

Metadata of the article that will be visualized in OnlineFirst

1 Article Title **A novel shape descriptor based on salient keypoints detection
for binary image matching and retrieval**

2 Article Sub-Title

3 ~~Please note: Images will appear in color online but will be printed in black and white.~~

3 Article Copyright ~~Springer Science+Business Media, LLC, part of Springer Nature~~

Year **2018**
(This will be the copyright line in the final PDF)

4 Journal Name Multimedia Tools and Applications

5 Family Name **Chatbri**

6 Particle

7 Given Name **Houssem**

8 Corresponding Author Suffix

9 Organization Dublin City University

10 Division Insight Centre for Data Analytics

11 Address Dublin, Ireland

12 e-mail houssem.chatbri@dcu.ie

13 Family Name **Kameyama**

14 Particle

15 Given Name **Keisuke**

16 Author Suffix

17 Organization University of Tsukuba

18 Division Faculty of Engineering, Information and Systems

19 Address Tsukuba, Japan

20 e-mail keisuke.kameyama@cs.tsukuba.ac.jp

21 Family Name **Kwan**

22 Particle

23 Given Name **Paul**

24 Author Suffix

25 Organization University of New England

26 Division School of Science and Technology

27 Address Amidale, NSW, Australia

28 e-mail paul.kwan@une.edu.au

29 Family Name **Little**

30 Particle

31 Given Name **Suzanne**

32 Author Suffix

33 Organization Dublin City University

34 Division Insight Centre for Data Analytics

35 Address Dublin, Ireland

36		e-mail	suzanne.little@dcu.ie
37		Family Name	O'Connor
38		Particle	
39		Given Name	Noel E.
40		Suffix	
41	Author	Organization	Dublin City University
42		Division	Insight Centre for Data Analytics
43		Address	Dublin, Ireland
44		e-mail	noel.oconnor@dcu.ie
45		Received	10 July 2017
46	Schedule	Revised	18 March 2018
47		Accepted	24 April 2018
48	Abstract	<p>We introduce a shape descriptor that extracts keypoints from binary images and automatically detects the salient ones among them. The proposed descriptor operates as follows: First, the contours of the image are detected and an image transformation is used to generate background information. Next, pixels of the transformed image that have specific characteristics in their local areas are used to extract keypoints. Afterwards, the most salient keypoints are automatically detected by filtering out redundant and sensitive ones. Finally, a feature vector is calculated for each keypoint by using the distribution of contour points in its local area. The proposed descriptor is evaluated using public datasets of silhouette images, handwritten math expressions, hand-drawn diagram sketches, and noisy scanned logos. Experimental results show that the proposed descriptor compares strongly against state of the art methods, and that it is reliable when applied on challenging images such as fluctuated handwriting and noisy scanned images. Furthermore, we integrate our descriptor in a content-based document image retrieval system using sketch queries as a step for query and candidate occurrence matching, and we show that it leads to a significant boost in retrieval performances.</p>	
49	Keywords	Shape descriptors - Salient keypoints - Image matching -	
	separated by ' - '	Sketch-based retrieval	
50	Foot note information		

A novel shape descriptor based on salient keypoints detection for binary image matching and retrieval

Housseem Chatbri¹ · Keisuke Kameyama² ·
Paul Kwan³ · Suzanne Little¹ · Noel E. O'Connor¹

Received: 10 July 2017 / Revised: 18 March 2018 / Accepted: 24 April 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract We introduce a shape descriptor that extracts keypoints from binary images and automatically detects the salient ones among them. The proposed descriptor operates as follows: First, the contours of the image are detected and an image transformation is used to generate background information. Next, pixels of the transformed image that have specific characteristics in their local areas are used to extract keypoints. Afterwards, the most salient keypoints are automatically detected by filtering out redundant and sensitive ones. Finally, a feature vector is calculated for each keypoint by using the distribution of contour points in its local area. The proposed descriptor is evaluated using public datasets of silhouette images, handwritten math expressions, hand-drawn diagram sketches, and noisy scanned logos. Experimental results show that the proposed descriptor compares strongly against state of the art methods, and that it is reliable when applied on challenging images such as fluctuated handwriting and noisy scanned images. Furthermore, we integrate our descriptor

Q2

✉ Housseem Chatbri
housseem.chatbri@dcu.ie
Keisuke Kameyama
keisuke.kameyama@cs.tsukuba.ac.jp
Paul Kwan
paul.kwan@une.edu.au
Suzanne Little
suzanne.little@dcu.ie
Noel E. O'Connor
noel.oconnor@dcu.ie

Q1

¹ Insight Centre for Data Analytics, Dublin City University, Dublin, Ireland

² Faculty of Engineering, Information and Systems, University of Tsukuba, Tsukuba, Japan

³ School of Science and Technology, University of New England, Armidale NSW, Australia

19 in a content-based document image retrieval system using sketch queries as a step for query
20 and candidate occurrence matching, and we show that it leads to a significant boost in
21 retrieval performances.

22 **Keywords** Shape descriptors · Salient keypoints · Image matching · Sketch-based retrieval

23 1 Introduction

24 Shape matching is a vibrant area of research on image analysis and retrieval due to
25 the numerous applications it allows [7]. Particularly, when dealing with binary images
26 where color and texture information are absent (e.g. silhouette images, scanned documents,
27 sketches, etc.), shape is the only available feature to be used for image representation and
28 matching [26].

29 Numerous methods have been presented for shape feature extraction in binary images
30 [54, 57]. Usually, images are subjected to contour detection or skeletonization before fea-
31 ture extraction in order to remove redundant information and reduce processing time [13].
32 Moreover, some methods select certain *keypoints* and use them to extract features [5, 35, 40].
33 In these cases, keypoints are selected based on their saliency or by using uniform sampling
34 from the shape contours.

35 Due to the absence of background information in binary images, keypoints are extracted
36 from the foreground pixels (i.e. regions, contours, or skeletons) and the background is omit-
37 ted. In this work, we introduce a shape descriptor that approaches the problem differently
38 by generating background information in binary images, and then involves it in feature
39 extraction. The main steps of the descriptor are the following:

- 40 – Keypoint extraction: An image transformation is used to generate background informa-
41 tion on the original binary image. Then, keypoints are extracted from the transformed
42 image using pixels' local area analysis.
- 43 – Keypoint selection: An objective measure of keypoint saliency is used to automatically
44 select the most important keypoints and filter out the redundant and sensitive ones.
- 45 – Feature representation: A feature vector is calculated for each keypoint by using the
46 distributions of contour points in the local area of the keypoint.

47 Our method, the binary salient keypoints (BSK) descriptor, is evaluated using silhouette
48 images of the Kimia 216 dataset [45] and the MPEG-7 CE-shape-1 part B dataset [6], hand-
49 written mathematical expressions of Zanibbi and Yu's dataset [56], hand-drawn diagram
50 sketches of Liang et al.'s dataset [28], and noisy scanned logo images of the Tobacco 800
51 dataset [60]. Experimental results on various types of images and a comparative evaluation
52 demonstrate that BSK is competitive compared with state of the art methods. Our code for
53 BSK is provided online.¹

54 The remainder of this paper is organized as follows: Section 2 reviews key methods
55 of shape matching. We present our descriptor in Section 3 and evaluate it in Section 4.
56 Concluding remarks and future work are presented in Section 5.

¹<https://github.com/hchatbri/bsk>

2 Related work

57

Research on shape matching has led to a large repository of shape descriptors that can be classified into methods using global and local features [57], graph-based methods [28], contour-based methods and skeleton based methods [13], in addition to methods using salient keypoints [5, 32, 35, 40].

Global methods extract features using the coarse information of the shape, and hence do not convey much information about the local details. Such methods include shape signatures [43], Fourier descriptors [58], and angular partitioning [10]. Global methods are robust against noise but on the detriment of representing fine details. On the other hand, other methods take into consideration the local region of the shape points, which makes them capable of capturing fine details of the shape. Such methods include curvature scale space (CSS) [32], shape contexts [5], and variations of local binary patterns [11, 18].

Graph-based methods represent features using graph structures in contrast to statistical methods which use statistical natures of appearances [28]. Advantages of graph-based methods are their ability to represent spatial and hierarchical relationships between the object parts [8], in addition to allowing partial matching. On the other hand, graph matching requires intensive computations and thus it is common to transform a graph into a numerical feature vector in order to speed up computations, which is done at the expense of some information loss [16, 27].

Contours and skeletons have been used as an intermediate representation before feature extraction. Contours are more robust against noise than skeletons, as skeletons tend to generate noisy branches and artifacts in presence of shape border perturbations [13]. On the other hand, skeletons are more suitable in applications that require the segmentation of the original object into its constituent parts for subsequent graph-based feature representation [3, 24, 44, 51].

