

2017-07-21

Handwritting: keyword spotting The Query by Example (QbE) case

Barlas, Georgios

Nova Science Publishers, Inc

<http://hdl.handle.net/11728/11652>

Downloaded from HEPHAESTUS Repository, Neapolis University institutional repository

1 Chapter 11

2 Handwritten keyword spotting The 3 Query by Example (QbE) case

4 *Georgios Barlas, Konstantinos Zagoris and Ioannis Pratikakis*

5 **Georgios Barlas, Konstantinos Zagoris and Ioannis Pratikakis**

6 Department of Electrical and Computer Engineering

7 Democritus University of Thrace

8 67100 Xanthi, Greece

9
10 Corresponding Author: ipratika@ee.duth.gr

12 Introduction

13 The traditional approach in document indexing usually involves an Optical Character
14 Recognition (OCR) step. Although OCR performs well in modern printed documents and
15 documents of high quality printing, in the case of handwritten documents OCR, several
16 factors affect the final performance like intense degradation, paper-positioning variations
17 (skew, translations, etc.) and writing styles variety.

18 Handwritten word spotting has attracted the attention of the research community in the
19 field of document image analysis and recognition since it appears to be a feasible solution
20 for indexing and retrieval of handwritten documents in the case that OCR-based methods
21 fail to deliver satisfactory results.

22 Handwritten keyword spotting (KWS) is the task of retrieving all instances of a given
23 query word in handwritten document image collections without involving a traditional OCR
24 step. There exist two basic variations for KWS approaches: (a) the Query by Example case
25 (QbE) where the query is a word image and (b) the Query by String case (QbS) where, as
26 the name implies, the query is a string. The study presented in this chapter will focus on
27 the QbE approach.

28 For a better understanding, QbE methods will be presented taking into account two dif-
29 ferent perspectives which relate to the use of segmentation and learning. The segmentation-

1 based methods are divided into 2 subcategories based upon the segmented entity which
2 could be either the word image or the textline. They are strongly dependent on the seg-
3 mentation step, so that to compare different methods regardless of segmentation errors,
4 many researchers do not implement a segmentation method but they use datasets where the
5 segments are given.

6 In the case of segmentation-free methods the whole image is tested against similarities
7 between the query image and the patches of the document image without segmenting it at
8 any level. The methods of this class, on the one hand bypass the step of segmentation but on
9 the other hand they cannot avoid searching for the words in parts of the image that may not
10 contain text. Therefore, segmentation-free methods avoid failures due to bad segmentation
11 but the running time increases considerably. It is worth-mentioning that the methods of this
12 class are not the trend.

13 Training-based methods are those that require training data at a particular stage of the
14 process. A common problem in these methods is the availability of training data. Further-
15 more, an extra weakness is that to apply such a method to a new word, usually ground
16 truthing work is required to obtain training data, which is quite time consuming and often
17 it has to be done totally manual.

18 Training - free are methods that as the name implies do not include any training stage
19 in the operational KWS pipeline. The training - free methods can be applied directly to
20 new word although, they usually require a particular configuration to be effective in the
21 corresponding text.

22 This chapter is structured as follows: Section "Segmentation-based Context" will
23 present the KWS methodologies that operate in a segmentation-based context wherein
24 methods based on training and methods that are independent of any training involvement
25 will be detailed. Both variations will be separately reviewed depending on the type of seg-
26 mentation which is used. In Section "Segmentation - Free Context", methodologies that do
27 account for a segmentation will be discussed with a particular focus on the use or not of
28 training. Section "Experimental Datasets and Evaluation Metrics" deals with an overview
29 of the current efforts for performance evaluation and a brief description of datasets that
30 were used in QbE KWS, while the Section "Conclusive Remarks" is dedicated to a fruitful
31 discussion which aims to identify the current trends of the QbE KWS.

32 **Segmentation-based Context**

33 In this section, segmentation-based methods are presented. Segmentation-based methods
34 have been categorized into training-based and training-free approach. Each category is then
35 subdivided into word image segmentation and textline segmentation context.

36 **Methods Based on Training**

37 **Word Image Segmentation Context**

38 In the work of Rodríguez-Serrano and Perronnin (2009), the method is based on a Semi-
39 Continuous - Hidden Markov Model (SC-HMM) coupled with a Gaussian Mixture Model
40 (GMM). SC-HMM is able to learn from a small set of samples. A segmentation algorithm

1 extracts sub-images that potentially represent words, employing state-of-the-art techniques
 2 based on projection profiles and clustering of gap distances. Then, a simple classifier using
 3 holistic features (per column, pixel count, Local Gradient Histogram (LGH)) is employed
 4 for performing a first rejection pass. The non-rejected word images are normalized with
 5 respect to slant, skew and text height, using standard techniques. Then, for each normalized
 6 word image, LGH features are computed by moving a window from left to right over the
 7 image and feature vectors are extracted at each position to build the feature vector sequence.
 8 Finally, using SC-HMM with GMM, a score is assigned to each feature vector sequence
 9 which is used to attribute the similarity with the query using a threshold. An overview of
 10 the methodology is shown in Figure 11.1.

11 The same framework was used in Rodríguez-Serrano et al. (2010) modified so that writers
 12 adaptation is achieved. For this purpose, a statistical adaptation technique was applied
 13 to change some of GMMs parameters at each document. Furthermore, SC-HMM was used
 14 in Rodríguez-Serrano and Perronnin (2012) to enrich features extraction, since in a left-to-
 15 right HMM, the states are ordered and the weights of the SC-HMMs can be viewed as a
 16 sequence of vectors. The distance between these vectors is computed using the Dynamic
 17 Time Warping (DTW) wherein Bhattacharyya measure is being used as local similarity. The
 18 use of SC-HMM in an unsupervised context was presented in Rodríguez-Serrano and Per-
 19 ronnin (2012) where character examples of existing fonts were used to create the training
 20 set.

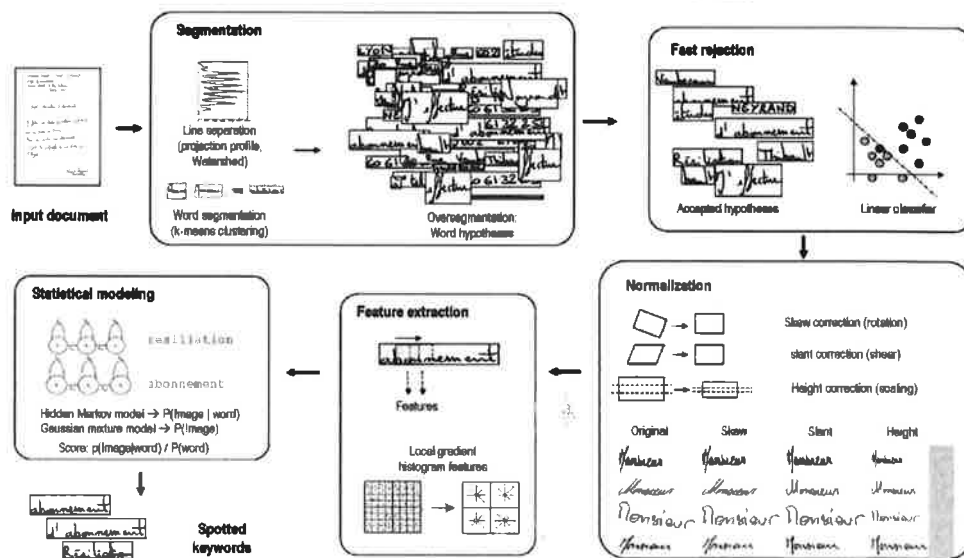


Figure 11.1. Overview of the word-spotting system presented in (Rodríguez-Serrano and Perronnin, 2009).

21 In the work in Almazán et al. (2012a), a preprocessing stage is initially applied using
 22 margins removal and anisotropic Gaussian filtering. Then, binarization and word segmen-
 23 tation is applied. In the core methodology, they use a hashing strategy based on Loci features

1 to prune word images and limit the candidate locations. A discriminative model is then ap-
2 plied to those locations. The discriminative learning relies upon a Support Vector Machine
3 (SVM) which sets the weights on the appearance features Histogram of Oriented Gradients
4 (HOG) to compute the final similarity.

5 Almazan et al. (2013) created a method that is both QbE and QbS and addresses
6 multi-writer WS. The QbE pipeline is based on Fisher Vectors (FV) computed over Scale-
7 Invariant Feature Transform SIFT descriptors. Then the FV are used to feed an SVM to
8 get the attribute scores. They report that any encoding method and classification algorithm
9 that transforms the input image into attribute scores could be used to replace them, but they
10 chose SVMs and FVs for simplicity and effectiveness.

