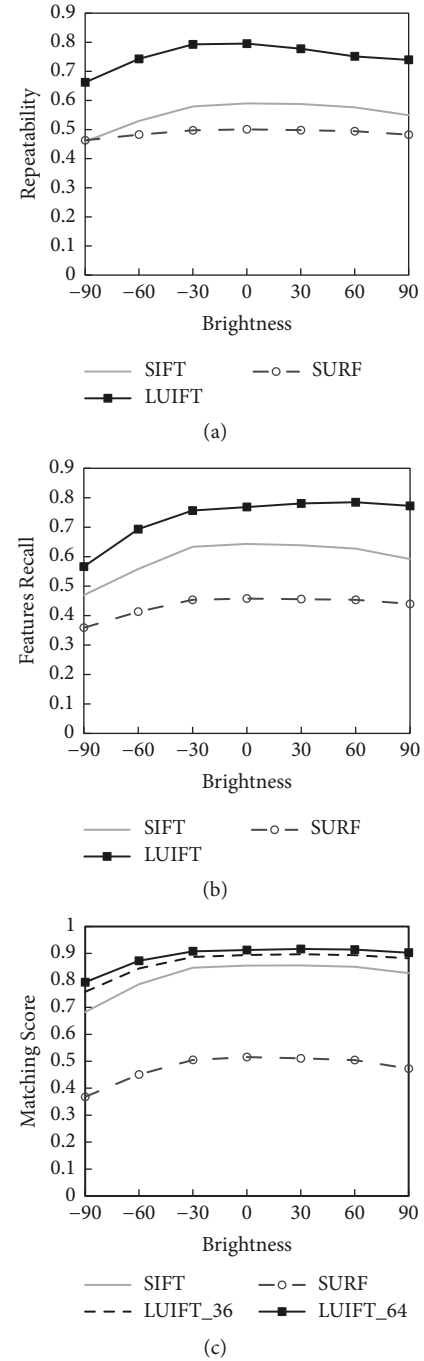


(a)

(b)

(c)

FIGURE 8: Performance of the tested methods on synthetic images rotated and scaled under the brightness changes. (a) Feature percentage that remains stable under brightness changes; (b) percentage of correct detected features with respect the original image; (c) feature descriptor performance.

Also the PHOS dataset [58] was used. The PHOS dataset contains 15 different scenes (see Figure 11) captured under different illumination conditions. Every scene of the dataset contains 15 different images: 9 images captured under different uniform illumination, varying the camera exposure between -4 and +4 from the original correctly exposed image (see Figure 12(a)); and 6 images under different degrees of nonuniform illumination, accomplished by adding a strong directional light source to uniform diffusive lights located around the objects (see Figure 12(b)).