Some descriptors extract a number of keypoints and generate a feature vector for each one of them. Keypoints can be extracted using uniform sampling from the shape contours without special consideration about the keypoints curvature or location, offering a way to extract keypoints without a bias [29]. This has been exploited in numerous descriptors [5, 17, 40, 49], yet it does not take into account keypoints' local characteristics that make some keypoints more important than others. In addition to binary images, uniform sampled keypoints have been used on grayscale images (e.g. magnetic resonance image (MRI) registration [36]) and they have been used on color images combined with other descriptors (e.g. combined with SIFT and segmentation patches in [21] for logo retrieval).

On the other hand, keypoints that are extracted based on their salience (e.g. corners, crossing points) are biologically plausible [39], although they can lead to false detection in regions of contour perturbations or texture [46]. High curvature points of the contour have been used for keypoint extraction [20, 35]. Early methods such as Curvature scale space (CSS) uses scale space filtering [53] to extract contour inflection points [1, 32]. Then, the contour deformation and merging of inflection points caused by scale space filtering are used for feature extraction. The CSS method is a global statistical method, designed to deal only with closed concave contours; Convex and complex shapes are poorly represented with the technique. In [23], Kopf et al. describe an attempt to extend the CSS technique and make it able to represent convex shapes. Their idea is to create a mapping of the original shape to a second shape, called mapped shape, where strong convex segments of the original shape become concave segments of the mapped shape, and significant curvatures in the original shape remain significant in the mapped shape. The mapping is done by enclosing the sketch

104 with a circle of radius R and locate the point P of the circle closest to each sketch pixel. The
105 sketch pixels are then mirrored on the tangent of the circle in P . The center of the circle is
106 the average position of sketch pixels.

107 Scale-space filtering has also been used to extract distinctive keypoints in intensity
108 images in the well-known SIFT descriptor [30]. However, it has been shown that SIFT key-
109 points are suboptimal compared to keypoints that are uniformly sampled from the shape
110 contours when using complex binary images such as historical hieroglyphs [41]. This result
111 is due to the absence of local changes of intensity in binary images that hinders scale-space
112 filtering from detecting distinctive keypoints and attributing them characteristic scales.
113 Scale-space has also been used for keypoint filtering [12], which proved to be effective but
114 on the expense of efficiency.

115 Curvature information has been also used for salient keypoints extraction in [35]. Here,
116 the salient points of a shape are defined as the higher curvature points along the shape
117 contour that are extracted using a noise-robust approach [34]. Then, each salient point is
118 represented with two values, the relative angular position of the salient point from the per-
119 spective of the shape centroid, and the *salience relevance* which characterizes the concavity
120 of the contour segmented around the salient point after applying a Gaussian filter to reduce
121 contour noise. Image matching corresponds an energy minimization function which give
122 the distance between the best pair of corresponding salient points.

123 In addition to high curvature points, other salient points have been used including end
124 points and branch points of the object's skeleton, and the vertices of the minimum enclosing
125 rectangle of the object [59]. For each salient point, features are calculated using a circular
126 layout of polar coordinates to calculate the distribution of some shape points which are
127 sampled using a maximum distance method. Finally, a feature vector is constructed using a
128 bag of words method.

129 Data-driven methods have been recently designed using deep learning [37, 50, 61].
130 Unlike the aforementioned methods, data-driven methods automatically learn salient fea-
131 tures using convolutional layers, in an attempt to mimic the way humans perceive shapes
132 and sketches [19]. Despite the success of such methods, they require large labeled datasets
133 for training and they usually need graphical processing units (GPUs) to alleviate compu-
134 tations. Due to these reasons, engineered features remain necessary for applications where
135 large labeled datasets are unavailable.

Q3136 2.1 Our contributions

137 Compared to the state of the art, our descriptor's main contributions are twofold:

- 138 – We demonstrate that the background of binary images, which before has not been con-
139 sidered enough for feature extraction, can be used to extract distinctive features. We
140 show that an image transform such as the distance transform (DT) can be used to
141 enable this. Our experiments show that extracting salient keypoints using this procedure
142 leads to improved robustness against noise that otherwise would easily corrupt object
143 contours.
- 144 – Our descriptor is modular and it proceeds in three main steps which are feature extrac-
145 tion, keypoint selection and feature representation. This is similar to frameworks of
146 widely-used color image descriptors (e.g. SIFT [30], SURF [4]). Consequently, we
147 adapt a framework that has been used for color images into the binary image domain.

3 The proposed descriptor

148

The binary salient keypoints descriptor (BSK) operates as follows: First, keypoints are extracted (Section 3.1). Then, a number of keypoints are selected among the extracted ones and the others are filtered out (Section 3.2). Finally, a feature vector is calculated for each keypoint (Section 3.3).

3.1 Keypoint extraction

153

In this step, a transformation is applied on the input binary image in order to generate background information. Then, points having specific characteristics in their local areas are used to extract keypoints.

For our image transformation, we use the distance transform (DT) [42]. DT generates a grayscale image where the intensity of each pixel corresponds to its distance from the nearest foreground pixel (Fig. 1c). Here, the distance between pixels is equal to their Manhattan distance as commonly used in DT implementations [31].

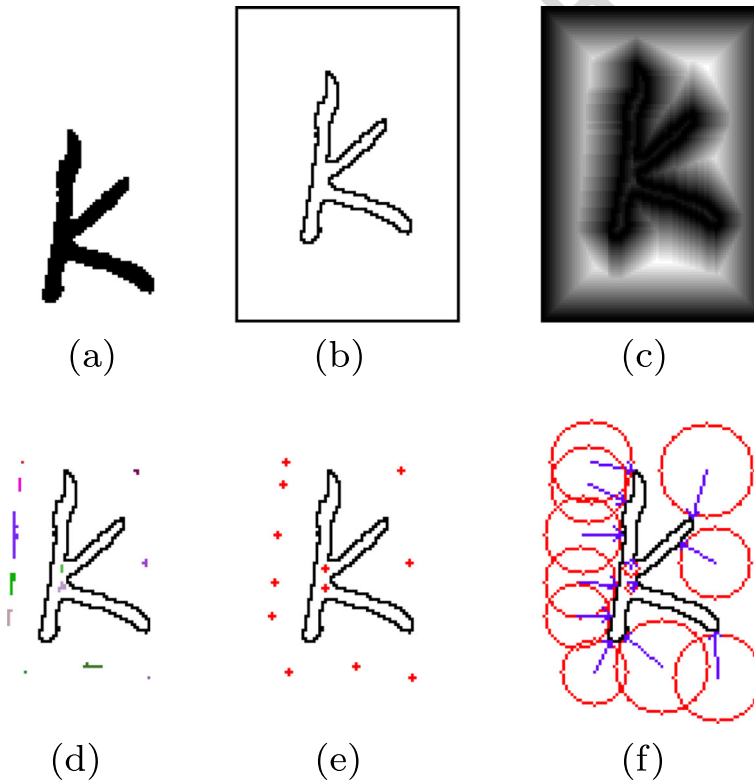


Fig. 1 Keypoint extraction steps: **a** Original binary image. **b** $W_F \times H_F$ image after normalization ($\alpha = 0.25$). **c** DT image. **d** Regions of equal maximal intensity highlighted in different colors **e** Keypoints ($k = 11$). **f** Keypoint vectors ($\alpha = 1$): Circle radii correspond to the keypoint distance from the nearest contour point, and arrows show the orientation of the vector delimited by the keypoint and its nearest contour point

161 Keypoints are extracted as follows: First, the original image (Fig. 1a) is normalized by
 162 applying contour detection and image translation (Fig. 1b). Then, background information is
 163 generated using DT (Fig. 1c). Before applying DT, a 1-pixel-width border frame is added to
 164 the normalized image in order to delimit the object so DT does not systematically generate
 165 maxima at the borders. Next, regions of equal maximal intensity are detected on the DT
 166 image using a $k \times k$ square window (Fig. 1d). A region of equal maximum intensity is the
 167 contiguous pixel "islands" that have higher intensities than their neighboring pixels. They
 168 correspond to the regions of highest intensity in Fig. 1c that are shown in different colors
 169 in Fig. 1d. Finally, the detected regions are represented using their centers of masses which
 170 are taken as keypoints (Fig. 1e).

171 Contour detection is used to produce a compact representation of the original image
 172 that reduces the number of foreground pixels but preserves the visual information [13, 55].
 173 Afterwards, keypoints can be extracted from regions inside and outside the object (Fig. 1e).

174 The dimensions (W_F, H_F) of the frame used before applying DT are calculated as
 175 follows:

$$W_F = (1 + a) W_{BB}, \quad H_F = (1 + a) H_{BB} \quad (1)$$

176 where W_{BB} and H_{BB} are the dimensions of the object's bounding box, and $a \geq 0$ is intro-
 177 duced to allow for a space between the object contours and the frame pixels in order to
 178 extract keypoints in these regions. The object is translated towards the center of the frame.
 179 In the present work, we set a empirically (Section 4.2), so that the bounding box is located
 180 in a good proximity from the foreground object (Fig. 1). A too small a would make the
 181 frame borders too close to the foreground object, which places the keypoints too close to
 182 the contour making them more vulnerable to contour noise, while a too large a would make
 183 the frame borders too far from the object contour, which increases the size of the feature
 184 extraction windows (Fig. 5) and hence puts more weight on global details of the object on
 185 the detriment of local details.