11 The study in Fernández et al. (2013) uses a previous method (Fernández et al., 2011)
12 which is extended in a way that the syntactic context in the document is used to infer
13 context. To achieve this, Markov Logic Networks (MLN) are used that are trained with
14 specific rules. The MLN can be considered as a collection of first order rules to each of
15 which it is assigned a real number, the weight. Each rule represents the rule in the domain,
16 while the weights indicate the strength of the rule.

17 The work in Aldavert et al. (2015) uses the Bag-of-Visual-Words (BoVW) model to
18 WS. The authors divide the procedure of creating BoVW in four basic steps: sampling,
19 description, encoding and pooling. In particular, they sampled densely the word images
20 using a fixed step at different scales. The description is derived from the HOG descriptor.
21 To encode the descriptors and create the codebook, the Locality-constrained Linear Coding
22 (LLC) algorithm was used. Finally, at pooling step, the Spatial Pyramid Matching (SPM)
23 technique applied so that spatial information could be used.

24 In the work at Sharma et al. (2015), they made experiments with Convolutional Neural
25 Networks (CNN). Only the classification layers were retrained to address the problem. The
26 CNN extracted a 4096-d feature for each word image which was achieved after discarding
27 the last fully connected layer and considering the activation of the units in the second fully
28 connected layer. For matching, standard Lp norms have been used.

29 **Textline Segmentation Context**

30 The system presented in Keaton et al. (1997) is composed by several modules. The “focus-
31 of attention” module concerns a cross-correlation testing between the query image and the
32 document for finding the candidate locations. The “preprocessing” module that consists
33 of estimation of word image zones at the upper, middle and lower level, filtering of stray
34 marks and skeletonization. The “feature extraction” module that concerns profile encoding
35 (20 Discrete Cosine Transform (DCT) coefficients from the profile extracted at each of the
36 three zones and cavity encoding which takes into account 2D spatial arrangement, as well
37 as the descender, and ascender information leading to a graph. Both encoding features are
38 combined to a new graph that contains the type, size and relative location of each feature,
39 which is considered as keyword signature graph. Finally, the keyword signature matching
40 is addressed by two distinct comparisons. First, a comparison is employed between the
41 profile encoding DCT coefficients wherein the resulting comparison is incorporated into
42 the graph as additional feature. In the sequel, the keyword signature graphs are compared
43 with probabilistic graph matching based on Bayesian evidential reasoning.

1 Training - Free Methods

2 Word Image Segmentation Context

3 In the paper Manmatha et al. (1996) the term “word spotting” was introduced for hand-
4 written documents as analogous to “word spotting” in speech processing. It was applied
5 on scanned greylevel document images. The steps of the algorithm are gaussian filtering,
6 subsampling to reduce image by half, binarization by thresholding and segmentation into
7 words. Then follows pruning taking into account (i) the aspect ratio and (ii) the size using
8 predefined thresholds. Finally, at matching stage, two different matching algorithms were
9 used based on the standard and the affine-corrected (SLH algorithm) Euclidean distance,
10 respectively.

11 In the sequel, Rath and Manmatha (2003b), was motivated by Kolcz et al. (2000) that
12 have used Dynamic Time Warping for matching in combination with a textline segmen-
13 tation method. In this approach, the textline segmentation has been replaced by a word
14 image segmentation. In the work of Rath and Manmatha (2003a), the feature set used for
15 the experiments was extended to 11 distinct features (4 projection profiles, 2 word pro-
16 files, background ink transitions, graylevel variance and 3 Gaussian smoothing variation
17 features) which were used as single or combined ones for matching with DTW. Never-
18 theless, the more recent work in Rath and Manmatha (2007) suggested only three distinct
19 features, namely, projection profiles, word profiles and background ink transitions which
20 were optimally matched when using DTW.

21 A study based on contours of words was presented in Adamek et al. (2007). They start
22 with local binarization of the image. To achieve smoothness at word outlines, they pre-
23 process the image applying morphological filtering. After binarization a word may split to
24 more than one connected component. To estimate the exact position of the word in the word
25 image a process based on horizontal and vertical projection histogram and a fixed threshold-
26 ing is applied. Then, after applying a series of heuristics rules, the connected components
27 of the word image are linked together to create a single component from which the contour
28 is extracted. They used Multiscale Convexity Concavity (MCC) representation for the con-
29 tour. MCC calculates the convexity and concavity along contour at different scales to create
30 2D matrix where rows correspond to scale level and columns to convexity or concavity. To
31 measure the matching between the contours the DTW algorithm were used. The distance
32 matrix of DTW is constructed by storing the distance between a pair of contour points
33 corresponding to the row and column of the MCC representation. The final dissimilarity
34 between contours is the normalized optimal path of the matrix. An alternative method that
35 was tested is the MCC-DCT where the 1D DCT is applied at MCC representation matrix
36 and the coefficients of DCT are combined to the final dissimilarity.

37 The work in Bhardwaj et al. (2008) presented an algorithm based on moment functions.
38 In the initial stage, they used horizontal and vertical profiles to segment the document into
39 lines and words, respectively. High order (up to 7) moment functions were used to extract
40 features and indexing each word image. They used cosine similarity metric for matching
41 and relevance feedback to improve the results.

42 A shape descriptor, the Compact Shape Portrayal Descriptor (CSPD) was presented in
43 Zagoris et al. (2011) which requires only 123 bits per word image. CSPD is based on five (5)
44 distinct features: (i) width height ratio, and the DCT coefficients of (ii) vertical projection,

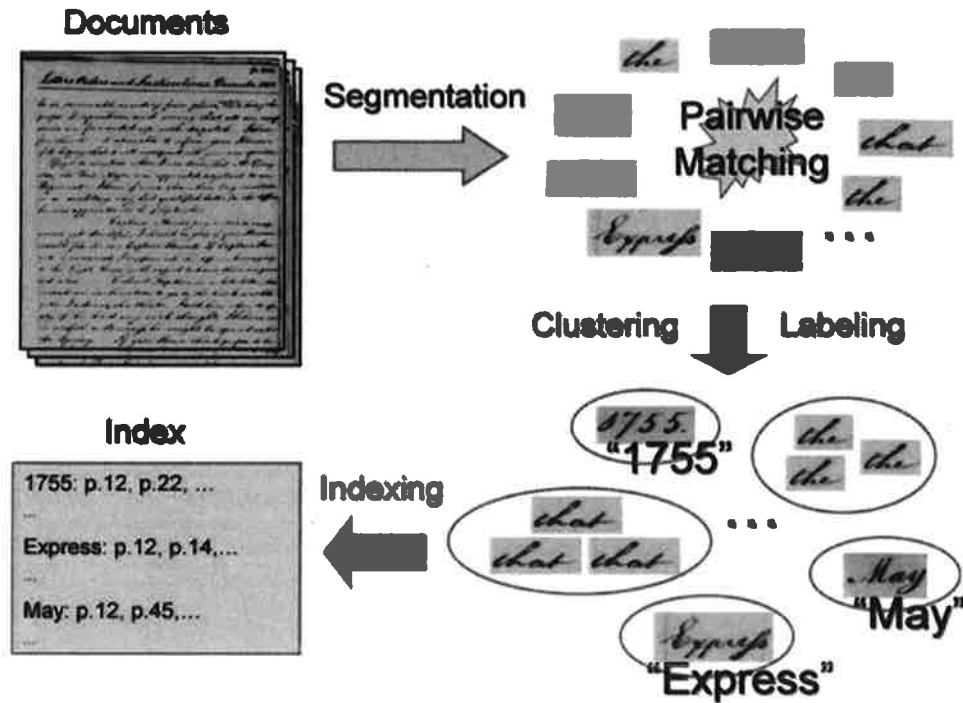


Figure 11.2. An illustration of the word spotting pipeline presented in (Rath and Manmatha, 2007).

1 (iii) horizontal projection, (iv) top shape projection and (v) bottom shape projection. For
 2 each feature Gustafson-Kessel fuzzy algorithm is applied for quantization to reduce the
 3 stored size of each descriptor. To refine the results, relevance feedback is applied. The
 4 user selects the best result and a training set is created. This training set is used to train an
 5 SVM to correct the results. The similarity measure that was used is a modified (weighted)
 6 Minkowski L_1 .

7 The study appeared in Can and Duygulu (2011) was motivated from a shape representa-
 8 tion method. They used already segmented word images provided with the datasets. The
 9 proposed methodology comprises 3 main steps. First, binarization is applied by using a
 10 threshold which is computed as the mean intensity value of the gray-scale image. Next,
 11 a contour extraction is used for each connected component in the binarized word image.
 12 Finally, a sequence of lines is created which is used as the descriptor for matching. The
 13 matching score is computed by first finding the distances between the line descriptors and
 14 then, summing all distances over the complete word image.