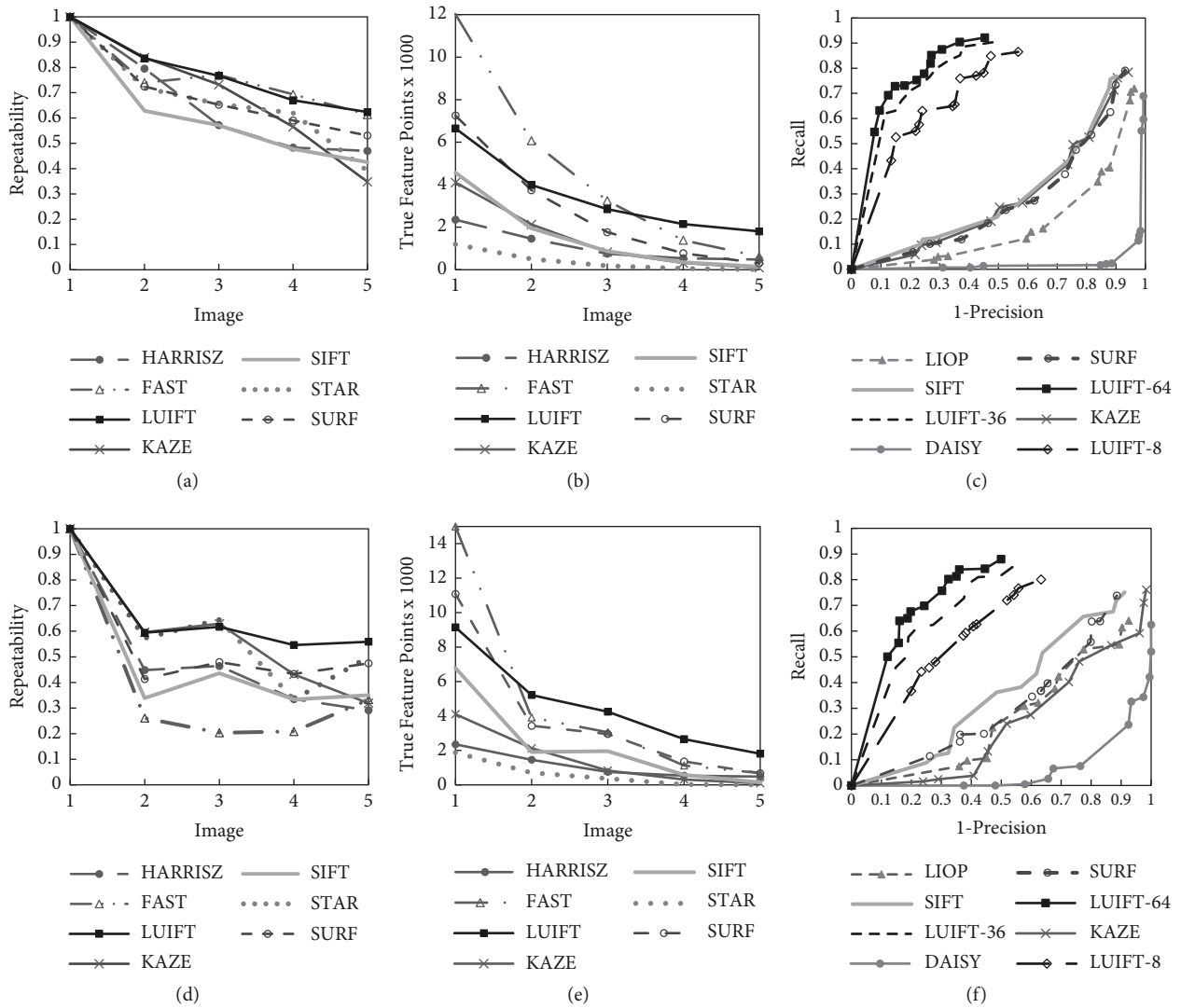FIGURE 9: OFFICE dataset: corridor and desktop scenes.