186 Regions of equal maximal intensity are detected using a $k \times k$ square window located at
 187 each DT image pixel. The parameter k affects the number of extracted local maxima. The
 188 larger k gets, the fewer keypoints are detected (Fig. 2). Therefore, parameter k controls the
 189 number of generated keypoints. In this paper, we set k empirically (Section 4.2.1), and we
 190 leave further investigation on setting k automatically for future work.

191 Due to using DT to generate background information, the extracted keypoints are in
 192 locus of symmetry between foreground pixels and thus characterize the object using its
 193 local symmetries. We anticipate the significance of such keypoints in shape representation

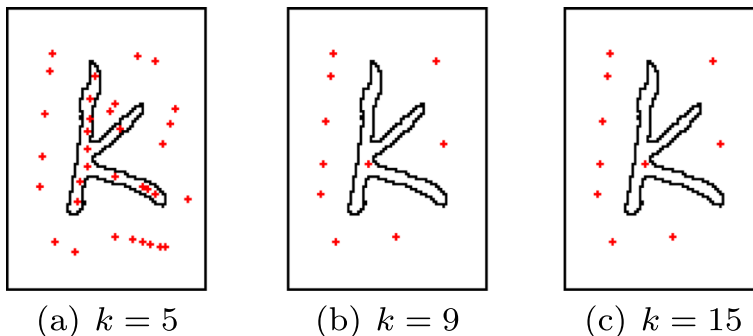


Fig. 2 Effect of the parameter k on the number of keypoints

due to the importance of symmetry as a characteristic of patterns that is exploited in human perception [52] and in computational image matching [25].

The complexity of the keypoint extraction step can be estimated as follows: The distance transform and regions of equal maximal intensity detection require two processes that browse the entire image pixels, hence they make a $2 \cdot O(n)$ complexity, with n here representing the number of image pixels. Then, keypoint detection in regions of equal maximal intensity make a complexity of $O(n)$.

3.2 Keypoint selection

The initial number of keypoints can be reduced by filtering out the redundant and sensitive keypoints. Redundant keypoints duplicate representing the same details of the image, and keypoints that are located very close to contours are sensitive to insignificant changes in image local details.

A measure of keypoint salience is introduced for keypoint ranking and selection. A salient keypoint is defined according to two aspects:

- It has few keypoints in its local area, and thus it is non-redundant.
- It is not located very close to foreground points, and thus it is robust against insignificant changes in image local details.

Formally, the salience $\gamma(i)$ of a keypoint K_i is calculated as follows:

$$\gamma(i) = \frac{d_i}{1 + n_i} \quad (2)$$

where d_i is the distance from keypoint K_i to its closest contour or frame border point, and n_i is the number of close keypoints. A keypoint K_j is considered close to K_i if it is located within a distance to K_i proportional to d_i .

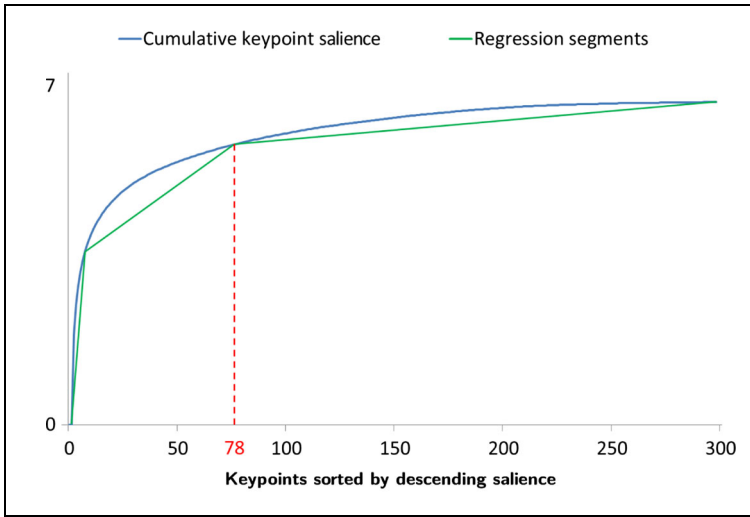
Our hypothesis for automatically selecting the most salient keypoints is as follows: We observe that the range of salience values commonly indicates three types of keypoints (Fig. 3c). The first type corresponds to few keypoints with extreme salience values, the second type corresponds to a larger number of keypoints with increasing redundancy, and the third type corresponds to keypoints with high redundancy and closeness to the contours or frame borders. Since keypoints of the third type are redundant and sensitive, they are filtered out.

In order to filter out keypoints of the third type, we calculate the cumulative keypoint salience $\Gamma(i)$ for a number i of keypoints ranked in their descending salience measures, as follows:

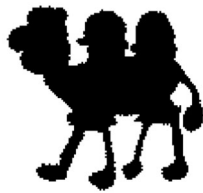
$$\Gamma(i) = \ln \left(\sum_{j=1}^i \gamma(j) \right) \quad (3)$$

Figure 3a shows a typical curve of Γ as a function of the number of accumulated keypoints. The curve of Γ can be roughly segmented into three parts corresponding to the types of keypoints. In order to find keypoints of each type, a two-dimensional search is used to detect the three segments that minimize the area between them and the curve of Γ . Then, keypoints corresponding to the first and second types are selected. In the literature, a similar strategy has been reported in [47] to automatically detect salient corner points in online sketches using scale-space filtering and digital ink attributes (e.g. pen speed, curvature).

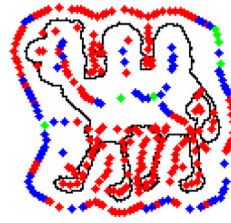
Figure 4 illustrates the benefit of automatic selection of keypoints using their salience scores. The top example shows matching between an image and its slightly different version



(a)

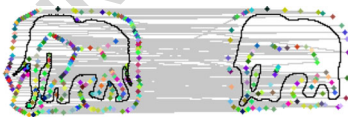


(b)

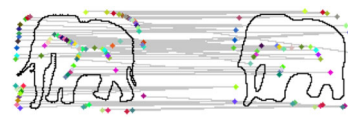


(c)

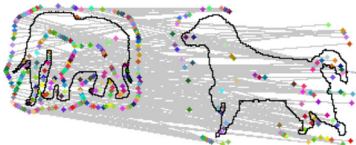
Fig. 3 Keypoint selection: **a** Curve approximation by three segments applied on image **b**, **c** keypoints of the first type in green, keypoints of second type in blue, and keypoints of third type in red. Automatic keypoint selection reduces the number of keypoints from 298 to 78



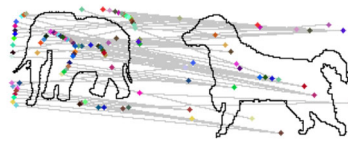
(a) $N_{Left} = 229, N_{Right} = 140,$
 $similarity = 98.95\%$



(b) $N_{Left} = 80, N_{Right} = 93,$
 $similarity = 98.95\%$



(c) $N_{Left} = 229, N_{Right} = 270,$
 $similarity = 98.30\%$



(d) $N_{Left} = 80, N_{Right} = 59,$
 $similarity = 97.76\%$

Fig. 4 Automatic keypoint selection reduces the number of keypoints while improving matching performances ($similarity$ is calculated according to (7))

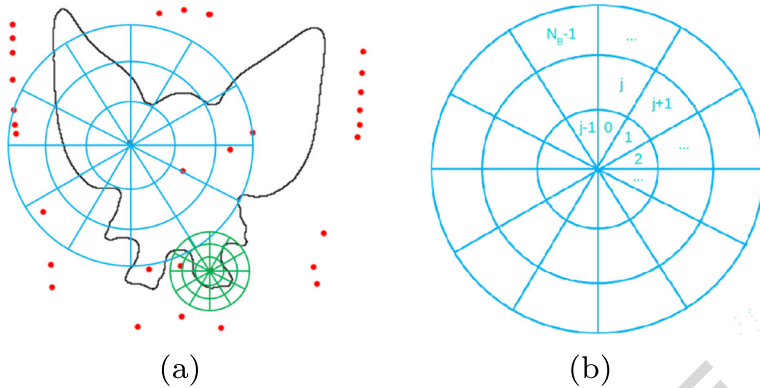


Fig. 5 Keypoint feature extraction using size-adaptive layouts

that is generated using a Gaussian filtering ($\sigma = 3$) followed by binarization [33], leading to remove the granularity of some local details. Using automatic keypoint selection does not affect the similarity between the two images, which shows that the filtered keypoints are not crucial for matching. On the other hand, the bottom example shows matching between two images belonging to different classes. Here, using automatic keypoint selection decreases the similarity, which shows that automatic keypoint selection has removed a significant number of keypoints causing false positives. In both cases, the reduction in the number of keypoints is considerable.