15 The basic premise in the work of Fornés et al. (2011) is that each word image is treated
 16 as a shape which can be represented using two models, namely the Blurred Shape Model
 17 (BSM) and the Deformable BSM. First, at preprocessing stage, segmented text lines are
 18 normalized by applying skew angle removal and slant correction. In the case of BSM,

1 the descriptor represents a probability distribution of the object shape. In the case of De-
2 formable BSM, every image is represented with two output descriptors, a vector which
3 contains the BSM value of each focus (equidistant distributed points) and the position of
4 each focus. The proposed matching technique lies upon the movement of focuses so that
5 its own BSM is maximized. It is shown that using Euclidean distance in both BSM and
6 deformable BSM methods outperforms the use of DTW.

7 Fernández et al. (2011) used Loci features (Glucksman, 1969) along eight directions.
8 They are computed on the skeleton of each word image which is achieved after a document
9 image binarization and word image segmentation step. The similarity is computed using
10 the Euclidean and Cosine distance.

11 The study in Diesendruck et al. (2012b) and Diesendruck et al. (2012a) was focusing on
12 building a search system for 1930-40 US Census data. The process starts with binarization,
13 morphological thinning and Hough transform to locate table lines since the documents are
14 in table format. Thus, the segmentation is based on table lines. Since, each cell contains one
15 word, the method is word-based segmented. Then, a signature vector was composed the first
16 10 coefficients of cosine transform of upper, lower and transition profiles. Since the data
17 set is quite large and the response time should be reasonable, hierarchical agglomerative
18 clustering with complete linkage is used to cluster the signature vectors.

19 The problem of sequential KWS was addressed in Fernández-Mota et al. (2014). In
20 sequential KWS the ordered sequence of indices is taken into account for finding similar
21 instances of words in a book. They experimented with descriptors that relate to a single
22 writer scenario (BSM, HOG, nrHOG) as well as descriptors that relate to a multiple writers
23 scenario (attribute-based approach).

24 The work in Howe (2013) is based on a part-structured modeling which aims to mini-
25 mize a deformation energy required to fit the query to the word image instance. The process
26 is initiated with a binarization. Then, skeletonization is applied to produce connected com-
27 ponents of a single-pixel width. The endpoints and junctions of the skeleton are used to
28 build a tree. An energy minimization of a function that comprises a deformation energy and
29 an observation energy term is finally addressed.

30 The method in Kovalchuk et al. (2014) is the winner of H-KWS 2014 competition. First,
31 they binarize the image by global thresholding and connected components are computed.
32 Then, pruning of connected components is followed based on heuristics that rely upon
33 properties of connected components. Using a regular grid of fixed size, they compute HOG
34 and LBP descriptors which result in a 250D vector. A max-pooling process is then applied
35 to the descriptor. The matching is made with L_2 distance.

36 In Wang et al. (2014) the authors initiate the process by applying the preprocessing step
37 presented in Wang et al. (2013). They use a graph representation model which is based on
38 the skeleton of each word image. In this graph, vertices are the structural interest points
39 and the strokes connecting them are the edges. The value of vertices corresponds to the
40 Shape Context descriptor while the value of edges corresponds to the length of stroke. The
41 computation of similarity between two word images concern the similarity of graphs for
42 each connected component existing in the query and the test word image which is used to
43 guide the DTW computation.

44 The work in Zagoris et al. (2014) is based on spatial information from word images.
45 First, gradient vectors are calculated. Because of the sensitivity of gradient to noise, an Otsu

1 like threshold is applied at gradient vectors. At the remaining points gradient orientation
2 is calculated. Next, a linear quantization of gradients orientation to a desirable amount of
3 levels follows. The quantization step also controls the amount of the final local points. After
4 quantization the corner points are characterized as initial keypoints (kP). The final points
5 are the dominant kPs according to Shannon entropy at an area where the kP is the center.
6 After the final points have been selected each area around kPs is divided to 9 subareas.
7 For each of these 9 areas using a voting system based on the weighted distance of each
8 point to kP, a 3-bin histogram is created. The combination of these 9 histograms to a 27-
9 bin histogram results in the descriptor, the Document-Specific Local Features (DSLFL). At
10 matching stage first, a normalization is applied for each word image. Then instead of a brute
11 force search, a Local Proximity Nearest Neighbour (LPNN) search is used by taking into
12 account the mean distance of each pixel from the mean center. Finally, Euclidean distance
13 is applied between the kPs and the results are presented in an ascending order. In Zagoris
14 et al. (2015), an extension of this work is presented using Relevance Feedback strategies
15 (CombSum, CombMin, Probabilistic model). It is reported that the optimal results are
16 achieved from CombMin model.

17 The goal in Papandreou et al. (2016) is to study the zoning features. Binarization and
18 deslanting are first applied in the query and candidate word images as pre-processing steps.
19 The zoning features are extracted after cutting the query word in vertical zones based on its
20 length and pixel distribution along the horizontal axis and adjusting these boundaries with
21 the corresponding zones in the candidate word image using DTW. In the sequel, the word
22 images are normalized and their features, which are based on pixel density, are extracted.
23 Finally, the final distance is the product of Euclidean distance of the two word images with
24 the distance provided from DTW.

25 In Mondal et al. (2016) the author introduces a new matching technique the Flexible
26 Sequence Matching (FSM) algorithm for KWS task. At the preprocessing stage, the docu-
27 ment images are first binarized by an adaptive technique (Gatos et al., 2006), after binariza-
28 tion along with border removal is applied to obtain proper text boundaries (Stamatopoulos
29 et al., 2010) and then a segmentation stage follows that partitions the documents in lines
30 or pieces of lines, up to words or parts of words, depending on the experiment to be con-
31 ducted. At features extraction stage grayscale and binary image are used to extract two
32 types of features, namely column-based features and Slit style HOG (SSHOG) (Terasawa
33 and Tanaka, 2009). Eight Column-based features, are extracted from the binary image.
34 Finally, at matching stage FSM is applied. FSM is similar in spirit to DTW but it is less
35 sensitive to local variations in the spelling of words and to local degradations effects within
36 the word image.

37 In a recent work Retsinas et al. (2016), three variances of the Projections of Oriented
38 Gradient (POG) descriptor (Retsinas et al., 2015) studied in the framework of the KWS
39 problem. The first variant, the global POG (gPOG) is slightly different from POG, for
40 which the main differences are: (i) it keeps different number of coefficients from DCT
41 and (ii) has 6 projections. The second variant k-segmented POG (IPOG), first segments
42 the word image to k overlapping images and then calculates the POG descriptor to each
43 of them. The third variant, fusion POG (fPOG), as the name implies is a fusion of gPOG
44 and IPOG descriptors. Finally, the Euclidean distance is used to attain the matching score.

1 It should be noted that at the preprocessing step, binarization, skew correction and height
2 normalization were applied.

3 **Textline Segmentation Context**

4 In the paper Kolcz et al. (2000), the approach is motivated by the success of dynamic-
5 programming based techniques for KWS in speech applications even when very limited
6 keyword models are present. It relies upon a line segmentation method to achieve distinct
7 textlines for each document image. The ink-density histogram and its Fast Fourier Trans-
8 form (FFT) spectrum is used to determine the textlines as well as the skew of the page. For
9 each textline, they extract the upper and the lower profile as well as transitional Features
10 (number of transitions between background and body in each column of pixels). Then
11 they used DTW to address matching in the KWS pipeline. They used heuristics to reduce
12 computational time which were similar to the ones used by the work in (Manmatha et al.,
13 1996).

14 Terasawa and Tanaka (2009) deals with language independent KWS method. They
15 have chosen a line-oriented approach because (i) word segmentation is impossible in some
16 languages and (ii) it can retrieve a hyphenated word that spans two lines. For each textline
17 image, a narrow sliding window is used for feature extraction. A variation of HOG features
18 is used, namely, the SSHOG features. Compared to the original HOG, the SSHOG uses a
19 narrow window and the computed gradient is signed. For feature matching, DTW is used.

20 In Wang et al. (2013) the authors use a coarse-to-fine strategy. They first remove the
21 noise with a smoothing filter and they apply textline segmentation based on Hough trans-
22 form. At coarse step, they apply a sliding window at the size of query word. They extract
23 4 textural features, namely, projection profile, upper and lower border and orientation dis-
24 tribution of skeleton pixels. Then they use DTW for the first three features and Chi-square
25 metric for orientation distribution. With an empirical threshold they chose the most similar
26 to the query regions, the regions of interest. The fine step is applied to the regions of interest.
27 In the fine step, morphological and topological properties are used. Morphological proper-
28 ties calculated using the Shape Context descriptor on selected interest points (branch points,
29 starting/ending points and high-curved points). Topological properties are obtained from a
30 skeletonized representation. The information of the properties is the input to a weighted
31 distance function for which Linear Discriminant Analysis is used to automatically get the
32 optimal weights.