FIGURE 10: Office dataset results. For the corridor scene set: (a) Feature detector *repeatability*, (b) correct feature points detected, and (c) feature descriptor *recall vs 1-precision* curve. For the desktop scene set: (d) feature detector *repeatability*; (e) correct feature points detected; and (f) feature descriptor *recall vs 1-precision* curve.
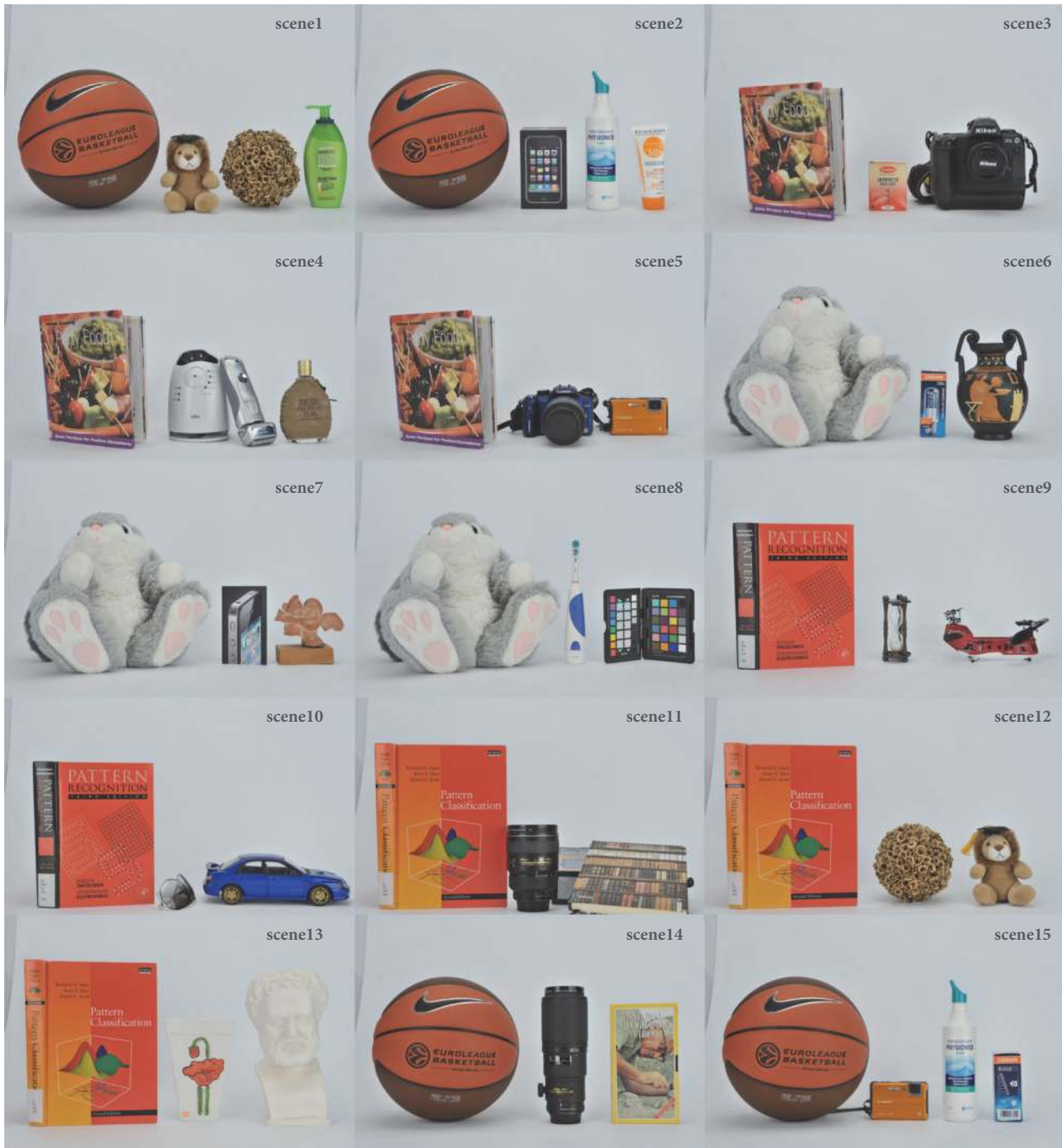
Figure 11: Scenes from the PHOS dataset. Each scene contains different types of objects.

Figure 13 shows the performance of the proposed LUIFT descriptor and the state-of-the-art methods on the PHOS dataset in terms of *repeatability* and the *recall vs 1-precision* curve. For the case of exposure variations, Figure 13(a) shows the average feature detector performance in terms of *repeatability*, and Figure 13(b) shows the average feature descriptor performance in terms of the *recall vs 1-precision* curve. For the case of nonuniform illumination variations, Figure 13(c) shows the average feature detector performance in terms of *repeatability*, and Figure 13(d) shows the average feature descriptor performance in terms of the *recall vs 1-precision* curve. The performance of the proposed LUIFT detector and descriptor is superior to that of all the tested methods, even for the LUIFT-8 descriptor.

(a)



(b)

Figure 12: Test images from scene15. (a) Exposure variations (EV); (b) nonuniform illumination variations.

Table 2: Computation time of the tested methods for the graffiti image ($800 \times 640$).

| Method | detected points | time (ms) |
|---|---|---|
| SIFT | 2916 | 802 |
| SURF | 4713 | 418 |
| LUIFT | 1719 | 1097 |

The performance of the tested methods for each scene set (including exposure and nonuniform illumination) are shown in Figure 14 in terms of the *recall vs 1-precision* curve. The proposed method outperforms the other descriptors in all the cases.

Finally, Table 2 shows computation time (ms) required by the tested methods for processing of the graffiti image ($800 \times 640$). As expected, the SURF descriptor is faster than all methods. That is because of the use of Haar-like filters and integral images to improve the processing time at expense of its performance. On the other hand, the SIFT descriptor is based on Laplacian of Gaussian approximations instead of getting second order derivatives, which are more computationally expensive. However, since Laplacian of Gaussian approximations are made by the difference of Gaussian images, there exist errors in feature location or losing features. Besides, the SIFT descriptor duplicate feature points if they get two prominent orientations, collecting more features than the proposed method. All the experiments were performed on a standard PC with Intel Xeon E5-1603 processor with 2.8GHz and 16GB of RAM.