The complexity of the keypoint selection step can be estimated as $2 \cdot O(n^2)$, with n here representing the number of initially extracted keypoints.

3.3 Feature representation and matching

The last step is to calculate a feature vector to each keypoint K_i . For this purpose, we use a scale-invariant circular layout which radius r_i is proportional to the distance between the keypoint K_i and its closest contour point (Fig. 5b):

$$r_i = \alpha \times d_i \tag{4}$$

where α is a heuristic. The idea is to set $\alpha > 1$ to allow taking into account the closest contour points in the smallest distance bins. Then, a histogram h_i is extracted by calculating the distribution of contour points in distance and angle bins, i.e. $h_i(j)$ holds the number of contour points that are located inside the feature window bin of index j (Fig. 5b). The distance between two histograms is expressed by the χ^2 statistic:

$$\chi^2(h_1, h_2) = \frac{1}{2} \sum_{j=0}^{N_B-1} \frac{[h_1(j) - h_2(j)]^2}{h_1(j) + h_2(j)} \tag{5}$$

where N_B is the number of bins in a keypoint histogram. Using the distance d_i to set the radius of the feature layout makes the descriptor scale-invariant.

255 The dissimilarity d between two images I_1 and I_2 is estimated by the cumulative
256 minimum distance between the images' keypoint histograms:

$$d(I_1, I_2) = \frac{1}{N_1} \sum_{i=0}^{N_1-1} \min_{0 \leq j < N_2} \left\{ \chi^2 \left(h_i^1, h_j^2 \right) \right\} \quad (6)$$

257 where N_1 and N_2 are the number of keypoints in images I_1 and I_2 . Because $d(I_1, I_2)$ is
258 asymmetric, we express the distance between two images I_1 and I_2 as follows:

$$D(I_1, I_2) = \frac{d(I_1, I_2) + d(I_2, I_1)}{2} \quad (D \in [0, 1]) \quad (7)$$

259 The smaller $D(I_1, I_2)$ is, the more similar I_1 and I_2 are.

260 The feature vector is translation-invariant due to using the object's bounding box for
261 image normalization. Scale-invariance is partly ensured in the keypoint filtering step (using
262 a radius d_i of the circular region used in the salience measure (Eq. 2) that changes with
263 the size of the image) and the feature representation step (since the radius of the feature
264 extraction circular window depends on each keypoint and also on the object's size), but
265 partly hindered by fixing parameter k making it scale-dependent. Rotation-invariance can
266 be ensured by using the orientation of the vector delimited by the keypoint and its nearest
267 contour point as a reference orientation (Fig. 1), or by using shifted matching of the key-
268 points' feature vectors. In case mirrored matching is necessary, it can be implemented by
269 mirroring one of the feature vectors and repeating the matching then taking the average.

270 The complexity of feature representation can be estimated as $O(m) \cdot O(n)$, with m here
271 representing the number of initially extracted keypoints, and n the number of contour points.
272 Feature matching requires $N_B \cdot O(m) \cdot O(n)$ with m and n are the number of keypoints in
273 images I_1 and I_2 respectively.

274 4 Experiments

275 4.1 Datasets

276 Evaluation is done using five datasets (Fig. 6): The Kimia 216 dataset [45] and the MPEG-
277 7 dataset [6] include silhouette images that are neat and which contain single component
278 objects. Zanibbi and Yu's dataset [56] contains handwritten mathematical expressions which
279 exhibit handwriting fluctuations and component displacement, which also appear in Liang
280 et al.'s dataset [28] of hand-drawn diagram sketches. The Tobacco 800 dataset [60] contains
281 logo images that are taken from scanned documents and they are the noisiest compared to
282 the other datasets. The datasets can be thought of as clusters' centers extracted from large
283 datasets that are typically used in data-driven approaches [19]. On the other hand, they are
284 fit to evaluate shape descriptors as they represent varied image classes and exhibit different
285 challenges (e.g. noise, handwriting fluctuations).

286 We used Kimia, Zanibbi and Yu, and Tobacco datasets as training datasets to empirically
287 set our parameters. The choice is made due to the reasonable sizes of these datasets, and
288 their characteristics relative to the remaining two datasets:

- 289 – Kimia 216 dataset can be considered a smaller subset of MPEG-7. It contains similar
290 classes and less image variations.



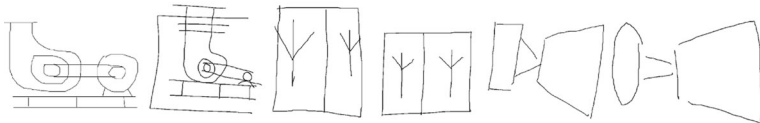
(a) Kimia’s dataset [41]: 216 images, 18 classes, and 12 instances



(b) MPEG-7 dataset [6]: 1400 images, 70 classes, and 20 instances

$$v = \frac{2v_+ + v_-}{v_+ + v_-} \quad v = \frac{2v_+ v_-}{v_+ + v_-} \left(\frac{1}{a^2} + \frac{1}{b^2} \right) \left(\frac{1}{a_2} + \frac{1}{b^2} \right)$$

(c) Zanibbi and Yu’s dataset [54]: 200 images, 20 classes, and 10 instances



(d) Liang et al. dataset [27]: 1086 images, 35 classes, and between 17 and 22 instances



(e) Tobacco 800 logo dataset [59]: 412 images, 35 classes, and between 1 and 68 instances

Fig. 6 Samples of the dataset images

- Zanibbi and Yu’s dataset is handwritten, which is the same main feature of Liang et al.’s dataset. 291
292
- Tobacco dataset is used for its significant noise. 293

4.2 Descriptor evaluation 294

Before evaluating the descriptor, we set its parameters as follows: The parameter for setting the normalization frame’s dimensions is set $a = 0.25$, which insures a scale-invariant frame with space between its borders and the object contours. A keypoint K_j is considered close to a keypoint K_i if the distance between them is equal or less than $\frac{d_i}{4}$, where d_i is the distance between keypoint K_i and its closest contour or frame border point. The radial and angular numbers of bins in the keypoint descriptor are set as 4 distance bins and 8 angle bins in order to make a trade-off between distinctiveness and robustness. A small

295
296
297
298
299
300
301

302 number of bins compromises the descriptor's distinctiveness, while a larger number of bins
 303 causes sensitivity to noise and fluctuations [41]. The constant for configuring the keypoint-
 304 dependent feature layout radius is set $\alpha = 1.5$ in order to insure taking into account the
 305 closest contour points in the smallest distance bins.

306 Evaluation is done using the *precision at n* metric [2], denoted $P@n$, which is calculated
 307 as follows:

$$P@n = \frac{| \{n \text{ retrieved images} \} \cap \{ \text{relevant images} \} |}{| \{n \text{ retrieved images} \} |} \quad (8)$$

308 Due to variations in the number of class instances, we set the number of retrieved images n
 309 as query-dependent and equivalent to the number of the query's class instances. This makes
 310 $P@n$ equal to *precision* and *recall*. The larger $P@n$ is, the better *precision* and *recall* the
 311 descriptor shows. In the following, we specify the parameter n in $P@n$ when it is fixed (e.g.
 312 results are shown for a single dataset with a fixed number of class instances).

313 4.2.1 Keypoint sampling evaluation

314 During keypoint extraction, the parameter k defines the size of the local maxima detection
 315 window and thus affects the number of extracted keypoints (Fig. 2). We evaluate the effect
 316 of this parameter on matching performances using the Kimia 216, Zanibbi and Yu, and
 317 Tobacco datasets as training datasets for this empirical setting.

318 Figure 7 shows curves of $P@n$ as a function of k . We observe that the matching per-
 319 formance eventually decreases when k increases, and that the best matching performances
 320 correspond to $k = 3$, which means that the best way is to keep a maximum number of key-
 321 points that will be later filtered during the keypoint selection step. According to the results
 322 of this experiment, we set $k = 3$ empirically and use it in subsequent experiments.