33 **Segmentation-free Context**

34 In this section, segmentation-free methods are presented. The methods of this section are
35 divided into training-based and training-free.

36 **Methods Based on Training**

37 The approach in Choisy (2007) is made to deal with KWS of isolated words on mail en-
38 velopes. It is character segmentation-free which relies on the dynamic creation of global

1 word models. This is achieved with the use of Non-Symmetric Half Plane - HMM (NSHP-
2 HMM) (Saon and Belaïd, 1997) which is a model hybrid of an HMM and a Markov Field
3 (MRF). Before applying NSHP-HMM, two preprocessing steps are applied (i) global slant
4 correction and (ii) a non-linear normalization that centers and normalizes the lower case
5 zone of the writing. The NSHP-HMM is trained at word level.

6 In the work of Rusinol et al. (2011), they used a BoVW model where each patch is
7 normalized by applying term frequency-inverse document frequency (tf-idf) model. BoVW
8 words model powered by SIFT descriptors which was further augmented by a word seg-
9 mentation. Then, Latent Semantic Indexing (LSI) (Deerwester et al., 1990) is used with
10 BoVW to retrieve relevant patches even if they do not contain the same exact features than
11 the query sub-image. Then Singular Value Decomposition (SVD) is applied to reduce the
12 descriptor dimension and to obtain a transformed space where patches having similar top-
13 ics but with different descriptors will lie close. At the retrieval stage cosine similarity were
14 used. In Rusinol et al. (2015), they have enhanced their preliminary version by including an
15 indexation scheme aimed to scale the proposed method to handle large datasets. The SVD
16 step replaced by Product Quantization (PQ) for this purpose. Also, a multi-length patch
17 representation is also introduced, which increases the retrieval performance by taking into
18 account the different possible lengths of the query words.

19 The study presented in Rusinol and Lladós (2012) comprises both fusion and relevance
20 feedback mechanisms. At fusion stage, three different fusion methods were tested, namely,
21 early fusion, combMAX and Borda count. For relevance feedback, three different methods
22 were tested, the Aocchio's algorithm (Cao et al., 2011), the Ide dec-hi method (Ide, 1971)
23 and the relevance score algorithm (Giacinto and Roli, 2004).

24 The work in Almazán et al. (2012b) is based on exemplar SVM for better representation
25 of features. Documents are represented with a grid of HOG descriptors, and a sliding
26 window approach is used to locate the document regions that are most similar to the query.
27 They use the Exemplar SVM framework to produce a better representation of the query
28 in an unsupervised way. Finally, the document descriptors are compressed with Product
29 Quantization (PQ) which has also the benefit of calculating the distance between the query
30 and the quantized document with the use of a look-up table.

31 The method presented in Dovgalecs et al. (2013) may operate with words or graph-
32 ics and is situated in the BoVW framework. In the offline stage, the BoVW is created
33 from densely sampled SIFT features. In the online stage, candidate zones detection works
34 by comparing query features with BoVW using chi-square distance. Then, the Longest
35 Weighted Profile (LWP) algorithm which enforce spatial ordering information characteris-
36 tics of words and graphical patterns alike, is used to compute the similarity score between
37 the query image and the candidate zones.

38 The study in Rothacker et al. (2013) is based on the use of Bag-of-Features with HMMs.
39 The method is divided into three parts. First, densely SIFT features are extracted in the
40 whole document image and 5% of those are used to create a codebook size of 4096 for
41 Bag-of-Features. Then, the Bag-of-Features representation feeds an HMM which encodes
42 the sequential visual appearance of features that are located in the query bounding box.
43 Finally, the document collection can be queried in a patch-based fashion where the output
44 is a map of probabilistic scores from which the query results can be retrieved. As shown
45 in the upper part of Figure 11.3 the document image representation is visualized. In the

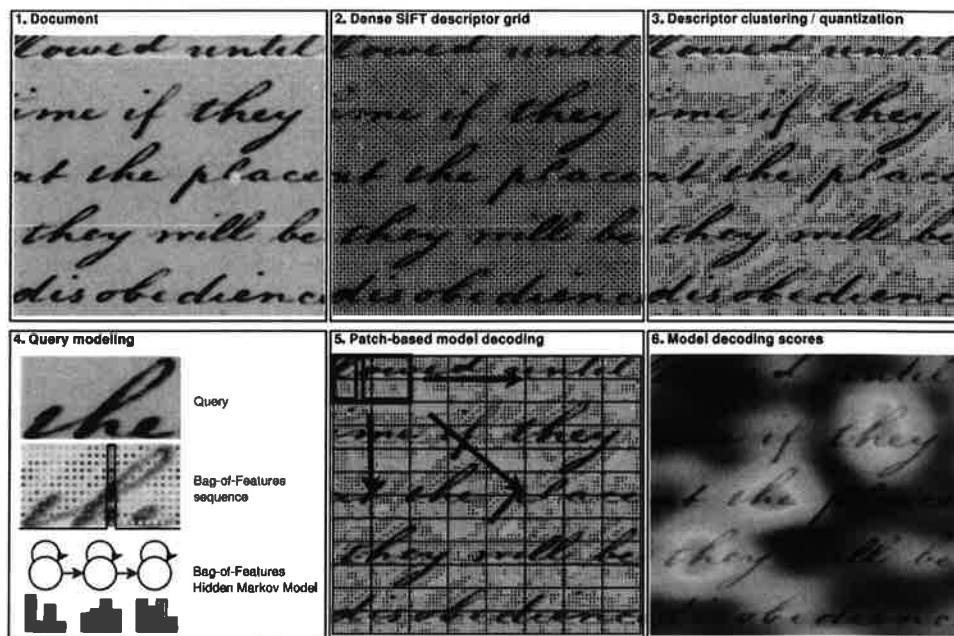


Figure 11.3. Overview of the segmentation-free word spotting method presented in (Rothacker et al., 2013).

- 1 lower part, the estimation of a query model and the patch-based decoding with respect to
- 2 that model are shown. Patch-based representations are evaluated at each grid point. The
- 3 scores obtained are visualized by interpolating them over the document image indicating
- 4 low to high responses with blue to red colors.

5 Training - Free Methods

6 The work in Leydier et al. (2009) focuses on medieval Latin manuscripts and is based on the
 7 observation that medieval Latin manuscripts have letters mainly composed of large vertical
 8 strokes. The algorithm has two main steps. In the first step, guides, gradients and Zones Of
 9 Interest (ZOI)s are extracted from the document image and the query word image. In the
 10 second step, cohesive matching is applied between the guides of the query image and those
 11 of the document image. For each match of guides, a check if ZOIs are matching, is applied,
 12 too. This work was enriched with a model that is the combination of an alphabet, a glyph
 13 book and a grammar as presented in Leydier et al. (2009). The model is used to create a
 14 character tree. The extra information from the character tree was used among the query
 15 word image and the document image for better extraction of ZOIs, guides and gradients.
 16 They also automated thresholds that were needed at ZOIs, guides and gradients extraction,
 17 but also at cohesive matching.

18 Zhang and Tan (2013) is motivated by the Heat Kernel Signature (HKS) which has been
 19 used for shape recognition. Actually, the Deformation and Light Invariant (DaLI) descriptor

1 was used which is applied by convolving Scale Invariant HKS (SI-HKS) with Gaussian
 2 kernels. To compute the similarity, a Delaunay Triangulation algorithm was applied to
 3 create a Triangular Mesh Structure (TMS) of keypoints detected in the word image and the
 4 document image, respectively. Finally, the similarity score is computed by building a score
 5 matrix which contains the optimal matching score and the optimal matching history.

6 The approach in Hast and Fornés (2016) is based on Putative Match Analysis (PUMA)
 7 (Hast and Marchetti, 2012), a technique that first introduced for matching aerial images.
 8 First, the input images are binarized using Otsu, and then smoothed with a Gaussian in order
 9 to find more key points. Then, four different kind of key points are detected in the word
 10 images, which basically detect lines, corners and blobs. Taking into account the detected
 11 keypoints, Fourier-based feature descriptors are computed. At the end, the matching is
 12 performed by an improved version of Random Sample Consensus (RANSAC), called as
 13 PUMA, which is able to allow a more relaxed matching among the word images.

14 Rabaev et al. (2016) focus on KWS by locating the query word in a document recur-
 15 sively in a scale-space pyramid. The proposed scheme does not depend on a specific choice
 16 of features. They experimented with the HOG descriptors, which have been shown to pro-
 17 vide good results. Chi-square distance is applied to compare HOG descriptors at all levels
 18 of the pyramid.