## 6. Conclusions

In this work, a robust phase-based descriptor for pattern recognition in degraded images using the scale-space monogenic signal and phase congruency approach was presented. With the help of computer simulation, the performance
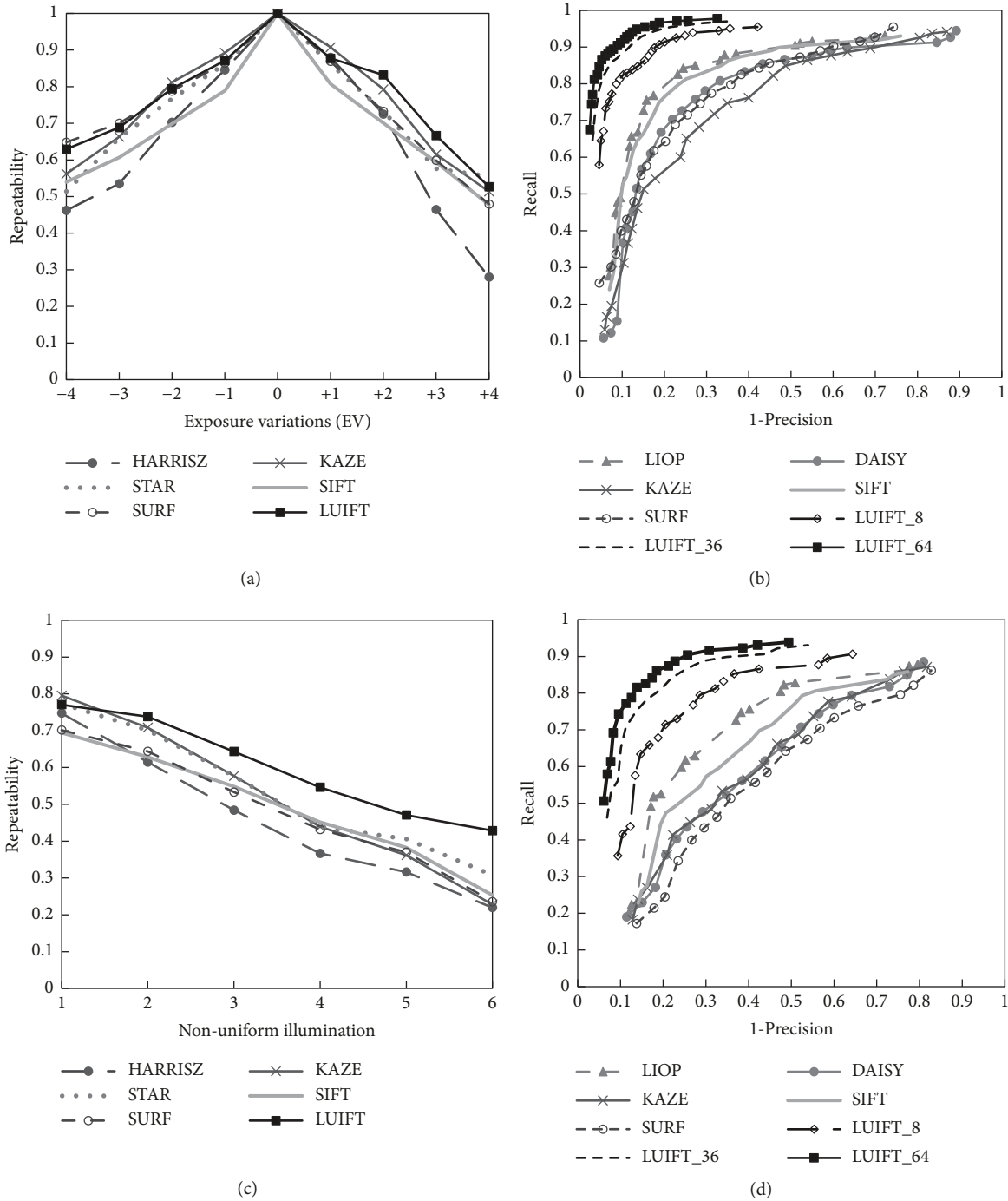
(a)