323 4.2.2 Keypoint distinctiveness evaluation

324 The distinctiveness of BSK's keypoints is assessed by comparison with equidistant sampling
 325 which is used in numerous descriptors, namely shape contexts [5]. We perform experiments

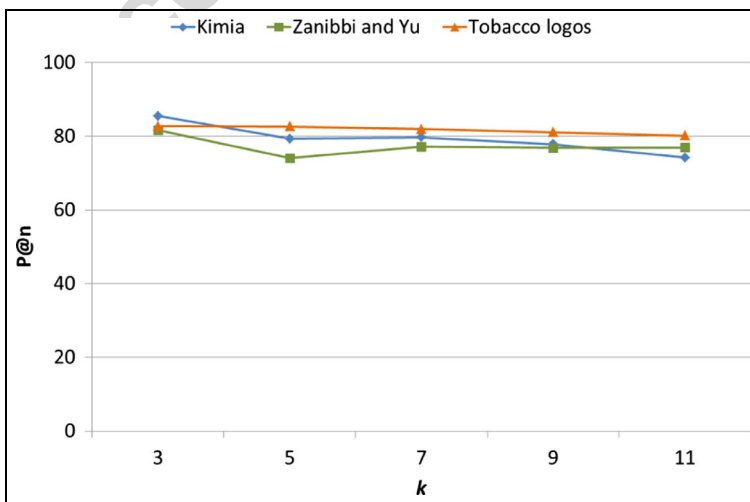


Fig. 7 Effect of varying the parameter k on $P@n$

of image retrieval using the Kimia 216 dataset where each image is used as a query and the average $P@12$ is calculated for all queries. We extract the same number of keypoints using BSK and shape contexts and perform matching using our keypoint matching steps (Section 3.3). In order to make the comparison between BSK keypoints and shape contexts fair, we introduced two modifications on the shape contexts: Features are extracted from equidistant keypoints from the contour and all the remaining contour points are considered when calculating the keypoint's histogram, unlike the original shape context descriptor where only the sampled keypoints are considered. In addition, scale-invariance is introduced by making the circular feature extraction layout's size adaptive to the shape by calculating the distance between each keypoint and its farther contour point, instead of using static log-polar layouts. Consequently, these modifications led to better results when compared with the original shape contexts considering only equidistant keypoints and using static log-polar layouts for feature extraction.

Figure 8 shows performances of BSK keypoints and shape contexts. For small numbers of extracted keypoints, using equidistant keypoints outperforms BSK keypoints. Then, starting from 40 keypoints, BSK outperforms shape contexts and the gap increases in correlation with the number of keypoints. In fact, using 40 BSK keypoints outperforms using 100 shape contexts. This result shows that our keypoints are distinctive and outperform the widely-used equidistant keypoint sampling scheme.

4.2.3 Keypoint selection evaluation

The keypoint selection step aims to reduce the number of keypoints by removing the redundant ones and the ones too close to the shape contour. Figure 9 shows retrieval performances expressed in $P@n$ as a function of the percentage of keypoints using the Kimia 216 dataset, Zanibbi and Yu's dataset, and Tobacco logos dataset. For the Kimia 216 dataset, performances increase when the percentage of keypoints increases. As for Zanibbi and Yu's and

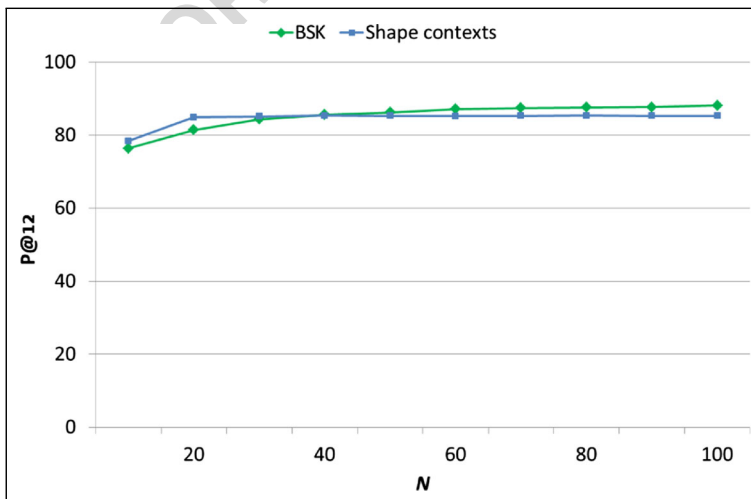


Fig. 8 $P@12$ as a function of the number of keypoints N for BSK and shape contexts on the Kimia 216 dataset

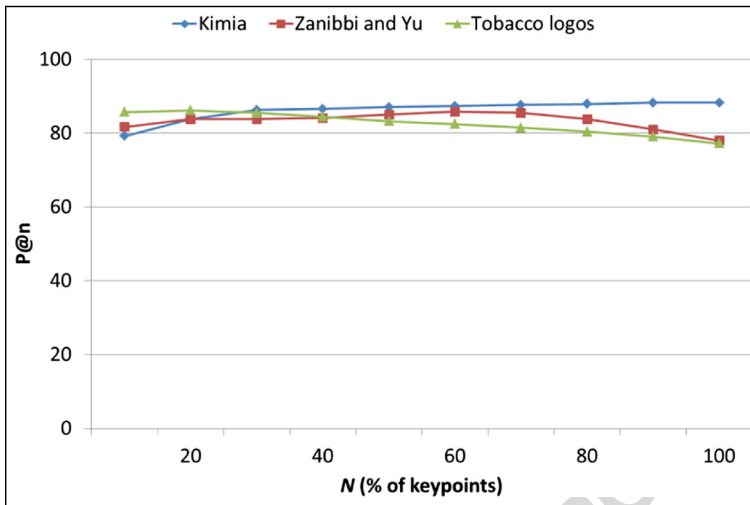


Fig. 9 $P@n$ as a function of the percentage of used keypoints relative to the total number of extracted keypoints using BSK

351 Tobacco logos datasets, optimal performances are obtained when not all of the keypoints
 352 are used (when 20% and 60% of keypoints are selected respectively).

353 Table 1 shows retrieval performances of BSK when all keypoints are used and when
 354 keypoint selection is performed. For all datasets, the reduction in number of keypoints is
 355 significant and roughly makes the third of total keypoints. In case of Zanibbi and Yu's and
 356 Tobacco logos datasets, matching performances improve. However, they decrease in case
 357 of Kimia 216 dataset. This result suggests that our keypoint salience-based selection is
 358 effective when the initial number of keypoints is relatively large (cases of Zanibbi and Yu's
 359 and Tobacco logos datasets). When the initial number of keypoints is relatively small (case
 360 of Kimia 216 dataset), the keypoint selection step would better be skipped. This can be done
 361 using a threshold on the initial number of keypoints.

362 4.3 Performance comparison with other descriptors

363 Tables 2, 3 and 4 show results of comparing our descriptor with other state of the art meth-
 364 ods. Results are shown according to metrics that are used in available published work; In
 365 case of Kimia's dataset, we calculate the retrieval performance metric reported in several
 366 published papers, that is the number of relevant retrieved images for each of the top 6 ranks

Table 1 $P@n$ and number of keypoints N using BSK with all keypoints and with selected keypoints

Implementation	All keypoints		Selected keypoints	
	$P@n$	N	$P@n$	N
Kimia 216	88.27 %	147	85.49 %	58
Zanibbi and Yu	78.0 %	1610	81.65 %	564
Tobacco logos	77.21 %	1203	82.74 %	379

Table 2 Performances using the Kimia 216 dataset [45]

Method	BSK	SRD [11]	SC [5]	PSSG [3]
$P@6$	93.83 %	86.96 %	92.12 %	99.22 %

and the percentage calculated by summing these numbers. We refer to it as $P@6$. In case of the MPEG-7 dataset, the $P@20$ metric is used. In case of Liang et al.'s dataset, comparison is done using the *mean average precision (MAP)* metric [2].

BSK yielded competitive performances in all the datasets we used. Results in Tables 2–4 show that BSK is effective in case of computer generated silhouette images of the Kimia 216 dataset and the MPEG-7 dataset, and hand-drawn sketch images of Liang et al.'s dataset that exhibit high sketching perturbations and drawing style variations. In addition, BSK reached $P@10 = 82.74\%$ and $P@n = 81.65\%$ on Zanibbi and Yu's dataset [56] and the Tobacco logos dataset [60] respectively, which compares well against the support region descriptor (SRD) [11] that gave $P@10 = 47.6\%$ and $P@n = 82.55\%$. This demonstrates the effectiveness of BSK in case of handwritten images of Zanibbi and Yu's dataset and in case of the noisy scanned images of the Tobacco logos dataset. It is worth mentioning that SRD is designed to be robust against noise by combining local and global features. Results on the MPEG-7 dataset are relatively lower due to two image variations: First, the dataset has significant scale variance that challenges our descriptor, which has some scale-invariance limitations due to fixing the parameter k during the keypoint extraction step. Second, the dataset has also a significant number of mirrored images, and we do not currently take this into account during the feature vector matching (5).