19 Experimental Datasets and Evaluation Metrics

20 Evaluation Metrics

21 The evaluation metrics that have been used for performance evaluation between different
 22 word spotting algorithms are inspired from information retrieval. Therefore, each retrieved
 23 item (word) is defined as relevant to original query (word - query) or not. Early reports
 24 on the KWS performance evaluation were simply taking the first n words and calculate the
 25 most basic retrieval metrics, the **Precision** and **Recall**.

$$26 \text{ Precision} = \frac{\{\text{relevant words}\} \cap \{\text{retrieved words}\}}{\{\text{retrieved words}\}} \quad (11.1)$$

$$27 \text{ Recall} = \frac{\{\text{relevant words}\} \cap \{\text{retrieved words}\}}{\{\text{relevant words}\}} \quad (11.2)$$

28 Precision is the fraction of retrieved words that are relevant to the search, while Recall
 29 is the fraction of the words that are relevant to the query. It is apparent that the above
 30 metrics are inversely related. To achieve a combined evaluation, the precision-recall curve
 31 is computed.

32 The **Precision - Recall Curve** is computed by the traditional 11-point interpolated av-
 33 erage precision approach (Manning et al., 2008), (Van Rijsbergen, 1979). For each query,
 34 the interpolated precision is measured at the 11 recall levels of 0.0, 0.1, 0.2, ..., 1.0.

Sometimes, the differences between the evaluation algorithms are very hard to observe
 especially, between very small performance results. Moreover, these graphs may not

Table 11.1. Descriptors, learning methods and similarity measures used from each method

	Method	Descriptors	Learning	Similarity
T r a i n i n g b a s e d	(Rodríguez-Serrano and Perronnin, 2009)	LGH	SC-HMM	Euclidean
	(Rodríguez-Serrano and Perronnin, 2012)	LGH	SC-HMM	DTW
	(Almazán et al., 2012a)	Loci, HOG	SVM	Dot product
	(Fernández et al., 2013)	Loci	MLN	Euclidean
	(Aldavert et al., 2015)	HOG	BoVW	Histogram Matching
	(Sharma et al., 2015)	Deep features	CNN	Lp-norm
	(Keaton et al., 1997)	DCT on profiles	Bayesian network	Graph matching
	(Choisy, 2007)	Column-wise binary patterns	NSHP-HMM	Posteriori Probability
	(Rusinol et al., 2011)	SIFT	BoVW	Histogram Matching
	(Almazán et al., 2012b)	HOG	Exemplar SVM	Euclidean
	(Dovgalecs et al., 2013)	SIFT	BoVW	Chi-square
	(Rothacker et al., 2013)	SIFT	BoVW, FSM	Histogram Matching
	T r a i n i n g f r e e	(Manmatha et al., 1996)	Profiles	
(Adamek et al., 2007)		MCC, DCT		DTW
(Bhardwaj et al., 2008)		Moments		Cosine
(Zagoris et al., 2011)		CSPD		Minkowski L1
(Can and Duygulu, 2011)		Sequence of lines		Line matching
(Fornés et al., 2011)		Deformable SVM		Euclidean, DTW
(Fernández et al., 2011)		Loci		Euclidean, Cosine
(Diesendruck et al., 2012b,a)		DCT on profiles		Euclidean
(Howe, 2013)		Endpoints and junctions of skeleton		Energy minimization
(Kovalchuk et al., 2013)		HOG, LBP		Euclidean
(Wang et al., 2014)		SC		DTW
(Zagoris et al., 2011)		DSLFP		Euclidean
(Papandreou et al., 2016)		Zoning features		Euclidean and DTW
(Mondal et al., 2016)		Column-based, SSHOG		FSM
(Retsinas et al., 2016)		POG, gPOG, IPOG, fPOG		Euclidean
(Kolcz et al., 2000)		Profiles		DTW
(Terasawa and Tanaka, 2009)		SSHOG		DTW
(Wang et al., 2013)		Profiles, SC		DTW
(Leydram et al., 2009)		ZOI		Cohesive matching
(Zhang and Tan, 2013)		DaLI		Minimum cost path between connected keypoints in a mesh grid
(Hast and Fornés, 2016)		Corners, blobs		PUMA
(Rabaev et al., 2016)		HOG		Chi-square

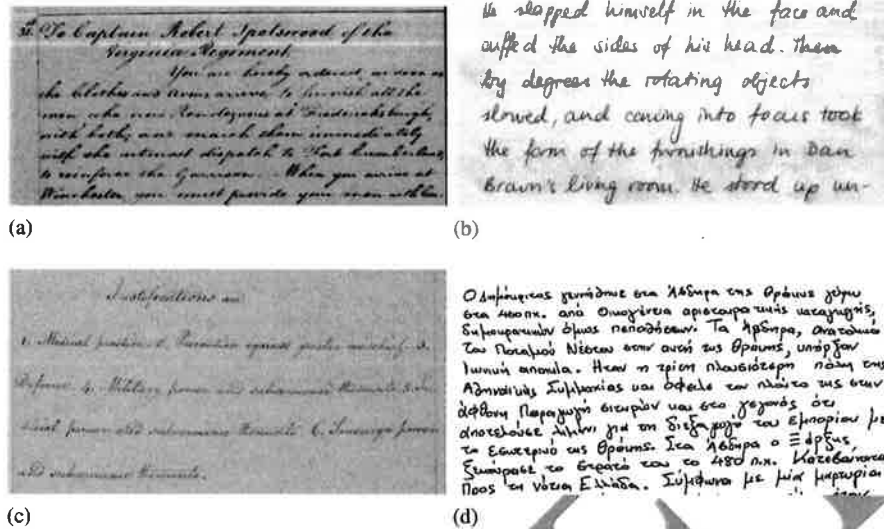


Figure 11.4. Samples from the most used datasets (a) GW (b) IAM (c) Bentham (d) Modern

1 contain all the desired information (Salton, 1992). Therefore, the need to evaluate the
 2 retrieval results with a single value is needed. The most common evaluation metric
 3 that can meet this requirement is the **Mean Average Precision (MAP)** (NIST, 2013;
 4 Chatzichristofis et al., 2011) which is defined as the average of the precision value obtained
 5 after each relevant word is retrieved:

$$7 \quad AP = \frac{\sum_{k=1}^n (P@k \times rel(k))}{\{\text{relevant words}\}} \quad (11.3)$$

8 where

9 **Precision at k items (P@k)** denoted as

$$11 \quad P@k = \frac{\{\text{relevant words}\} \cap \{k \text{ retrieved words}\}}{\{k \text{ retrieved words}\}} \quad (11.4)$$

12 with the function of relevance denoted as follows:

$$14 \quad rel(k) = \begin{cases} 1 & \text{if word at rank } k \text{ is relevant} \\ 0 & \text{if word at rank } k \text{ is not relevant} \end{cases}$$

15 Finally, the **Mean Average Precision (MAP)** is calculated by averaging the AP for all the
 16 queries, denoted as:

$$18 \quad MAP = \sum_{q=1}^Q AP_q \quad (11.5)$$

19 where Q is the total number of queries.

20 It is worth to note that in the experiments for the segmentation-free case, a resulting
 21 bounding box may not match exactly with the word bounding box from ground-truth

1 corpora. Thus, a correct match is registered when the relative overlapping area is over a
 2 certain threshold. For the sake of consistency, in every segmentation-free experiment the
 3 overlapping area is defined as:

$$4 \quad OA = \frac{A \cap B}{A \cup B} \quad (11.6)$$

6 where OA is the overlapping area, A the resulting bounding box and B is the ground-truth.

7 The challenging nature of KWS in handwritten documents has motivated the organi-
 8 zation of three dedicated international competitions in conjunction with the International
 9 Conference on Frontiers of Handwriting Recognition (ICFHR) and the International
 10 Conference on Document Analysis and Recognition (ICDAR). In particular, the ICFHR
 11 2014 Handwritten Keyword Spotting Competition (ICFHR-2014) (Pratikakis et al., 2014),
 12 the ICDAR 2015 Competition on Keyword Spotting for Handwritten Documents (ICDAR-
 13 2015) (Puigcerver et al., 2015) and the ICFHR 2016 Handwritten Keyword Spotting
 14 Competition (ICFHR-2016) (Pratikakis et al., 2016) have been the venues where research
 15 groups have been competed in two different KWS scenarios, namely, segmentation-free
 16 and segmentation-based. Table 11.2, 11.3 and 11.4 shows the results for the ICFHR-2014,
 17 ICDAR-2015 and ICFHR-2016 competitions, respectively.