(b)



(c)



(d)

Figure 13: Performance of the tested methods on the PHOS dataset in terms of *repeatability* and the *recall vs 1-precision* curves. (a and b) Exposure variation results; (c and d) nonuniform illumination variation results.

of the proposed method was compared with that of the state-of-the-art methods. The proposed method shows a superior performance under illumination variations, and noise degradations. Besides, the obtained results on typical dataset for evaluation of feature detection and matching performance are competitive with those obtained with the state-ofthe-art descriptors. The performance of the proposed method can be further improved by including into the

design pyramidal scale decomposition. Since the proposed method is inherently local, its fast GPU implementation is straightforward.

## Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this article.
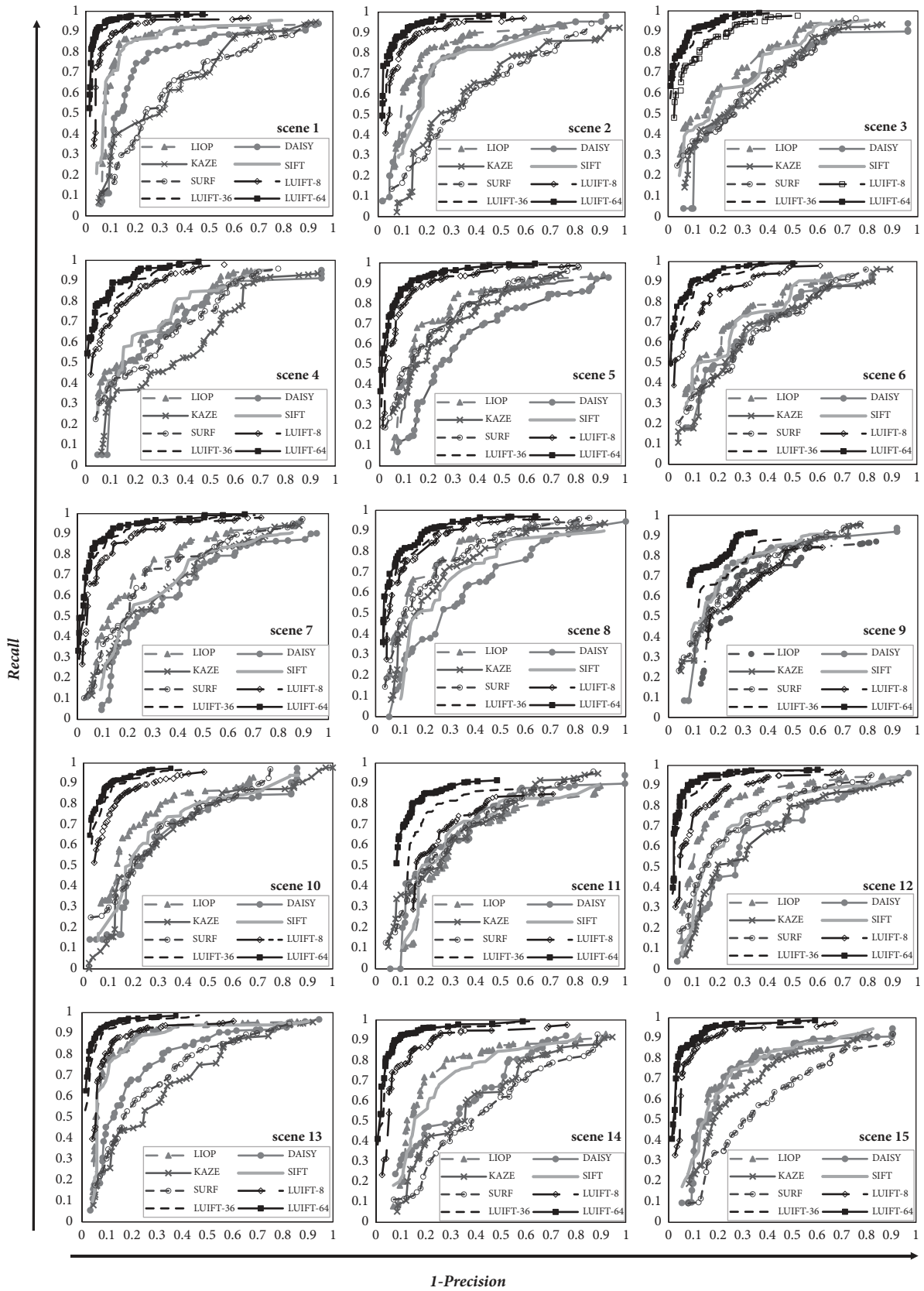
FIGURE 14: Performance of the tested methods on the PHOS dataset in terms of the *recall vs 1-precision* curve.
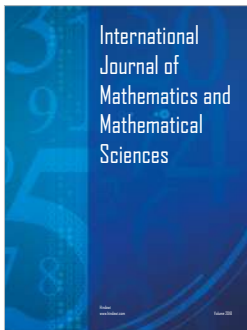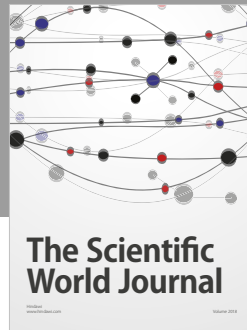
## Acknowledgments