Although BSK does not show supremacy over all other descriptors, results showed that it compares strongly against various types of methods. BSK outperformed shape context that uses equidistant sampling (Table 2) and other salience-based keypoint descriptors such as the contour salience descriptor (CS) [15] (Table 3), the minimal spanning tree (MST), Laplacian spectrum with geometry (LS+G) [16] and the LS+G [16] descriptors (Table 4). BSK also compares well against methods that combines local and global features such as SRD (Table 2), and graph-based methods such as MST, LS+G, and TPG [28] (Table 4). On the other hand, BSK was outperformed by PSSG [3] on the Kimia dataset [45] and TSDIZ [20] and SSD+GF [35] on the MPEG-7 dataset [6]. In case of PSSG [3], the use of skeleton pruning makes PSSG robust against contour noise, which explains the improved performances on the Kimia dataset. PSSG [3] skeletons, on the other hand, are vulnerable to shape ambiguity [38], but this problem is minor in the Kimia dataset and does not affect PSSG. As for TSDIZ [20] and SSD+GF [35], we explain the results by their invariance to scale, since SSD+GF [35] is based on the multiscale tensor scale transform, and SSD+GF [35] uses a scale-invariant salience detection that analyzes the curvature of contour points.

We further evaluate BSK's robustness against contour noise by comparing it with other keypoint-based descriptors. For this purpose, we generated noisy versions of the Kimia dataset images [45]. First, we removed the small contour perturbations using a Gaussian filter ($\sigma = 1$) followed by binarization [33]. Then, we produced 10 sets of images with

Table 3 Performances using the MPEG-7 dataset [6]

Method	BSK	CS [15]	TSDIZ [20]	SSD+GF [35]
$P@10$	75.48 %	36 %	81 %	85 %

Table 4 Performances using Liang et al.'s dataset [28]

Method	BSK	MST [27]	LS+G [16]	TPG [28]
<i>MAP</i>	83.83 %	29.8 %	50.9 %	61.6 %

404 contour noise levels from 10% to 100%. The noise is generated by a random removal of
 405 a percentage of contour pixels. Using the noisy sets of images, BSK is compared against
 406 shape contexts [5] (as used in Section 4.2.2) and a similar descriptor that uses the Harris
 407 detector [22], as a widely-used corner detector. For the three descriptors, the same number
 408 of keypoints are selected, which is equal to the number of salient keypoints selected auto-
 409 matically by BSK. For shape contexts, a similar number of keypoints are selected by using
 410 uniform sampling. As for the Harris-based descriptor, the keypoints are selected according
 411 to their descending Harris detector response.

412 Figure 10 shows examples of keypoints extracted using the three descriptors for a neat
 413 image with smooth contours and its noisy version after generating 50% contour noise. We
 414 observe that BSK and uniform sampling, used in shape contexts, produce keypoints that are
 415 sparse and cover all image details, while keypoints produced by the Harris detector tend to
 416 be localized on few corners. In order to estimate the effect of noise on keypoint location
 417 shifting, we calculate the keypoint *average shift* as follows:

$$averageshift = \frac{1}{N_2} \sum_{j=1}^{N_2} \min_i |\vec{K}_i \vec{K}_j| \tag{9}$$

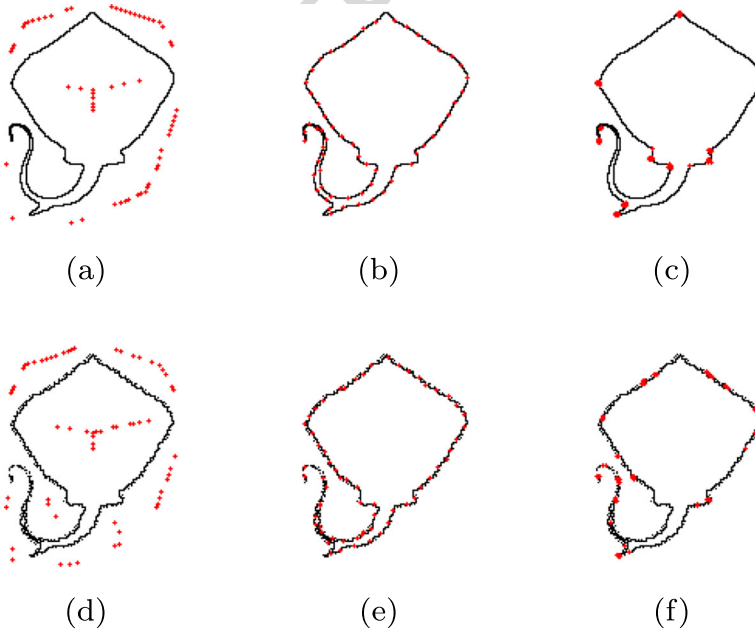


Fig. 10 Keypoints extraction methods. Top: results for a neat image (60 keypoints are extracted automatically). Bottom: results for an image with 50% contour noise (62 keypoints are extracted automatically). Right to left: BSK, equidistant sampling, and Harris detector

where N_1 and N_2 are the numbers of keypoints in the neat image and its noisy version respectively, and $\overrightarrow{K_i K_j}$ is equal to the Euclidean distance between a keypoint K_i of the neat image and a keypoint K_j in the noisy image. By taking the minimum value of $|\overrightarrow{K_i K_j}|$, we find keypoint K_i that corresponds to the previous location of keypoint K_j before a shift caused by noise occurs.

For the examples in Fig. 10, the rounded values of *average shift* for BSK, uniform sampling and Harris detector are 3 pixels, 2 pixels, and 15 pixels respectively (fixing the number of salient keypoints does not change the behavior of *average shift*). This shows that BSK and uniform sampling are more robust than Harris detector, since their keypoints do not shift much when exposed to contour noise. Accordingly, BSK and uniform sampling outperform the Harris detector in terms of retrieval performances, as shown in Fig. 11. BSK also outperforms uniform sampling although BSK's *average shift* value is slightly larger than uniform sampling's *average shift* value. This result provides an evidence that extracting keypoints from the background, instead of the contours, is a good strategy to reduce the effect of noise.

4.4 Evaluation in content-based document image retrieval

We integrate BSK in a document image retrieval system reported by Chatbri et al. [14]. This system takes input in the form of sketched mathematical expressions, and outputs a ranked list of document images that contain the user' query. This is done by using a finding the connected components of the document image that are similar to the connected components of the query using contour points distribution histograms, and then locate the ones that have a spatial arrangement similar to the query. BSK is integrated as a last step that further compares the query with the detected occurrences in the database images.

Table 5 shows a performance comparison including Chatbri et al.'s original system against when BSK is integrated, in addition to another content-based retrieval system by Zanibbi and Yu [56]. Performances are expressed with two metrics: $P - Recall$ expresses the system's ability to find the correct document image (i.e. document *page*), and $A - Recall$

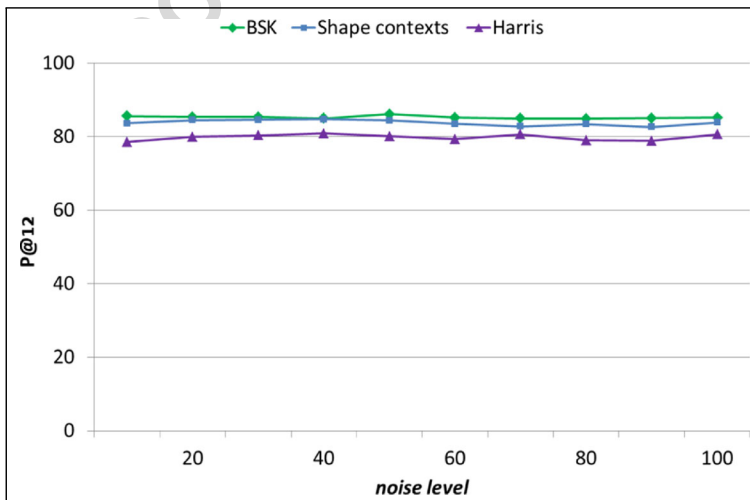


Fig. 11 $P@12$ as a function of *noise level* for BSK, shape contexts [5] and the Harris detector [22] using noisy versions of the Kimia dataset images [45]

Table 5 Average values of *P-Recall* and *A-Recall* calculated for $n = 1, 5, 10$

Method	n	Printed queries		Handwritten queries	
		<i>P-Recall</i>	<i>A-Recall</i>	<i>P-Recall</i>	<i>A-Recall</i>
Chatbri et al. [14]	1	100%	94.28%	40.0%	27.83%
	5	100%	96.78%	63.5%	51.15%
	10	100%	96.78%	73.5%	57.92%
Chatbri et al. [14] + BSK	1	92.5%	89.29%	54.0%	47.84%
	5	100%	96.29%	70.0%	59.89%
	10	100%	96.78%	75.0%	62.43%
Zanibbi and Yu [56]	1	.	90%	38.6%	26.7%
	5	.	90%	54.9%	39.8%
	10	.	90%	63.2%	43.3%

444 expresses the system's ability to find the correct *area* of the query's occurrence inside
 445 the document image [14, 56]. The metrics are calculated for the top- n retrieved document
 446 images ranked by occurrence similarity with the query.

447 According to the results, BSK improves retrieval performances especially when hand-
 448 written queries are used. Improvement reaches 20% of *A-Recall* when the top-1 images
 449 are retrieved. On the other hand, performances drop in case of printed queries.