Table 11.2. Experimental results for the ICFHR-2014 Competition

Method	Segmentation-based				Segmentation-free			
	BENTHAM		MODERN		BENTHAM		MODERN	
	mAP	P@5	mAP	P@5	mAP	P@5	mAP	P@5
(Kovalchuk et al., 2014)	0.524	0.738	0.338	0.388	0.416	0.609	0.263	0.539
(Almazan et al., 2013)	0.513	0.724	0.323	0.706	-	-	-	-
(Howe, 2013)	0.462	0.718	0.282	0.569	0.363	.556	0.163	0.417
(Leydier et al., 2009)	-	-	-	-	0.205	0.335	0.087	0.234
(Pantke et al., 2013)	-	-	-	-	0.337	0.543	0.091	0.245

Table 11.3. Experimental results for the ICDAR-2015 Competition

Method	Segmentation-based		Segmentation-free	
	mAP	P@5	mAP	P@5
Pattern Recognition Group TU Dortmund University	0.424	0.406	0.276	0.343
(Almazan et al., 2013)	0.300	0.342	0.082	0.109

Table 11.4. Experimental results for the ICFHR-2016 Competition

Method	Segmentation-based		Segmentation-free	
	Botany	Konzil.	Botany	Konzil.
	mAP	mAP	mAP	mAP
Computer Vision Center (CVCDAG) Universitat Autònoma de Barcelona, Spain	75.77	77.91	0	0
Pattern Recognition (PRG) TU Dortmund University, Germany	89.69	96.05	15.89	52.20
(Kovalchuk et al., 2014)	50.64	71.11	37.48	61.78
Visual Information and Interaction (QTOB) Uppsala University, Sweden	54.95	82.15	-	-

Table 11.5. Segmentation type and Datasets used from each method.

	Method	Segmentation type		Dataset	
		Line	Word		
T r a i n i n g b a s e d	(Rodríguez-Serrano and Perronnin, 2009), (Rodríguez-Serrano and Perronnin, 2012), (Rodríguez-Serrano et al., 2010)		X	French	
	(Rodríguez-Serrano and Perronnin, 2012)		X	GW, French, IFN/ENIT	
	(Almazán et al., 2012a), (Almazan et al., 2013)		X	CB	
	(Fernández et al., 2013)		X	CB	
	(Aldavert et al., 2015)		X	GW, Bentham, Modern	
	(Sharma et al., 2015)		X	IAM	
	(Keaton et al., 1997)	X		AIS	
	(Choisy, 2007)	-	-	HMS	
	(Rusinol et al., 2011, 2015)	-	-	GW	
	(Almazán et al., 2012b)	-	-	GW	
	(Dovgalecs et al., 2013)	-	-	GW	
	(Rothacker et al., 2013)	-	-	GW	
	T r a i n i n g f r e e	(Manmatha et al., 1996)		X	MUMUND, Hudson
		(Rath and Manmatha, 2003a,b, 2007)		X	GW
(Adamek et al., 2007)				GW	
(Bhardwaj et al., 2008)			X	GW, IAM	
(Zagoris et al., 2011)			X	GW, Greek	
(Can and Duygulu, 2011)			X	GW, OTM	
(Fornés et al., 2011)			X	GW	
(Fernández et al., 2011), (Fernández-Mota et al., 2014)				CB	
(Diesendruck et al., 2012b,a)			X	USC	
(Howe, 2013)			X	GW	
(Kovalchuk et al., 2014)			X	GW	
(Wang et al., 2014)			X	GW, CB	
(Zagoris et al., 2014)			X	GW, Bentham, Modern	
(Papandreou et al., 2016)			X	GW, Bentham	
(Mondal et al., 2016)			X	GW, Japanese	
(Kretsinas et al., 2016)			X	Bentham, Modern	
(Kolcz et al., 2006)		X		AIS	
(Terasawa and Tanaka, 2009)		X		GW, Japanese	
(Wang et al., 2013)		X		CITERE	
(Lecun et al., 2007, 2009)		-	-	GW	
(Zhang and Tan, 2013)		-	-	GW	
(Hast and Fornés, 2016)		-	-	CB	
(Rabaev et al., 2016)		-	-	GW, CG, Arabic	

1 “CB” stands for a collection of 50 pages from handwritten marriage licenses from the
2 Barcelona Cathedral written in 1617’. URL: <http://dag.cvc.uab.es/the-esposalles-database/>

3 The “Bentham” dataset is part of the H-KWS 2014 contest’s dataset. It consists of high
4 quality (approximately 3000 pixels width and 4000 pixels height) handwritten manuscripts.
5 The documents are written by Jeremy Bentham (1748-1832) himself as well as by Ben-
6 tham’s secretarial staff over a period of sixty years.

1 The “Modern” dataset is also part of the H-KWS 2014 contest’s dataset. It consists of
2 modern handwritten documents from the ICDAR 2009 Handwritten Segmentation Contest.
3 These documents originate from several writers that were asked to copy a given text. They
4 do not include any non-text elements (lines, drawings, etc.) and are written in four (4)
5 languages (English, French, German and Greek).

6 A dataset that comprises 1539 pages of modern off-line handwritten English text written
7 by 657 different writers is denoted as “IAM”. URL: [http://www.fki.inf.unibe.ch/databases](http://www.fki.inf.unibe.ch/databases/iam-handwriting-database)
8 [/iam-handwriting-database](http://www.fki.inf.unibe.ch/databases/iam-handwriting-database)

9 The Archives of the Indies in Seville (AIS), is a repository that represents the official
10 communication between the Spanish Crown and its New World colonies and spans approx-
11 imately four centuries (i.e. 15th-19th).

12 At Choisy (2007) a dataset that consists of French handwritten mails collection was
13 used for Handwritten Mail Sorting (HMS) task, wherein 1522 mail pages are manually
14 labelled.

15 At Manmatha et al. (1996) as dataset two single pages were used. One obtained from
16 the DIMUND document server, thus denoted as “DIMUND” and the other single page was
17 taken from a collection in the library of the University of Massachusetts. This page is a letter
18 written by James S. Gibbons to Erasmus Darwin Hudson and it is denoted as “Hudson”.

19 An Ottoman dataset denoted as “OTM” comprises documents written with a commonly
20 encountered calligraphy style called Riqqa, which was used in official documents. Consists
21 of 257 words in three pages of text. URL: <http://courses.washington.edu/otap/>

22 US Census forms from 1930 and 1940 comprise a dataset denoted as “USC”.

23 Scanned images of the Japanese manuscript “The diary of Matsumae Kageyu” by
24 Akoku Raishiki comprise a dataset denoted as “Japanese”.

25 Letters written by different French philosophers constituting 4 collections are denoted
26 as “CITERE”. There are 11 pages containing approximately 2000 words, where 51 words
27 were used as queries. URL: <http://citere.hypotheses.org/>

28 The Cairo Geniza (CG) collection that consists of 12 document images dated to the 10th
29 century. This collection exhibits smeared characters, bleed through, and stains. The page
30 size is about 1650 × 2330 pixels, and the collection contains 1371 words of 921 different
31 transcriptions. URL: <http://www.genizah.org/>

32 A collection of 10 pages of Islamic manuscripts from Harvard University denoted as
33 “Arabic” consists of documents that are dated from 12th to 15th centuries. The ground truth
34 for this collection is given in terms of word-parts. Since word-parts are relatively small, for
35 the experiments, 5117 largest word-parts were chosen with 929 different transcriptions. The
36 page size is 1600 pixels. URL: <http://ocp.hul.harvard.edu/ihp/>

37 **Conclusive Remarks**

38 The major difference between segmentation-based and segmentation-free methods is the
39 different search space (distinct word images versus the whole document image) where they
40 operate. This is the basis of each disadvantage or advantage that each approach entails.

41 The main advantage of the segmentation-based methods is the retrieval speed. The
42 ability to know the words boundaries inside the document provide profound advantages

1 with respect to an efficient retrieval performance. Therefore, segmentation-based methods
2 are mainly based on word segmentations as there is only one method i.e. Keaton et al.
3 (1997) which alternatively uses textline segmentation.

4 On the other hand, if the document is very noisy or very complex to apply a word seg-
5 mentation method then a segmentation-free methodology is the most appropriate approach.
6 Unfortunately, current segmentation-free KWS methods seem not to be appealing since
7 the indexing storage size and the retrieval computation are very costly even when they are
8 dealing with a medium size (100 pages) datasets. Moreover, the complicated issue of man-
9 aging directly the whole document image is the main reason that few works deal with the
10 segmentation-free approach.