## References

[1] T. Deselaers, D. Keysers, and H. Ney, "Features for image retrieval: an experimental comparison," *Information Retrieval*, vol. 11, no. 2, pp. 77–107, 2008.

[2] S. Liu and X. Bai, "Discriminative features for image classification and retrieval," *Pattern Recognition Letters*, vol. 33, no. 6, pp. 744–751, 2012.

[3] D. Fortun, P. Bouthemy, and C. Kervrann, "Optical flow modeling and computation: A survey," *Computer Vision and Image Understanding*, vol. 134, pp. 1–21, 2015.

[4] S. Tang, M. Andriluka, and B. Schiele, "Detection and tracking of occluded people," *International Journal of Computer Vision*, vol. 110, no. 1, pp. 58–69, 2014.

[5] A. Jain, A. A. Ross, and K. Nandakumar, *Introduction to biometrics*, Springer Science & Business Media, 2011.

[6] B. Zitová and J. Flusser, "Image registration methods: a survey," *Image and Vision Computing*, vol. 21, no. 11, pp. 977–1000, 2003.

[7] P. Moreels and P. Perona, "Evaluation of features detectors and descriptors based on 3D objects," *International Journal of Computer Vision*, vol. 73, no. 3, pp. 263–284, 2007.

[8] T. Tuytelaars and K. Mikolajczyk, "Local invariant feature detectors: a survey," *Foundations and Trends in Computer Graphics and Vision*, vol. 3, no. 3, pp. 177–280, 2007.

[9] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of the 7th IEEE International Conference on Computer Vision (ICCV '99)*, vol. 2, pp. 1150–1157, Kerkyra, Greece, September 1999.

[10] E. Tola, V. Lepetit, and P. Fua, "A fast local descriptor for dense matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '08)*, pp. 1–8, Anchorage, Alaska, USA, June 2008.

[11] H. Bay, T. Tuytelaars, and L. van Gool, "SURF: speeded up robust features," in *Computer Vision-ECCV2006*, vol. 3951 of *Lecture Notes in Computer Science*, pp. 404–417, Springer, Berlin, Germany, 2006.

[12] P. F. Alcantarilla, A. Bartoli, and A. J. Davison, "KAZE Features," in *Computer Vision – ECCV 2012*, vol. 7577 of *Lecture Notes in Computer Science*, pp. 214–227, Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.

[13] Y. Verdie, . Kwang Moo Yi, P. Fua, and V. Lepetit, "TILDE: A Temporally Invariant Learned DEtector," in *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5279–5288, Boston, MA, USA, June 2015.

[14] F. Attneave, "Some informational aspects of visual perception," *Psychological Review*, vol. 61, no. 3, pp. 183–193, 1954.

[15] K. Mikolajczyk and C. Schmid, "Scale & affine invariant interest point detectors," *International Journal of Computer Vision*, vol. 60, no. 1, pp. 63–86, 2004.

[16] K. Mikolajczyk, T. Tuytelaars, C. Schmid et al., "A comparison of affine region detectors," *International Journal of Computer Vision*, vol. 65, no. 1-2, pp. 43–72, 2005.

[17] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615–1630, 2005.