450 4.5 Discussion

451 The proposed descriptor is able to extract distinctive keypoints as demonstrated by compar-
 452 ison with similar numbers of shape context keypoints extracted using equidistant contour
 453 points sampling on the Kimia 216 dataset. In fact, BSK is able to outperform shape contexts
 454 using significantly fewer keypoints. This is further proven when BSK outperforms methods
 455 that detect salient points in the image contour using the other datasets. An interesting direc-
 456 tion motivated by these results is to combine BSK keypoints with salient keypoints of the
 457 contour for the sake of better distinctiveness.

458 Experiments on challenging images, such as fluctuated handwritten mathematical
 459 expressions of Zanibbi and Yu's dataset and hand-drawn diagram sketches of Liang et al.'s
 460 dataset, demonstrate the reliability of BSK, as it outperforms largely other methods. Meth-
 461 ods used for comparison include graph-based descriptors which are known for their high
 462 matching performances and ability to perform partial matching. The reliability of BSK is
 463 further demonstrated when assessed on the noisy scanned images of the Tobacco logos
 464 dataset.

465 The keypoint selection based on keypoint salience is effective in reducing the number
 466 of keypoints without significantly compromising the descriptor's distinctiveness. However,
 467 the performances improve when the initial number of keypoints is relatively large. For this
 468 purpose, a threshold on the initial number of keypoints can be used to activate or skip the
 469 salience-based keypoint selection.

470 BSK is adequate to be used for applications of image retrieval from document image
 471 databases. This is shown by the performance improvement it leads to when integrated in a
 472 standard document image retrieval system.

473 Finally, BSK is currently not suitable for real-time applications. In order to become so,
 474 the following can be done:

- Use parallel computing: Currently, all the procedures are executed sequentially due to limited memory. Otherwise, several steps of the algorithm can be made faster, including a parallel implementation of the distance transform [9], faster connected components extraction, and feature extraction for each keypoint in parallel. 475
476
477
478
- Resize the images to a reasonably smaller scale, which will speed up the steps aforementioned, and use integral images [48] to speed up local processes. 479
480

5 Conclusions and future work 481

In this paper, we introduced a descriptor for binary image matching using image salient keypoints. The proposed binary salient keypoints descriptor (BSK) generates background information in binary images, then extracts keypoints using pixels that have specific characteristics in their local areas. A measure of keypoint salience is used for automatically selecting the most salient keypoints and filtering out the redundant and sensitive ones. 482
483
484
485
486

The proposed descriptor has been evaluated using five public datasets of silhouette images, handwritten mathematical expressions, hand-drawn diagram sketches, and scanned logo images. Experimental results and comparison with state of the art methods demonstrated that BSK has competitive matching performances when applied on various types of images, including challenging images of fluctuated handwriting and noisy scanned images. Furthermore, BSK's integration in a content-based document image retrieval system leads to improving the system's performances considerably. 487
488
489
490
491
492
493

BSK paves the way for future research on salient keypoints detection in the background of binary images, as an unconventional new way of binary image analysis. In addition, it can be improved by tuning its keypoint extraction, filtering, and feature representation modular stages. We identify areas of future work as follows: 494
495
496
497

- Scale-invariance can be improved by setting parameter k automatically. For instance, instead of fixing the value of k according to empirical results on a number of datasets, k can be set according to each image taking into account its characteristics that can lead to produce more keypoints (e.g. scale, texture). On the other hand, the feature vector matching equation (5) can be modified to implement mirror matching. 498
499
500
501
502
- On the other hand, it would be interesting to make the parameters of BSK set in an evolutionary or data-driven way. For instance, one can try using a genetic algorithm where the genetic representation uses BSK's parameters and the fitness function is a performance metric (e.g. P@n) in a training dataset. 503
504
505
506
- We introduce a specific definition of keypoint salience that is based on the proximity between keypoints and the distance between a keypoint and the object contour. Alternatively, other definition of keypoint salience can be defined for specific applications and compared. 507
508
509
510

In addition to image matching, it would be interesting to investigate applying our descriptor in similar applications such as image registration, particularly magnetic resonance (MR) images where the shape information provides significant features [36]. Moreover, it would be interesting to apply our descriptor in color images, in combination with other descriptors. For instance, BSK would be fit to embed in the logo retrieval framework proposed in [21], preceded by edge detection, and combined with color images features (SIFT and segmentation patches). 511
512
513
514
515
516
517

518 **Acknowledgments** This work has emanated from a research grant in part from the Monbukagakusho
 519 Scholarship sponsored by the Japanese Government, in part from the Irish Research Council (IRC) under
 520 Grant Number GOIPD/2016/61, and in part from Science Foundation Ireland (SFI) under Grant Number
 521 SFI/12/RC/2289 (Insight Centre for Data Analytics). The authors would also like to thank Dr. Richard
 522 Zanibbi for providing his dataset [56], Dr. Mathieu Delalandre and Dr. Alireza Alaei for their assistance with
 523 extracting logos from the Tobacco 800 dataset, and Dr. Shuang Liang for her assistance with her dataset [28].

524 References

- 525 1. Abbasi S, Mokhtarian F, Kittler J (1999) Curvature scale space image in shape similarity retrieval.
 526 *Multimed Syst* 7(6):467–476
- 527 2. Baeza-Yates R, Ribeiro-Neto B (1999) *Modern information retrieval*, vol 463. ACM Press, New York
- 528 3. Bai X, Latecki LJ (2008) Path similarity skeleton graph matching. *IEEE Trans Pattern Anal Mach Intell*
 529 30(7):1282–1292
- 530 4. Bay H, Tuytelaars T, Van Gool L (2006) Surf: speeded up robust features. In: *European conference on*
 531 *computer vision (ECCV)*, pp 404–417
- 532 5. Belongie S, Malik J, Puzicha J (2002) Shape matching and object recognition using shape contexts. *IEEE*
 533 *Trans Pattern Anal Mach Intell* 24(4):509–522
- 534 6. Bober M (2001) MPEG-7 visual shape descriptors. *IEEE Trans Circ Syst Vid Technol* 11(6):716–719
- 535 7. Breuß M (2013) *Innovations for shape analysis: models and algorithms*. Springer Science & Business
 536 Media
- 537 8. Bunke H, Riesen K (2012) Towards the unification of structural and statistical pattern recognition.
 538 *Pattern Recogn Lett* 33(7):811–825
- 539 9. Cao T-T, Tang K, Mohamed A, Tan T-S (2010) Parallel banding algorithm to compute exact distance
 540 transform with the gpu. In: *ACM SIGGRAPH Symposium on interactive 3d graphics and games*. ACM,
 541 pp 83–90
- 542 10. Chalechale A, Naghdy G, Mertins A (2005) Sketch-based image matching using angular partitioning.
 543 In: *IEEE Transactions on systems, man and cybernetics, part A: systems and humans*
- 544 11. Chatbri H, Kameyama K, Kwan P (2013) Sketch-based image retrieval by size-adaptive and noise-robust
 545 feature description. In: *International conference on digital image computing: techniques and applications*
 546 (DICTA). IEEE, pp 1–8
- 547 12. Chatbri H, Davila K, Kameyama K, Zanibbi R (2015) Shape matching using keypoints extracted from
 548 both the foreground and the background of binary images. In: *International Conference on image*
 549 *processing theory, tools and applications (IPTA)*. IEEE, pp 205–210
- 550 13. Chatbri H, Kameyama K, Kwan P (2015) A comparative study using contours and skeletons as shape
 551 representations for binary image matching. *Pattern Recognition Letters*
- 552 14. Chatbri H, Kameyama K, Kwan P (2015) Towards a segmentation and recognition-free approach
 553 for content-based document image retrieval of handwritten queries. In: *Asian Conference on pattern*
 554 *recognition (ACPR)*. IAPR
- 555 15. da S Torres R, Falcao AX (2007) Contour salience descriptors for effective image retrieval and analysis.
 556 *Image Vis Comput* 25(1):3–13
- 557 16. Demirci FM, van Leuken RH, Veltkamp RC (2008) Indexing through laplacian spectra. *Comput Vis*
 558 *Image Underst* 110(3):312–325
- 559 17. Donoser M, Riemenschneider H, Bischof H (2010) Efficient partial shape matching of outer contours,
 560 281–292
- 561 18. Dubey SR et al (2016) Multichannel decoded local binary patterns for content-based image retrieval.
 562 *IEEE Trans Image Process* 25(9):4018–4032
- 563 19. Eitz M, Hays J, Alexa M (2012) How do humans sketch objects? *ACM Trans Graph* 31(4):44–1
- 564 20. Fernanda AA, Paulo AV, da S Torres MR, Falcão AX (2010) Shape feature extraction and description
 565 based on tensor scale. *Pattern Recogn* 43(1):26–36
- 566 21. Fu J, Wang J, Lu H (2010) Effective logo retrieval with adaptive local feature selection. In: *ACM*
 567 *International conference on multimedia (ACM MM)*. ACM, pp 971–974
- 568 22. Harris C, Stephens M (1988) A combined corner and edge detector. In: *Alvey vision conference*, pp
 569 147–151
- 570 23. Kopf S, Haenselmann T, Effelsberg W (2005) Enhancing curvature scale space features for robust shape
 571 classification. In: *International conference on multimedia and expo (ICME)*. IEEE, p 4–pp