11 Concerning features extraction, it is worth noting, that the majority of the works are
12 relying on some form of gradients like SIFT, HOG, LGH, SSHOG, POG. Although, the
13 initial approaches were using profile features, recent works understood the spatial texture
14 information is more robust than shape especially for the handwritten documents. Lastly,
15 some very recent works use features that are obtained using Deep Neural Networks which
16 are called deep features.

17 For the training-based methods, the BoVW has been extensively used as a standalone
18 learning component or combined with other models like the HMMs. The connection with
19 the HMMs was motivated by the use of HMMs for handwritten transcription using a mod-
20 elling inherent to the way a human makes a transcription. The recent advent of CNNs has
21 started to appear in the KWS context (Sharma et al., 2015).

22 Until recently, the most common used dataset was the George Washington dataset for
23 which, there was not a common evaluation protocol and each researcher employing its own
24 subset and query set. Fortunately, the recent KWS competitions have set the ground for a
25 concise performance evaluation framework.

26 Finally, in some works (Zagoris et al., 2011, 2014; Rusinol et al., 2011; Bhardwaj et al.,
27 2008) the KWS pipeline is coupled with a relevance feedback mechanism which introduces
28 the user in the retrieval loop, thus, improving the final retrieval performance.

1 Bibliography

- 2 Adamek, T., O'Connor, N. E., and Smeaton, A. F. (2007). Word matching using single
3 closed contours for indexing handwritten historical documents. *International Journal of*
4 *Document Analysis and Recognition (IJ DAR)*, 9(2-4):153–165.
- 5 Aldavert, D., Rusinol, M., Toledo, R., and Lladós, J. (2013). Integrating visual and textual
6 cues for query-by-string word spotting. In *Document Analysis and Recognition (ICDAR),*
7 *2013 12th International Conference on*, pages 511–515. IEEE.
- 8 Aldavert, D., Rusiñol, M., Toledo, R., and Lladós, J. (2015). A study of bag-of-visual-words
9 representations for handwritten keyword spotting. *International Journal on Document*
10 *Analysis and Recognition (IJ DAR)*, 18(3):223–234.
- 11 Almazán, J., Fernández, D., Fornés, A., Lladós, J., and Valveny, E. (2012a). A coarse-to-
12 fine approach for handwritten word spotting in large scale historical documents collec-
13 tion. In *Frontiers in Handwriting Recognition (ICFHR), 2012 International Conference*
14 *on*, pages 455–460. IEEE.
- 15 Almazán, J., Gordo, A., Fornés, A., and Valveny, E. (2012b). Efficient exemplar word
16 spotting. In *Bmyc*, volume 1, page 3.
- 17 Almazan, J., Gordo, A., Fornés, A., and Valveny, E. (2013). Handwritten word spotting with
18 corrected attributes. In *Proceedings of the IEEE International Conference on Computer*
19 *Vision*, pages 1017–1024.
- 20 Bhardwaj, A., Jose, D., and Govindaraju, V. (2008). Script independent word spotting in
21 multilingual documents. In *IJCNLP*, pages 48–54.
- 22 Can, E. F. and Duygulu, P. (2011). A line-based representation for matching words in
23 historical manuscripts. *Pattern Recognition Letters*, 32(8):1126–1138.
- 24 Cao, H., Govindaraju, V., and Bhardwaj, A. (2011). Unconstrained handwritten docu-
25 ment retrieval. *International Journal on Document Analysis and Recognition (IJ DAR)*,
26 14(2):145–157.
- 27 Chatzichristofis, S. A., Zagoris, K., and Arampatzis, A. (2011). The trec files: the (ground)
28 truth is out there. In *Proceedings of the 34th international ACM SIGIR conference on*
29 *Research and development in Information Retrieval*, pages 1289–1290. ACM.

- 1 Choisy, C. (2007). Dynamic handwritten keyword spotting based on the nshp-hmm. In
2 *Document Analysis and Recognition, 2007. ICDAR 2007. Ninth International Conference*
3 *on*, volume 1, pages 242–246. IEEE.
- 4 Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990).
5 Indexing by latent semantic analysis. *Journal of the American society for information*
6 *science*, 41(6):391.
- 7 Diesendruck, L., Marini, L., Kooper, R., Kejriwal, M., and McHenry, K. (2012a). Digitiza-
8 tion and search: A non-traditional use of hpc. In *E-Science (e-Science), 2012 IEEE 8th*
9 *International Conference on*, pages 1–6. IEEE.
- 10 Diesendruck, L., Marini, L., Kooper, R., Kejriwal, M., and McHenry, K. (2012b). A frame-
11 work to access handwritten information within large digitized paper collections. In *E-*
12 *Science (e-Science), 2012 IEEE 8th International Conference on*, pages 1–10. IEEE.
- 13 Dovgalecs, V., Burnett, A., Tranouez, P., Nicolas, S., and Heutte, L. (2013). Spot it! find-
14 ing words and patterns in historical documents. In *Document Analysis and Recognition*
15 *(ICDAR), 2013 12th International Conference on*, pages 1039–1043. IEEE.
- 16 Fernández, D., Lladós, J., and Fornés, A. (2011). Handwritten word spotting in old
17 manuscript images using a pseudo-structural descriptor organized in a hash structure.
18 In *Iberian Conference on Pattern Recognition and Image Analysis*, pages 628–635.
19 Springer.
- 20 Fernández, D., Marinai, S., Lladós, J., and Fornés, A. (2013). Contextual word spotting
21 in historical manuscripts using markov logic networks. In *Proceedings of the 2nd In-*
22 *ternational Workshop on Historical Document Imaging and Processing*, pages 36–43.
23 ACM.
- 24 Fernández-Mota, D., Manmatha, R., Fornés, A., and Lladós, J. (2014). Sequential word
25 spotting in historical handwritten documents. In *Document Analysis Systems (DAS),*
26 *2014 11th IAPR International Workshop on*, pages 101–105. IEEE.
- 27 Fornés, A., Frinken, V., Fischer, A., Almazán, J., Jackson, G., and Bunke, H. (2011). A
28 keyword spotting approach using blurred shape model-based descriptors. In *Proceedings*
29 *of the 2011 workshop on historical document imaging and processing*, pages 83–90.
30 ACM.
- 31 Gatos, B., Pratikakis, I., and Perantonis, S. J. (2006). Adaptive degraded document image
32 binarization. *Pattern recognition*, 39(3):317–327.
- 33 Giacinto, G. and Roli, F. (2004). Instance-based relevance feedback for image retrieval. In
34 *NIPS*, pages 489–496.
- 35 Glucksman, H. A. (1969). Classification of mixed-font alphabets by characteristic loci.
36 Technical report, DTIC Document.

- 1 Hast, A. and Fornés, A. (2016). A segmentation-free handwritten word spotting approach
2 by relaxed feature matching. In *Document Analysis Systems (DAS), 2016 12th IAPR*
3 *Workshop on*, pages 150–155. IEEE.
- 4 Hast, A. and Marchetti, A. (2012). Putative match analysis: a repeatable alternative to
5 ransac for matching of aerial images. In *International Conference on Computer Vision*
6 *Theory and Applications, VISAPP2012, Rome, Italy, 24-26 February, 2012*, pages 341–
7 344. SciTePress.
- 8 Howe, N. R. (2013). Part-structured inkball models for one-shot handwritten word spotting.
9 In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*,
10 pages 582–586. IEEE.
- 11 Ide, E. (1971). New experiments in relevance feedback. *The SMART retrieval system*, pages
12 337–354.
- 13 Keaton, P., Greenspan, H., and Goodman, R. (1997). Keyword spotting for cursive docu-
14 ment retrieval. In *Document Image Analysis, 1997.(DIA'97) Proceedings., Workshop on*,
15 pages 74–81. IEEE.
- 16 Kolcz, A., Alspector, J., Augusteijn, M., Carlson, R., and Popescu, G. V. (2000). A line-
17 oriented approach to word spotting in handwritten documents. *Pattern Analysis & Appli-*
18 *cations*, 3(2):153–168.
- 19 Kovalchuk, A., Wolf, L., and Dershowitz, N. (2014). A simple and fast word spotting
20 method. In *Frontiers in Handwriting Recognition (ICFHR), 2014 14th International*
21 *Conference on*, pages 3–8. IEEE.
- 22 Leydier, Y., Lebourgeois, F., and Emptoz, H. (2007). Text search for medieval manuscript
23 images. *Pattern Recognition*, 40(12):3552–3567.
- 24 Leydier, Y., Ouji, A., LeBourgeois, F., and Emptoz, H. (2009). Towards an omnilingual
25 word retrieval system for ancient manuscripts. *Pattern Recognition*, 42(9):2089–2105.
- 26 Manmatha, R., Han, C., and Riseman, E. M. (1996). Word spotting: A new approach to
27 indexing handwriting. In *Computer Vision and Pattern Recognition, 1996. Proceedings*
28 *CVPR '96, 1996 IEEE Computer Society Conference on*, pages 631–637. IEEE.
- 29 Manning, C. D., Raghavan, P., Schütze, H., et al. (2008). *Introduction to information*
30 *retrieval*, volume 1. Cambridge university press Cambridge.
- 31 Mondal, T., Ragot, N., Ramei, J.-Y., and Pal, U. (2016). Flexible sequence matching
32 technique: An effective learning-free approach for word spotting. *Pattern Recognition*,
33 60:596–612.
- 34 NIST, T. (2013). Trec nist. [Online]. Available:
35 <http://trec.nist.gov/pubs/trec16/appendices/measures.pdf>.