[18] F. Tang, S. H. Lim, N. L. Chang, and H. Tao, "A novel feature descriptor invariant to complex brightness changes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '09)*, pp. 2631–2638, IEEE, June 2009.

[19] M. Heikkilä, M. Pietikäinen, and C. Schmid, "Description of interest regions with local binary patterns," *Pattern Recognition*, vol. 42, no. 3, pp. 425–436, 2009.

[20] R. Gupta, H. Patil, and A. Mittal, "Robust order-based methods for feature description," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, (CVPR '10)*, pp. 334–341, June 2010.

[21] Z. Wang, B. Fan, and F. Wu, "Local intensity order pattern for feature description," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV '11)*, pp. 603–610, Barcelona, Spain, November 2011.

[22] J. Wensveen and B. Wick, "Eye, brain, and vision," *Optometry and Vision Science*, vol. 72, no. 10, p. 773, 1995.

[23] E. Gladilin and R. Eils, "On the role of spatial phase and phase correlation in vision, illusion, and cognition," *Frontiers in Computational Neuroscience*, vol. 9, 2015.

[24] G. Papari and N. Petkov, "Edge and line oriented contour detection: state of the art," *Image and Vision Computing*, vol. 29, no. 2-3, pp. 79–103, 2011.

[25] S. M. Smith and J. M. Brady, "SUSAN, a new approach to low level image processing," *International Journal of Computer Vision*, vol. 23, no. 1, pp. 45–78, 1997.

[26] E. Rosten and T. Drummond, "Machine learning for high-speed corner detection," in *Computer Vision—ECCV 2006: Proceedings of the 9th European Conference on Computer Vision, Graz, Austria, May 7–13, 2006, Part I*, vol. 3951 of *Lecture Notes in Computer Science*, pp. 430–443, Springer, Berlin, Germany, 2006.

[27] E. Simo-Serra, C. Torras, and F. Moreno-Noguer, "DaLI: deformation and light invariant descriptor," *International Journal of Computer Vision*, vol. 115, no. 2, pp. 136–154, 2015.

[28] K. M. Yi, E. Trulls, V. Lepetit, and P. Fua, "LIFT: Learned Invariant Feature Transform," in *Computer Vision – ECCV 2016*, vol. 9910 of *Lecture Notes in Computer Science*, pp. 467–483, Springer International Publishing, Cham, 2016.

[29] G. Levi and T. Hassner, "LATCH: Learned arrangements of three patch codes," in *Proceedings of the 2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1–9, Lake Placid, NY, USA, March 2016.

[30] C. Harris and M. Stephens, "A combined corner and edge detector," in *Proceedings of the 4th Alvey Vision Conference*, pp. 147–151, 1988.

[31] H. P. Moravec, "Visual mapping by a robot rover," in *Conference on Artificial Intelligence*, pp. 598–600, 1979.

[32] B. P. R, "Rotationally invariant image operators," in *Proceedings of the in International Joint Conference on Pattern Recognition*, pp. 579–583, 1987.

[33] F. Bellavia, D. Tegolo, and C. Valenti, "Improving Harris corner selection strategy," *IET Computer Vision*, vol. 5, no. 2, pp. 87–96, 2011.

[34] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.

[35] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR '05)*, vol. 1, pp. 886–893, June 2005.

[36] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-up robust features (surf)," *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, 2008.

[37] M. Agrawal, K. Konolige, and M. R. Blas, "Censure: Center surround extremas for realtime feature detection and matching," in *Conference on Computer Vision*, pp. 102–115, 2008.

[38] A. Alahi, R. Ortiz, and P. Vandergheynst, "FREAK: fast retina keypoint," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR '12)*, pp. 510–517, June 2012.

[39] M. Calonder, V. Lepetit, C. Strecha, and P. Fua, "BRIEF: binary robust independent elementary features," in *Computer Vision—ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part IV*, vol. 6314 of *Lecture Notes in Computer Science*, pp. 778–792, Springer, Berlin, Germany, 2010.

[40] S. Leutenegger, M. Chli, and R. Y. Siegwart, "BRISK: binary robust invariant scalable keypoints," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV '11)*, pp. 2548–2555, Barcelona, Spain, November 2011.