24. Laiche N, Larabi S, Ladraa F, Khadraoui A (2014) Curve normalization for shape retrieval. *Signal Process Image Commun* 29(4):556–571 572
25. Lee S (2013) Symmetry-driven shape description for image retrieval. *Image Vis Comput* 31(4):357–363 574
26. Liang S, Sun Z (2008) Sketch retrieval and relevance feedback with biased svm classification. *Pattern Recogn Lett* 29(12):1733–1741 575
27. Liang S, Li R-H, Baciu G (2011) A graph modeling and matching method for sketch-based garment panel design. In: *International conference on cognitive informatics & cognitive Computing (ICCI CC)*. IEEE, pp 340–347 577
28. Liang S, Luo J, Wenyn L, Wei Y (2015) Sketch matching on topology product graph. *IEEE Trans Pattern Anal Mach Intell* 37(8):1723–1729 580
29. Ling H, Jacobs DW (2007) Shape classification using the inner-distance. *IEEE Trans Pattern Anal Mach Intell* 29(2):286–299 582
30. Lowe DG (2004) Distinctive image features from scale-invariant keypoints. *Int J Comput Vis* 584
31. Meijster A, Roerdink JBTM, Hesselink WH (2000) A general algorithm for computing distance transforms in linear time. In: *Mathematical morphology and its applications to image and signal processing*. Springer, pp 331–340 585
32. Mokhtarian F, Abbasi S, Kittler J (1996) Robust and efficient shape indexing through curvature scale space. In: *British machine and vision conference (BMVC)*, vol 96 588
33. Otsu N (1975) A threshold selection method from gray-level histograms. *Automatica* 11(285-296):23–27 590
34. Pedrosa GV, Barcelos CAZ (2010) Anisotropic diffusion for effective shape corner point detection. *Pattern Recogn Lett* 31(12):1658–1664 591
35. Pedrosa GV, Batista MA, Barcelos CAZ (2013) Image feature descriptor based on shape salience points. *Neurocomputing* 593
36. Qi W, Zou C, Yuan Y, Hongbing L, Yan P (2013) Image registration by normalized mapping. *Neurocomputing* 101:181–189 595
37. Qian Y, Yang Y, Song Y-Z, Xiang T, Hospedales T (2015) Sketch-a-net that beats humans [arXiv:1501.07873](https://arxiv.org/abs/1501.07873) 597
38. Ren Z, Yuan J, Meng J, Zhang Z (2013) Robust part-based hand gesture recognition using kinect sensor. *IEEE Trans Multimed* 15(5):1110–1120 599
39. Richards W, Hoffman DD (1985) Codon constraints on closed 2D shapes. *Comput Vis Graph Image Process* 31(3):265–281 600
40. Roman-Rangel E, Marchand-Maillet S (2014) Hoosc128: a more robust local shape descriptor. In: *Pattern recognition*, volume 8495 of *Lecture notes in computer science*. Springer, pp 172–181 603
41. Roman-Rangel E, Marchand-Maillet S (2015) Shape-based detection of maya hieroglyphs using weighted bag representations. *Pattern Recogn* 48(4):1161–1173 604
42. Rosenfeld A, Pfaltz JL (1966) Sequential operations in digital picture processing. *J ACM (JACM)* 13(4):471–494 606
43. Roy Davies E (2004) *Machine vision: theory, algorithms, practicalities*. Elsevier 607
44. Saha PK, Borgefors G, Sanniti di Baja G (2015) A survey on skeletonization algorithms and their applications. *Pattern Recognition Letters* 609
45. Sebastian T, Klein P, Kimia B (2001) Recognition of shapes by editing shock graphs. In: *International conference on computer vision (ICCV)*, vol 1. IEEE, pp 755–755 610
46. Sebe N, Qi T, Loupias E, Lew MS, Huang TS (2003) Evaluation of salient point techniques. *Image Vis Comput* 21(13):1087–1095 611
47. Sezgin TM, Davis R (2007) Scale-space based feature point detection for digital ink. In: *ACM SIGGRAPH 2007 courses*. ACM, p 36 612
48. Shafait F, Keysers D, Breuel TM (2008) Efficient implementation of local adaptive thresholding techniques using integral images. In: *Document recognition and retrieval XV* 616
49. Shu X, Xiao-Jun W (2011) A novel contour descriptor for 2D shape matching and its application to image retrieval. *Image Vis Comput* 29(4):286–294 618
50. Song J, Song Y-Z, Xiang T, Hospedales T, Ruan X (2016) Deep multi-task attribute-driven ranking for fine-grained sketch-based image retrieval. In: *British Machine vision conference (BMVC)*, vol 3 619
51. Sundar H, Silver D, Gagvani N, Dickinson S (2003) Skeleton based shape matching and retrieval. In: *Shape modeling international*. IEEE, pp 130–139 622
52. Tyler CW (2002) *Human symmetry perception and its computational analysis*. Psychology Press 623
53. Witkin AP (1984) Scale-space filtering: a new approach to multi-scale description. In: *International conference on acoustics, speech, and signal processing (ICASSP)*, vol 9. IEEE, pp 150–153 624
54. Yang M, Kpalma K, Ronsin J (2008) A survey of shape feature extraction techniques. *Pattern Recogn*, 43–90 625

- 631 55. Yang X, Liu H, Latecki LJ (2012) Contour-based object detection as dominant set computation. *Pattern*
632 *Recogn* 45(5):1927–1936
- 633 56. Zaniibbi R, Yu L (2011) Math spotting: retrieving math in technical documents using handwritten query
634 images. In: *International Conference on document analysis and recognition (ICDAR)*
- 635 57. Zhang D, Guojun L (2004) Review of shape representation and description techniques. *Pattern Recogn*
636 37(1):1–19
- 637 58. Zhang D, Lu G (2002) A comparative study of fourier descriptors for shape representation and retrieval.
638 In: *Asian Conference on computer vision (ACCV)*, pp 646–651
- 639 59. Zhao Peng, Guoqin Wu, Yijuan Lu, Xianwen Wu, Yao Sheng (2016) A novel hand-drawn sketch
640 descriptor based on the fusion of multiple features. *Neurocomputing* 213:66–74
- 641 60. Zhu G, Doermann D (2007) Automatic document logo detection. In: *International conference on*
642 *document analysis and recognition (ICDAR)*, pp 864–868
- 643 61. Zhu F, Xie J, Fang Y (2016) Learning cross-domain neural networks for sketch-based 3d shape retrieval.
644 In: *AAAI Conference on artificial intelligence*. AAAI Press, pp 3683–3689



Houssem Chatbri is a Postdoctoral Researcher at the Insight Centre for Data Analytics, Dublin City University in Ireland. He graduated from University of Tsukuba in Japan in 2016. His research interests include multimedia retrieval, image and video analysis and computer vision.



Keisuke Kameyama is a Professor of computer science at the Faculty of Engineering, Information and Systems, University of Tsukuba in Japan. His research interest include pattern recognition, machine learning, signal processing, computational intelligence and multimedia retrieval. He graduated from Tokyo Institute of Technology in 1991.



Paul Kwan is a Professor of computer science at the University of New England in Australia. He received a BSc and an MSc degree in computer science from Cornell University and University of Arizona (USA) in 1986 and 1988, respectively, then a PhD degree in 2003 from University of Tsukuba, Japan. His research interests include artificial intelligence, computer vision and image processing, complex systems and computational modelling, and bioinformatics.



Suzanne Little is a Lecturer at the School of Computing, Dublin City University, Ireland. She completed her PhD at the University of Queensland, Australia. Her research interests include machine learning, computer vision, multimedia analytics, semantic search and data integration in various applications such as security, technology enhanced learning, biomedical, multimedia archives (news), autonomous vehicles, internet of things and smart communities.



Noel E. O'Connor is a Professor in the School of Electronic Engineering, the Director of Information Technology and the Digital Society Hub, and a Principal Investigator at the Insight Centre for Data Analysis, Dublin City University (DCU), Ireland. obtained his PhD from Dublin City University in 1998. His research interests include machine learning, image and video analysis and computer vision with various applications such as autonomous vehicles, e-learning smart cities an security.

UNCORRECTED PROOF

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES:

- Q1. Please check affiliations if captured and presented correctly.
- Q2. Housseem Chatbri was captured as the corresponding author. Please check if correct.
- Q3. Please check section heads if assigned to its appropriate levels.
- Q4. Please provide significance of bold emphasis of tables 2–5.