- 1 Pantke, W., Märgner, V., Fecker, D., Fingscheidt, T., Asi, A., Biller, O., El-Sana, J., Saabni,
2 R., and Yehia, M. (2013). Hadara—a software system for semi-automatic processing of
3 historical handwritten arabic documents. In *Archiving Conference*, volume 2013, pages
4 161–166. Society for Imaging Science and Technology.
- 5 Papandreou, A., Gatos, B., and Zagoris, K. (2016). An adaptive zoning technique for word
6 spotting using dynamic time warping. In *Document Analysis Systems (DAS), 2016 12th*
7 *IAPR Workshop on*, pages 387–392. IEEE.
- 8 Pratikakis, I., Zagoris, K., Gatos, B., Louloudis, G., and Stamatopoulos, N. (2014). Icfhr
9 2014 competition on handwritten keyword spotting (h-kws 2014). In *Frontiers in Hand-*
10 *writing Recognition (ICFHR), 2014 14th International Conference on*, pages 814–819.
11 IEEE.
- 12 Pratikakis, I., Zagoris, K., Gatos, B., Puigcerver, J., Toselli, A. H., and Vidal, E. (2016).
13 Icfhr2016 handwritten keyword spotting competition (h-kws 2016). In *Frontiers in*
14 *Handwriting Recognition (ICFHR), 2016 15th International Conference on*, pages 613–
15 618. IEEE.
- 16 Puigcerver, J., Toselli, A. H., and Vidal, E. (2015). Icdar2015 competition on keyword
17 spotting for handwritten documents. In *Document Analysis and Recognition (ICDAR),*
18 *2015 13th International Conference on*, pages 1176–1180. IEEE.
- 19 Rabaev, I., Kedem, K., and El-Sana, J. (2016). Keyword retrieval using scale-space pyra-
20 mid. In *Document Analysis Systems (DAS), 2016 12th IAPR Workshop on*, pages 144–
21 149. IEEE.
- 22 Rath, T. M. and Manmatha, R. (2003a). Features for word spotting in historical manuscripts.
23 In *Document Analysis and Recognition, 2003. Proceedings. Seventh International Con-*
24 *ference on*, pages 218–222. IEEE.
- 25 Rath, T. M. and Manmatha, R. (2003b). Word image matching using dynamic time warping.
26 In *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer*
27 *Society Conference on*, volume 2, pages II–II. IEEE.
- 28 Rath, T. M. and Manmatha, R. (2007). Word spotting for historical documents. *Internat-*
29 *ional Journal of Document Analysis and Recognition (IJ DAR)*, 9(2-4):139–152.
- 30 Retsinas, G., Gatos, B., Stamatopoulos, N., and Louloudis, G. (2015). Isolated character
31 recognition using projections of oriented gradients. In *Document Analysis and Recogni-*
32 *tion (ICDAR), 2015 13th International Conference on*, pages 336–340. IEEE.
- 33 Retsinas, G., Louloudis, G., Stamatopoulos, N., and Gatos, B. (2016). Keyword spotting
34 in handwritten documents using projections of oriented gradients. In *Document Analysis*
35 *Systems (DAS), 2016 12th IAPR Workshop on*, pages 411–416. IEEE.
- 36 Rodríguez-Serrano, J. A. and Perronnin, F. (2009). Handwritten word-spotting using hidden
37 markov models and universal vocabularies. *Pattern Recognition*, 42(9):2106–2116.

- 1 Rodríguez-Serrano, J. A. and Perronnin, F. (2012). A model-based sequence similarity
2 with application to handwritten word spotting. *IEEE Transactions on Pattern Analysis
3 and Machine Intelligence*, 34(11):2108–2120.
- 4 Rodríguez-Serrano, J. A. and Perronnin, F. (2012). Synthesizing queries for handwritten
5 word image retrieval. *Pattern Recognition*, 45(9):3270–3276.
- 6 Rodríguez-Serrano, J. A., Perronnin, F., Sánchez, G., and Lladós, J. (2010). Unsuper-
7 vised writer adaptation of whole-word hmms with application to word-spotting. *Pattern
8 Recognition Letters*, 31(8):742–749.
- 9 Rothacker, L., Rusinol, M., and Fink, G. A. (2013). Bag-of-features hmms for
10 segmentation-free word spotting in handwritten documents. In *Document Analysis and
11 Recognition (ICDAR), 2013 12th International Conference on*, pages 1305–1309. IEEE.
- 12 Rusinol, M., Aldavert, D., Toledo, R., and Lladós, J. (2011). Browsing heterogeneous docu-
13 ment collections by a segmentation-free word spotting method. In *Document Analysis
14 and Recognition (ICDAR), 2011 International Conference on*, pages 63–67. IEEE.
- 15 Rusinol, M., Aldavert, D., Toledo, R., and Lladós, J. (2015). Efficient segmentation-free
16 keyword spotting in historical document collections. *Pattern Recognition*, 48(2):545–
17 555.
- 18 Rusinol, M. and Lladós, J. (2012). The role of the users in handwritten word spotting appli-
19 cations: query fusion and relevance feedback. In *Frontiers in Handwriting Recognition
20 (ICFHR), 2012 International Conference on*, pages 55–60. IEEE.
- 21 Salton, G. (1992). The state of retrieval system evaluation. *Information processing &
22 management*, 28(4):441–449.
- 23 Saon, G. and Belaid, A. (1997). High performance unconstrained word recognition system
24 combining hmms and markov random fields. *International Journal of Pattern Recogni-
25 tion and Artificial Intelligence*, 11(05):771–788.
- 26 Sharma, A. et al. (2015). Adapting off-the-shelf cnns for word spotting & recognition. In
27 *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*,
28 pages 986–990. IEEE.
- 29 Stamatopoulos, N., Gatos, B., and Georgiou, T. (2010). Page frame detection for dou-
30 ble page document images. In *Proceedings of the 9th IAPR International Workshop on
31 Document Analysis Systems*, pages 401–408. ACM.
- 32 Terasawa, K. and Tanaka, Y. (2009). Slit style hog feature for document image word spot-
33 ting. In *Document Analysis and Recognition, 2009. ICDAR'09. 10th International Con-
34 ference on*, pages 116–120. IEEE.
- 35 Van Rijsbergen, C. (1979). Information retrieval 2nd edition butterworths.
- 36 Wang, P., Eglin, V., Garcia, C., Largeton, C., Lladós, J., and Fornés, A. (2014). A novel
37 learning-free word spotting approach based on graph representation. In *Document Anal-
38 ysis Systems (DAS), 2014 11th IAPR International Workshop on*, pages 207–211. IEEE.

- 1 Wang, P., Eglin, V., Garcia, C., Largeton, C., and McKenna, A. (2013). A comprehensive
2 representation model for handwriting dedicated to word spotting. In *Document Analy-*
3 *sis and Recognition (ICDAR), 2013 12th International Conference on*, pages 450–454.
4 IEEE.
- 5 Zagoris, K., Ergina, K., and Papamarkos, N. (2011). Image retrieval systems based on
6 compact shape descriptor and relevance feedback information. *Journal of Visual Com-*
7 *munication and Image Representation*, 22(5):378–390.
- 8 Zagoris, K., Pratikakis, I., and Gatos, B. (2014). Segmentation-based historical handwrit-
9 ten word spotting using document-specific local features. In *Frontiers in Handwriting*
10 *Recognition (ICFHR), 2014 14th International Conference on*, pages 9–14. IEEE.
- 11 Zagoris, K., Pratikakis, I., and Gatos, B. (2015). A framework for efficient transcription
12 of historical documents using keyword spotting. In *Proceedings of the 3rd International*
13 *Workshop on Historical Document Imaging and Processing*, pages 9–14. ACM.
- 14 Zhang, X. and Tan, C. L. (2013). Segmentation-free keyword spotting for handwritten doc-
15 uments based on heat kernel signature. In *Document Analysis and Recognition (ICDAR),*
16 *2013 12th International Conference on*, pages 827–831. IEEE.