[41] A. V. Oppenheim and J. S. Lim, "The importance of phase in signals," *Proceedings of the IEEE*, vol. 69, no. 5, pp. 529–541, 1981.

[42] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction, and functional architecture in the cat's visual cortex," *The Journal of Physiology*, vol. 160, pp. 106–154, 1962.

[43] M. C. Morrone, J. Ross, D. C. Burr, and R. Owens, "Mach bands are phase dependent," *Nature*, vol. 324, no. 6094, pp. 250–253, 1986.

[44] M. C. Morrone and R. A. Owens, "Feature detection from local energy," *Pattern Recognition Letters*, vol. 6, no. 5, pp. 303–313, 1987.

[45] M. C. Morrone and D. C. Burr, "Feature detection in human vision: a phase-dependent energy model," *Proceedings of the Royal Society B Biological Science*, vol. 235, no. 1280, pp. 221–245, 1988.

[46] B. Robbins and R. Owens, "2D feature detection via local energy," *Image and Vision Computing*, vol. 15, no. 5, pp. 353–368, 1997.

[47] P. Kovesi, "Image features from phase congruency," *Videre: Journal of Computer Vision Research*, vol. 1, no. 3, pp. 1–26, 1999.

[48] P. Kovesi, "Phase congruency: a low-level image invariant," *Psychological Research*, vol. 64, no. 2, pp. 136–148, 2000.

[49] P. Kovesi et al., "Edges are not just steps," in *Conference on Computer Vision*, vol. 8, pp. 22–28, Melbourne, 2002.

[50] M. Felsberg and G. Sommer, "A New Extension of Linear Signal Processing for Estimating Local Properties and Detecting Features," in *Mustererkennung 2000*, Informatik aktuell, pp. 195–202, Springer Berlin Heidelberg, Berlin, Heidelberg, 2000.

[51] A. G. Tescher, J. Diaz-Escobar, and V. Kober, "A robust HOG-based descriptor for pattern recognition," in *Proceedings of the SPIE Optical Engineering + Applications*, p. 99712A, San Diego, California, United States.

[52] M. Felsberg and G. Sommer, "The monogenic signal," *IEEE Transactions on Signal Processing*, vol. 49, no. 12, pp. 3136–3144, 2001.

[53] M. Felsberg and G. Sommer, "The monogenic scale-space: a unifying approach to phase-based image processing in scale-space," *Journal of Mathematical Imaging and Vision*, vol. 21, no. 1, pp. 5–26, 2004.

[54] M. Felsberg, Low-level image processing with the structure multivector, 2002.

[55] G. Carneiro and A. D. Jepson, "Phase-Based Local Features," in *Computer Vision — ECCV 2002*, vol. 2350 of *Lecture Notes in Computer Science*, pp. 282–296, Springer Berlin Heidelberg, Berlin, Heidelberg, 2002.

[56] M. Muja and D. G. Lowe, "Fast approximate nearest neighbors with automatic algorithm configuration," in *Proceedings of the 4th International Conference on Computer Vision Theory and Applications (VISAPP '09)*, vol. 1, pp. 331–340, February 2009.

[57] V. H. Diaz-Ramirez, K. Picos, and V. Kober, "Target tracking in nonuniform illumination conditions using locally adaptive correlation filters," *Optics Communications*, vol. 323, pp. 32–43, 2014.

[58] V. Vonikakis, D. Chrysostomou, R. Kouskouridas, and A. Gasteratos, "A biologically inspired scale-space for illumination invariant feature detection," *Measurement Science and Technology*, vol. 24, no. 7, 2013.

Advances in
Operations Research

Advances in
Decision Sciences

Journal of
Applied Mathematics

**The Scientific
World Journal**

Journal of
Probability and Statistics

International
Journal of
Mathematics and
Mathematical
Sciences

Journal of
Optimization

International Journal of
Engineering
Mathematics

International Journal of
Analysis

# Hindawi

Submit your manuscripts at
www.hindawi.com

Journal of
Complex Analysis

Advances in
Numerical Analysis

Mathematical Problems
in Engineering

International Journal of
Differential Equations

Discrete Dynamics in
Nature and Society

International Journal of
Stochastic Analysis

Journal of
Mathematics

Journal of
Function Spaces

Abstract and
Applied Analysis

Advances in
Mathematical Physics